

Colecția ***UNIVERSITARIA***

***Bogdan IONESCU***

***ANALIZA și PRELUCRAREA  
SECVENTELOR VIDEO***  
***Indexarea automată după conținut***



***Bogdan IONESCU***

---

***ANALIZA și PRELUCRAREA  
SECVENTELOR VIDEO***  
***Indexarea automată după conținut***



---

București, 2009

**Copyright © 2009, Bogdan IONESCU**  
Toate drepturile sunt rezervate autorului.

*Adresă: Editura TEHNICĂ  
Str. Olari, nr. 23, sector 2  
Bucureşti, România  
cod 024056*

**www.tehnica.ro**

Referenți științifici:

**Prof. univ. dr. ing. Constantin VERTAN**  
**Conf. univ. dr. ing. Mihai CIUC**

Lucrare editată cu sprijinul  
Consiliului Național al Cercetării Științifice  
din Învățământul Superior

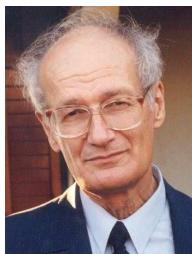
**Descrierea CIP a Bibliotecii Naționale a României**  
**IONESCU, BOGDAN**  
**Analiza și prelucrarea sevențelor video: indexarea**  
**automată după conținut / Bogdan Ionescu. – București:**  
Editura Tehnică, 2009  
Bibliogr.  
ISBN 978-973-31-2354-5

621.397.42  
681.772.7

---

## Prefață

---



Fostul meu student și actualul coleg dr. Bogdan Ionescu m-a rugat să-i prefațez această primă carte, ceea ce fac cu placere - el fiind dintr-o pleiadă de tineri care aduc cinste școlii românești. Bogdan Ionescu și-a terminat de curând (un an) un doctorat în cotutelă (România - Franța) în domeniul relativ nou și foarte important al prelucrării semnalelor multidimensionale, și a avut șansa unui subiect a cărui stringență crește pe zi ce trece: indexarea după conținut a bazelor de date video (mai pe românește: a bibliotecilor de filme - în cazul lui, de filme artistice de animație).

Ce vrea să zică asta - indexarea după conținut - cititorul va găsi în primul capitol, dar sunt tentat să zic și aici, în aceste rânduri, câteva cuvinte: problema nu e chiar nouă. Cu ceva zeci de ani în urmă am aflat că pe alte meleaguri oamenii se ocupau, pentru cuvinte, cu alcătuirea unor asemenea dicționare. Cele alfabetice, pe care le avem și noi, îți explică ce vrea să zică un cuvânt pe care îl ai dar al cărui sens nu îl știi; dar sunt și probleme de alt fel: acolo era un exemplu de întâmplare în academia spaniolă - un vorbitor nu-și aducea aminte cum se cheamă un om născut pe vapor (noi n-avem cuvânt pentru acest concept). Ne trebuie dicționare care să ne ducă de la concept la cuvânt. Despre unele popoare primitive se zice că aveau zeci de cuvinte pentru a denumi diferite tipuri de nori; noi n-avem, dar am putea eventual descrie forma lor, mișcarea lor, ca să precizăm la care ne referim când vrem să povestim o întâmplare concretă.

Într-o bibliotecă de un miliard de cărți, cu câte 500 de pagini fiecare și cu 2000 de semne pe pagină avem nevoie doar de 50 de cifre binare pentru a identifica orice literă, ceea ce mi se pare extrem de puțin - la îndemâna

umanului: le cuprindem cu ochiul dintr-o privire, pe un rând. Oare nu e posibil să avem căi/o cale de a ajunge la ”obiectul” dorit dintr-o colecție vastă, cunoscându-l prin calitățile sale (făcute cumva măsurabile: da-nu, roșu-albastru-galben-verde, o valoare întreagă între 1 și 100, 17 grade de turtire a unui cerc în elipsă, etc.)? ”Obiectele” de care vorbeam pot fi entități foarte complexe: o imagine, o secvență de film mut, entități ”multimodale” (vorbă, sunete, imagini, text, etc.). Parcă suntem tentați a zice da. Dar acum vine partea dificilă a problemei, și în același timp frumoasă prin efortul de creație pe care ni-l cere (aspectul care ne provoacă, ne desfide, englezul ar zice ”challenging”): pe de o parte, în cazul concret al unei colecții de un tip dat (de pietre, de găze, de filme), care sunt atributele, cum le definim ca să caracterizăm cât mai compact și mai corect, acea colecție; pe de altă parte, în fața unui obiect din colecție, cum măsurăm *automat*, adică nu prin intervenția omului (în cazul acesta avem nevoie de un specialist în domeniu!), aceste atributे.

Fără acest mic amănunt aici, ”automat”, suntem pierduți fiindcă operația manuală de adnotare cu atrbute a obiectelor este consumatoare de timp în aşa măsură că ne face întreprinderea lipsită de sens.

În momentul de față al scurtei noastre istorii de câteva sute de ani, suntem în pericol de a fi ”înecată în informații” care pe de o parte multe ne sunt vitale și pe de alta, în ansamblul lor ne copleșesc, fără a putea ajunge la cele de care avem nevoie suntem ca însăsatul din pustiu peste care năvălește marea. Indexarea automată după conținut ne poate salva.

La laboratorul nostru din Politehnica bucureșteană, aceste preocupări sunt de dată mai veche (aș menționa aici preoccupările prof. Constantin Vertan în timpul unor stagii în Franța și apoi aici), dar tomul lui Bogdan Ionescu este prima carte dedicată acestui subiect, și în particular indexării video, și cred că trebuie să o salutăm cu entuziasm fiind sosită într-un moment când e nevoie de ea. Sper că o vor urma altele și că subiectul va atrage și pe alți tineri cercetători spre binele nostru al tuturor. Felicitări autorului pentru munca asiduă depusă și calitatea lucrării rezultate.

Prof. univ. dr. ing. Vasile BUZULOIU  
București 17 noiembrie 2008

---

## Cuvântul autorului

---

Indexarea automată după conținut a datelor este un domeniu ce câștigă din ce în ce mai mult teren, datorită necesității crescânde de exploatare a volumelor mari de date. Dacă, nu până demult, puteam vorbi de o lipsă informațională, progresul tehnologic a făcut ca în zilele noastre să ne confruntăm cu o adevărată explozie de informație.

Din acest amalgam informațional, un interes aparte îl au *informațiile multimedia*, ce sunt definite ca fiind o combinație de tipuri de conținut, printre care cele mai uzuale sunt: textul, sunetul și imaginile.

În societatea modernă, informația multimedia face parte din viața noastră cotidiană și imi este greu să-mi imaginez că vom mai putea vreodată să ne lipsim de ea. De exemplu, telefonul portabil a devenit indispensabil și ne însotește pretutindeni, acesta fiind un adevărat centru multimedia în miniatură. Prin intermediul acestuia, putem accesa informațiile multimedia din rețeaua Internet, putem folosi mesageria electronică, putem înregistra, stoca, reda și distribui filme sau imagini în orice moment. Fiecare persoană, a devenit astfel, cu voie sau fără voie, un ”consumator” de date multimedia.

Motivată în principal de un interes comercial, dezvoltarea infrastructurii de stocare și transmisie a datelor a dus la apariția unei noi probleme, și anume: *Cum accesăm informația multimedia utilă dintr-un vast amalgam de date? Cum facem să găsim aceea informație pe care o dorim?* Problema ar fi una simplă în cazul a câtorva date, dar când o astfel de colecție poate conține la ordinul a sute de mii de documente video, de exemplu, care la rândul lor conțin sute de mii de imagini, problema pare imposibil de rezolvat.

Soluția actuală existentă este dată de *sistemele de indexare după conținut* a bazelor de date. Conceptul de indexare este definit ca fiind procesul de adnotare a informației existente într-o colecție de date, prin adăugarea de

informații suplimentare despre conținutul acesteia, informații numite și *indici* de conținut. Pe baza indicilor, sistemul poate grupa datele în funcție de similaritate, în categorii, subcategorii și aşa mai departe. De exemplu, dacă dispunem de o bază de documente video, ideal, în urma indexării automate, acestea pot fi regrupate în funcție de gen în: filme, muzică, desene animate, știri, etc., sau la un nivel de detaliu mai ridicat, în subcategorii precum: film de ficțiune, dramă, documentar, etc. În acest fel, căutarea informației dorite este restrânsă la căutarea într-o subcategorie din care aceasta face parte, reducând astfel timpul de căutare și totodată imbunătățind precizia căutării.

Pe de altă parte, procesul de indexare nu este opțional, ci este strict necesar într-o colecție mare de date. În acest caz, o informație care nu a fost indexată este practic inexistentă pentru utilizator, cu toate că aceasta este prezentă în bază. Să luăm exemplul simplu al unui sistem de indexare a fișierelor, prezent în orice sistem de operare. Acesta ordonează datele în funcție de nume, tipul conținutului, data creării etc., în directoare și subdirectoare. Dacă un anumit fișier nu a fost indexat, cu toate că acesta se află fizic pe dispozitivul de stocare, acesta este transparent pentru utilizator, fiind imposibil de localizat.

Sistemele de indexare, pe parcursul evoluției, au trecut de la o abordare sintactică a procesului de adnotare la o abordare semantică, cum este cazul sistemelor actuale. Diferența dintre acestea este una semnificativă. *Adnotarea sintactică* se limitează la caracterizarea conținutului datelor cu atrbute numerice de nivel scăzut, precum măsuri statistice, diversi parametri, etc. Din păcate, o astfel de abordare este implicit adresată unui public avizat în domeniu, căutarea informației necesitând cunoștințe tehnice. Pe de altă parte, *adnotarea semantică* are ca scop descrierea conținutului datelor într-un mod cât mai apropiat de modul de percepție uman. Astfel, localizarea datelor devine naturală și accesibilă publicului larg, fiind ghidată de un limbaj textual. De exemplu, căutarea filmelor în funcție de valorile vitezei medii de deplasare a obiectelor în scenă nu este evidentă, pe când o căutare în funcție de conținutul de acțiune (redus, ridicat) este pe înțelesul tuturor.

Această lucrare vine să adreseze tocmai această problematică a indexării automate după conținut a datelor multimedia, punând accentul pe secvențele de imagini, domeniu de mare actualitate în străinătate în acest moment, dar încă la începuturi în România.

Această lucrare propune un studiu bibliografic detaliat al literaturii de specialitate din acest domeniu, abordând direcțiile fundamentale de analiză și prelucrare a secvențelor de imagini în contextul sistemelor de indexare după conținut. Astfel, sunt prezentate atât aspecte teoretice (principii și metode), cât și exemple concrete (sisteme, aplicații), punând la dispoziția cititorului

o bibliografie mai mult decât generoasă (peste 280 de citări ale unor articole din reviste și conferințe internaționale de specialitate). Cartea este adresată atât începătorilor în domeniul prelucrării și analizei de imagini și video, cât și celor deja experimentați, constituind un ghid de bună practică și totodată un sistem de indexare a realizărilor semnificative din domeniu.

Ideea scierii acestui manuscris, a apărut în urmă cu mai bine de cinci ani, odată cu demararea tezei mele de doctorat realizată în cotutelă, pe de-o parte la laboratorul LAPI - Laboratorul de Analiza și Prelucrarea Imaginilor din Universitatea "Politehnica" din București, sub îndrumarea Domnului Profesor Vasile Buzuloiu, cât și la laboratorul, la vremea respectivă, LAMII - Laboratoire d'Automatique et de Micro-Informatique Industrielle din Université de Savoie, sub îndrumarea Domnului Profesor Patrick Lambert. Tematica abordată a constat în studiul și dezvoltarea unui sistem de indexare automată după conținut a secvențelor de animație din cadrul Festivalului Internațional al Filmului de Animăție de la Annecy, echivalentul în domeniul animației al festivalului de film de la Cannes. Studiul bibliografic și cercetarea detaliată realizată cu această ocazie, precum și faptul că doar o parte din acestea au putut fi valorificate în teza de doctorat (din motive obiective de spațiu), m-au condus spre ideea unei posibile redactări ulterioare a unei cărți dedicate.

Această idee avea să se concretizeze după susținerea tezei de doctorat, când am participat la competiția de granturi de Resurse Umane, organizată de CNCSIS - Consiliul Național al Cercetării Științifice din Învățământul Superior, programul RP de stimulare a revenirii în țară a tinerilor cercetători români. Proiectul propus venea să continue natural cercetarea realizată în străinătate până în acel moment, și anume propunea extinderea studiului indexării spre baze de date generice de documente video, precum și dezvoltarea unei aplicații software de adnotare și navigare virtuală în baza de date. Obținerea grantului RP-2 mi-a permis actualizarea studiului bibliografic realizat anterior, îmbunătățirea acestuia, precum și dezvoltarea de noi direcții de studiu.

Astfel, rezultatele cercetării până în acest moment s-au concretizat în șapte capitole. În primul capitol am detaliat problematica indexării după conținut a datelor multimedia, punând accentul pe metodele de analiză și adnotare de conținut a secvențelor de imagini, ce fac subiectul acestei lucrări. De asemenea, am realizat o trecere în revistă a tehnicilor de indexare a imaginilor, sunetului, secvențelor de imagini și, respectiv, video.

Capitolul al doilea abordează o problemă de prelucrare a secvențelor de imagini ce este premergătoare adnotării propriu-zise a conținutului, dar totodată necesară, și anume segmentarea temporală a secvenței, atât în unități sintactice (plane video), cât și semantice (scene video). Segmentarea

temporală, prin detecția schimbărilor de plan, permite înțelegerea structurii temporale a secvenței, necesară în etapele ulterioare de prelucrare, indiferent că este vorba de o indexare sintactică sau de nivel semantic superior.

Capitolul al treilea propune o analiză a metodelor de caracterizare a informației fundamentale a secvențelor de imagini și anume mișcarea. Pornind de la studierea problematicii estimării mișcării la nivel de imagine, am realizat o trecere în revistă a diverselor direcții de studiu abordate de metodele de analiză și caracterizare a conținutului de mișcare din secvență.

Capitolul al patrulea abordează o altă informație reprezentativă a secvențelor de imagini, ce joacă un rol important în percepția vizuală, și anume, conținutul de culoare. Pornind de la modalitățile clasice de reprezentare a culorilor folosind spațiile de culoare, și ajungând până la o descriere perceptuală cu ajutorul teoriei culorilor, am realizat o trecere în revistă a modalităților de caracterizare a conținutului de culoare, atât static, la nivel de imagine, cât și dinamic, la nivel de secvență de imagini.

Capitolul al cincilea propune un studiu al metodelor de rezumare automată a conținutului secvențelor de imagini, atât statică (în imagini) cât și dinamică (în mișcare). Rezumarea de conținut joacă un rol important pentru indexare, deoarece permite pe de-o parte reducerea drastică a timpului vizualizării datelor dintr-o bază mare de date, cât și reducerea redundanței informaționale pentru alte etape de prelucrare.

Capitolul al săselea face trecerea dintre nivelul sintactic de adnotare și cel semantic, prin abordarea tehniciilor de formalizare cu concepte fuzzy a datelor numerice de nivel scăzut.

În final, capitolul al șaptelea prezintă un studiu al tehniciilor de clasificare nesupervizată a datelor (automată), cât și al tehniciilor de clasificare supervizată (ce folosesc o etapă de învățare). Tehnicile de clasificare prezintă un real interes pentru procesul de indexare, deoarece pe baza atributelor de conținut, determinate în etapa de adnotare, acestea pot grupa datele în colecții de date omogene.

Pentru mai multe detalii referitor la aspecte aplicative ale prelucrării și analizei secvențelor de imagini în contextul indexării, cititorul poate consulta site-ul proiectului de indexare RP-2, și anume: <http://alpha.imag.pub.ro/VideoIndexingRP2/>.

Sper sincer ca această lucrare să constituie un ajutor și o referință pentru cei interesați de problemele prelucrării secvențelor de imagini, și că alții o să-mi urmeze exemplul și o să ducă mai departe cercetarea românească din acest domeniu.

S.l. univ. dr. ing. Bogdan IONESCU  
București 30 noiembrie 2008

---

## Mulțumiri

---

Această lucrare nu s-ar fi concretizat fără suportul grantului de cercetare CNCSIS - Consiliului Național al Cercetării Științifice din Învățământul Superior, Resurse Umane, RP-2 (2007-2009), intitulat "Dezvoltarea de Metode de Indexare Semantică după Conținut a Bazelor de Documente Video: Aplicații la Navigare, Căutare și Rezumare Automată a Conținutului"<sup>1</sup>. În acest sens, ţin să mulțumesc *Domnului Președinte CNCSIS Profesor Ioan Dumitrușche*, *Domnului Vicepreședinte Profesor Mihai Gîrțu* și *Domnului Director Profesor Adrian Curaj*, inițiatorii programului RP de reintegrare. De asemenea, ţin să mulțumesc *Doamnei Director Adjunct Magdalena Crîngășu* și *Doamnei Consilier Adriana Rotar*, pentru ajutorul acordat cât și pentru informațiile prețioase oferite pe durata desfășurării proiectului.

Țin să mulțumesc laboratorului LAPI - Laboratorul de Analiza și Pre-lucrarea Imaginilor, din Universitatea "Politehnica" din București, și astfel *Domnului Profesor Vasile Buzuloiu*, pentru acceptarea mea în colectivul de cercetare, pentru prietenia arătată de-a lungul timpului cât și pentru încadrarea prețioasă acordată pe parcursul formării mele profesionale. Mulțumesc colegilor mei profesori, *Constantin Vertan* și *Mihai Ciuc*, pentru ajutorul important, pentru atmosfera cordială din cadrul laboratorului precum și pentru modelul de conduită arătat.

Mulțumesc în mod special *Domnului Profesor Adrian Badea* și *Domnului Profesor Corneliu Burileanu*, pentru prietenia acordată, pentru ajutorul prețios, pentru sugestiile valoroase și pentru suportul constant de-a lungul formării mele științifice.

Țin să mulțumesc *Domnului Profesor Nicolae Vasiliu* pentru ajutorul

---

<sup>1</sup>vezi site-ul proiectului "<http://alpha.imag.pub.ro/VideoIndexingRP2/>".

acordat publicării acestei lucrări și pentru sfaturile prețioase. De asemenea, mulțumesc *Domnului Profesor Ilie Prisecaru* pentru suportul acestuia și pentru ajutorul acordat.

Mulțumesc de asemenea *Domnului Profesor Teodor Petrescu* și *Domnului Profesor Dan Stoichescu* pentru sprijinirea activității mele de cercetare în cadrul Facultății de Electronică, Telecomunicații și Tehnologia Informației și a Catedrei de Electronică Aplicată și Ingineria Informației, precum și pentru ajutorul acordat.

Vreau să mulțumesc laboratorului LISTIC - Laboratoire d'Informatique, Systèmes, Traitement de l'Information et de la Connaissance, Annecy, Franța, și astfel *Domnului Profesor Philippe Bolon* pentru co-finanțarea tezei mele de doctorat realizată în domeniul analizei și prelucrării secvențelor de imagini. De asemenea, mulțumesc călduros *Domnului Profesor Patrick Lambert* și *Domnului Profesor Didier Coquin* pentru încadrarea mea pe parcursul tezei de doctorat, pentru sfaturile acordate cât și pentru suportul constant al acestora pe toată durata cercetării efectuate la Annecy.

De asemenea, ţin să mulțumesc tuturor colaboratorilor externi ce au susținut proiectul de indexare video pe care l-am inițiat:

- *Domnul Profesor Daniel Bouillot*, IMUS - Institut de Management de l'Université de Savoie și CITIA - Cité de l'Image en Mouvement,
- *Domnul Profesor Patrick Lambert* și *Domnul Profesor Philippe Bolon* - LISTIC, Polytech'Savoie, Annecy-Franța,
- *Domnul Profesor Robert Laganière* - VIVA - The Video, Image, Vision and Autonomous Systems Research Laboratory, Ottawa-Canada,
- *Domnul Emmanuel Quillet* și *Domnul Director Hervé Lièvre* CERIMES - Centre de Ressources et d'Information sur les Multimédias pour L'Enseignement Supérieur din Ministère Enseignement Supérieur et Recherche Français.

*Mulțumesc în mod special prietenei mele Monica care m-a sprijinit în tot ce am întreprins până în prezent și care a avut răbdarea să corecteze acest manuscris.*

Cu această ocazie, ţin să mulțumesc călduros *Doamnei Eugenia Burcea* pentru tot sprijinul prețios acordat și *Doamnei Director Diana Cocârlă* ce mi-a sugerat redactarea acestei cărți.

Nu în ultimul rând, vreau să mulțumesc călduros Editurii Tehnice și astfel *Domnului Director Roman Chirilă* pentru acceptarea publicării acestei cărți, pentru finanțarea a o parte din costurile de publicare, precum și pentru munca deloc neglijabilă depusă pentru aducerea la viață a manuscrisului acestei cărți.

---

## Remerciements

---

Ces travaux ne se seraient véritablement concrétisés sans le soutien assuré par le Grant de recherche du CNCSIS, le Conseil National de la Recherche Scientifique de l'Enseignement Supérieur de la Roumanie, Ressources Humaines, RP-2 (2007-2009), intitulé "Le Développement de Méthodes d'Indexation Sémantique du Contenu des Documents Vidéo: Application à la Navigation, Recherche et Résumé Automatique du Contenu"<sup>2</sup>. A ce titre là, je tiens à remercier *M. Ioan Dumitrache*, *M. Mihai Gîrțu* et *M. Adrian Curaj*, les promoteurs du programme RP. Je tiens également à remercier *Mme Adriana Rotar* et *Mme Magdalena Crîngășu* pour leur aide et pour les informations précieuses fournies pendant le déroulement du projet.

Je tiens à remercier le laboratoire LAPI de l'Université "Politehnica" de Bucarest, Laboratoire d'Analyse et Traitement d'Images, et particulièrement *M. Vasile Buzuloiu*, pour m'avoir accueilli au sein de son équipe de recherche, pour son amitié et son encadrement précieux pendant ma formation professionnelle. Je remercie également mes collègues professeurs, *M. Constantin Vertan* et *M. Mihai Ciuc*, pour leur aide précieuse, pour la bonne ambiance qu'ils ont su créer au sein du laboratoire et pour l'exemple de conduite qu'ils ont été.

Je remercie particulièrement à *M. Adrian Badea* et *M. Corneliu Burileanu*, pour leur amitié, leur aide précieuse, leurs nombreux conseils et leur soutien constant tout au long de ma formation scientifique.

Je remercie *M. Nicolae Vasiliu* pour m'avoir aidé dans la publication de ce livre et pour ses précieux conseils. Egalement, je remercie *M. Ilie Prisecaru* pour son soutien et son aide.

---

<sup>2</sup>voir la page web du projet <http://alpha.imag.pub.ro/VideoIndexingRP2>.

Je veux également adresser tous mes remerciements à *M. Teodor Petrescu* et *M. Dan Stoichescu* pour leur aide et pour avoir encouragé et soutenu mon activité de recherche à la Faculté d'Electronique, Télécommunications et Technologie de l'Information et au département d'Electronique Appliquée et d'Ingénierie de l'Information.

Je tiens à remercier le laboratoire LISTIC, Laboratoire d'Informatique, Systèmes, Traitement de l'Information et de la Connaissance d'Annecy, ainsi que *M. Philippe Bolon*, directeur du LISTIC, pour le cofinancement de ma thèse de doctorat portant sur l'analyse et le traitement de séquences d'image. Je remercie également chaleureusement *M. Patrick Lambert* et *M. Didier Coquin* pour leur encadrement pendant ma thèse, leurs nombreux conseils et leur soutien constant tout au long de mes stages de recherche à Annecy.

J'adresse tous mes remerciements aux collaborateurs étrangers qui ont soutenu le projet d'indexation vidéo que j'ai monté à Bucarest:

- *M. Daniel Bouillot* de l'IMUS - Institut de Management de l'Université de Savoie et CITIA - Cité de l'Image en Mouvement,
- *M. Patrick Lambert* et *M. Philippe Bolon* du LISTIC, Polytech'Savoie, Annecy-France,
- *M. Robert Laganière* de VIVA - The Video, Image, Vision and Autonomous Systems Research Laboratory, Ottawa-Canada,
- *M. Emmanuel Quillet* et *M. Hervé Lièvre* du CERIMES - Centre de Ressources et d'Information sur les Multimédias pour l'Enseignement Supérieur du Ministère Enseignement Supérieur et Recherche Français.

*Je remercie tout particulièrement mon amie Monica pour son soutien, pour la patience dont elle a fait preuve pour corriger ce manuscrit et pour avoir été constamment proche de moi.*

A cette occasion, je remercie chaleureusement *Mme Eugenia Burcea* pour son aide très précieuse et *Mme Diana Cocârță* qui m'a suggéré de réaliser ce livre.

Enfin, je tiens à remercier chaleureusement la Maison d'édition "Editura Tehnică" ainsi que *M. Roman Chirilă* pour avoir accepté la publication de ces travaux, pour le financement d'une partie des frais et aussi pour l'aide apportée pour donner vie à ce manuscrit.

---

## Cuprins

---

<b>1 Conceptul de indexare după conținut</b>	<b>1</b>
1.1 Definirea conceptului de indexare . . . . .	2
1.2 Sistemele de indexare de imagini . . . . .	5
1.3 Sistemele de indexare a sunetului . . . . .	7
1.4 Sistemele de indexare a secvențelor de imagini . . . . .	8
1.4.1 Prințipiu adnotării de conținut . . . . .	10
1.4.2 Adnotarea semantică a conținutului . . . . .	14
1.4.3 Sistemul de navigare în baza de date . . . . .	20
1.4.4 Sistemul de căutare în baza de date . . . . .	22
1.5 Sistemele de indexare video . . . . .	26
1.6 Concluzii . . . . .	28
<b>2 Segmentarea temporală</b>	<b>29</b>
2.1 Structura temporală a unei secvențe . . . . .	30
2.2 Descompunerea în plane video . . . . .	33
2.2.1 Detecția de "cuts" . . . . .	33
2.2.2 Detecția de "fades" . . . . .	45
2.2.3 Detecția de "dissolves" . . . . .	50
2.2.4 Evaluarea detecției tranzițiilor video . . . . .	56
2.2.5 Constituirea planelor video . . . . .	59
2.3 Detecția scenelor video . . . . .	61
2.3.1 Tehnici de clasare automată a scenelor . . . . .	63
2.3.2 Tehnici de descompunere în scene . . . . .	65
2.3.3 Aplicații ale analizei scenelor video . . . . .	68
2.4 Concluzii . . . . .	71

<b>3 Analiza mișcării</b>	<b>73</b>
3.1 Estimarea mișcării . . . . .	76
3.1.1 Metodele diferențiale . . . . .	80
3.1.2 Metodele parametrice . . . . .	83
3.1.3 Metodele stohastice . . . . .	85
3.1.4 Metodele de estimare pe blocuri de pixeli . . . . .	87
3.1.5 Fluxul video MPEG . . . . .	97
3.2 Analiza mișcării camerei video . . . . .	99
3.2.1 Analiza mișcării camerei în domeniul comprimat . . . . .	101
3.2.2 Analiza mișcării în domeniul spațio-temporal . . . . .	103
3.3 Concluzii . . . . .	105
<b>4 Analiza de culoare</b>	<b>107</b>
4.1 Spațiile de culoare . . . . .	109
4.1.1 Sisteme de culori primare . . . . .	110
4.1.2 Sisteme pe bază de luminanță-crominanță . . . . .	115
4.1.3 Sisteme perceptuale . . . . .	117
4.1.4 Sisteme de axe independente . . . . .	123
4.2 Conținutul de culoare la nivel de imagine . . . . .	124
4.2.1 Analiza pe bază de histogramă . . . . .	125
4.2.2 Analiza pe baza denumirii culorilor . . . . .	129
4.2.3 Analiza senzației induse de culoare . . . . .	133
4.3 Conținutul de culoare în secvențele de imagini . . . . .	138
4.4 Concluzii . . . . .	141
<b>5 Rezumarea automată de conținut</b>	<b>143</b>
5.1 Construcția rezumatelor statice . . . . .	146
5.1.1 Clasificarea metodelor existente . . . . .	147
5.1.2 Mecanismul de extragere a imaginilor cheie . . . . .	153
5.2 Construcția rezumatelor dinamice . . . . .	162
5.2.1 Informația conservată de rezumat . . . . .	164
5.2.2 Procesul de generare a rezumatului dinamic . . . . .	167
5.3 Metodele de evaluare a rezumatelor . . . . .	170
5.3.1 Analiza descriptivă a rezultatului . . . . .	171
5.3.2 Utilizarea unei măsuri matematice . . . . .	171
5.3.3 Testele de evaluare . . . . .	173
5.4 Concluzii . . . . .	175
<b>6 Formalizarea fuzzy</b>	<b>177</b>
6.1 Introducerea conceptului de incertitudine . . . . .	178
6.2 Logica booleană și logica fuzzy . . . . .	181

6.3	Formalizarea pe baza regulilor fuzzy . . . . .	184
6.3.1	Variabilele fuzzy . . . . .	185
6.3.2	Principiul inferenței fuzzy . . . . .	187
6.4	Avantajele reprezentării fuzzy . . . . .	194
6.5	Aplicabilitatea sistemelor fuzzy . . . . .	195
6.6	Concluzii . . . . .	199
<b>7</b>	<b>Clasificarea după conținut a datelor</b>	<b>201</b>
7.1	Clasificarea nesupervizată a datelor . . . . .	204
7.1.1	Etapele clasificării nesupervizate . . . . .	205
7.1.2	Metodele existente de clasificare nesupervizată . . . . .	206
7.2	Clasificarea supervizată . . . . .	217
7.2.1	Etapele clasificării supervizate . . . . .	218
7.2.2	Metodele existente de clasificare supervizată . . . . .	220
7.3	Concluzii . . . . .	233
	<b>Bibliografie</b>	<b>235</b>



# CAPITOLUL 1

---

## Conceptul de indexare după conținut

---

**Rezumat:** *Un rol important în societatea modernă îl are informația multi-media (imagini, sunet, text, video). Datorită exploziei tehnologice cu care ne confruntăm, volumul informațional a devenit foarte mare. Problema actuală nu este lipsa de informație, ci găsirea informației utile într-un vast amalgam de informații. Soluția este dată de sistemele de indexare după conținut a datelor. În acest capitol vom face o trecere în revistă a literaturii de specialitate din domeniu, prezentând din perspectiva proprie diversele metode și tehnici folosite de sistemele actuale de indexare după conținut a imaginilor, a sunetului, a secvențelor de imagini, precum și video. Prezentarea se va focaliza pe tehniciile de analiză și adnotare a conținutului secvențelor de imagini, ce fac subiectul acestei cărți.*

Dacă în urmă cu aproximativ un deceniu, cantitatea de informație multimedia disponibilă era foarte redusă, iar accesarea acesteia se realiza dificil și inefficient, în zilele noastre putem vorbi despre o *explozie informațională*. Accesul la date, fie că este vorba de imagini, sunet, text sau video, a devenit strict necesar și face parte integrantă din viața noastră de zi cu zi.

Evoluția tehnologică a dispozitivelor de achiziție și prelucrare a datelor (calculatoare personale, medii de stocare, dispozitive de redare și captură audio-video) cât și a infrastructurii de transmisie de date (protocole de transmisie fără fir: WiFi, BlueTooth, rețele LAN de mare viteză, telefonia

multimedia 3G) au dus la creșterea exponențială a volumului informațional prin simplificarea stocării și prelucrării acestuia. Astfel *că problema cu care ne confruntăm acum, nu este lipsa de informație, ci, dimpotrivă imposibilitatea de a selecționa din volumul informațional imens disponibil, informația utilă căutată.*

Pentru a răspunde acestei noi provocări, cercetările actuale din domeniile prelucrării de semnal și a vederii asistate de calculator, au dus la elaborarea a ceea ce numim *sisteme de indexare după conținut* sau CBRS ("Content-Based Retrieval Systems").

## 1.1 Definirea conceptului de indexare

Conceptul de indexare este definit ca fiind **procesul de adnotare** a informației existente într-o colecție de date, prin adăugarea de informații suplimentare relative la conținutul acesteia [Kyungpook 06], informații numite și *indici* de conținut. Această etapă este necesară accesării colecției de date, deoarece permite catalogarea automată în funcție de conținut a datelor.

Într-o colecție de date suficient de vastă, putem spune că datele care nu au fost adnotate sunt practic inexistente pentru utilizator. Un exemplu simplu de sistem de indexare este însuși sistemul de fișiere al oricărui calculator personal. Acesta ne furnizează datele aflate pe diversele medii de stocare (disc dur, memorie externă, etc.) sub formă de fișiere ce sunt indexate după nume, extensie, dată, etc. Să ne imaginăm situația în care un fișier a fost omis din această listă de indici, cu toate că el este prezent fizic pe suportul de stocare, acesta va fi invizibil și inaccesibil pentru utilizatorul de rând.

Procesul de adnotare a datelor este văzut din două perspective: pe de-o parte există *adnotarea manuală*, iar pe de altă parte *adnotarea automată*. Gradul de complexitate al adnotării este direct proporțional cu nivelul de detaliu dorit pentru accesarea datelor. Dacă se dorește ca utilizatorul să poată accesa datele folosind criterii mai complexe, ca de exemplu căutarea unei anumite secvențe video pentru care nu se cunoaște nici numele, nici extensia fișierului, dar totuși utilizatorul dispune de informații referitoare la conținutul vizual al acesteia, în această situație, procesul de indexare va fi mult mai complex, necesitând înțelegerea de către calculator a conținutului datelor.

Astfel, în cazul unei indexări după criterii complexe de conținut, adnotarea manuală este foarte dificil de realizat, deoarece necesită un număr important de operatori umani. Aceștia ar trebui să "răsfoiască" manual întregul conținut al bazei de date pentru definirea indicilor de conținut. Luând în calcul că o astfel de colecție poate conține milioane de înregistrări, timpul

necesar indexării manuale devine semnificativ. În acest moment, cercetările existente în domeniu se focalizează pe dezvoltarea de algoritmi de adnotare automată a conținutului, mai ales în cazul datelor ce necesită un timp important pentru vizualizare, ca de exemplu documentele video.

Cu toate că adnotarea conținutului datelor este soluția optimală pentru a accesa informația utilă dintr-o vastă colecție de date, aceasta nu este și suficientă. Adnotarea în sine nu oferă decât o serie de date suplimentare, putem spune, de nivel scăzut, care deseori sunt inaccesibile utilizatorului neavizat.

Pentru a accesa baza de date, utilizatorul trebuie să disponă de o interfață grafică software prin care să poată accesa sau vizualiza ușor datele, fie pe baza indicilor, fie în mod direct. Aceasta trebuie să aibă o funcționalitate naturală și intuitivă. Sistemul care permite utilizatorului să vizualizeze conținutul bazei de date poartă numele de **sistem de navigare**.

Pe de altă parte, accesul la date presupune un proces de căutare. Utilizatorul trebuie să mai disponă, pe lângă sistemul de navigare, de utilitar software care să-i permită căutarea informațiilor dorite în baza de date. Căutarea se realizează prin formularea de cereri de căutare sau "queries". Pentru ușurință, o astfel de cerere trebuie să fie exprimată într-un limbaj natural, apropiat de limbajul uman, cum ar fi de exemplu "caută filme de acțiune" sau "caută imaginile ce conțin peisaje". Sistemul care răspunde acestor cerințe poartă numele de **sistem de căutare**.

În Figura 1.1 am prezentat diagrama simplificată de funcționare a unui sistem de indexare.

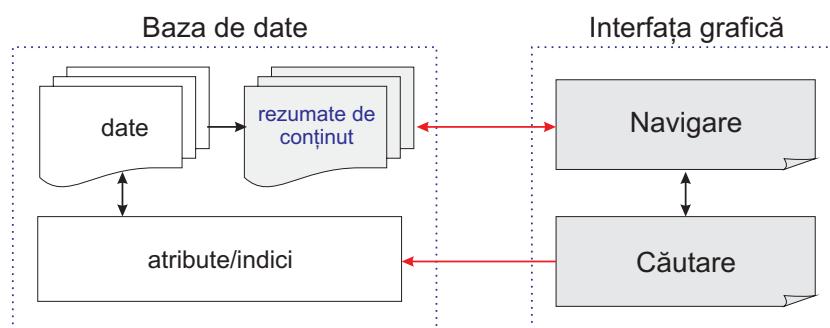


Figura 1.1: Prințipiu de funcționare al unui sistem de indexare după conținut.

Astfel, pentru realizarea indexării după conținut, în primă etapă sistemul analizează datele din baza de date în vederea generării indicilor/atributelor de conținut. Aceștia pot fi: fie date de *nivel semantic scăzut*, precum măsuri

statistice, parametri numerici (de exemplu în cazul documentelor video: histograme de culoare, câmpuri vectoriale de mișcare, etc.), fie date simbolice de *nivel semantic ridicat* (de exemplu în cazul imaginilor statice: obiecte de interes, percepția culorilor, recunoaștere text "încrustat" în imagine, etc.). În final, fiecare înregistrare va avea asociată o colecție de astfel de atrbute ce vor caracteriza conținutul acesteia.

Pentru etapa de adnotare a bazei de date, timpul disponibil nu este critic, în sensul că aceasta este efectuată o singură dată în mod "offline", în momentul creării bazei de date (utilizatorul nu este conectat la baza de date). Tot în această etapă, optional, se pot genera *descrieri compacte ale conținutului* datelor, precum scurte rezumatate pentru secvențele video sau pasaje de text reprezentative pentru documentele textuale. Aceste descrieri vor fi folosite ulterior de sistemul de navigare pentru a facilita accesul la date, utilizatorul putând apela la aceste informații pentru a evita accesarea integrală a conținutului acestora. Totuși, la cererea utilizatorului, sistemul de navigare poate furniza un acces direct la baza de date.

Sistemul de căutare va permite utilizatorului să localizeze datele dorite prin formularea de cereri de căutare. Acestea, după cum am menționat anterior, trebuie formulate într-un limbaj natural apropiat de modul de percepție umană. Pentru a fi înțelese de sistem, cererile de căutare sunt mai întâi convertite în indici folosind același mecanism ca și în cazul adnotării inițiale a bazei de date. Mai departe, căutarea propriu-zisă se efectuează prin compararea indicilor de căutare cu cei deja existenți în baza de date. Rezultatele căutării vor fi acele date ale căror indici sunt cei mai apropiati din punct de vedere al unuia sau a mai multor *criterii de similaritate* de indicii de căutare (de exemplu, pentru indici numerici este vorba de măsuri de distanță, baze de reguli, etc.).

Rezultatele sunt puse la dispoziția utilizatorului prin intermediul sistemului de navigare. Frecvent, în această etapă, pentru a ameliora performanțele sistemului, se permite interacția cu utilizatorul prin ceea ce numim "feedback"<sup>1</sup> al sistemului. Astfel, utilizatorul își va exprima gradul de satisfacție față de rezultatele obținute, sistemul autoinstruindu-se în funcție de preferințele acestuia, pentru a furniza datele cele mai relevante pentru căutare.

În concluzie, sistemele de indexare actuale sunt rezultatul necesității accesării după conținut a datelor multimedia, ce sunt din ce în ce mai diverse. Tehnicile existente au fost adaptate la tipul de date ce trebuie indexate, astfel întâlnim:

---

<sup>1</sup>conceptul de feedback este definit în general ca fiind procesul prin care o proporție din semnalul de ieșire al unui sistem este trecut la intrarea sistemului. În contextul indexării după conținut, "feedback"-ul sistemului constă în informația furnizată de utilizator pentru ameliorarea performanțelor acestuia.

- sisteme **CBIR** ("Content-Based Image Retrieval"): dezvoltate pentru accesarea bazelor de imagini statice (imagini medicale, fotografii, picturi, etc.);
- sisteme **CBAR** ("Content-Based Audio Retrieval"): dezvoltate pentru accesarea bazelor de înregistrări audio (muzică, voce, sunete, etc.);
- sisteme **CBISR** ("Content-Based Image Sequence Retrieval"): dezvoltate pentru indexarea secvențelor de imagini, înțelegând prin aceasta, orice înșiruire temporală de imagini statice sau orice document video în absența informației sonore;
- sisteme **CBVR** ("Content-Based Video Retrieval"): dezvoltate pentru indexarea bazelor de documente audio-vizuale (știri, sport, filme, etc.).

Pentru a înțelege mai bine conceptul de indexare a secvențelor video precum și a tehniciilor folosite de sistemele existente, tematică ce face subiectul acestei cărti, în cele ce urmează vom prezenta succint caracteristicile generale ale sistemelor CBIR și CBAR. Problematica adnotării conținutului de imagine și respectiv sunet influențează în mod direct tehniciile de indexare a documentelor video, acestea din urmă fiind informații audio-vizuale.

## 1.2 Sistemele de indexare de imagini

Unele dintre primele sisteme de indexare au fost sistemele de indexare a bazelor de date de imagini statice sau CBIRS ("Content-Based Image Retrieval Systems"). Necesitatea acestora a fost dată de creșterea semnificativă a numărului de imagini ce trebuiau accesate și stocate, în special în domeniile prioritare precum medicină (baze de date de imagini medicale: radiografii, tomografii, etc.) sau domeniul militar (baze de imagini satelitare), dar și de interes pur comercial, cum ar fi comercializarea bazelor de date de imagini destinate domeniului Prepress<sup>2</sup>.

Global tehniciile de adnotare a conținutului imaginilor sunt orientate către trei axe principale de analiză, și anume: *analiza de culoare*, *analiza formelor* și *analiza de textură* [Smeulders 00].

---

<sup>2</sup>termenul Prepress este folosit în industria tipografică și de publicare pentru a desemna procesele și procedurile software (procesare imagini, text) cât și hardware (creare film tipografic, pregătire tipar) necesare pregătirii manuscrisul sau a creației artistice pentru imprimarea finală pe suportul fizic (hârtie).

**Analiza de culoare** este una dintre tehniciile de adnotare a imaginilor cel mai frecvent folosită, deoarece însuși sistemul vizual uman este bazat pe prelucrarea informației de culoare (unde luminoase de diverse frecvențe). Astfel, sistemele de indexare a imaginilor ce folosesc exclusiv analiza culorilor sunt folosite în aplicații practice în care distribuția de culoare este trăsătura fundamentală a imaginii, precum indexarea picturilor [Lay 04] sau a fotografiilor [Flickner 95]. Culorile sunt analizate folosind diverse spații de culoare, de la cele clasice, precum sistemul RGB (Roșu-Verde-Albastru), până la sisteme perceptuale în care culorile sunt structurate pentru a fi în concordanță cu percepția vizuală umană, precum sistemul HSV (Tintă-Saturație-Intensitate) sau Lab (Luminozitate și diferențe cromatice) [Mojsilovic 00] (un studiu detaliat al spațiilor de culoare este prezentat în Secțiunea 4.1).

**Analiza formelor** se folosește de proprietățile geometrice ale obiectelor conținute în imagine pentru a caracteriza scena. Această analiză presupune detecția în prealabil a acestora folosind tehnici de segmentare bazate pe contur sau pe regiuni de pixeli. Succesul adnotării este astfel direct condiționat de calitatea segmentării în obiecte din imaginii. Mai mult, caracteristicile obiectelor din scenă, obținute în urma analizei formelor, nu trebuie să fie dependente de unghiul sub care a fost prelevată imaginea. Un anumit obiect trebuie caracterizat în același fel chiar dacă a fost imortalizat din unghiuri diferite. Metodele de adnotare existente propun pentru a soluționa această problemă folosirea de descriptori de formă invariante la transformările geometrice ce pot surveni în imagine [Rivlin 95]. De asemenea, se încearcă și soluționarea problemei suprapunerii obiectelor de interes, suprapunere survenită în urma proiecției scenei reale 3D, în spațiul 2D al imaginii. Aceasta poate duce la recunoașterea eronată a obiectului de interes [Schmid 97].

**Analiza de textură** este de asemenea des folosită, deoarece permite caracterizarea proprietăților materialelor prezente în imagine, ca de exemplu caracterizarea texturii pielii pentru detecția persoanelor sau a fețelor prezente în imagine. Metodele existente de analiză de textură folosesc, fie metode clasice de caracterizare a texturii: matrice de co-ocurență, parametri fractali, etc. [Gimel'farb 96], fie abordări mai complexe, cum ar fi analiza Markoviană [Choi 98].

Tendința actuală a sistemelor de indexare de imagini constă în dezvoltarea de metode mixte ce folosesc colaborarea celor trei modalități de analiză, și anume, culoare-formă-textură, profitând astfel de avantajele oferite de fiecare dintre cele trei surse de informație. Ne limităm în această lucrare la

prezentarea a doar câtorva generalități ale sistemelor CBIR. Pentru un studiu bibliografic detaliat al tehniciilor de indexare de imagini existente, cititorul se poate raporta la lucrarea de sinteză propusă în [Smeulders 00].

Dezvoltarea sistemelor de indexare după conținut a imaginilor a dus la apariția a două noi provocări în cercetarea științifică din domeniu. Prima dintre ele este cunoscută sub numele de **paradigma semantică**<sup>3</sup> ("semantic gap") ce este enunțată ca fiind: *discrepanța dintre informațiile extrase în mod automat din imagine și semnificația semantică pe care le-o putem atribui acestora.* Astfel, tehniciile de indexare de imagini duc deseori la rezultate corecte din punct de vedere al algoritmilor de calcul, dar care au o semnificație semantică redusă pentru utilizator. A doua problemă este **paradigma senzorială** ("sensor gap") ce este enunțată ca fiind: *discrepanța care există între informațiile prezente în scena reală 3D și informațiile furnizate de imagine, imagine ce reprezintă o proiecție discretă 2D obținută în momentul înregistrării scenei* [Smeulders 00].

În concluzie, datorită acestor limitări informaționale și tehnologice, sistemele de indexare de imagini existente au tendința să furnizeze rezultate ce nu sunt întotdeauna în conformitate cu realitatea, prezentă în scena reală, a cărei reprezentare o constituie imaginea.

### 1.3 Sistemele de indexare a sunetului

O altă categorie de sisteme de indexare sunt sistemele de indexare a sunetului sau CBARS ("Content-Based Audio Retrieval Systems"). În acest caz, datele indexate sunt documentele audio, precum înregistrările de voce sau înregistrările muzicale.

În general, pentru adnotarea conținutului audio, datele sunt analizate temporal folosind două niveluri de detaliu. Un prim nivel de analiză este *nivelul cadrelor*. Un cadru audio este definit ca fiind o fereastră temporală de analiză de durată redusă. Un al doilea nivel de analiză este *nivelul segment*, un segment audio fiind o fereastră temporală de lungă durată sau chiar întreaga secvență audio ("long-term clip level"). Atributele de conținut specifice datelor audio sunt calculate atât în domeniul *temporal* cât și *frecvențial*. Metodele folosite sunt în mare parte metode clasice ale domeniului de prelucrare de semnal sau de prelucrare a vocii. Proprietățile audio sunt exprimate cu parametri de nivel scăzut specifici, precum: volum, număr de treceri prin zero ale semnalului, tonalitate, parametri spectrali etc. [Wang 00].

---

<sup>3</sup>termenul de semantică este definit global ca fiind relația între semne și lucrurile la care acestea se referă sau înțelesul lor ("denotata").

O primă direcție larg abordată de sistemele de indexare audio o constituie indexarea vocii [Naphade 02]. În acest caz, pentru a ușura analiza, metodele de adnotare a conținutului sunt aplicate în condiții simplificate, cum ar fi de exemplu absența zgomotului de fond sau folosirea unui dicționar de termeni predefinit [Research 05].

O altă direcție de studiu este clasificarea automată a sunetelor în categorii predefinite, precum: voce, muzică, scene de violență, etc. Metodele existente folosesc cu predilecție două abordări: analiza bazată pe reguli și analiza bazată pe modele [Naphade 02].

De mare interes comercial s-a dovedit a fi clasificarea automată a genurilor muzicale, folosită pentru căutarea în bazele de documente muzicale disponibile în format electronic, precum mp3, ogg<sup>4</sup>, etc. Acestea sunt de regulă indexate după gen (blues, clasic, jazz), numele artistului și titlul piesei muzicale. Metodele de indexare folosite utilizează în general algoritmi de clasificare, precum metoda K-means, rețelele neuronale sau sistemele expert (pentru o descriere detaliată a tehniciilor de clasificare a datelor, vezi Capitolul 7). Clasificarea este efectuată în acest caz după parametri specifici, precum: timbrul muzical, armonicitate, melodicitate și ritm [Scaringella 06].

Interesul pentru sistemele de indexare audio, putem spune că este în principal unul comercial, datorat creșterii exponențiale a volumului de documente muzicale "tranzacționate" în format electronic. Pe de altă parte, tehniciile avansate de adnotare semantică a conținutului audio, precum recunoașterea vocii sau a vorbitorului, își au aplicație în principal în sistemele de indexare multimodală complexe (imagine-sunet-text), precum sistemele de indexare video, fiind rar dezvoltate individual.

## 1.4 Sistemele de indexare a secvențelor de imagini

Înlocuirea suportului magnetic de stocare a secvențelor video cu stocarea în format digital, a orientat sistemele de indexare a imaginilor către secvențele de imagini. Sistemele de indexare a secvențelor de imagini sau CBISRS ("Content-Based Image Sequence Retrieval Systems") sunt la origine extensia temporală a sistemelor de indexare a imaginilor CBIR. În cazul CBISR, datele prelucrate nu sunt imagini statice independente, ci serii temporale de

---

<sup>4</sup>Ogg Vorbis este un format de compresie audio comparabil cu celealte formate existente pentru stocarea și redarea audio digitală. Diferența față de acestea constă în faptul ca nu este patentat, putând fi astfel folosit liber, fără a necesita drept de autor. Aceasta poate încapsula atât informație video cât și text.

imagini sau pe scurt *imagini în mișcare* (dinamice). Deseori se întâmplă ca secvențele de imagini să fie numite artificial filme sau documente video. Diferența dintre acestea constă în faptul că documentele video conțin în plus informația audio. În acest context, o secvență de imagini poate fi definită ca fiind informația spațio-temporală a unui document video.

În cazul sistemelor de indexare a secvențelor de imagini, paradigma senzorială (vezi Secțiunea 1.2) este mai puțin sesizabilă datorită informațiilor suplimentare furnizate de acestea, dintre care cea mai importantă este **informația temporală**. Pe de altă parte, particularitățile sistemelor de indexare a secvențelor de imagini, vor aduce o serie de noi dificultăți și cerințe de prelucrare.

O primă problemă este *cantitatea importantă de date* ce trebuie prelucrate. La o frecvență de 25 de imagini pe secundă (de exemplu standardul PAL<sup>5</sup>) o secvență de 10 minute conține nu mai puțin de 15000 de imagini statice, numite și cadre ("frames"). Astfel, o singură secvență de imagini poate fi echivalentul mai multor baze de date de imagini statice. În cazul unei baze de secvențe de imagini, volumul de date este absolut impresionant, deoarece, în mod ușual, aceasta poate conține mii de secvențe. Volumul important de date antrenează după sine dificultatea accesului și a prelucrării acestora.

Pe de altă parte, pe lângă informația spațială furnizată de imagini, secvențele de imagini mai conțin *informație temporală*. Dacă într-un sistem de indexare a imaginilor, două imagini ce conțin aceleași obiecte, pot fi considerate ca fiind similare ca și conținut, într-un sistem de indexare a secvențelor de imagini, două secvențe ce conțin aceleași obiecte pot avea un conținut complet diferit, acest lucru datorându-se aspectului dinamic sau temporal. Astfel, comportamentul obiectelor și evoluția temporală a acestora în scenă sunt informații esențiale pentru înțelegerea conținutului unei sevențe de imagini și drept urmare, pentru procesul de indexare a acestora.

O altă specificitate a secvențelor de imagini este *structura ierarhică* a informației. Într-o secvență, imaginile sunt grupate în ceea ce numim *plane video*. Acestea constituie unitatea sintactică de bază a secvenței și sunt definite ca fiind grupuri de imagini ce au proprietatea de unitate temporală, spațială și de acțiune. La un nivel semantic superior, conținutul secvenței este structurat în unități semantice, precum scenele sau episoadele (structura temporală a secvențelor de imagini este detaliată pe larg în Capitolul 2).

Tinând cont de diversitatea informațională furnizată de secvențele de imagini, tehniciile existente de prelucrare au trebuit să se adapteze pentru a tine cont de aspectul temporal și structural al datelor.

---

<sup>5</sup>PAL este prescurtarea pentru Phase Alternating Line ce reprezintă un standard de condare a semnalului TV color, pe 625 de linii, la frecvența de 50Hz.

În cele ce urmează, vom prezenta pe larg metodele de prelucrare folosite de sistemele de indexare a secvențelor de imagini existente, evidențiind cele trei părți componente de prelucrare ale acestora și anume: *sistemul de adnotare a conținutului*, *sistemul de navigare* în baza de date și *sistemul de căutare*.

### 1.4.1 Principiul adnotării de conținut

După cum am menționat în capitolul introductiv, procesul de adnotare al conținutului datelor constă în *crearea atributelor* sau a indicilor ce permit sistemului înțelegerea automată a conținutului bazei de date. Aceste informații sunt în general proprietăți ale secvenței extrase la nivel de pixel, la nivel de regiuni de pixeli, la nivel de imagine sau grup de imagini.

Indexarea de imagini, fiind un caz particular, simplificat, al indexării secvențelor de imagini, face ca metodele specifice de prelucrare a secvențelor să aibă la bază informațiile folosite la analiza imaginilor statice, și anume: analiza de culoare, textură și formă. Acestea caracterizează proprietățile spațiale ale imaginilor.

Specificitatea secvențelor de imagini constă pe de-o parte în analiza temporală a *evoluției atributelor* extrase la nivel de imagine, precum și în analiza *structurii temporale* a secvenței și a *informației dinamice de miscare*. În Figura 1.2 am sintetizat diferențele surse de informații ce intervin în procesul de adnotare al conținutului secvențelor de imagini.

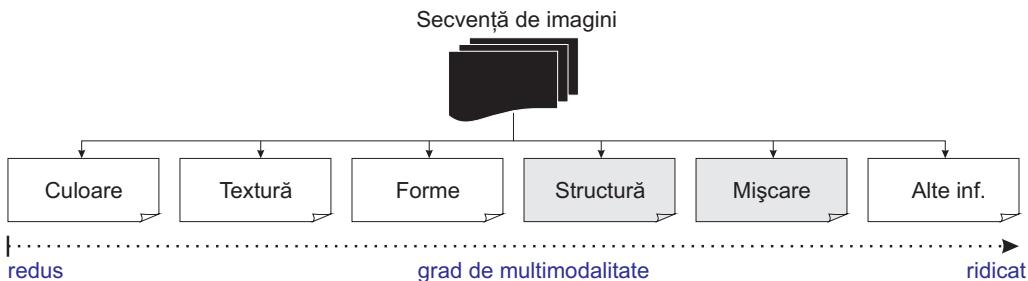


Figura 1.2: Sursele de informații exploatațe de sistemele de indexare a secvențelor de imagini (elementele marcate cu culoarea gri sunt specifice secvențelor).

În cele ce urmează, vom face o trecere în revistă a metodelor și a tehnicilor de indexare ce folosesc fiecare dintre modalitățile menționate anterior.

**Culoarea** este unul dintre atributele cele mai importante pentru caracterizarea conținutului unei imagini. Analiza conținutului de culoare permite

găsirea similarităților vizuale dintre secvențele de imagini. În particular, pentru a cuantifica similaritatea locală sau globală, marea parte a metodelor existente folosesc măsuri statistice, precum *histogramele color*. O astfel de abordare multi-rezoluție este propusă în [Calic 02b] unde pentru a realiza indexarea după conținut, histogramele color sunt calculate pe imagini cu diverse rezoluții. Sistemul propune astfel utilizatorului un număr variabil de niveluri de detaliu, iar o măsură de pertinență este calculată pentru a controla gradul de degradare al imaginilor.

Totuși, măsurile statistice bazate pe histogramă nu conțin informație spațială. Pentru a rezolva această problemă, s-au propus diferite soluții. De exemplu, [Chen 99] folosește histograme augmentate în care fiecărui pixel îi se adaugă informații statistice suplimentare referitoare la similaritatea față de pixelii vecini (medie statistică, entropie, varianță, etc.).

Metodele bazate pe histogramă nu permit nici analiza evoluției temporale, specifică secvențelor de imagini. Pentru a depăși aceste inconveniente, unele metode propun caracterizarea conținutului secvenței pe baza studiului evoluției temporale a unor vectori de caracteristici locale de culoare. O astfel de abordare este propusă în [Zhong 97], unde secvențele de imagini sunt caracterizate prin proprietăți specifice obiectelor de interes: culori specifice, dimensiuni, poziție și traiectorie de deplasare.

Drept alte metode de analiză a conținutului de culoare în vederea indexării, putem menționa metodele ce folosesc arbori de decizie fuzzy pentru extragerea regulilor de indexare [Detyniecki 03], sau modele ale distribuției de rapoarte de culoare ("color ratio models") ce sunt calculate folosind informația de contur [Adjero 01] (o descriere detaliată a metodelor de analiză a culorii în contextul indexării după conținut este prezentată în Capitolul 4).

**Textura** în cazul secvențelor de imagini este folosită cu precădere pentru caracterizarea proprietăților materialelor prezente în scenă sau a obiectelor de interes. Metodele existente folosesc în general aceiași parametri de textură ca în cazul indexării imaginilor statice [Vertan 04]. În literatura de specialitate întâlnim puține tehnici de adnotare a conținutului ce se folosesc exclusiv de analiza de textură, marea parte a metodelor existente folosind abordări mixte în care se regăsește și informația de textură. Cu toate acestea, ca exemplu, putem menționa analiza de textură folosită în [Chang 99] pentru segmentarea cadrelor, sau metoda propusă în [Bouthemy 98] unde câmpul vectorial de mișcare al secvenței este văzut din prisma unei texturi, iar proprietățile temporale ale acestuia sunt caracterizate cu ajutorul parametrilor specifici, precum matricea de co-ocurență.

**Formele.** Parametrii ce caracterizează forma obiectelor de interes, ca și în

cazul sistemelor de indexare a imaginilor statice, sunt analizați în domeniul spațial al imaginii. Specificitatea secvențelor de imagini constă în faptul că obiectele se deplasează în scenă, deplasare ce se traduce în spațiul imaginii prin transformări geometrice progresive ale obiectului vizat. Astfel, *invarianța la transformări geometrice* este una dintre proprietățile fundamentale ce trebuie să îndeplinească descriptorii de formă ce vor servi la indexarea conținutului secvenței. Descriptorii de formă cel mai frecvent folosiți sunt momentele invariante și descriptorii Fourier. În [Mehtre 97] eficiența descriptorilor bazați pe contur (Fourier, UFF - "UNL Fourier Features", etc.) este comparată cu cea a descriptorilor bazați pe regiuni de pixeli (momente invariante, momente Zernike). Invarianța acestora este imbunătățită prin propunerea unor abordări mixte între diversele tipuri de descriptori: momente invariante și descriptori Fourier, sau momente invariante și caracteristici UFF. O altă direcție de studiu este analiza evoluției temporale a formelor obiectelor, des întâlnită și în metodele de urmărire temporală a obiectelor ("object tracking"), cum ar fi de exemplu metoda bazată pe modele multi-rezoluție a contururilor active propusă în [Mazière 00].

**Structura temporală.** Dacă analiza de culoare, textură și formă sunt specifice atât imaginilor statice cât și secvențelor de imagini, structura temporală este o informație specifică doar secvențelor de imagini și documentelor video. În etapa de montaj a secvenței, planele video sunt concatenate pentru a defini ceea ce are să fie evoluția temporală a evenimentelor secvenței (vezi Capitolul 2). Metodele de adnotare a conținutului folosite de sistemele de indexare a secvențelor de imagini sunt bazate pe *segmentarea temporală* în plane, pe *extragerea de imagini cheie* ("key frames"<sup>6</sup>) și pe *analiza similarității* între unitățile temporale ale secvenței: scene, episoade, etc. Modalitatea structurală în care a fost constituită secvența oferă informații prețioase relative la conținutul semantic al acesteia. Astfel, o secvență de acțiune va avea o densitate ridicată de plane video de scurtă durată, în timp ce o secvență a unui reportaj TV, este foarte probabil să conțină doar câteva plane video. Ca metode existente de analiză a modului în care este structurată secvența, putem menționa folosirea de modele probabiliste precum lanțuri Markov ascunse [Ferman 99] sau adnotarea conținutului pe baza analizei și evaluării scenelor video [Vendriga 01].

**Mișcarea.** Analiza de mișcare este o etapă de analiză naturală în cazul secvențelor de imagini, deoarece mișcarea reprezintă însăși esența unei sec-

---

<sup>6</sup>o imagine cheie este o imagine considerată ca fiind reprezentativă pentru conținutul unității structurale din care face parte (plan, scenă, segment, etc.).

vențe. În sistemele de indexare a secvențelor de imagini, caracterizarea conținutului de mișcare este realizată de regulă pe baza estimării câmpului vectorial de mișcare. Vectorii de deplasare sunt fie estimați la nivel de imagine, fie recuperati direct din fluxul video comprimat MPEG (vezi Secțiunea 3.1.5). O primă direcție de studiu o constituie adnotarea *spațio-temporală*. Aceasta include segmentarea, urmărirea de obiecte și caracterizarea mișcării în cadrul anumitor pasaje de interes ale secvenței. Ca exemplu putem menționa sistemul propus în [Dagtas 00], unde mișcarea spațio-temporală a obiectelor este folosită pentru a determina și caracteriza evenimentele importante ale secvenței; sau sistemul VideoQ propus în [Chang 98], dedicat exclusiv caracterizării globale a mișcării obiectelor.

O altă direcție de studiu pentru adnotarea conținutului de mișcare o constituie analiza mișcării camerei video. Metodele existente clasifică diversele mișcări globale în: mișcări translational, de rotație, mărire sau micșorare ("zoom in", "zoom out"), basculare, etc. (vezi Secțiunea 3.2). În această categorie putem menționa metoda propusă în [Lee 01] unde diversele mișcări ale camerei video sunt clasate folosind modele de mișcare predefinite și rețele neuronale, sau abordarea din [Fablet 02] ce folosește modele Gibbs pentru a prezenta derivele semnalului spațio-temporal al imaginii.

În ceea ce privește alte categorii de abordări, putem menționa caracterizarea informației de mișcare folosind traectoria obiectelor sau a regiunilor de interes din imagine [Hsu 02], sau metoda inedită prezentată în [Zeng 02] ce propune transpunerea informației temporale în domeniul spațial al imaginii prin constituirea de hărți de mișcare.

**Alte direcții de studiu** folosesc diferite surse de informație obținute în mod indirect. Una dintre abordările cel mai frecvent folosite pentru procesul de adnotare al conținutului este detecția și analiza prezenței persoanelor în scenă. Aceasta este realizată pe baza detecției caracteristicilor specifice acestora, precum: prezența culorii pielii în imagine, prezența feței, a ochilor, etc. [Acosta 02]. Metodele de localizare a feței sunt bazate pe tehnici de clasificare supervizată, precum rețelele neuronale sau modelele Markov ascunse [Ben-Yacoub 99]. În aceeași categorie putem include și abordările bazate pe detecția anumitor obiecte de interes din scenă, ce sunt considerate ca fiind reprezentative pentru conținutul secvenței, cum ar fi detecția prezenței mașinilor propusă în [Schneiderman 00].

O altă informație prețioasă furnizată de secvențele de imagini este textul "încrustat" în imagine, text ce poate corespunde adnotărilor textuale, genericului, subtitrărilor, scorului în secvențele sportive, numărului de înmatriculare al mașinilor, diverselor indicatoare, etc. Metodele de adnotare a conținutului ce folosesc această informație sunt bazate pe recunoașterea au-

tomată a caracterelor sau OCR ("Optical Character Recognition"<sup>7</sup>). Un exemplu este sistemul propus în [Kim 00b] unde regiunile din imagine ce conțin text încrustat sunt mai întâi izolate folosind o clasificare pe bază de rețele neuronale, iar mai departe literele sunt segmentate și identificate.

#### 1.4.2 Adnotarea semantică a conținutului

Metodele existente de adnotare a conținutului secvențelor de imagini se împart în două mari categorii:

- **metode de adnotare sintactică**, ce sunt utilizate de prima generație de sisteme de indexare, precum cele enumerate în paragrafele anterioare,
- **metode de adnotare semantică**, ce reprezintă noua direcție de analiză folosită de marea parte a sistemelor de indexare actuale.

Adnotarea sintactică este definită generic ca fiind adnotarea ce se referă la *relațiile dintre unitățile de nivel scăzut constituente ale secvenței și modul de constituire a structurii acesteia*. Aceasta se poate realiza pe baza atributelor de nivel scăzut extrase din secvență, precum parametri statistici calculați la nivel de pixel sau regiuni de pixeli, proprietăți geometrice ale obiectelor, structura temporală a secvenței sau vectori de mișcare. De regulă, indicii obținuți în urma procesului de adnotare sunt valori numerice ce descriu atributele enumerate mai sus dar și relațiile sintactice ce pot exista între acestea. Extrași la acest nivel de percepție, indicii sintactici sunt accesibili doar pentru utilizatorul avizat. De exemplu, căutarea unei secvențe de imagini care să conțină 30% mișcare de translație și 20% mișcare de rotație, nu constituie o formulare prea relevantă pentru utilizator.

În contrast cu adnotarea sintactică, adnotarea semantică a conținutului propune o descriere perceptuală la un nivel similar cu nivelul de percepție uman. Informațiile de nivel scăzut obținute în urma analizei sintactice pot fi convertite în concepte lingvistice folosind informații "a priori" despre conținutul secvenței. Totuși, obținerea unei descrieri semantice de conținut necesită înțelegerea completă a conținutului secvenței, astfel că pentru aceasta se preferă o abordare multimodală (imagine-sunet-text).

Un sistem semantic este definit generic ca fiind *orice sistem ce implică o colecție de simboluri (vocabularul sistemului), reguli ce permit constituirea de propoziții, reguli de desemnare și reguli de validare*. În cazul sistemelor

---

<sup>7</sup>recunoașterea automată a caracterelor reprezintă procesul mecanic sau electronic de traducere a imaginilor ce conțin scris de mâna, scris de mașină sau text imprimat (de regulă rezultate în urma procesului de scanare) în text editabil de către calculator.

de indexare, termenul de "semantic" își conservă acest sens. Acesta se traduce prin *codarea interpretării datelor pentru a servi unei aplicații specifice* [Smeulders 00]. Astfel, sistemele de indexare semantică implică existența unui *set de simboluri și reguli* ce permit interpretarea lingvistică a anumitor evenimente sau proprietăți ale secvențelor de imagini.

Adnotarea semantică a conținutului a fost abordată pentru prima oară în sistemele de indexare a imaginilor, dar aceasta era greu de realizat deoarece proprietățile semantice ale scenei sunt dificil de extras dintr-o simplă imagine statică. Datorită informațiilor suplimentare furnizate de secvențele de imagini (informația spațio-temporală și de mișcare), analiza semantică devine mai naturală în acest caz. De exemplu, dacă luăm cazul unei imagini ce surprinde un jucător de fotbal, singurele caracteristici ce reies din analiza imaginii sunt fizionomia acestuia și prezența sa în scenă. Pe de altă parte, dacă dispunem de secvența ce îl surprinde pe jucător, putem determina dacă acesta va marca golul, modul în care acesta joacă, despre ce meci este vorba, etc., informații semantice esențiale pentru înțelegerea conținutului secvenței.

Pentru a înțelege mai bine diferența dintre cele două categorii de adnotări, în Figura 1.3 am ilustrat un exemplu concret de adnotare sintactică și respectiv semantică în cazul unei secvențe de fotbal (axa orizontală reprezintă axa temporală, secvența este rezumată în doar câteva imagini reprezentative).

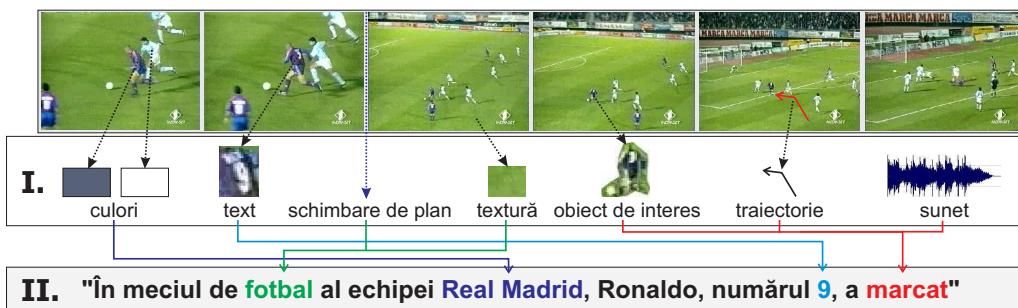


Figura 1.3: Diferența dintre adnotarea sintactică (punctul I.) și semantică (punctul II.). Săgețile colorate indică gradul de implicare al parametrilor de nivel scăzut în construirea descrierii semantice.

Astfel, în acest caz, adnotarea sintactică ne va furniza doar informații relative la scenă și la proprietățile acesteia, precum culoare, prezență text, textură, traiectoria obiectelor în mișcare, ritmul de desfășurare al acțiunii, etc. Pe de altă parte, adnotarea semantică va da sens acestor informații: în mod ideal textura verde va indica că este vorba de un meci de fotbal, culorile jucătorilor (obiecte în mișcare) vor dezvăluia echipele, recunoașterea

numerelor de pe tricou va identifica jucătorii, segmentarea obiectului de interes, urmărirea acestuia și prezența zgomotului specific vor indica marcarea golului. Punând cap la cap toate informațiile, sistemul va ”înțelege” că este vorba despre un meci de fotbal al echipei Real Madrid în care jucătorul cu numărul 9, Ronaldo, marchează.

Paradigma senzorială enunțată în cazul sistemelor de indexare a imaginilor este mai puțin pronunțată în cazul sistemelor de indexare a secvențelor de imagini, acest lucru datorându-se în principal informațiilor suplimentare ce facilitează înțelegerea conținutului secvenței. Cu toate acestea, paradigma semantică, de asemenea prezintă în sistemele de indexare a imaginilor, ia ampioare în cazul secvențelor de imagini datorită *lipsei de corelație dintre informația pe care o putem recupera din conținutul datelor și interpretarea care i-o atribuim* [Smeulders 00].

Astfel că, un sistem de indexare semantică eficient, trebuie să reunească următoarele trăsături importante [Naphade 02]:

- în primul rând este *capacitatea de analiză semantică* pe baza cererilor de căutare formulate de utilizator (vezi Secțiunea 1.4.4),
- un sistem eficient trebuie să fie *multimodal*, reunind și armonizând metode de analiză ce folosesc diversele modalități ale secvenței: imagine, text, sunet, etc.,
- relațiile existente între atrbutele de nivel scăzut și perceptia lor semantică trebuie *rezumate în mod eficient* pentru ca sistemul să fie capabil să ofere utilizatorului o descriere semantică coerentă.

Tendința actuală a sistemelor de indexare a secvențelor de imagini către analiza semantică a fost motivată și de atenția acordată relativ noului standard de compresie video și anume standardul MPEG-7<sup>8</sup> [Wang 00]. Noul standard video încearcă să introducă în procesul de codare, direct în fluxul de date, informații semantice referitoare la conținutul secvenței. Astfel, în momentul indexării, acestea vor putea fi recuperate direct din fluxul MPEG, eliminând astfel procesul de adnotare.

Pentru o descriere mai amănunțită a sistemelor de indexare semantică, cititorul se poate raporta la studiile prezentate în [Naphade 02] și [Snoek 05b]. În cele ce urmează vom prezenta obiectivele sistemelor de indexare semantică precum și dificultățile impuse de analiza semantică a conținutului.

---

<sup>8</sup>standardul MPEG-7 este un standard de descriere a conținutului multimedia. Acesta folosește descrieri suplimentare atașate conținutului video clasic MPEG, pentru a facilita indexarea automată după conținut. Standardul MPEG-7 este denumit formal și Interfață de Descriere a Conținutului Multimedia.

### A. Obiectivele sistemelor de indexare semantică

Global, obiectivele sistemelor de indexare semantică a secvențelor de imagini pot fi structurate pe patru direcții de studiu [Naphade 02], și anume:

- analiza **structurilor de nivel înalt** ale secvenței,
- clasificarea după **gen**,
- analiza **dependentă de domeniul** de aplicație,
- analiza **independentă de domeniul** de aplicație.

**Analiza structurilor de nivel înalt** prezente în secvență presupune detectia și analiza diferitelor pasaje cu semnificație semantică, precum scenele de dialog, spoturile publicitare, diverse evenimente etc. Un exemplu este metoda propusă în [Hauptmann 98] unde frecvența schimbărilor de plan și prezența imaginilor constante, negre, sunt folosite pentru a detecta pasajele publicitare în secvențele de știri. O abordare similară, dar de această dată aplicată pentru detecția pasajelor publicitare în filme, folosește viteza de schimbare a contururilor și analiza amplitudinii vectorilor de mișcare pentru a detecta scenele de acțiune [Lienhart 97]. Un alt exemplu este metoda propusă în [Alatan 01] ce detectează scenele de dialog folosind analiza sunetului și detecția și localizarea prezenței fețelor în imagine. La cel mai jos nivel semantic de analiză se găsesc metodele de detecție a evenimentelor de interes, precum metoda multimodală propusă în [Haering 00]. Aceasta se folosește de analiza de culoare, textură și mișcare într-un clasificator pe bază de rețele neuronale, pentru a identifica pasajele de vânătoare în secvențele documentare.

**Clasificarea după gen.** Un alt obiectiv al sistemelor de indexare semantică este clasarea secvențelor după gen sau tip. Datorită interesului comercial prezentat de acestea, genurile cel mai des vizate sunt știrile, reportajele sportive, filmele și spoturile publicitare. Un exemplu este metoda propusă în [Truong 00a] ce folosește ca atrbute durata planelor video, procentul de apariție al diferitelor tranziții video precum și parametri de culoare, pentru a clasa secvențele de imagini în genurile: animație, publicitate, videoclipuri, știri și sport. O metodă generică, aplicabilă pentru orice tip de secvențe, este propusă în [Kobla 00]. Aceasta folosește ca informație gradul de repetitivitate al mișcării, prezența textului și prezența mișcărilor specifice ale camerei video sau de obiecte. Genul secvenței este mai departe determinat folosind arbori de decizie. O altă abordare interesantă este prezentată în [Colombo 99] unde secvențele publicitare sunt clasate în funcție de percepția conținutului,

în patru sub-genuri: genul practic, critic, utopic și animat. Informațiile folosite pentru aceasta sunt: saturăția culorilor, prezența în imagine a liniilor orizontale, mișcarea și respectiv statistica tranzițiilor video.

**Analiza dependentă de domeniul de aplicație.** În acest caz, metodologia de analiză este adaptată domeniului de aplicație. Astfel, informațiile sunt extrase pe baza expertizei domeniului vizat și sunt specifice tipului de secvențe analizate (vezi studiul bibliografic prezentat în [Snoek 05b]). Fiecare gen de secvență prezintă o serie de particularități ce pot face identificarea și analiza conținutului acestora mai ușoare. De exemplu, întotdeauna un desen animat va avea culori pastelate și nenaturale, sau într-o secvență a unui meci de biliard, culoarea predominantă va fi culoarea verde. Dezavantajul acestui tip de abordare este dat de faptul că nu va putea fi aplicată în alte domenii decât cel vizat. Ca exemplu de analiză dependentă de domeniul de aplicație, putem menționa analiza secvențelor sportive de basket din [Saur 97], a secvențelor medicale din [Fan 04] sau a filmelor de animație din [Ionescu 08].

**Analiza independentă de domeniul de aplicație.** Aceasta este direcția de studiu cea mai dificilă abordată de sistemele de indexare semantică actuale. În prezent nu există multe sisteme care să fie independente de domeniul de aplicație. Cercetările existente încearcă să dezvolte metode de adnotare și clasificare automată a conținutului secvențelor generice, fără a avea cunoștință de proveniența acestora și fără a folosi cunoștințe "a priori" despre conținutul acestora. Ca exemplu de astfel de sisteme putem menționa folosirea de reprezentări probabiliste ale conținutului multimedia propusă în [Naphade 01a], sau folosirea abordărilor de clasificare semantică [Qian 99]. Această direcție rămâne totuși o provocare pentru sistemele de indexare actuale, progresul științific actual din domeniu neputând oferi o soluție definitivă.

## B. Dificultățile ridicate de analiza semantică

Sistemele de indexare semantică existente răspund mai mult sau mai puțin exigențelor actuale de indexare a secvențelor de imagini. Succesul acestora depinde în principal de modalitatea în care se depășesc problemele impuse de adnotarea semantică. În cele ce urmează vom prezenta diversele puncte critice ale adnotării semantice [Fan 04].

**Problema concordanței.** O primă dificultate a adnotării semantice este concordanța dintre analiza de nivel scăzut și descrierea semantică. Aceasta

deinde de eficiență și de calitatea parametrilor numerici folosiți pentru inferență semantică. Pentru a putea face distincția între diversele concepte semantice ce pot fi atribuite datelor analizate, diversitatea parametrilor numerici utilizați trebuie să fie suficient de vastă. Marea parte a sistemelor existente de indexare a secvențelor de imagini folosesc ca informație de plecare pentru extragerea atributelor de conținut, descompunerea secvenței în plane video sau în obiecte semantice (de exemplu scene, pasaje de interes, etc.) [Fan 01]. Pe de altă parte, acest tip de atribute de nivel scăzut sunt dificil de asociat conceptelor semantice, în acest caz o implementare complet automată fiind greu realizabilă [Erol 00].

**Problema modelării.** O a doua dificultate a adnotării semantice este modelarea conceptelor semantice. Datorită paradigmii semantice (vezi Secțiunea 1.2), sistemele actuale nu sunt capabile să furnizeze un acces la baza de date care să fie în totalitate semantic. Diferite soluții au fost totuși propuse pentru a soluționa sau reduce această problemă. Putem menționa metodologia dezvoltată în [Adames 02] ce folosește informații "a priori" din domeniul de aplicație pentru a genera regulile perceptuale aferente descrierii semantice. De asemenea, o altă soluție propusă este folosirea conceptului de "feedback" al sistemului pentru a exploata interacția utilizator-sistem în vederea ameliorării performanțelor indexării [Cox 00]. Tot aici putem menționa metodele de tip "machine learning"<sup>9</sup> ce exploatează corelațiile ascunse ce există între datele multimodale [Barnard 03] (pentru o descriere detaliată a metodelor de clasificare supervizată, vezi Secțiunea 7.2).

**Problema clasificării.** Clasificarea semantică a datelor este de asemenea una dintre dificultățile întâlnite în sistemele de indexare semantică. Metodele existente se împart în metode bazate pe sisteme de reguli constituite pe baza expertizei domeniului de aplicație [Alatan 01] și metode statistice [Wang 01]. Prima categorie de metode se limitează la a folosi doar reguli perceptuale, fără a lua în calcul relațiile existente între informațiile furnizate de diversele modalități ale secvenței. În contrast, metodele statistice permit exploatarea relațiilor ascunse dintre date, dar performanțele acestora sunt dependente de eficacitatea parametrilor aleși, precum și de etapa de antrenare a clasificatorilor folosiți.

**Problema selecției.** O altă dificultate este selecția atributelor și reduce-

---

<sup>9</sup>"machine learning" este un sub-domeniu al inteligenței artificiale ce se ocupă cu proiectarea și dezvoltarea de algoritmi și tehnici ce permit calculatorului să simuleze procesul de învățare.

rea dimensiunii spațiului parametrilor. Intuitiv, putem spune că folosirea unui număr cât mai mare de atribute de nivel scăzut implică o mai bună putere de discriminare, și prin urmare, o adnotare semantică mai eficientă. Totuși, creșterea numărului de atribute implică creșterea radicală a timpului de antrenare al clasificatorilor, precum și a redundanței datelor. Astfel că este foarte important ca din ”marea de atribute” disponibile să alegem doar pe acele care sunt cele mai eficiente pentru indexare.

**Problema organizării.** Organizarea bazei de date și accesul la conținut sunt două aspecte ale sistemelor de indexare cel puțin la fel de importante ca cele enumerate anterior. Din păcate, domeniul bazelor de date și cel al vederii asistate de calculator, nu interacționează încă suficient pentru a propune structurarea bazei de date în funcție de necesitățile de indexare semantică [Smeulders 00]. În cazul ideal, baza de date ar trebui concepută încă din momentul constituirii în aşa fel încât să permită ca accesul la date să se realizeze în mod intuitiv, folosind criterii apropiate de modul de percepție umană [Benitez 01].

### 1.4.3 Sistemul de navigare în baza de date

Accesul la conținutul unei baze de secvențe de imagini se realizează dificil, în primul rând datorită volumului important de informații conținute de aceasta (număr foarte mare de imagini ce trebuie vizualizate). Vizualizarea în parte a fiecărei secvențe este aproape imposibilă. O astfel de bază poate conține milioane de secvențe, astfel că timpul necesar vizualizării poate fi de ordinul anilor. Pentru a facilita accesul la informație, sistemele de indexare a secvențelor de imagini<sup>10</sup> se folosesc de reprezentări compacte de conținut (rezumate). Rezumatele sunt puse la dispoziția utilizatorului prin intermediul unui pachet de utilitate software ce are ca rol principal să permită utilizatorului vizualizarea rapidă și eficientă a conținutului bazei de date. Acestea constituie **sistemul de navigare** în baza de date.

O primă modalitate de construcție a rezumatelor constă în rezumarea conținutului secvenței pe baza structurii temporale (de exemplu, folosind descompunerea în plane). În acest caz, este posibilă construcția a două tipuri diferite de rezumat, și anume: *rezumatul în imagini* (static), care la bază este definit ca fiind o colecție de imagini reprezentative pentru conținutul secvenței, și *rezumatul în mișcare* (dinamic), ce este definit ca fiind o colecție de pasaje reprezentative ale secvenței (pentru o descriere detaliată a tehnici-

---

<sup>10</sup> și nu numai, aceasta fiind valabilă și pentru sistemele de indexare video și respectiv audio.

cilor de rezumare automată de conținut, cititorul se poate raporta la Capitolul 5).

Rezumatele de conținut permit utilizatorului să-și facă rapid o idee globală asupra conținutului secvenței. Astfel, rezumatul static permite reprezentarea conținutului vizual al secvenței în doar câteva imagini, ce sunt ușor accesibile utilizatorului prin sistemul de navigare (timpul de vizualizare fiind neglijabil). Pe de altă parte, rezumatul dinamic aduce un plus de informație la nivelul acțiunii prezente în secvență, informație ce nu este disponibilă în rezumatul static. Totuși, fiind el însuși o secvență, în funcție de nivelul de detaliu furnizat, timpul necesar vizualizării acestuia este mai ridicat decât în cazul rezumatului static, dar net inferior timpului de vizualizare integrală a secvenței.

Ca exemple de sisteme de indexare a secvențelor de imagini ce folosesc un sistem de navigare bazat pe rezumate, putem menționa sistemul propus în [Zhu 05], unde conținutul secvențelor este rezumat folosind "imagini cheie" ("key-frames"<sup>11</sup>). Imaginile astfel reținute, sunt prezentate utilizatorului sub forma unei broșuri. Un alt exemplu este sistemul propus în [Houten 03], unde conținutul secvențelor este prezentat sub formă de colecții de fragmente ("patches"), unde fragmentele sunt definite ca fiind pasaje ale secvenței de aceeași natură semantică (de exemplu scene de dialog, interviuri, etc.).

Un caz particular de rezumare a conținutului o constituie vizualizarea imaginilor reprezentative ale secvenței într-un spațiu 3D, în care primele două dimensiuni sunt date de spațiul imaginii, iar a treia dimensiune o constituie axa temporală a secvenței. De exemplu, în [Vogl 99] secvențele sunt rezumate cu serii temporale de "imagini cheie" pe care utilizatorul le poate vizualiza într-un mediu virtual interactiv. Un sistem similar, numit tunelul temporal ("time tunnel") este prezentat în [Electric 05]. "Imaginile cheie" ale secvenței sunt vizualizate stratificat în funcție de evoluția temporală a cadrelor secvenței (vezi Figura 1.4.a).

O altă modalitate de vizualizare compactă a conținutului secvențelor de imagini este pe baza reprezentării acesteia ca o structură ierarhică. De exemplu, sistemul propus în [Eidenberger 04] folosește un sistem de navigare interactivă constituit pe baza descriptorilor de conținut furnizați de standardul video MPEG-7 [Jeannin 01]. Astfel, acesta pune la dispoziția utilizatorului două modalități de reprezentare arborescentă a conținutului secvenței, și anume: o reprezentare a conținutului de plane video și o reprezentare a structurii temporale. Pe fiecare nivel ierarhic, datele sunt structurate sub formă de hărți cu auto-organizare ("Self-Organizing Maps"). Un alt exemplu de reprezentare ierarhică este cea folosită de sistemul ViBE [Chen 06]. În

---

<sup>11</sup>vezi explicația de la pagina 12.

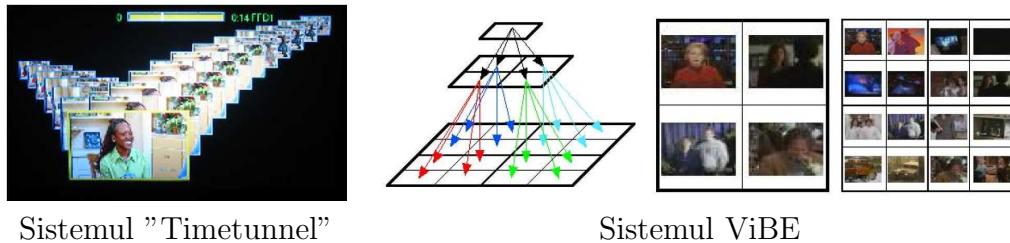


Figura 1.4: Exemple de modalități de vizualizare a conținutului secvențelor de imagini, folosite de sistemul de navigare în baza de date.

aceasta, planele video sunt vizualizate sub formă de structuri arborescente de ”imagini cheie”, ce sunt clasate într-o serie de categorii preudo-semantice determinate ”a priori”. Datele sunt prezentate utilizatorului sub formă piramidală, pe mai multe niveluri de detaliu (vezi Figura 1.4.b).

O categorie aparte de sisteme de navigare sunt sistemele disponibile ”online” pe Internet. Acestea se folosesc de interfață grafică a programului de navigare sau ”Internet browser”. Ca exemple putem menționa sistemul Vimix propus în [Yao 01] în care conținutului secvențelor este ierarhizat pe baza limbajului XML<sup>12</sup>. Un alt sistem este sistemul BIBS [Rowe 01] ce propune o organizare ierarhică liniară a secvenței precum și sincronizarea vizuală a pasajelor secvenței cu adnotări textuale de conținut.

#### 1.4.4 Sistemul de căutare în baza de date

Scopul unui sistem de indexare după conținut este de a furniza utilizatorului posibilitatea de *a căuta și accesa* simplu și eficient informațiile prezente în baza de date. Pentru aceasta, sistemul folosește mai multe etape de prelucrare. În primă fază sunt create *adnotările de conținut* ale datelor, ce permit gruparea acestora în funcție de similaritatea conținutului (vezi Secțiunea 1.4.1). În continuare, *accesul la baza* de date este efectuat cu ajutorul sistemului de navigare (vezi Secțiunea anterioară).

Datorită numărului important de informații disponibile în baza de date, aceste două etape nu sunt încă suficiente pentru a accesa în mod eficient conținutul bazei. O ultimă cerință este posibilitatea de căutare a datelor dorite. Sistemul de căutare permite utilizatorului să localizeze datele prin

---

<sup>12</sup>XML este acronimul pentru ”Extensible Markup Language” și reprezintă un limbaj de uz general folosit în principal pentru a furniza informații suplimentare la conținutul datelor prin crearea unui limbaj de tip ”markup”. Acesta este considerat ca fiind un limbaj extensibil, deoarece permite utilizatorului să-și definească propriile elemente.

interrogarea sistemului. Această interogare se face pe bază de cereri de căutare sau ”queries”. Pentru o căutare optimală și eficientă, sistemul trebuie să permită formularea cererilor de căutare într-un limbaj intuitiv și natural, pe de-o parte accesibil pentru utilizatorul neavizat, dar și inteligibil pentru sistem.

În general, sistemul de căutare funcționează în felul următor:

- **formularea cererii:** mai întâi utilizatorul concepe cererea de căutare. Aceasta poate fi formulată, fie pe baza unui exemplu a ceea ce caută, fie folosind o descriere textuală a conținutului datelor căutate, sau pe baza unei descrieri grafice schematiche a proprietăților datelor căutate.
- **conversia în descriptori sintactici sau semantici:** mai departe sistemul de căutare traduce cererea utilizatorului în attribute sintactice de nivel scăzut sau în attribute semantice de conținut (în funcție de tipul sistemului de indexare). Mecanismul folosit este similar cu cel folosit de sistem în etapa de adnotare a conținutului bazei de date.
- **căutarea propriu-zisă:** căutarea se realizează prin compararea atributelor cererii de căutare cu cele deja stocate în baza de date. Folosind diverse măsuri de distanță și similaritate între attribute, sistemul va căuta datele ce sunt cele mai apropiate de criterile formulate (de exemplu, distanță minimă în cazul parametrilor numerici). Rezultatele obținute vor fi vizualizate în sistemul de navigare.
- **interacția cu utilizatorul:** în mod optional, sistemul poate interacționa cu utilizatorul (”feedback”) pentru a-și ameliora algoritmii de căutare. În acest caz, utilizatorul este încurajat să-și exprime gradul de satisfacție cu privire la rezultatele furnizate de căutare. Aceste date vor servi ca exemple pentru antrenarea ulterioară a sistemului.

Calitatea unui sistem de căutare depinde de mai mulți factori. Mai întâi este vorba de calitatea atributelor ce caracterizează datele și de puterea discriminatorie a măsurilor de similaritate folosite pentru compararea acestora. Totuși, rezultatele căutării sunt dependente în mare măsură de modul în care a fost formulată cererea de căutare. Dacă sistemul nu este capabil să înțeleagă criteriile utilizatorului, căutarea eşuează încă din faza incipientă.

Astfel, calitatea căutării este dependentă mai întâi de nivelul de cunoaștere de către utilizator a datelor căutate. În acest caz întâlnim mai multe situații posibile [Maillet 03]:

- utilizatorul știe cu siguranță că secvența căutată se află în baza de date. În acest caz, ținta este unică iar utilizatorul va fi capabil să formuleze

eficient cererea de căutare. Căutarea se va repeta până în momentul în care secvența va fi găsită.

- utilizatorul caută o anumită secvență, dar nu este sigur că aceasta este prezentă în baza de date. În acest caz, sistemul de indexare trebuie să îi pună la dispoziție algoritmi de căutare preciși și eficienți, pentru ca acesta să decidă rapid dacă secvența este cu adevărat prezentă în baza de date.
- utilizatorul caută o secvență folosind un exemplu sau pe baza unuia sau a mai multor evenimente de interes prezente în aceasta. În această situație, sistemul de indexare trebuie să îl ghidzeze pe utilizator pe tot parcursul căutării. Sistemul trebuie să interacționeze cu utilizatorul pentru a filtra rezultatele obținute ("feedback"). De asemenea, sistemul trebuie să prezinte eficient rezultatele căutării, de regulă folosind reprezentări și descrieri compacte de conținut, pentru ca utilizatorul să se decidă rapid dacă este vorba de secvența căutată, sau în caz contrar, să repete căutarea cu alte criterii.

În al doilea rând, calitatea căutării este condiționată și de modalitatea în care utilizatorul formulează cererea de căutare [Fan 04]:

- **folosirea unui exemplu:** în acest caz, cererea este formulată folosind un model al datelor [Tong 01]. De exemplu, utilizatorul caută toate secvențele ce sunt asemănătoare cu o anumită secvență de care acesta dispune. Similaritatea va fi tradusă pe baza conținutului multimodal al secvenței în: similaritate de culoare, de tehnici de mișcare, a acțiunii, a obiectelor de interes prezente, etc. Această modalitate de căutare se dovedește a fi mai puțin eficientă în situația în care utilizatorul nu cunoaște domeniul respectiv, caz în care acesta nu va fi capabil să indice un bun exemplu pentru căutare.
- **folosirea reprezentărilor textuale:** în acest caz, cererea de căutare este formulată prin sintetizarea caracteristicilor datelor dorite sub formă textuală [Benitez 01]. Acest tip de formulare a căutării, fiind textuală, este apropiată de modalitatea de percepție umană. De exemplu, utilizatorul caută toate "serialele TV", sau la un nivel semantic superior toate "filmele dramatice". Principalul inconvenient la acest tip de căutare este lipsa de sens a anumitor adnotări textuale automate, generate de sistemul de indexare, ce pot conduce la rezultate eronate ale căutării.

- **folosirea sistemului de navigare:** în acest caz, utilizatorul poate folosi direct sistemul de navigare în baza de date pentru a efectua căutarea [Smith 99]. Acest tip de căutare este adecvat utilizatorilor neavizați, ce nu au cunoștință asupra conținutului bazei de date și nici asupra criteriilor de căutare. Totuși, principala constrângere a acestei modalități de căutare o constituie faptul că în general, datele din baza de date nu sunt structurate pe criterii semantice de conținut. Astfel că timpul de localizare a informației dorite poate deveni important.

În cele ce urmează vom prezenta câteva dintre sistemele de căutare folosite de sistemele de indexare a secvențelor de imagini actuale. Marea parte a acestor sisteme, analizează proprietățile conținutului secvențelor folosindu-se de relațiile temporale existente între atributele de conținut, relații puse în evidență pentru prima oară în [Allen 83]. Ca exemple de astfel de sisteme putem menționa sistemul SMOOTH [Kosch 01], GOALGLE și News RePortal [Snoek 05a], în care căutarea este efectuată pe bază de criterii semantice sau temporale (vezi Figura 1.5).



Figura 1.5: Exemple de motoare de căutare: sistemul ”Goalgle” de căutare a secvențelor de gol în secvențele de fotbal și sistemul ”News RePortal” de căutare a stîrilor.

Un alt exemplu este sistemul SoccerQ [Chen 05] ce permite căutarea după criterii semantice a secvențelor sportive în funcție de prezența anumitor evenimente de interes. Căutarea este realizată pe trei niveluri structurale: la nivel de secvență, la nivel de segment și la nivel de variabilă (prin variabilă înțelegând proprietăți ale obiectelor, ca de exemplu numele echipelor). Cererile de căutare sunt formulate într-un limbaj natural, precum: ”caută toate secvențele ce conțin offside”. Alte abordări folosesc predefinirea limbajului de formulare a cererii de căutare, precum sistemul pro-

pus în [Donderler 04]. O abordare diferită este prezentată în [Liu 02a], unde relațiile spațiale și temporale existente între obiectele din scenă sunt modelate prin relații textuale între simboluri, folosind sirurile 3D ("3D-strings"). Astfel, problema căutării este transformată într-o problemă de similaritate între simboluri textuale.

În concluzie, sistemele de căutare în baza de date sunt dependente de domeniul de aplicație (sport, film, etc.), fiind adaptate acestuia. Criteriile de căutare sunt construite pe baza expertizei domeniului respectiv. O posibilă soluție pentru constituirea unui sistem general valabil, aplicabil în cazul oricărei categorii de secvențe de imagini, constă în reunirea "experienței" sistemelor de căutare dezvoltate pentru fiecare domeniu aplicativ în parte, într-un sistem global. În acest fel, se va profita de metodologiile de căutare cele mai performante ale fiecărei categorii de secvențe.

## 1.5 Sistemele de indexare video

Sistemele de indexare video după conținut sau CBVRS ("Content-Based Video Retrieval Systems") sunt extensia naturală a sistemelor de indexare a secvențelor de imagini și a sistemelor de indexare a sunetului, deoarece datele prelucrate sunt în acest caz *documente audio-vizuale*. Un document video este definit ca fiind o secvență de imagini care conține și informație audio (coloana sonoră). Din această cauză, frecvent în literatura de specialitate, cuvântul video este utilizat în mod abuziv pentru a desemna secvențele de imagini. Pe de altă parte, definirea unei frontiere precise între cele două tipuri de sisteme, de indexare de secvențe și respectiv de indexare video, nu este întotdeauna un lucru ușor, deoarece deseori acestea se întrepătrund. În cele ce urmează, ne vom referi prin sistem de indexare video la orice sistem de indexare ce ia în calcul imaginea și sunetul.

În absența informației audio, sistemele de indexare a secvențelor de imagini sunt un caz particular al sistemelor de indexare video. Astfel, toate metodele utilizate de acestea sunt aplicabile și în cazul sistemelor de indexare video. Descrierea sistemelor de indexare video din acest subcapitol nu va relua în discuție metodele deja prezentate în Secțiunea 1.4, în cadrul sistemelor de indexare a secvențelor de imagini, ci ne vom focaliza atenția doar asupra metodelor specifice prelucrării informației audio-vizuale.

Majoritatea sistemelor de indexare video nu folosesc abordări cu adevărat multimediale, în care informația vizuală și audio este prelucrată simultan, în strânsă corelație. În realitate, cele două modalități sunt mai întâi prelucrate în mod separat, iar rezultatele sunt apoi fuzionate pentru a obține descrierea conținutului audio-vizual [Naphade 02]. De exemplu, analiza de imagine

poate fi folosită pentru a realiza segmentarea în plane a secvenței iar sunetul poate fi folosit ulterior pentru a cataloga conținutul video, cum este cazul sistemului propus în [Wang 00].

Astfel, sistemele ce folosesc integrarea analizei multimodale, imagine-sunet, sunt de două tipuri [Snoek 05a]:

- sisteme ce folosesc o abordare pe *bază de reguli* definite prin expertiza domeniului de aplicație [Babaguchi 02],
- sisteme ce folosesc *abordări statistice* pe bază de algoritmi de clasificare a datelor [Han 02b].

Prima categorie de sisteme analizează independent fiecare modalitate video, iar rezultatele sunt fuzionate la final folosind o clasificare pe bază de reguli. În această categorie putem include sistemul propus în [Babaguchi 02], în care sunetul este folosit mai întâi pentru a identifica o serie de cuvinte specifice anumitor evenimente de interes din înregistrările meciurilor de fotbal (de exemplu: momentul golului, aclamațiile publicului, etc.), iar informația vizuală este folosită pentru catalogarea conținutului scenei. Pentru a doua categorie de sisteme, printre abordările statistice cel mai des folosite, putem menționa: metodele ce folosesc rețele Bayes dinamice [Naphade 01b], arbori de decizie [Zhou 02] sau clasificări de tip ”Support Vector Machines” [Lin 02] (vezi Secțiunea 7.2.2).

Analiza multimodală folosită de sistemele de indexare video necesită metode eficiente de fuziune între diferitele surse de informație. Dificultatea unei astfel de prelucrări constă în următoarele aspecte:

- **sincronizarea datelor:** este necesară pentru a omogeniza informațiile provenite de la diferitele surse de informație. Soluția cea mai frecvent adoptată constă în conversia tuturor datelor la un sistem unic de referință, cum ar fi axa temporală. Datele sunt astfel sincronizate în funcție de momentul temporal de producere al acestora [Snoek 05b].
- **alegerea modelului adecvat:** constă în introducerea de informații suplimentare, ce nu sunt disponibile în momentul exact al producerii evenimentului semantic analizat, acestea fiind prelevate din documentul video, fie dinaintea acestuia (sens negativ al timpului) sau de după producerea lui (sens pozitiv al timpului).
- **redundanța parametrilor folosiți:** adnotarea multimodală a conținutului folosește o serie de parametri de conținut ce sunt calculați pentru fiecare modalitate a documentului video. Problema care apare este corelația puternică dintre aceștia, astfel că o etapă de decorelare și selecție este strict necesară [Wang 00].

În concluzie, sistemele actuale de indexare a documentelor video folosesc fuziunea diferitelor tipuri de atribute extrase la nivel de imagine, sunet și text încrustat în imagine. Scopul final este adnotarea conținutului în vederea accesării ulterioare a acestuia la un nivel semantic apropiat de percepția umană.

## 1.6 Concluzii

În acest capitol am prezentat conceptul de indexare după conținut a datelor multimedia. Astfel am făcut o trecere în revistă a metodelor și tehnicilor folosite în sistemele actuale de indexare a imaginilor (CBIR), a sunetului (CBAR), a secvențelor de imagini (CBIRS) precum și a documentelor video (CBVR).

Global, analizând tendințele cercetării din acest domeniu, putem spune că există două direcții distințe. Pe de-o parte sunt sistemele ce folosesc **analiza sintactică a conținutului**. În acest caz, adnotarea bazei de date este realizată cu descriptori statistici de nivel scăzut ce sunt calculați folosind informații precum: culoare, formă, textură, sunet, mișcare, text, etc. Acești descriptori sunt de regulă măsuri numerice complexe, dificil accesibile pentru un utilizator neavizat în domeniu.

Pe de altă parte, întâlnim sistemele ce folosesc **analiza semantică**. Acestea propun descriptori perceptuali ai conținutului datelor, de regulă obținuți pe baza descrierilor sintactice, prin adăugarea expertizei manuale a domeniului vizat. Descriptorii semantici sunt exprimați cu ajutorul unui dicționar de simboluri similar vocabularului folosit de limbajul uman.

Sistemele actuale de indexare după conținut tind să evolueze exclusiv spre *descrierea semantică automată* a conținutului datelor în încercarea de simplificare a problematicii accesării bazelor de date multimedia. Totuși pentru a atinge acest nivel de descriere a datelor, trebuie să depășește o serie de probleme, de la problema selecției atributelor cele mai reprezentative până la paradigma semantică (lipsa de corespondență între parametrii matematici de care dispunem și interpretarea lor semantică).

Sistemele actuale nu au reușit încă să adopte o soluție definitivă pentru aceste probleme, soluția provizorie fiind simplificarea prelucrării prin adaptarea metodelor la domeniul de aplicație. Majoritatea sistemelor de indexare semantică folosesc astfel informații "a priori" despre conținutul datelor, fapt ce a dus la adaptarea metodelor la diversele tipuri de date.

## CAPITOLUL 2

---

### Segmentarea temporală

---

**Rezumat:** Structura temporală a unei secvențe de imagini este similară cu modul în care este structurată o carte, unde diversele capitole sunt înălțuite pentru a constitui narativitatea. În acest capitol vom discuta problematica descompunerii temporale a secvenței în plane video, descompunere ce permite înțelegerea structurală a conținutului acesteia. Astfel, vom analiza diversele metode de detecție a tranzițiilor video, atât abrupte, cât și graduale. De asemenea, vom prezenta și problematica descompunerii secvenței în unități structurale de nivel semantic superior planelor, precum scenele video. Dacă planele video pot fi considerate ca fiind unitățile sintactice de bază ale secvenței, atunci scenele pot fi văzute ca fiind unitățile semantice, acestea permitând o înțelegere mai profundă a conținutului secvenței.

Segmentarea temporală a unei secvențe de imagini este definită ca fiind procesul de divizare al acesteia în *unități structurale* de bază, numite și **plane video**. Fiind o etapă premergătoare înțelegerei structurale a secvenței, marea majoritate a tehnicilor existente de analiză a conținutului secvențelor de imagini, folosesc ca punct de plecare segmentarea în plane [Lienhart 01b].

Procesul de decupare în plane poate fi văzut și din prisma procesului de editare al secvenței ce are loc în studio, la momentul montajului. Planele video sunt concatenate folosind diverse tehnici specifice sau efecte speciale (de exemplu tranzitii video) pentru a da naștere secvenței finale, proces ce este

numit în literatura de specialitate și ”final cut”. În acest sens, segmentarea temporală poate fi percepută ca fiind procesul invers editării secvenței ce are loc în studio.

## 2.1 Structura temporală a unei secvențe

Din punct de vedere al structurii temporale, o secvență de imagini poate fi reprezentată pe mai multe niveluri ierarhice. Acestea sunt ilustrate în Figura 2.1, astfel:

- **nivelul imagine**: reprezintă nivelul structural cu gradul de granularitate cel mai mare (cel mai detaliat nivel) și este reprezentat de toate imaginile conținute în secvență.
- **nivelul planelor video**: corespunde imaginilor secvenței ce au fost filmate între două porniri consecutive ale camerei video. Secvența de imagini astfel obținută are proprietatea de continuitate vizuală (vezi [Corridoni 95]).

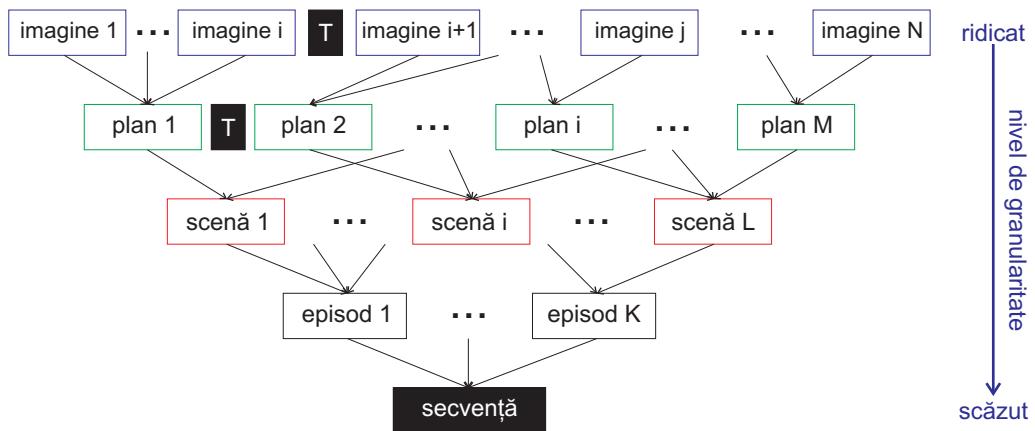


Figura 2.1: Structura ierarhică a unei secvențe de imagini ( $T$  reprezintă o tranziție video).

- **nivelul scenelor**: corespunde grupurilor de plane video ce sunt corelate din punct de vedere al conținutului semantic. Acestea trebuie să respecte regula celor trei unități: unitate de loc, unitate de timp și unitate de acțiune [Corridoni 95].

- **nivelul episoadelor:** corespunde grupurilor de scene ce sunt similare din punct de vedere al acțiunii globale (de exemplu, episoadele unei serii TV) [Bimbo 99].
- **nivelul secvenței:** este nivelul structural cu gradul de granularitate cel mai mic și este reprezentat de secvența însăși.

Marea parte a metodelor de analiză a secvențelor de imagini prelucrează secvența la nivel de plan video. Celelalte niveluri ierarhice, precum scenele sau episoadele, sunt folosite cu predilecție de sistemele de indexare semantică, deoarece detecția acestora presupune o analiză perceptuală de conținut.

Într-o secvență, planele video sunt concatenate pe baza **tranzitiei video** (vezi Figura 2.1). O tranzitie video este un efect vizual folosit pentru a lega imaginea de sfârșit a unui plan, de imaginea de început a planului următor. În funcție de tipul transformărilor 2D aplicate imaginilor, tranzitiiile video existente se împart în cinci clase:

- **clasa de identitate:** tranzitiiile din această categorie nu modifică imaginile planelor video și nici nu adaugă imagini suplimentare (vezi [Lienhart 01b]). În această categorie se află doar tranzitiile de tip "cut", numite și tranzitii abrupte. Un "cut" produce o discontinuitate vizuală în secvență, deoarece planele vecine sunt alipite în mod direct (vezi Figura 2.2).
- **clasa spațială:** din această categorie fac parte tranzitiiile ce aplică imaginilor planelor transformări spațiale [Lienhart 01b] (vezi Figura 2.2). Ca exemple putem menționa efectele de tip "wipes" în care o imagine este înlocuită progresiv de o alta folosind o margine de o anumită formă, efectele de tip "mattes" care de regulă sunt folosite pentru a combina imaginea din planul principal cu imaginea de fundal sau efectele de tip "page turns" în care noua imagine este descoperită simulând răsfoirea paginii unei cărți.
- **clasa cromatică:** în acest caz, imaginile planelor video sunt modificate prin transformări de culoare [Lienhart 01b]. Ca exemple putem menționa tranzitiiile de tip "fade" și "dissolve" (vezi Figura 2.2). Un "fade" este o tranzitie ce permite, fie dizolvarea progresivă a unei anumite imagini într-o imagine constantă, de regulă neagră, ceea ce numim "fade-out", fie apariția progresivă a unei imagini dintr-o imagine constantă, proces numit "fade-in". O tranzitie de tip "dissolve" este în general definită de superpoziția unui efect "fade-out" peste un efect "fade-in", suprapunere ce are ca efect vizual dizolvarea unei imagini în alta.

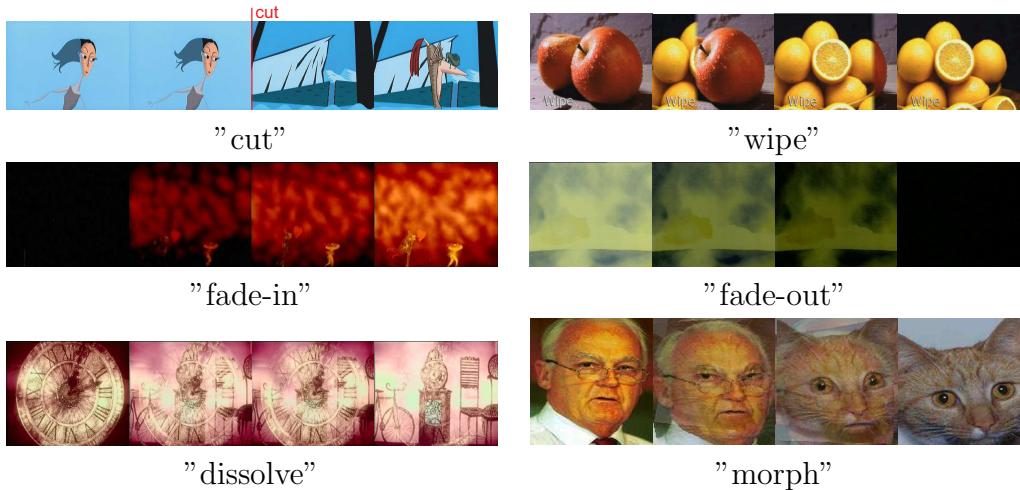


Figura 2.2: Exemple de tranzitii video (pentru fiecare tranzitie au fost prezentate doar cateva imagini reprezentative, axa orizontală reprezentând axa temporală, sursă imagini [Folimage 06] [Wikipedia 08] [Morphing 08]).

- **clasa spațio-cromatică:** tranzitiiile video din această categorie sunt o combinație a clasei spațiale și cromatice, imaginile planelor fiind modificate atât prin transformări spațiale, cât și cromatice [Lienhart 01b]. În această categorie se regăsesc toate efectele de tip "morphing"<sup>1</sup> (vezi Figura 2.2). Cu toate acestea, anumite transformări din clasa cromatică pot fi încadrate și în această categorie, un exemplu fiind transformările de tip "dissolve" ce înglobează mișcări ale camerei video.
- **clasa temporală:** reprezintă o categorie aparte de tranzitii video. În anumite situații, tranzitia de la un plan video la altul se face temporal folosind o mișcare 3D a camerei video (vezi Secțiunea 3.2). De exemplu, camera video filmează un obiect de interes, iar apoi se translatează și se focalizează pe un punct de interes îndepărtat, din fundalul imaginii. Astfel, anumite mișcări 3D ale camerei video, cu toate că nu sunt tranzitii video propriu-zise în sensul definiției enunțate anterior, au rolul de a face legătura între două momente distințe ale secvenței (două plane diferite), putând astfel fi considerate drept tranzitii.

Din punct de vedere al duratei, tranzitiiile video se împart în două categorii, astfel întâlnim *tranzitii abrupte* sau "cuts" și *tranzitii graduale*, precum

<sup>1</sup>"morphing" este un efect special ce presupune metamorfozarea unei imagini în alta prin tranzitii uniforme și constante.

”fade”, ”dissolve”, ”matte” etc. Dintre toate tranzitiiile existente, cel mai frecvent folosite sunt tranzitiiile abrupte de tip ”cut”, deoarece sunt simplu de utilizat și nu introduc întârzieri în derularea conținutului secvenței. Cu o frecvență de apariție cu cel puțin un ordin de măsură mai mic sunt tranzitiiile graduale, dintre care, cele mai întâlnite sunt de tip ”fade” și ”dissolve”.

Folosirea unui anumit tip de tranzitie pentru a face legătura între două plane nu este aleatorie. Fiecare tranzitie are o semnificație semantică ce este adaptată conținutului secvenței. De exemplu, folosirea frecventă a tranzitiiilor de tip ”cut” are ca efect creșterea dinamismul secvenței [Colombo 99]. Pe de altă parte, tranzitiiile de tip ”dissolve” și ”fade” sunt folosite frecvent pentru a schimba timpul sau locul acțiunii [Lienhart 01b], în timp ce folosirea unei tranzitii ”fade-out” urmată de o tranzitie ”fade-in”, introduce un moment de pauză în derularea secvenței și este folosită de regulă pentru a trece la un alt capitol al acțiunii.

## 2.2 Descompunerea în plane video

În cele ce urmează, ne vom limita la trecerea în revistă a algoritmilor și tehniciilor de detectie a tranzitiiilor cel mai frecvent analizate de sistemele de indexare semantică a secvențelor de imagini. Astfel, vom prezenta detectia tranzitiiilor de tip ”cut”, ”fade-in”, ”fade-out” și respectiv ”dissolve”.

Pentru mai multe informații cu privire la algoritmii folosiți pentru detectarea altor tipuri de tranzitii video, cititorul se poate raporta la lucrările [Bimbo 99] (”wipes” și ”mattes”), [Song 02] (”wipes”), [Ren 03] (tranzitii temporale și alte tipuri de tranzitii) sau [Hanjalic 02] (abordare generică a problematicii detectiei tranzitiiilor video folosind metode statistice).

### 2.2.1 Detectia de ”cuts”

După cum am precizat și în paragrafele anterioare, tranzitiiile de tip ”cut” sunt tranzitiiile cel mai frecvent folosite în secvențele de imagini. Acestea sunt definite ca fiind concatenarea directă a două plane video adiacente din punct de vedere temporal,  $P_1(x, y, t)$  și  $P_2(x, y, t)$ , unde  $(x, y)$  reprezintă coordonatele spațiale ale imaginii iar  $t$  coordonata temporală.

Din punct de vedere matematic, secvența ce rezultă în urma concatenării,  $S(x, y, t)$ , este dată de relația următoare [Lienhart 01b]:

$$S(x, y, t) = (1 - u(t - t_{cut})) \cdot P_1(x, y, t) + u(t - t_{cut}) \cdot P_2(x, y, t) \quad (2.1)$$

unde  $t_{cut}$  reprezintă momentul de timp ce corespunde primei imagini de după

tranzită „cut” iar  $u()$  este funcția treaptă unitate:

$$u(t) = \begin{cases} 1 & \text{dacă } t \geq 0 \\ 0 & \text{altfel} \end{cases} \quad (2.2)$$

Definit în acest fel, un „cut” are proprietatea de a introduce o *discontinuitate vizuală* în fluxul secvenței. Metodele existente de detecție folosesc diverse abordări pentru a măsura tocmai această discontinuitate. Cu toate acestea, în general, se pot identifica anumite etape comune de prelucrare.

O primă etapă constă în parametrizarea cadrelor secvenței prin extragerea de *parametri specifici* tranzităilor de tip „cut”, ca de exemplu histograme de culoare, medii spațiale, etc. Mai departe, variația temporală a parametrilor extrași este evaluată între imaginile la momentele  $k$  și  $k + l$ , unde  $k = 0, \dots, N_{sec}$  reprezintă indicele cadrului curent,  $N_{sec}$  este numărul total de imagini al secvenței iar  $l \geq 1$  reprezintă pasul temporal de analiză<sup>2</sup>. Pentru aceasta se folosesc *măsuri de distanță*, sau mai general *măsuri de similaritate*.

Valorile de discontinuitate astfel obținute sunt folosite pentru a localiza tranzită prin compararea acestora cu un anumit prag  $T$ , operație numită și „thresholding”. Dacă valoarea discontinuității se dovedește a fi superioară pragului  $T$ , atunci este foarte probabil ca un „cut” să fi avut loc între imaginile  $k$  și  $k + l$ .

Principalele puncte critice ale unei astfel de abordări au fost bine evidențiate în [Hanjalic 02], astfel:

- o primă problemă este **puterea discriminatorie a parametrilor**. Performanța algoritmului de detecție este total dependentă de puterea de discriminare a parametrilor aleși. În cazul în care parametrii nu sunt reprezentativi pentru conținut, atunci valorile funcției de discontinuitate pentru tranzită de tip „cut” vor fi similare cu cele obținute pentru celealte cadre ale secvenței. Separarea acestora va fi astfel imposibilă.
- o a doua problemă este alegerea **măsurii de distanță**. Aceasta trebuie aleasă astfel încât să furnizeze valori neglijabile pentru imaginile similare, imagini ce aparțin aceluiași plan, și respectiv valori importante pentru imaginile foarte diferite ce separă o tranzită de tip „cut”.
- în ultimul rând este dificultatea alegерii **pragului de detecție**. Astfel, alegerea unui prag  $T$  prea mic va duce la mărirea numărului de false detecții, deoarece pe lângă valorile de discontinuitate specifice unui

---

<sup>2</sup>metodele de analiză și prelucrare a secvențelor de imagini se raportează de regulă la numărul total de cadre al secvenței și mai puțin la durata totală a acesteia. Astfel că, indicele unei imagini va reprezenta numărul cadrului în secvență și nu indicele temporal.

”cut”, va lăsa să treacă și alte astfel de variații mai mici datorate altor factori (de exemplu, mișcării camerei video, mișcării de obiecte). Pe de altă parte, un prag prea ridicat va mări numărul de tranzitii ce vor trece nedetectate.

O problemă stringată ce este responsabilă pentru marea parte a falselor detecții, este disimilaritatea imaginilor același plan video. Aceasta este cauzată în principal de trei situații: mișcarea camerei video, mișcarea obiectelor din scenă sau fluctuațiile intensității luminoase a imaginii. Pentru a compensa aceste situații, o soluție constă în folosirea și a altor surse de informații, pe lângă valorile de discontinuitate. De exemplu, compensarea mișcării poate fi folosită pentru a anula efectul mișcării globale a camerei video [Marichal 98], divizarea imaginii în mai multe regiuni și medierea funcției de discontinuitate în acestea poate fi folosită pentru a reduce influența mișcării obiectelor [Ionescu 06a], detectia blițului aparatului foto poate fi folosită pentru a elimina falsele detecții provocate de acesta [Heng 99] (îndeosebi în secvențele de știri), etc.

În literatura de specialitate există o vastă diversitate de metode de detectie de ”cut” ce ameliorează sau corectează influența factorilor enumerate mai sus, precum metodele propuse în [Bimbo 99], [Lienhart 01b], [Fernando 01], [Hanjalic 02] sau [Ren 03]. În funcție de natura parametrilor folosiți pentru a măsura discontinuitatea vizuală specifică acestui tip de tranzitie, metodele existente se împart în cinci categorii principale, astfel:

- metode bazate pe *analiza intensității pixelilor* din imagine,
- metode bazate pe *analiza contururilor* din imagine,
- metode bazate pe *analiza de mișcare*,
- metode ce analizează conținutul secvenței în formatul original al datelor video, cum ar fi *fluxul comprimat MPEG*.
- alte metode ce folosesc *alte surse de informație*.

În cele ce urmează, vom face o trecere în revistă a tehnicielor și algoritmilor de detectie folosite de metodele din fiecare dintre cele cinci categorii enumerate.

### **Metode bazate pe analiza intensității pixelilor**

Metoda cea mai simplă de evaluare a discontinuității vizuale produsă de o tranzitie de tip ”cut”, constă în folosirea diferenței, la nivel de intensitate a pixelilor, dintre imaginile consecutive la momentele  $k$  și  $k + l$ , unde  $l$

reprezintă pasul de analiză. De exemplu, în [Otsuji 91] un ”cut” este detectat în momentul în care numărul de pixeli,  $N_{pixels}$ , ce se schimbă de la o imagine la alta depășește un anumit prag  $T$ , astfel  $N_{pixels} \geq T$ , unde  $N_{pixels}$  este dat de relația următoare:

$$N_{pixels} = \frac{1}{NM} \sum_{x=1}^X \sum_{y=1}^Y D_{k,k+l}(x, y) \quad (2.3)$$

unde  $X \cdot Y$  reprezintă dimensiunea imaginii iar  $D_{k,k+l}(x, y)$  este definit ca fiind:

$$D_{k,k+l}(x, y) = \begin{cases} 1 & \text{dacă } |I_k(x, y) - I_{k+l}(x, y)| > T_1 \\ 0 & \text{altfel} \end{cases} \quad (2.4)$$

unde  $I_k(x, y)$  reprezintă imaginea la momentul  $k$  iar  $T_1$  este un alt prag estimat în funcție de gradul de schimbare al intensității la nivel de pixel.

Principala problemă a acestui gen de abordare constă în sensibilitatea preciziei detecției la prezența zgomotului în imagine sau a mișcărilor globale ale camerei video, ce au ca efect obținerea de valori semnificative pentru  $N_{pixels}$ . Dintre metodele ce folosesc tehnici similare, putem enumera:

- [Zhang 93] ce propune reducerea zgomotului prin filtrarea mediană a imaginilor înaintea efectuării diferenței dintre acestea,
- [Boreczky 98] ce clasifică distanțele dintre pixeli folosind modele Markov ascunse,
- [Kobla 99] ce folosește ca măsură de disimilaritate distanța Euclidiană calculată în spațiul de culoare YUV și RGB.

O altă cauză ce influențează valoarea distanței dintre imagini, sunt transformările geometrice ce pot surveni în imagine. Pentru a diminua influența acestora, unele metode de detecție apelează la calculul histogramelor de intensitate a pixelilor. Acestea sunt calculate, fie pe imaginea de nivele de gri, fie folosind informația de culoare.

Metoda cea mai des întâlnită constă în măsurarea discontinuității vizuale produse între imaginile la momentele  $k$  și  $k + l$  pe baza sumei distanței dintre binii<sup>3</sup> histogramelor celor două imagini [Yeo 95]. Pentru a reduce influența schimbărilor de intensitate luminoasă ce pot surveni în imagine, [Furht 95] propune folosirea de histograme calculate în spațiul de culoare HSV (H-nuanță, S-saturație și V-intensitate). Astfel, intensitatea luminoasă este separată de

---

<sup>3</sup>histogramele sunt funcții ce contabilizează valorile din anumite grupuri de valori. Acestea sunt numite și clase, dar în contextul histogramei sunt cunoscute sub numele de ”bini” (containeri ce acumulează valori la aceeași rată cu frecvența clasei).

informația de culoare. [Arman 93a] propune calcularea histogramelor folosind doar componente H și S ce formează suprafața 2D,  $HC$ . Aceasta este folosită mai departe la estimarea discontinuității vizuale astfel:

$$D(k, k + l) = \sum_{x=1}^X \sum_{y=1}^Y |d_{k,k+l}(x, y)| \times \Delta_H \times \Delta_C \quad (2.5)$$

unde  $d_{k,k+l}(x, y)$  reprezintă diferența dintre binii de coordonate  $(x, y)$  (apartenând suprafeței  $HC$ ) pentru imaginile la momentele  $k$  și  $k + l$  iar  $\Delta_H$ ,  $\Delta_C$  sunt pașii de cuantizare a componentelor  $H$  și  $S$ .

Alte abordări bazate pe calculul histogramelor încearcă să îmbunătățească invarianta detecției folosind diverse măsuri de distanță ce sunt calculate în spații de culoare precum: HSV, YIQ, Lab, Luv, etc. [Lienhart 01b] (un studiu detaliat al spațiilor de culoare existente este prezentat în Secțiunea 4.1). De exemplu, [Shen 97] propune folosirea distanței Hausdorff<sup>4</sup> multi-nivel pentru a calcula distanța între histograme, [Drew 00] propune ca măsură de similaritate intersecția histogramelor calculată folosind distanța dintre culorile proiectate în spațiul CbCr și rb, [Kim 02] propune calculul histogramelor în spațiul de culoare YUV, [Ma 01] folosește intersecția între histograme și diferența între culorile medii ale blocurilor de pixeli, etc.

Pentru a reduce influența mișcării obiectelor sau a zgometului prezent în imagine, alte metode calculează histogramele la nivel de bloc de pixeli și nu pentru întreaga imagine. În [Nagasaki 92], imaginile la momentele  $k$  și  $k + l$  sunt divizate în 16 blocuri de pixeli iar histogramele,  $H_{k,i}$ , și respectiv  $H_{k+l,i}$ , sunt calculate pentru blocurile de pixeli  $b_i(k)$  și  $b_i(k + l)$ , unde  $i = 1, \dots, 16$  reprezintă indicele blocului. Pentru a compara histogramele astfel obținute se folosește testul  $\chi^2$ , astfel:

$$D(i) = \sum_{j=0}^{63} \frac{[H_{k,i}(j) - H_{k+l,i}(j)]^2}{H_{k+l,i}(j)} \quad (2.6)$$

unde  $j$  reprezintă indicele unui bin al histogramei, în acest caz  $j = 0, \dots, 63$ .

Un alt exemplu este metoda propusă în [Ionescu 07c], unde imaginile sunt divizate în patru regiuni egale (cadrane) în scopul reducerii influenței mișcării obiectelor care intră sau ies din scenă, asupra valorilor histogramelor. Pentru fiecare cadran al imaginii este calculată o histogramă color. Funcția de discontinuitate vizuală este calculată ca fiind media distanțelor Euclidiene obținute între cele patru histograme ale imaginilor la momentele  $k$  și respectiv

---

<sup>4</sup>vezi explicația de la pagina 172.

$k + 2$ , astfel:

$$D(k, k + 2) = \frac{1}{4} \sum_{q=1}^4 d_E^q(k) \quad (2.7)$$

unde  $q$  reprezintă indicele cadranului,  $q \in \{1, 2, 3, 4\}$ , iar  $d_E^q(k)$  este distanța Euclidiană:

$$d_E^{q,2}(k) = \sum_{j=1}^{216} [H_{k+2,q}(j) - H_{k,q}(j)]^2 \quad (2.8)$$

unde  $j$  reprezintă indicele culorii iar  $H_{k,q}(j)$  reprezintă histograma color a cadranului  $q$  din imaginea la momentul  $k$ .

Mai mult, valorile funcției de discontinuitate astfel obținute sunt derivate pentru a reduce fluctuațiile cauzate de mișcările globale ale camerei video sau de diferențe succesive de culoare. Prințipiu de funcționare este ilustrat în Figura 2.3. Un "cut" produce o discontinuitate între două imagini succese, discontinuitate ce se traduce printr-o valoare semnificativă a funcției de discontinuitate. Pe de altă parte, imaginile ce preced și respectiv, ce succed tranzitiei, sunt foarte similare, valorile funcției de discontinuitate fiind în acest caz reduse. Astfel, o tranzitie de tip "cut" are o semnătură particulară a valorilor funcției de discontinuitate, și anume o succesiune de valori de tip *SHS* (*S*-valoare redusă, *H*-valoare ridicată) ce poate fi pusă în evidență cu ajutorul derivatei temporale.

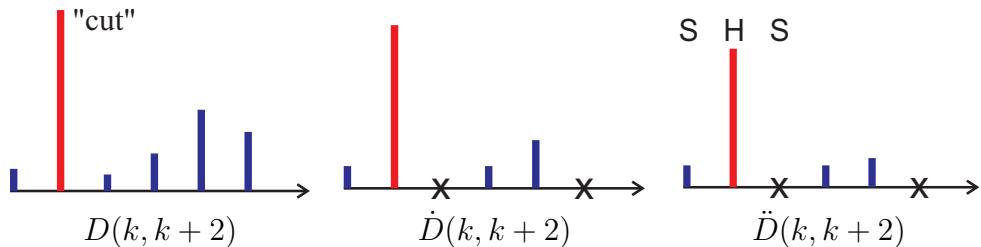


Figura 2.3: Folosirea derivatei temporale pentru accentuarea tranzițiilor de tip "cut" (axa  $oX$  corespunde axei temporale, valorile negative ale derivatei, marcate cu simbolul  $\times$ , sunt anulate deoarece conțin informație redundantă pentru detecție).

Un studiu comparativ al diferitelor metode de detecție de "cut" bazate pe calculul histogramelor color, pe fluxul video MPEG și respectiv pe estimarea mișcării, este propus în [Gargi 00]. Astfel, în contextul prezentat, cea mai eficientă metodă de detecție se dovedește a fi o metodă bazată pe histogramă, ce folosește intersecția histogramelor în spațiul de culoare al lui

Munsell (MTM, vezi Secțiunea 4.1). Din punct de vedere al complexității de calcul, aceasta are însă o complexitate moderată, raportat la celelalte metode testate. Acest lucru este totuși valabil și în general, metodele bazate pe analiza intensității pixelilor se dovedesc încă a fi metodele cele mai robuste de detecție.

### Metode bazate pe analiza de contur

Acestea folosesc pentru detecție analiza contururilor obiectelor prezente în imagine. Un "cut" produce, pe lângă discontinuitatea vizuală a fluxului secvenței, și o discontinuitate structurală la nivel de imagine. Astfel, contururile obiectelor prezente în imaginea ce precede un "cut" nu se vor mai regăsi în imaginea ce succede tranziția. Metodele existente exploatează tocmai această proprietate.

Algoritmii propuși în [Zabih 95] și [Zabih 99] folosesc pentru detectie calculul raportului de schimbare a conturului, ECR ("Edge Change Ratio"). Aceasta este definit pentru imaginile la momentele  $k$  și  $k + l$  în felul următor:

$$ECR_{k+l} = \max \left( \frac{X_k^{out}}{\sigma_k}, \frac{X_{k+l}^{in}}{\sigma_{k+l}} \right) \quad (2.9)$$

unde  $\sigma_k$  reprezintă numărul de puncte de contur existente în imaginea la momentul  $k$  iar  $X_k^{out}$  și  $X_{k+l}^{in}$  reprezintă numărul de puncte de contur ce au dispărut din imaginea la momentul  $k$ , și respectiv, numărul de puncte de contur ce au apărut în imaginea la momentul  $k + l$ .

Pentru a îmbunătății invarianța raportului ECR la prezența mișcării, [Zabih 95] propune folosirea compensării de mișcare ce este calculată folosind distanța Hausdorff<sup>5</sup>. Punctele de contur analizate în imaginea curentă la momentul  $k$ , ce sunt apropriate de punctele de contur din imaginea următoare analizată la momentul  $k+l$ , nu vor fi luate în calcul decât dacă distanța Hausdorff depășește 6 pixeli. Un alt exemplu este metoda propusă în [Kim 02] ce folosește pentru detecție histograme color calculate în spațiul de culoare YUV, precum și evaluarea raportului de potrivire al contururilor, EMR ("Edge Matching Rate"). În [Lienhart 00] sunt propuse mai multe metode pentru detecția tranzițiilor video, metode ce folosesc diverse informații, precum informații de contur, histograme și analiza de mișcare.

Un studiu comparativ al performanțelor metodelor bazate pe analiza de contur precum și a celor bazate pe histogramă este propus în [Lienhart 01a] și [Lupatini 98]. Astfel, testele efectuate arată că metodele bazate pe analiza contururilor sunt mai puțin eficiente și necesită un timp de calcul mult mai

---

<sup>5</sup>vezi explicația de la pagina 172.

ridicat decât metodele bazate pe histogramă. Totuși, metodele bazate pe contur au avantajul de a putea fi utilizate în același timp și pentru detecția tranzițiilor video graduale, precum tranzitiiile de tip "fade" sau "dissolve" [Lienhart 01b].

### Metode bazate pe analiza de mișcare

Pentru detecție, metodele bazate pe mișcare pornesc de la ipoteza că o tranzitie de tip "cut" produce în secvență o discontinuitate a mișcării. Detecția este astfel bazată pe estimarea câmpului vectorial de mișcare (un studiu detaliat al tehniciilor de estimare a mișcării este prezentat în Secțiunea 3.1). Estimarea mișcării este în general realizată folosind metode bazate pe blocuri de pixeli ("block-based"), în principal datorită faptului că acestea oferă un bun compromis între complexitatea de calcul și precizia rezultatelor. Procedeul constă mai întâi în împărțirea imaginilor analizate în blocuri disjuncte de pixeli. Mai departe, pentru fiecare bloc de pixeli,  $b_i(k)$ , de indice  $i$ , din imaginea la momentul  $k$ , se caută blocul cel mai similar,  $b_{i,j}(k+l)$ , de indice  $j$ , din imaginea la momentul următor  $k+l$  ( $l$  este pasul de analiză), astfel  $b_i(k) \approx b_{i,j}(k+l)$ .

Similaritatea între blocurile de pixeli se traduce prin minimizarea unei funcții de cost,  $F_c()$ , ce poate fi distanța Euclidiană între pixeli, eroarea pătratica medie, eroarea absolută, etc. Astfel, eroarea minimală este dată de ecuația următoare:

$$\tilde{D}_{k,k+l}(i) = \min|_{j=1,\dots,N_{cand}} F_c(b_i(k), b_{i,j}(k+l)) \quad (2.10)$$

unde  $\tilde{D}_{k,k+l}(i)$  reprezintă valoarea minimă a funcției de cost,  $N_{cand}$  reprezintă numărul de blocuri  $b_{i,j}(k+l)$  din imaginea la momentul  $k+l$  ce sunt "andidate" pentru a fi similare cu blocul curent,  $b_i(k)$ . Acestea definesc ceea ce se numește fereastră de căutare.

Dacă imaginile la momentele  $k$  și  $k+l$  sunt imagini vecine în interiorul unui plan video, atunci valorile funcției de cost sunt mici sau chiar apropiate de zero. Pe de altă parte, dacă imaginea la momentul  $k$  este imaginea ce precede un "cut", valorile funcției de cost vor fi importante. Aceasta se datorează faptului că între cele două imagini există o diferență vizuală semnificativă, iar blocurile de pixeli ale imaginii  $k$  nu vor fi regăsite în imaginea  $k+l$ .

De exemplu, în [Shahraray 95] imaginile sunt divizate în 12 blocuri disjuncte de pixeli iar compensarea mișcării este realizată folosind ca funcție de cost diferența între intensitățile pixelilor. Valorile  $\tilde{D}_{k,k+l}()$  astfel obținute sunt ordonate și normalizate între 0 și 1 fiind notate cu  $d_{k,k+l}^s()$ . Măsura de

discontinuitate între imaginile la momentele  $k$  și  $k + l$  este calculată pe baza unui anumit set de ponderi,  $c_i$ , în felul următor:

$$DSC(k, k + l) = \sum_{i=1}^{12} c_i \cdot d_{k,k+l}^s(i) \quad (2.11)$$

Mai departe, tranzițiile de tip "cut" sunt detectate prin filtrarea valorilor lui  $DSC(k, k + l)$  cu un prag predefinit.

O altă abordare similară este metoda propusă în [Porter 00] ce folosește ca funcție de cost corelația dintre blocurile de pixeli, calculată de această dată în domeniul frecvențial.

Metoda propusă în [Hanjalic 02] folosește pentru detecție informații precum statistică duratei tranzițiilor, compensarea mișcării și amplitudinea diferențelor temporale ce survin în secvență. Alte abordări cu o complexitate de calcul mai ridicată, folosesc estimarea fluxului optic iar măsurile de similaritate între imagini sunt calculate pe baza vectorilor de mișcare și a deplasărilor survenite în imagine [Zhong 96] [Lupatini 98].

Din punct de vedere global, metodele de detecție a tranzițiilor de tip "cut" bazate pe analiza mișcării sunt mai puțin eficiente și precise decât metodele bazate pe histogramă [Gargi 00]. Estimarea mișcării este o operație cu o complexitate de calcul ce poate fi importantă, fiind mult mai complexă decât evaluarea unor histograme color [Lienhart 01b]. Cu toate acestea, metodele bazate pe mișcare rămân o soluție pentru secvențele ce au un conținut preponderent dinamic, caz în care metodele bazate pe contur și respectiv pe analiza intensității pixelilor, tind să obțină un număr ridicat de false detecții, datorat în principal diferențelor succesive dintre imaginile în mișcare.

### Metode de detecție în domeniul comprimat

Metodele de detecție din această categorie exploatează direct informația din domeniul comprimat al fluxului video MPEG, ca de exemplu coeficienții transformatei Cosinus Discrete (DCT - "Discrete Cosine Transform"<sup>6</sup>).

De exemplu, [Fernando 01] propune detecția tranzițiilor de tip "cut" în secvențele codate MPEG-2. Metoda propusă folosește analiza numărului de predicții ale macro-blocurilor în cadrele de tip B (compresie bidirecțională ce folosește informație atât din imaginile anterioare, cât și din cele ulterioare imaginii curente).

---

<sup>6</sup>"Discrete Cosine Transform" sau DCT reprezintă transformata cosinus discretă. Aceasta permite reprezentarea unei mulțimi finite de valori într-un spațiu de reprezentare frecvențial, sub forma unei sume de funcții cosinus ce oscilează pe diferite frecvențe.

Metoda propusă în [Arman 93b] folosește pentru detecție sub-blocuri de pixeli, de dimensiune  $8 \times 8$ , codate DCT, ce sunt alese din  $n$  regiuni conexe ale imaginii curente  $k$ . Pentru toate aceste blocuri sunt reținuți, în mod aleator, doar 64 din coeficienții DCT ai componentei alternative. Astfel, fiecare imagine va fi reprezentată în domeniul comprimat de un vector de coeficienți,  $V_k = (c_1, c_2, \dots, c_{64})$ , unde  $c_i$  reprezintă coeficientul de indice  $i$ ,  $i = 1, \dots, 64$ . Similaritatea imaginilor la momentele  $k$  și  $k + l$  este calculată pe baza produsului scalar normalizat al vectorilor  $V_k$ , dat de relația:

$$\Psi_{k,k+l} = \frac{V_k \cdot V_{k+l}}{|V_k| \cdot |V_{k+l}|} \quad (2.12)$$

Un "cut" este detectat în momentul în care  $1 - |\Psi| > T$ , unde  $T$  reprezintă pragul de discontinuitate. Metode similare ce folosesc fluxul MPEG sunt propuse și în [Zhang 94] sau [Meng 95].

Metodele de detecție ce folosesc direct informațiile furnizate de fluxul MPEG nu necesită o etapă prealabilă de decompresie a datelor, etapă ce este necesară tuturor celorlalte metode ce folosesc o analiză la nivel de imagine. Coeficienții fluxului MPEG conțin suficientă informație pentru a detecta discontinuitățile produse de tranzițiile de tip "cut". Complexitatea de calcul în acest caz este foarte redusă, lucru ce facilitează implementarea în timp real<sup>7</sup> a algoritmilor de detecție. Totuși, deseori precizia metodelor bazate pe fluxul MPEG este inferioară celorlalte metode datorită incoerenței vectorilor de mișcare (vezi Secțiunea 3.1.5). Soluția de compromis constă în decompresia datelor până la un anumit nivel de detaliu.

Ca tendință generală, este posibil ca noile standarde de codare, precum standardul de compresie video MPEG-7, să înglobeze informații referitoare la structura temporală a secvenței, precum distribuția de plane, scene, etc., lucru ce va face ca etapa de segmentare temporală să nu mai fie necesară [Wang 00].

### Alte metode

În această categorie putem încadra metodele ce folosesc pentru detecție alte surse de informație decât cele enumerate anterior. De exemplu, metoda propusă în [Boreczky 98] transformă problema detecției într-o problemă de clasificare. Astfel, aceasta propune segmentarea temporală în plane video pe baza modelelor Markov ascunse (HMM - "Hidden Markov Models"). Diferitele stări ale HMM sunt folosite pentru a modeliza diferențele tipuri de segmente ale secvenței. Pentru mai multe detalii referitoare la folosirea modelelor Markov

---

<sup>7</sup>vezi explicația de la pagina 102.

în segmentarea temporală, cititorul se poate raporta la studiul [Wang 00]. O altă metodă ce tratează problematica segmentării, generic, independent de conținutul secvenței este propusă în [Hanjalic 02]. Aceasta este o abordare statistică ce folosește pentru segmentare minimizarea probabilității erorii medii de detecție.

O abordare inedită a problemei detecției tranzițiilor de tip "cut" este propusă în [Guimaraes 03]. În prima fază, fiecare imagine a secvenței este rezumată cu o singură linie de pixeli, și anume, diagonala principală a imaginii. Liniile de pixeli astfel obținute pentru întreaga secvență sunt juxtapuse pentru a forma o singură imagine, numită și ritm vizual al secvenței. Tranzițiile de tip "cut" apar în ritmul vizual sub formă de tranziții verticale. Pentru detecția acestora sunt folosite metode de detectie de contur și de morfologie matematică.

Aceste alte tipuri de abordări, relativ particulare sau inedite, au fost în general testate și aplicate în cazuri particulare. Din această cauză, nu dispunem de suficiente teste comparative pentru a trage o concluzie relativă la performanța acestora în comparație cu celelalte metode existente. Totuși, prin aceste abordări se încearcă valorificarea și a altor surse de informație decât cele clasice: imagine-contur-mișcare.

### **Problematica estimării pragului de detectie**

Marea majoritate a metodelor de măsurare a discontinuităților produse de tranzițiile de tip "cut" folosesc unul sau mai multe praguri pentru detectie. Noțiunea de *similaritate* între două imagini se rezumă în final la compararea unei măsuri de distanță cu un anumit prag, calculat sau fixat "a priori". Dacă valoarea acesteia depășește valoarea pragului, atunci vorbim despre *imagini disimilare*, în caz contrar, imaginile sunt considerate ca fiind similare. Această operație este cunoscută în literatura de specialitate sub numele de "thresholding".

Astfel, alegerea adecvată a pragurilor este esențială pentru precizia detectiei. După cum am menționat la începutul acestui capitol, un prag ales prea mic va avea ca efect creșterea numărului de false detectii, iar un prag prea ridicat va conduce la un număr mare de tranziții ce vor trece nedetectate. Pentru un studiu bibliografic complet relativ la metodele de estimare a pragului de detectie, cititorul se poate raporta la lucrările [Lienhart 01b] și [Hanjalic 02].

Primele abordări de estimare a pragurilor de detectie erau bazate pe *metode euristice*, ca cele propuse în: [Otsuji 91], [Nagasaki 92], [Arman 93b]. Astfel, pragurile erau alese "a priori" ca rezultat al expertizei manuale a datelor folosite. Alte abordări propuneau o *analiză statistică* a distribuției valo-

rilor funcției de discontinuitate folosită la detecție. De exemplu, în [Zhang 93] este propusă modelizarea acestei distribuții cu funcții Gaussiene de medie  $\mu$  și varianță  $\sigma^2$ . Astfel, pragul de detecție,  $T$ , era definit ca fiind:

$$T = \mu + r \cdot \sigma \quad (2.13)$$

unde parametrul  $r$  corespunde unei probabilități de falsă detecție ce este fixată "a priori".

Un alt exemplu este abordarea propusă în [Ionescu 06a]. Aceasta pornește de la observația că discontinuitatea produsă de tranzițiile de tip "cut" este mai puțin frecventă decât alte discontinuități produse de variații de culoare sau de mișcare. Astfel, în această ipoteză, pragul  $T$  definit în ecuația 2.13 va avea o valoare prea scăzută ce va duce la un număr important de false detecții. Pentru a corecta acest lucru, metoda propusă calculează pragul în două etape. În prima etapă sunt selectate maximele locale ale funcției de discontinuitate ce sunt superioare mediei globale a acesteia. Pragul de detecție este calculat în a doua etapă ca fiind valoarea medie a acestor maxime pentru întreaga secvență, asigurând astfel o valoare apropiată de valoarea optimă.

Aceste două tipuri de abordări, euristică și statistică, au ca rezultat determinarea unui prag global de detecție ce este același pentru întreaga secvență.

O altă metodă constă în calcularea pragurilor în mod *adaptiv*. [Yeo 95] propune calcularea pragului de detecție,  $T$ , în funcție de informația temporală a secvenței. Sevența este împărțită în ferestre temporale de analiză de dimensiune  $N$ . Astfel, un "cut" este detectat în mijlocul ferestrei curente de analiză, dacă funcția de discontinuitate,  $D(k, k+1)$ , calculată între imaginile la momentele  $k$  și respectiv  $k+1$ , îndeplinește următoarele condiții:

$$D(k, k+1) = \max|_{i=-\frac{N}{2}, \dots, \frac{N}{2}} \{D(k+i, k+1+i)\} \quad (2.14)$$

$$D(k, k+1) \geq \alpha \cdot D_{smax} \quad (2.15)$$

unde  $D_{smax}$  reprezintă a doua valoare maximală a funcției de discontinuitate obținută pentru fereastra temporală curentă de dimensiune  $N$ , iar  $\alpha$  este un parametru determinat în funcție de forma funcției de discontinuitate. Pentru mai multe detalii referitor la metodele de calcul adaptiv al pragului de detecție, cititorul se poate raporta la [Gargi 00] și [Truong 00b].

O altă abordare a problemei estimării pragului de detecție sunt *metodele mixte* ce combină metodele adaptive cu abordările statistice. De exemplu, [Hanjalic 97] propune modelarea cu funcții Gausiene a distribuției valorilor de discontinuitate în ferestre temporale de analiză. Parametrul  $\alpha$ , din ecuația 2.15, este determinat pe baza analizei probabilității "a priori" de falsă detecție propusă în [Zhang 93].

Calculul *pragului optimal* pentru detecție reprezintă o altă direcție de studiu. Metodele din această categorie sunt inspirate din teoria detecției statistice. Acestea folosesc informații statistice despre distribuția tranzițiilor de tip "cut" obținute în urma expertizei manuale a unui număr semnificativ de secvențe de imagini. Regula de detecție a discontinuității unui "cut" este calculată prin minimizarea erorii de detecție [Vasconcelos 00].

### 2.2.2 Detectia de "fades"

După cum am menționat în partea introductivă a acestui capitol, o tranzitie de tip "fade" este o tranzitie video graduală ce corespunde efectului optic prin care imaginea de interes apare progresiv dintr-o imagine constantă, de regulă neagră. Acest efect este numit și "fade-in". Procesul invers, respectiv de dispariție progresivă a imaginii curente către un fond constant, este numit "fade-out" (vezi Figura 2.2). Deseori, cele două tipuri de "fade" sunt folosite împreună, unul după altul, pentru a forma o secvență de tip "fade-out" - "fade-in". În acest caz, cele două tranzitii se comportă ca o singură tranzitie globală numită și grup de "fade".

Secvența de imagini ce constituie o tranzitie de tip "fade" de durată  $T$ , notată  $F(x, y, t)$ , în care  $(x, y)$  reprezintă spațiul imaginii iar  $t$  dimensiunea temporală, este definită matematic ca fiind transformarea intensităților pixelilor secvenței  $S_1(x, y, t)$  printr-o funcție monotonă  $f(t)$  [Lienhart 01b]. Astfel, aceasta este dată de relația:

$$F(x, y, t) = f(t) \cdot S_1(x, y, t) \quad (2.16)$$

unde  $0 \leq t \leq T$ .

Tranzitiiile de tip "fade-in" folosesc funcții monotone,  $f(t)$ , cu proprietatea că  $f(0) = 0$  și  $f(T) = 1$ . Similar, pentru o tranzitie de tip "fade-out", funcția  $f(t)$  trebuie să îndeplinească condițiile  $f(0) = 1$  și respectiv  $f(T) = 0$ . În cele mai multe cazuri, funcția  $f(t)$  este aleasă ca fiind liniară și este definită în felul următor:

$$f_{fade-in}(t) = \frac{t}{T} \quad (2.17)$$

$$f_{fade-out}(t) = 1 - \frac{t}{T} \quad (2.18)$$

Comparate cu varietatea de metode disponibile pentru detecția tranzitiiilor abrupte de tip "cut", metodele existente pentru detecția de "fade" sunt mai puțin numeroase. Acest lucru se datorează pe de-o parte complexității algoritmilor de detecție precum și a faptului că frecvența acestora în secvență este net inferioară tranzitiiilor de tip "cut". Mai mult, tranzitiiile graduale

pot fi approximate fără pierderi semnificative pentru segmentarea temporală, cu tranzitii abrupte, pe când reciproca nu este valabilă.

Metodele de detectie de "fade" existente sunt orientate către trei axe principale de studiu, și anume:

- metode bazate pe *analiza intensității pixelilor*,
- metode bazate pe *analiza contururilor*,
- alte metode ce folosesc *alte surse de informație* decât cele clasice enumerate anterior.

Pentru un studiu bibliografic complet al literaturii de specialitate, cititorul se poate raporta la lucrările [Lienhart 01b], [Hanjalic 02] sau [Ren 03]. În cele ce urmează, ne vom limita la o prezentare succintă a particularităților metodelor din fiecare categorie.

### Metode bazate pe analiza intensității pixelilor

Una dintre primele metode de detectie a fost propusă în [Zhang 93] și se baza pe folosirea a două praguri. Aceasta modifică o metodă de detectie a tranzitiei de tip "cut" bazată pe analiza distanței între histograme. Astfel, un "cut" era detectat dacă disimilaritatea între două imagini succesive era superioară unui prag  $\tau_{start}$ . Dacă la un moment  $k$ , disimilaritatea între două imagini succesive era superioară unui alt prag,  $\tau_{cand}$  ( $\tau_{cand} < \tau_{start}$ ), dar totodată inferioară pragului inițial  $\tau_{start}$ , atunci imaginea  $k$  era considerată ca o potențială imagine de început a unei tranzitii graduale. Mai departe, această imagine era comparată cu imaginile următoare, iar distanța între acestea acumulată. În momentul în care cumulul distanțelor depășea pragul  $\tau_{start}$ , dar în același timp diferențele individuale dintre imaginile succesive rămâneau inferioare pragului  $\tau_{cand}$ , atunci era detectată o tranzitie de tip "fade".

Pe durata unei tranzitii de tip "fade", una dintre informațiile caracteristice acesteia o reprezintă schimbarea intensității luminoase a imaginii. Astfel, alte abordări studiază evoluția dispersiei intensității luminoase globale a imaginii. În condiții de ergodicitate<sup>8</sup>, dispersia secvenței unei tranzitii de tip "fade" poate fi exprimată ca:

$$\sigma(F(x, y, t)) = f(t) \cdot \sigma(S_1(x, y)) \quad (2.19)$$

---

<sup>8</sup>un proces aleator este considerat ca fiind ergodic dacă momentele statistice sunt egale cu momentele temporale.

unde  $F(\cdot)$ ,  $f(\cdot)$  și  $S_1(\cdot)$  au aceeași semnificație ca în ecuația 2.16 iar  $\sigma(\cdot)$  reprezintă dispersia. Se poate observa că dispersia secvenței unui "fade" respectă monotonia funcției  $f(t)$ .

Metoda de detecție propusă în [Lienhart 99a] localizează mai întâi în secvență imaginile monocromatice pentru care valoarea varianței intensității luminoase este apropiată de zero. Aceste imagini particulare sunt posibile imagini de început sau finale ale unor tranziții de tip "fade-in" și respectiv "fade-out". Mai departe, detecția este realizată prin analiza creșterii progressive a intensității luminoase și a dispersiei acesteia în sensul pozitiv al axei temporale. Liniaritatea este verificată prin evaluarea erorii de aproximare a valorilor acestora cu dreptele de regresie ce le modelează.

O abordare similară, bazată pe analiza varianței intensității luminoase, este propusă în [Alattar 97]. Aceasta propune o detecție preliminară pe baza analizei punctelor de extrem negativ ale derivatei a două a valorilor varianței intensității luminoase a imaginii. Detecția unui "fade" este ulterior confirmată dacă valoarea derivatei întâi a intensității luminoase medii este constantă între două puncte de extrem negativ.

Metoda propusă în [Truong 00a] vine să combine metodele propuse în [Lienhart 99a] și [Alattar 97]. Astfel, în primă fază a detecției vor fi selectate imaginile monocromatice. Dintre acestea sunt reținute doar imaginile ce furnizează valori apropiate de extremele negative ale derivatei secunde a varianței intensității luminoase din imagine. Un "fade" este mai departe detectat dacă sunt satisfăcute criteriile următoare:

- derivata întâi a valorilor intensității medii rămâne constantă și nu își schimbă semnul,
- valoarea medie a pantei derivatei întâi trebuie să fie superioară unui anumit prag,
- valoarea varianței intensității luminoase pentru prima și ultima imagine a tranziției trebuie de asemenea să fie superioare unui anumit prag.

Una dintre principalele surse de erori în cazul metodelor de detecție de "fade" este prezența mișcării. Aceasta face ca intensitatea luminoasă să aibă variații neliniare ce nu mai respectă monotonia funcției  $f(t)$  (vezi ecuația 2.16). În [Fernando 99], pentru a reduce influența mișcării, detecția este realizată folosind atât măsuri statistice ale intensității pixelilor, cât și măsuri statistice ale semnalului de crominanță. Analiza este efectuată în spațiul de culoare YCbCr (Y-luminanță și Cb, Cr-diferențe cromatice). Media semnalului de crominanță,  $C = \frac{C_b + C_r}{2}$ , se dovedește a fi mai puțin sensibilă la prezența mișcării decât media semnalului de luminanță  $Y$ . Aceste două

informații sunt folosite pentru a defini parametrul  $R(k)$  ce reprezintă raportul de schimbare incrementală a mediei semnalului de luminanță relativ la semnalul de crominanță la momentul  $k$ :

$$R(k) = \begin{cases} \frac{\Delta_k^Y}{\Delta_k^C} & \text{dacă } k < L_1 \text{ sau } k \geq (L_1 + T) \\ \frac{|C_0 - m_{k+1}^Y + (L_1 - k) \cdot \Delta_k^Y|}{|C_0 - m_{k+1}^C + (L_1 - k) \cdot \Delta_k^C|} & \text{dacă } L_1 \leq k < (L_1 + T) \end{cases} \quad (2.20)$$

unde  $\Delta_k^Y$  și  $\Delta_k^C$  reprezintă schimbările incrementale ale mediei semnalului  $Y$  și respectiv ale semnalului  $C$ ,  $m_{k+1}^Y$  și  $m_{k+1}^C$  sunt valorile medii ale lui  $Y$  și  $C$  calculate la momentul  $k+1$ ,  $L_1$  reprezintă momentul de început al tranziției,  $T$  este durata totală a acesteia, iar  $C_0$  reprezintă nivelul semnalului video pentru începutul tranzиiei.

Detectia este realizată mai departe pe baza analizei valorilor lui  $R(k)$ . Pe durata unui "fade", acestea trebuie să rămână aproximativ constante, ceea ce face ca valorile diferenței  $|R(k) - R(k - 1)|$  să fie apropiate de zero.

O abordare diferită, cu o complexitate de calcul mai redusă, a principiului metodei din [Fernando 99], este propusă în [Ionescu 05a] unde detectia este realizată pe baza analizei colaborative a informației de culoare și a intensității luminoase. Astfel, detectia este realizată prin analiza evoluției temporale a trei parametrii, și anume:

- valoarea medie a componentei de intensitate luminoasă a imaginii,  $\bar{Y}$ , calculată în spațiul de culoare  $YCbCr$ ,
- varianța globală în imagine a componentei  $Y$ ,  $\sigma^2(Y)$ ,
- valoarea absolută a diferenței dintre valorile medii pe imagine ale componentelor  $Cb$  și  $Cr$ ,  $|\bar{Cb} - \bar{Cr}|$ .

Varianța  $\sigma^2(Y)$  este folosită pentru a detecta începutul unei tranzиii de tip "fade-in" și respectiv sfârșitul unei tranzиii de tip "fade-out", deoarece în acest caz, varianța are valori apropiate de zero (fiind vorba de imagini constante). Detectia propriu-zisă este realizată prin localizarea evoluției crescătoare ("fade-in") și respectiv descrescătoare ("fade-out") a parametrilor  $\bar{Y}$  și  $|\bar{Cb} - \bar{Cr}|$ . Parametrul  $|\bar{Cb} - \bar{Cr}|$  se dovedește a fi mai puțin sensibil la prezența mișcării decat  $\bar{Y}$ , care la rândul său se dovedește a fi mai eficient în absența mișcării, astfel că cei doi parametri sunt folosiți disjunctiv.

### Metode bazate pe analiza de contur

O altă direcție de studiu este analiza bazată pe informația de contur. Metodele din această categorie se folosesc de ipoteza că pe durata unei tranzиii

de tip "fade", contururile obiectelor din imagine fie dispar ("fade-out"), fie apar gradual ("fade-in").

O măsură cantitativă a schimbării de contur este dată de raportul ECR ("Edge Change Ratio") propus în [Zabih 95], măsură folosită și la detecția tranzițiilor de tip "cut" (vezi Secțiunea 2.2.1). Pe durata unei tranziții de tip "fade-in" numărul de pixeli de contur ce apar în imagine, notat  $ECR_{in}$ , este superior numărului de pixeli de contur ce dispar, notat  $ECR_{out}$ , astfel  $ECR_{in} > ECR_{out}$ . Similar, pentru un "fade-out" întâlnim situația inversă, și anume  $ECR_{in} < ECR_{out}$ .

Dacă tranzițiile de tip "cut" corespundeau valorilor maxime ale funcției ECR, tranzиile de tip "fade" precum și alte tranzиii graduale sunt caracterizate de intervale de valori semnificative ale ECR [Zabih 99]. Marea majoritate a metodelor de detecție existente bazate pe analiza de contur folosesc abordări similare, ce se bazează pe contabilizarea punctelor de contur ce apar sau dispar din scenă [Yu 97] [Lupatini 98].

Global, metodele bazate pe analiza conturului sunt mai puțin eficiente decât cele bazate pe analiza intensității pixelilor [Lienhart 01b]. Acest lucru se datorează în mare parte sensibilității ridicate a detecției la prezența mișcării sau a diverselor efecte de culoare. Acestea conduc la schimbarea substanțială a distribuției conturilor din imagine și astfel la eșuarea detecției.

### Alte metode

În această categorie putem menționa mai întâi *abordările mixte*, ce fuzionează diferitele surse de informație disponibile, precum metodele ce folosesc informații statistice ale intensității pixelilor, dar care realizează analiza în domeniul comprimat al coeficienților DCT [Bimbo 99]. Pe de altă parte putem menționa abordările inovantive ce încearcă exploatarea de noi surse de informație pentru a ameliora punctele slabe ale metodelor clasice existente. Astfel putem menționa:

- utilizarea informației de mișcare pentru a reduce influența acesteia și astfel numărul de false detecții [Porter 01],
- utilizarea ritmului vizual bazat pe histogramă pentru a reduce influența zgomotului din imagine [Guimaraes 03],
- exploatarea de parametri calculați în domeniul frecvențial al coeficienților FFT ("Fast Fourier Transform") ai imaginii [Miene 01],
- utilizarea de abordări statistice pentru detecția generică a tranzиilor graduale [Heng 01].

În general, fiecare metodă propusă prezintă o serie de avantaje cât și de inconveniente, care în anumite situații vor permite ameliorarea detecției raportat la metodele clasice, dar vor da greș în altele. De exemplu, metoda propusă în [Guimaraes 03], bazată pe ritmul vizual, detectează poziția în secvență cât și durata tranzițiilor de tip "fade" cu o precizie de o imagine, pe când în majoritatea abordărilor clasice, respectarea duratei tranziției nu este un criteriu de calitate. Pe de altă parte, metoda propusă eșuează în cazul în care tranzițiile prezintă diferențe de luminozitate între imaginea de început și cea de final.

### 2.2.3 Detectia de "dissolves"

Un alt tip de tranziție graduală des întâlnită în secvențele de imagini sunt tranzițiile de tip "dissolve". Un "dissolve" reprezintă efectul obținut prin transformarea graduală progresivă a unei imagini în alta. Transformarea este realizată la nivel de intensitate a pixelilor și nu prin transformări geometrice (vezi Figura 2.2).

Din punct de vedere matematic, secvența unei tranziții de tip "dissolve",  $D(x, y, t)$ , de durată  $T$ , unde  $(x, y)$  reprezintă spațiul imaginii iar  $t$  este dimensiunea temporală, este definită pe baza secvențelor  $S_1(x, y, t)$  și  $S_2(x, y, t)$  în felul următor:

$$D(x, y, t) = f_1(t) \cdot S_1(x, y, t) + f_2(t) \cdot S_2(x, y, t) \quad (2.21)$$

unde  $0 \leq t \leq T$ . În funcție de forma funcțiilor  $f_1(t)$  și  $f_2(t)$  folosite, întâlnim mai multe tipuri de "dissolve" [Lienhart 01b].

Cel mai frecvent folosite sunt tranzițiile de tip "cross-dissolve" (vezi Figura 2.4). Acestea sunt construite ca suprapunerea unei tranziții de tip "fade-out" cu o tranziție de tip "fade-in", astfel funcțiile  $f_1(t)$  și  $f_2(t)$  fiind cele utilizate de tranzițiile de tip "fade":

$$f_1(t) = 1 - \frac{t}{T}, \quad f_2(t) = \frac{t}{T} \quad (2.22)$$

unde  $0 \leq t \leq T$ .

O altă categorie de "dissolve" mai puțin utilizată sunt tranzițiile "dissolve" aditive ("additive dissolve", vezi Figura 2.4). Acestea pot fi văzute ca fiind suma unei tranziții de tip "fade-out" cu o tranziție de tip "fade-in". În acest caz, cele două funcții  $f_1(t)$  și  $f_2(t)$  sunt definite astfel:

$$f_1(t) = \begin{cases} 1 & \text{dacă } t \leq c_1 \\ \frac{T-t}{T-c_1} & \text{altfel} \end{cases} \quad (2.23)$$

$$f_2(t) = \begin{cases} \frac{t}{c_2} & \text{dacă } t \leq c_2 \\ 1 & \text{altfel} \end{cases} \quad (2.24)$$

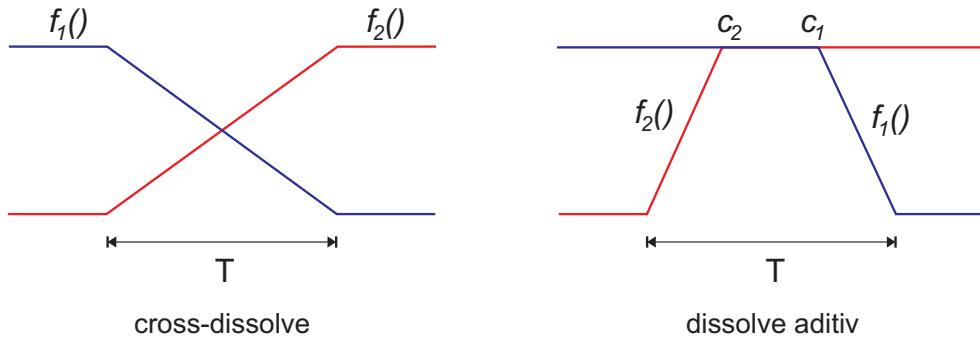


Figura 2.4: Funcțiile de scalare folosite de tranzițiile de tip ”dissolve” (axa orizontală corespunde axei temporale).

unde  $c_1, c_2 \in (0; T)$ ,  $c_2 < c_1$  și  $0 \leq t \leq T$ .

Din punct de vedere al diferenței vizuale dintre secvențele  $S_1(x, y, t)$  și respectiv  $S_2(x, y, t)$ , ce constituie tranzitia ”dissolve”, acestea se împart în trei categorii [Lienhart 01b], astfel:

- secvențele  $S_1$  și  $S_2$  au *distribuții de culoare suficient de diferite* astfel încât tranzitia ”dissolve” rezultată să fie confundată cu una sau mai multe tranzitii de tip ”cut”,
- secvențele  $S_1$  și  $S_2$  au *distribuții de culoare similară* ce nu sunt detectabile folosind metode de detecție a discontinuității vizuale bazate pe histogramă. Pe de altă parte, diferențele la nivel structural ale obiectelor prezente în scenă sunt importante, fiind detectabile cu metode bazate pe analiza contururilor,
- secvențele  $S_1$  și  $S_2$  au *distribuția de culoare și structura obiectelor din scenă similară*. Acest tip de ”dissolve” este de fapt un caz particular al unui efect de ”morphing”<sup>9</sup>.

Similar cu metodele de detecție a tranzitiei de tip ”fade”, metodele existente de detecție a tranzitiei de tip ”dissolve” se împart în trei categorii [Lienhart 01b]:

- metode bazate pe *analiza intensității pixelilor* din imagine,
- metode bazate pe *analiza contururilor*,
- alte metode ce folosesc *alte surse de informație* decât cele menționate anterior.

<sup>9</sup>vezi explicația de la pagina 32.

Pentru un studiu bibliografic complet al literaturii de specialitate, cititorul se poate raporta la lucrările [Bimbo 99] [Lienhart 01b], [Hanjalic 02], [Ren 03] sau [Su 05a]. În cele ce urmează, vom face o trecere în revistă a tehnicielor de detectie reprezentative pentru fiecare categorie de metode.

### **Metode bazate pe analiza intensității pixelilor**

Dacă pe durata unei tranziții de tip ”cross-dissolve”, mișcarea obiectelor sau a camerei video este neglijabilă, atunci pornind de la ecuația 2.19 putem obține următoarea relație:

$$\frac{\partial D(x, y, t)}{\partial t} = \frac{S_2(x, y) - S_1(x, y)}{T} \quad (2.25)$$

unde  $0 \leq t \leq T$  iar  $T$  reprezintă durata tranziției. Astfel, o posibilă metodă de detectie o reprezintă localizarea în secvență a tuturor schimbărilor liniare ale intensității pixelilor, metodă propusă în [Hampapur 95]. Această abordare, conform chiar ipotezei de plecare, este foarte sensibilă la prezența zgomotului în imagine cât și la prezența mișcării.

Metoda propusă în [Gu 97] vine cu o serie de îmbunătățiri menite să amelioreze invarianța detecției. Detectia este realizată de această dată în domeniul comprimat al fluxului MPEG pe baza coeficientelor DCT ai informației de luminanță. Pentru imagini succesive, este calculat procentul de blocuri de pixeli pentru care diferența absolută între coeficientii DCT are o evoluție specifică unei tranziții de tip ”dissolve”. Un ”dissolve” va fi detectat în cazul în care valorile astfel obținute rămân superioare unui anumit prag pe durata a 10 până la 60 de imagini (durata medie maximală a unui ”dissolve” calculată la o frecvență de cadre de 30 imagini/s).

O abordare ce nu folosește direct modelul matematic al unui ”dissolve” o constituie analiza comportamentului temporal al pixelilor din imagine într-un anumit spațiu de caracteristici. De exemplu, [Nam 00] aproximează evoluția temporală a fiecărui coeficient DCT într-o fereastră de durată  $L$ , cu funcții ”B-spline”<sup>10</sup>. Dacă pe durata unui ”dissolve” valoarea dispersiei funcției temporale de evoluție a intensității pixelilor prezintă valori semnificative, eroarea de aproximare a evoluției coeficientilor DCT cu funcții ”B-spline”, trebuie să prezinte valori reduse. Aceasta se datorează efectului ”sintetic” pe care îl are tranziția ”dissolve” în evoluția temporală a secvenței. [Nam 00]

---

<sup>10</sup>În analiza numerică, o funcție ”spline” este o funcție polinomială pe porțiuni folosită la interpolarea datelor. ”B-spline” este un caz particular, fiind o funcție ”spline” care pentru anumite valori date ale gradului polinomial, gradului de uniformitate și ale domeniului de definiție, oferă un suport minimal al funcției.

definește eroarea de aproximare ca fiind:

$$e(x, y, t) = \frac{1}{L+1} \sum_{i=t-L/2}^{t+L/2} [D(x, y, i) - D_{\text{spline}}(x, y, i)]^2 \quad (2.26)$$

unde  $D_{\text{spline}}(x, y, t)$  este aproximarea de tip "B-spline" iar  $L = 31$ . Comportamentul acestei funcții va fi diferit pentru diversele tipuri de pasaje ale secvenței. Astfel, pe durata unui "dissolve",  $e(x, y, t)$  prezintă valori mici dar varianța dintre imagini este semnificativă. Pe de altă parte, un pasaj ce conține mișcări de obiecte sau mișcări ale camerei video, va conduce atât la o valoare semnificativă a varianței cât și a erorii.

Dacă presupunem că cele două secvențe,  $S_1(x, y, t)$  și  $S_2(x, y, t)$ , ce constituie tranzitia de tip "dissolve" sunt procese aleatoare ergodice statistic independente, atunci varianța secvenței tranzitiei poate fi exprimată astfel:

$$\text{Var}\{D(x, y, t)\} = f_1^2(t) \cdot \text{Var}\{S_1(x, y)\} + f_2^2(t) \cdot \text{Var}\{S_2(x, y)\} \quad (2.27)$$

unde varianțele secvențelor  $S_1$  și  $S_2$  sunt independente de timp.

În cazul unui "cross-dissolve" (vezi ecuația 2.22), ecuația anterioară devine:

$$\text{Var}\{D(x, y, t)\} = \frac{(T-t)^2}{T^2} \cdot \text{Var}\{S_1(x, y)\} + \frac{t^2}{T^2} \cdot \text{Var}\{S_2(x, y)\} \quad (2.28)$$

și mai departe:

$$\text{Var}\{D(x, y, t)\} = c \cdot (t-a)^2 - b \quad (2.29)$$

unde

$$a = \frac{T \cdot \text{Var}\{S_1(x, y)\}}{\text{Var}\{S_1(x, y)\} + \text{Var}\{S_2(x, y)\}} \quad (2.30)$$

$$b = \frac{\text{Var}\{S_1(x, y)\} \cdot \text{Var}\{S_2(x, y)\}}{\text{Var}\{S_1(x, y)\} + \text{Var}\{S_2(x, y)\}} \quad (2.31)$$

$$c = \frac{\text{Var}\{S_1(x, y)\} + \text{Var}\{S_2(x, y)\}}{T^2} \quad (2.32)$$

Calculând mai departe derivata întâi și secundă a varianței vom obține relațiile următoare:

$$\frac{\partial \text{Var}\{D(x, y, t)\}}{\partial t} = 2 \cdot c \cdot (t-a) \quad (2.33)$$

$$\frac{\partial^2 \text{Var}\{D(x, y, t)\}}{\partial t^2} = 2 \cdot c \quad (2.34)$$

Astfel, se observă că pe parcursul unei tranzitii de tip ”dissolve”, evoluția temporală a varianței intensității pixelilor are un comportament parabolic. Acest lucru are ca implicație valori apropiate de zero ale derivatei secunde, calculată înainte și după ”dissolve”, și valori constante și pozitive pe durata tranzitiei.

Această ipoteză a fost folosită pentru prima oară în metoda propusă în [Alattar 93] ce exploata prezența a două puncte de extrem negativ în derivata secundă, puncte ce corespundeau începutului și respectiv sfârșitului tranzitiei. Alte exemple sunt metodele propuse în [Fernando 99], [Gu 97] sau [Truong 00b]. De exemplu, metoda din [Truong 00b] folosește pentru detecție următoarele considerații:

- derivata întâi a varianței intensității pixelilor pe parcursul unui ”cross-dissolve” trebuie să fie o funcție monoton crescătoare ce pornește cu o valoare negativă și se oprește într-o valoare pozitivă,
- varianța intensității pixelilor pentru cele două secvențe  $S_1$  și  $S_2$  trebuie să fie superioară unui anumit prag  $\tau_{min}$ ,
- durata unui ”dissolve” nu trebuie să depășească un anumit interval de valori,  $[T_{min}, T_{max}]$ , determinat experimental.

O abordare diferită este propusă în [Su 05b]. În prima fază, pentru fiecare imagine analizată la momentul  $k$  se constituie o nouă imagine pseudo-binară,  $B_k^L(x, y)$ , unde  $L + 1$  reprezintă durata minimă a unei tranzitii de tip ”dissolve”. Pixelii din imagine ce prezintă o creștere sau respectiv o diminuare progresivă a valorii intensității, în intervalul temporal  $[k - L; k]$ , sunt marcați în imaginea  $B_k^L(x, y)$  cu valoarea 1, fiind pixeli activi. În mod similar, ceilalți pixeli sunt marcați cu 0 fiind pixeli inactivi. Pentru detecție se folosește un indicator al cantitatii de pixeli activi din imaginea pseudo-binară. Dacă acesta depășește valoarea unui anumit prag statistic, atunci detecția este validată.

### Metode bazate pe analiza de contur

O altă categorie de metode o reprezintă metodele ce folosesc pentru detecție informația de contur. Acestea se folosesc de ipoteza conform căreia pe durata unei tranzitii de tip ”dissolve”, contururile obiectelor din imaginea de început a tranzitiei dispar progresiv, în timp ce contururile obiectelor din imaginea finală a tranzitiei apar progresiv. Astfel, contrastul imaginii va avea o evoluție descrescătoare odată cu apropierea de mijlocul tranzitiei.

O măsură cantitativă a gradului de schimbare al contururilor din imagine este raportul ECR definit în [Zabih 95] pentru a servi la detecția tranzitilor

de tip "cut" și "fade" (vezi ecuația 2.9). Unele metode existente de detecție, precum metoda propusă în [Zabih 99], folosesc măsuri similare pentru a contabiliza numărul de pixeli de contur ce apar și respectiv ce dispar din imagine. Totuși aceste abordări sunt foarte sensibile la prezența mișcării ce duce la modificarea contururilor și astfel la creșterea numărului de false detecții.

Metoda din [Lienhart 99a] propune pentru detecție o nouă mărime numită contrastul de contur sau  $EC$  ("Edge-Based Contrast"). Aceasta este o măsură matematică ce pune în evidență punctele de contur ce contrastează legătura între punctele de "contur slabă" și "cele semnificative". Dacă presupunem că  $K(x, y, t)$  reprezintă harta de contururi a imaginii curente analizate la momentul  $t$ , iar pragurile  $\tau_w$  și respectiv  $\tau_s$  definesc punctele de contur slabă ( $K(x, y, t) < \tau_w$ ) și respectiv semnificative ( $K(x, y, t) > \tau_s$ ), atunci raportul  $EC$  este dat de relația următoare:

$$EC(K) = 1 + \frac{s(K) - w(K) - 1}{s(K) + w(K) + 1} \quad (2.35)$$

unde  $s(K)$  și  $w(K)$  reprezintă numărul total de puncte de contur semnificative și respectiv slabă prezente în  $K(x, y, t)$ . Astfel, tranzițiile de tip "dissolve" au o semnătură particulară în evoluția temporală a valorilor  $EC$  și anume prezența unui minim local mărginit de variații abrupte ale valorilor  $EC$ .

Global, metodele bazate pe analiza de contur sunt mai puțin eficiente decât abordările ce folosesc analiza intensității pixelilor, acestea fiind mai vulnerabile la prezența mișcării sau a zgomotului în imagine.

### Alte metode

În această categorie putem menționa metoda propusă în [Lienhart 01a] ce vine cu o abordare diferită a problematicii detecției tranzițiilor de tip "dissolve". Astfel, în loc să localizeze tranzițiile pe baza analizei unei funcții de similaritate între cadre, detecția se realizează pe bază de comparații cu modele de "dissolve" predefinite. Pentru aceasta, un sintetizor automat de "dissolve" generează un număr foarte mare de astfel de tranziții folosind ca parametri durata tranziției precum și poziția imaginii centrale. Tranzițiile obținute sunt folosite pentru a antrena în mod iterativ, pe baza metodei "bootstrap"<sup>11</sup>, un clasificator optimal precum rețelele neuronale sau "sup-

---

<sup>11</sup>"bootstrap" este o metodă de inferență statistică. Aceasta estimează proprietățile unui estimator prin măsurarea acestora pe baza mostrelor unei distribuții aproximative ale acestuia.

port vector machines”<sup>12</sup>, ce vor servi ulterior la detectarea tranzițiilor din secvență.

Un alt exemplu este metoda propusă în [Boccignone 00] ce efectuează detectia direct în fluxul comprimat MPEG. Pentru aceasta sunt propuși doi parametri, și anume: parametrul  $D$  ce reprezintă procentul de blocuri de pixeli ce prezintă o diferență importantă între coeficienții DCT, și respectiv parametrul  $\sigma_{mv}$  ce caracterizează gradul de aleatoricitate al vectorilor de mișcare. Detectia este realizată prin analiza evoluției temporale a celor două mărimi. Un ”dissolve” este caracterizat de valori ale lui  $D$  și  $\sigma_{mv}$  ce depășesc un anumit prag determinat în mod automat.

Alte abordări ale detectiei de ”dissolves” se folosesc de modele Markov ascunse [Boreczky 98], de analiza similarității cadrelor în domeniul frecvențial al coeficienților FFT [Miene 01] sau de estimarea de mișcare cu metode bazate pe blocuri de pixeli [Porter 01].

#### 2.2.4 Evaluarea detectiei tranzițiilor video

Am vorbit în capitolele anterioare de metodele de detectie a tranzițiilor video, precum și de performanțele acestora. La complexitatea procesului de detectie al algoritmilor existenți se adaugă o problemă conexă, ce nu este imediat evidentă. Aceasta o reprezintă **evaluarea** preciziei procesului de detectie.

Evaluarea detectiei constă în principiu în calcularea a o serie de *erori de detectie* pe baza a ceea ce numim ”realitate de teren” (”groundtruth”<sup>13</sup>). O ”realitate de teren” în acest caz, este o *segmentare de referință* a secvenței ce presupune marcarea pozițiilor reale ale tranzițiilor video și etichetarea acestora prin analiza manuală, cadru cu cadru, a secvenței. Acestea constituie datele de referință pentru validare. Tranzițiile video rezultate din procesul de detectie sunt ulterior comparate cu ”realitatea de teren” pentru a evalua erorile de detectie și astfel precizia algoritmului.

Procesul de evaluare este o problemă complexă datorită volumului mare de date conținut într-o secvență de imagini, date ce trebuie analizate manual pentru constituirea segmentării de referință. Mai mult, anumite tranziții video în practică pot fi interpretabile, lucru ce duce la existența mai multor segmentări de referință, în funcție de modul de percepție al persoanei care le-a realizat.

În ceea ce privește procesul de validare, sunt vizate două situații de eroare ale algoritmului de detectie, și anume:

---

<sup>12</sup>”Support Vector Machine” sau SVM reprezintă o colecție de metode ”înrudite” de învățare supervizată, ce sunt folosite pentru clasificarea datelor. Acestea fac parte din categoria clasificatorilor liniari generalizați (vezi Secțiunea 7.2.2).

<sup>13</sup>vezi explicatia de la pagina 170.

- algoritmul nu a detectat anumite tranziții prezente în secvență, această situație este o **eroare de nedetectare**,
- algoritmul a detectat în mod eronat anumite pasaje ale secvenței drept tranziții, această situație fiind o **eroare de falsă detectie**.

Pe baza acestor două situații posibile, metodele de evaluare existente propun diverse măsuri de eroare. Astfel, o primă abordare constă în calcularea erorii de detectie,  $E_D$ , ce corespunde procentului de tranziții nedetectate, precum și a erorii de falsă detectie,  $E_{FD}$ , ce corespunde procentului de false detectii. Acestea sunt date de relațiile următoare:

$$E_D = \frac{N_t - GD}{N_t} \quad (2.36)$$

$$E_{FD} = \frac{FD}{N_t} \quad (2.37)$$

unde  $N_t$  reprezintă numărul total de tranziții prezente în secvență,  $GD$  reprezintă numărul de tranziții ce au fost detectate corect iar  $FD$  reprezintă numărul de false detectii. Folosind cele două măsuri,  $E_D$  și  $E_{FD}$ , *performanța algoritmului de detectie este maximală când acestea sunt minime*.

Abordarea cea mai frecvent întâlnită, constă în calculul erorilor de tip "precision" și "recall". Acestea sunt date de relațiile următoare:

$$Precision = \frac{GD}{GD + FD} \quad (2.38)$$

$$Recall = \frac{GD}{N_t} \quad (2.39)$$

unde  $N_t$ ,  $GD$  și  $FD$  au aceeași semnificație ca mai sus.

Eroarea de tip "precision" este într-un fel o măsură cantitativă a numărului de false detectii. Aceasta este maximală (valoare 1) pentru  $FD = 0$  și deci în cazul în care nu au avut loc false detectii. Pe de altă parte, eroarea de tip "recall" este o măsură a numărului de detectii corecte, aceasta având valoarea maximală pentru  $GD = N_t$  (adică toate tranzițiile prezente au fost detectate). Folosind cele două măsuri, "precision" și "recall", *performanța algoritmului de detectie este maximală atunci când cele două mărimi au valori maxime*.

Indiferent de metoda adoptată, validarea detectiei trebuie efectuată pe baza evaluării atât a detectiilor corecte cât și a falselor detectii. Un algoritm ce detectează toate tranzițiile prezente în secvență nu este neapărat un algoritm eficient, acesta putând furniza adițional un număr foarte mare de false detectii, și vice-versa.

O soluție pentru a reprezenta simultan cele două erori o constituie prezentarea rezultatelor sub formă de curbe de precizie. Pentru aceasta, algoritmul de detecție este rulat de mai multe ori pentru valori diferite ale parametrilor de reglaj (de exemplu, valori diferite ale pasului de analiză, valori diferite ale pragului de detecție, etc.). Erorile de detecție astfel obținute, "precision" și "recall", sunt reprezentate grafic ca puncte în spațiul bidimensional format de acestea. Procesul se poate repeta similar și pentru alți algoritmi de detecție. În final, metoda cea mai eficientă va fi metoda ce va furniza punctul cel mai apropiat de colțul din dreapta sus al graficului, prezentând valorile cele mai ridicate ale "precision" și "recall".

În Figura 2.5 am ilustrat un exemplu de curbe de precizie obținute pentru mai mulți algoritmi de detecție de "cut" (marcați cu simboluri diferite), precum și pentru un același algoritm executat cu parametri de reglaj diferenți (marcat cu simboluri de aceeași culoare). Astfel, pe baza curbei de precizie obținute putem decide care sunt algoritmii de detecție cei mai preciști, precum și care sunt seturile de parametri optimali ce conduc la rezultatele cele mai performante (marcați în Figura 2.5 cu cercul roșu).

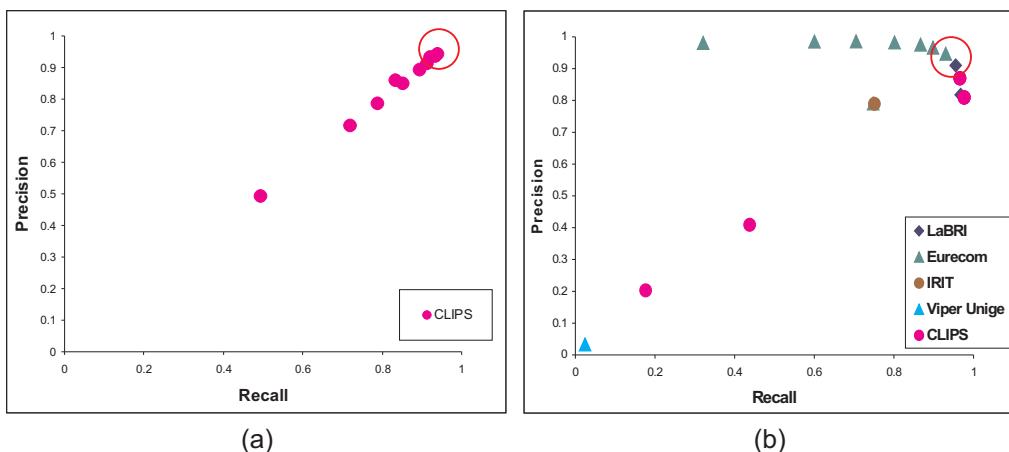


Figura 2.5: Exemplu de curbe de precizie (sursă [ARGOS 06]): (a) o singură metodă rulată pentru mai multe valori ale parametrilor de reglaj, (b) mai multe metode și reglaje diferențiale ale parametrilor.

În concluzie, punctul cheie al procesului de evaluare îl constituie disponibilitatea unei segmentări de referință. Datorită vastei diversității de secvențe de imagini disponibile, este aproape imposibilă constituirea manuală a unei referințe pentru fiecare secvență în parte. Mai mult, precizia algoritmului de detecție este dependentă de genul secvenței folosite, astfel, o secvență cu un conținut bogat în efecte vizuale va fi mai probabil să declanșeze un număr

mai mare de false detecții decât o secvență cvasi-uniformă. Soluția constă în validarea detecției folosind o bază de secvențe de test, special constituită, ce dispune de "realitate de teren". Din păcate, datorită constrângerilor de "drept de autor" cu care se confruntă fiecare secvență, o astfel de bază nu este încă disponibilă pentru publicul larg, aproape fiecare metodă existentă fiind testată pe o bază de secvențe particulară.

Totuși, notabile sunt eforturile depuse de o serie de campanii de evaluare a algoritmilor de prelucrare video, precum campania ARGOS - "Campagne d'Evaluation d'Outils de Surveillance de Contenus Vidéo" [ARGOS 06] sau campania TRECVID - "Video Retrieval Evaluation" [Trecvid 08], ce au ca scop standardizarea procesului de evaluare prin constituirea unei baze unice de test cât și a "realității de teren" aferente.

### 2.2.5 Constituirea planelor video

După cum am menționat în partea introductivă a acestui capitol, o secvență de imagini este constituită prin concatenarea planelor video pe baza tranzițiilor video. Astfel, planele video constituie *unitatea structurală de bază* a secvenței, descompunerea în plane stând la baza marii majorități a metodelor existente de analiză și prelucrare a secvențelor de imagini.

Până în acest punct al lucrării, în scopul realizării segmentării temporale a secvenței, am prezentat diversele tehnici de detecție existente a tranzițiilor video cel mai frecvent folosite, și anume *tranzițiile abrupte* de tip "cut" și *tranzițiile graduale* de tip "fade" și "dissolve". Localizarea individuală a tranzițiilor în secvență nu este însă suficientă pentru a obține descompunerea în plane, pentru aceasta trebuind luate în calcul o serie de aspecte practice [Ionescu 05a].

Detectia tranzițiilor de tip "cut", raportată la detectia celoralte tranziții, este cea mai probabilă să returneze un număr important de false detecții. Acest lucru se datorează faptului că toate celelalte tranziții graduale sunt surse de discontinuitate ale fluxului vizual, fiind de regulă confundate cu una sau mai multe tranziții de tip "cut". Anumite efecte sau schimbări de culoare, precum blițul camerei foto, produc de asemenea o schimbare vizuală importantă. Acestea nu sunt schimbări de plan dar este foarte probabil să fie detectate drept schimbări abrupte de tip "cut". Evitarea acestor situații constă în folosirea adițională de algoritmi de detecție specifici acestor efecte.

Trebuie ținut cont și de faptul că diversele tranziții video detectate nu sunt sincronizate între ele, algoritmii de detecție fiind în general independenți unul de altul. Astfel, pentru a constitui segmentarea temporală în plane video a secvenței o primă etapă constă în sincronizarea temporală a tranzițiilor detectate prin ordonarea cronologică a acestora. Planele video vor fi definite

astfel ca fiind intervalele continue de imagini cuprinse între două tranziții succesive ce țin cont de următoarele posibile reguli [Ionescu 05a] (vezi Figura 2.6):

- toate tranzițiile de tip ”cut” detectate în intervalul de timp ce corespunde unei tranziții graduale sunt eliminate,

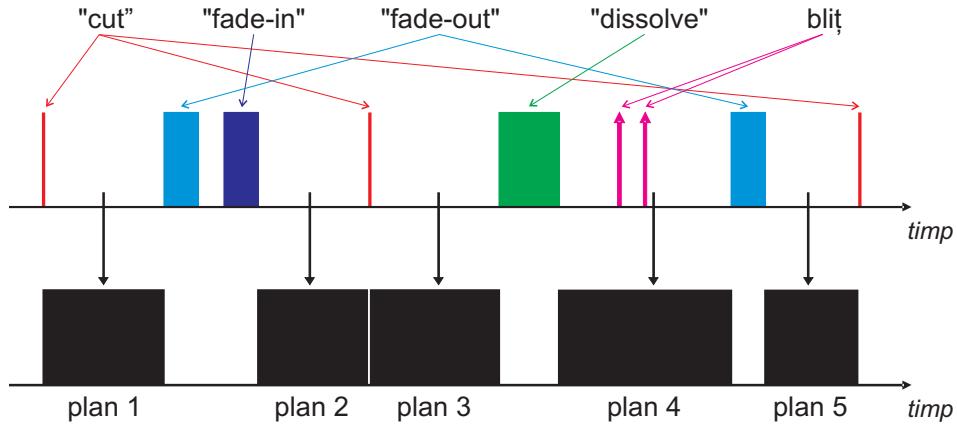


Figura 2.6: Constituirea planelor video ( fiecare tip de tranziție este reprezentat cu o culoare diferită).

- imaginile tranzițiilor video nu vor face parte din plan, fiind imagini nesemnificative pentru conținutul acestuia,
- efectele vizuale cu proprietatea că imaginea de început este similară cu imaginea de sfârșit, (de exemplu blițul aparatului foto din secvențele de știri [Heng 99] sau efectele de tip ”short color change” din filmele de animație [Ionescu 07c]), nu produc o schimbare de plan,
- planele video cu o durată inferioară unui anumit prag  $T_{plan} \approx 5$  imagini (determinat empiric) pot fi eliminate deoarece nu conțin informație esențială pentru secvență, fiind aproape neperceptibile,
- planele video cuprinse între două tranziții de tip ”fade-out” și ”fade-in”, sunt eliminate dacă durata lor este inferioară unui prag  $T_{fade}$ , acestea conținând doar imagini constante și de regulă negre. Dacă această inserție de imagini constante are o durată importantă, atunci poate fi considerată ca fiind un plan, având în acest caz o semnificație semantică și anume introducerea unui moment de pauză în desfășurarea acțiunii.

## 2.3 Detectia scenelor video

Segmentarea temporală a unui film clasic furnizează în general între 600 și 1500 de plane video. În diversele metode de prelucrare, pentru a reduce redundanța temporală, strategia cel mai des folosită constă în rezumarea fiecarui plan cu una sau mai multe imagini reprezentative. Astfel, se poate obține o reprezentare compactă a secvenței în cel puțin 600 până la 1500 de imagini, în funcție de numărul de plane. Volumul informațional este însă în acest caz încă mult prea ridicat pentru a permite o analiză și vizualizare eficientă a conținutului secvenței. Astfel, de multe ori este necesară o reprezentare structurală a conținutului secvenței pe un nivel superior planelor video. O astfel de reprezentare o constituie descompunerea secvenței în **scene**.

Detectia elementelor structurale de nivel semantic superior, precum scenele, nu este folosită doar pentru a evalua conținutul secvenței, acestea permitând și accesarea secvenței la un nivel semantic. În general, durata și frecvența de apariție a scenelor constituie elemente importante pentru studiul ritmului de desfășurare al acțiunii, sau al tehnicilor de realizare a secvenței. De exemplu, o scenă constituită din plane de scurtă durată implică un conținut de acțiune important. Pe de altă parte, o scenă în care durata planelor descrește progresiv, are ca efect creșterea suspansului acțiunii.

În literatura de specialitate, scenele sunt numite frecvent și *unități ale narativului* sau *paragrafe video*. Folosind terminologia specifică domeniului producției de film, putem defini scenele ca *fiind un ansamblu redus de plane video ce sunt unificate de locul acțiunii sau de anumite evenimente de interes* [Beaver 94].

În limbaj științific, o scenă se traduce printr-un grup de plane video ce prezintă caracteristici semantice comune. Inspirat din teoria teatrului clasic, conținutul unei scene trebuie să respecte regula celor trei unități: *unitatea de timp, unitatea de loc și unitatea de acțiune* [Corridoni 95]. În Figura 2.7 am ilustrat un exemplu de repartiție în scene pentru un extras dintr-un film abstract de animație [Ionescu 05b]. Astfel, o scenă  $X$  va conține planele de tip  $X$ , unde  $X \in \{A, B, C, D, E, F\}$ . Pentru fiecare plan video am figurat doar imaginile de început și respectiv de sfârșit. Se poate observa că o scenă va conține acele plane video, vecine temporal, cu un conținut similar, în care sunt prezente aceleași personaje și a căror acțiune se desfășoară în același loc.

Conceptul teoretic prin care filmul este constituit ca o combinație armonioasă de elemente sintactice (plane) legate între ele printr-un meta-limbaj<sup>14</sup>

---

<sup>14</sup>în domeniul lingvistic sau al logicii, un ”meta-limbaj” este un limbaj folosit pentru a caracteriza sau pentru a face afirmații cu privire la alte limbaje.

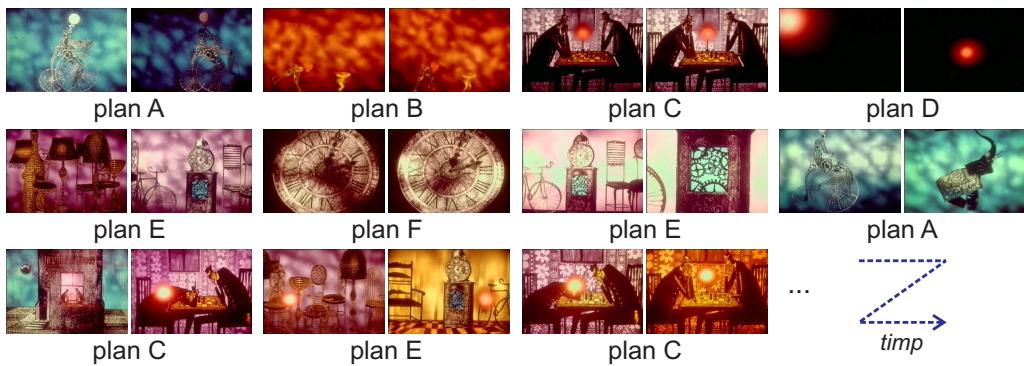


Figura 2.7: Exemplu de descompunere în scene (axa orizontală corespunde axei temporale, imagini din filmul de animație "Cœur de Secours" [CICA 06]).

poartă numele de paradigmă de montaj [Corridoni 95]. În acest sens, putem spune că, dacă planele video sunt considerate ca fiind unitățile sintactice ale secvenței, folosite cu predilecție pentru înțelegerea structurală a conținutului, atunci scenele constituie unitățile semantice ale secvenței, folosite pentru înțelegerea la nivel perceptual a conținutului secvenței.

După cum am precizat în capitolul introductiv al acestei cărți, tehniciile existente de prelucrare și analiză a imaginilor, nu sunt încă suficient de performante pentru a permite o înțelegere semantică completă a conținutului secvențelor de imagini, lucru ce este strict necesar pentru o detecție corectă a scenelor video. Cu toate acestea, tehniciile existente de detecție încearcă să "trișeze" prin transpunerea problemei de la un nivel semantic de analiză, la un nivel de analiză reprezentat prin mărimi numerice de nivel scăzut (nivel sintactic). Astfel, conceptele semantice de unitate temporală, loc și acțiune vor fi exprimate în termeni de similaritate între diversi parametri de culoare, textură, formă, mișcare, etc. Revenirea de la nivelul sintactic la cel semantic se va face în acest caz pe baza expertizei "a priori" a domeniului de aplicație.

Metodele existente de analiză a scenelor video sunt orientate spre două direcții de studiu distințte, și anume:

- pe de-o parte sunt metodele de *clasare automată a scenelor* în funcție de conținutul semantic al acestora, cum ar fi de exemplu separarea scenelor de dialog, scenelor de vânătoare, etc. Clasele predefinite folosite sunt de regulă specifice fiecărui domeniu (stiri, film, etc.),
- pe de altă parte, sunt metodele ce au ca scop *descompunerea în scene* sau care detectează schimbările de scenă.

Pentru un studiu bibliografic complet al problematicii de segmentare în elemente structurale de nivel semantic, cititorul se poate raporta la lucrările [Bimbo 99], [Kang 01] sau [Snoek 05b]. În cele ce urmează vom face o trecere în revistă a principalelor caracteristici ale tehnicilor folosite de metodele din cele două categorii enunțate.

### 2.3.1 Tehnici de clasare automată a scenelor

Obiectivul metodelor de clasare automată a conținutului scenelor este de a grupa diversele pasaje ale secvenței în categorii semantice predefinite. Aceasta poate fi realizată pe mai multe niveluri semantice, după cum urmează [Wang 00]:

- **nivel semantic de bază:** la acest nivel pasajele secvenței sunt grupate în clase elementare, precum: scene filmate în exterior/interior, scene de acțiune, scene fără acțiune, etc.,
- **nivel semantic intermediar:** la acest nivel sunt detectate scenele ce conțin evenimente de interes ușual, precum scenele de dialog dintre personaje, scenele ce se derulează în anumite locații specifice, scenele unor evenimente artistice, etc.,
- **nivel semantic propriu-zis:** la acest nivel, pentru localizarea scenelor este necesară înțelegerea conținutului de acțiune al secvenței. Scenele căutate conțin în acest caz evenimente complexe, de exemplu se caută scena ”uraganului apărut în Florida în 2004” sau ”celebrării Noului An la baza turnului Eiffel”.

Metodele propuse pentru clasarea automată a scenelor sunt foarte diverse. De exemplu, [Alatan 01] propune pentru detectarea scenelor de dialog din filme, folosirea informațiilor obținute în urma detectiei și localizării fețelor cu modele Markov ascunse. Scenele de violență sunt detectate în [Nam 98] pe baza analizei activității spațio-temporale a secvenței, prezenței în imagine a urmelor de sânge și a flăcărilor, precum și a analizei schimbărilor de energie ale semnalului audio.

Abordarea propusă în [Saraceno 98] folosește atât informația audio cât și vizuală pentru a clasa scenele în categoriile următoare: scene de dialog, scene de narativă, scene de acțiune și scene generice. Pentru aceasta, secvențele sunt mai întâi descompuse în plane video și audio. Planele audio sunt detectate pe baza analizei energiei semnalului audio. Acestea sunt încadrate într-o dintre categoriile următoare: pasaje fără sunet, pasaje de vorbire, pasaje de muzică sau zgromot.

Descompunerea în plane video este realizată pe baza informației de culoare. Folosind o metodă de cuantificare vectorială, în planele obținute anterior sunt localizate modele de blocuri de pixeli ce sunt specifice scenelor vizate. Mai departe, planele sunt grupate în funcție de caracteristicile audio-vizuale comune în grupuri omogene.

Detectia propriu-zisă a scenelor vizate este realizată pe bază de reguli. De exemplu, pe parcursul unei scene de dialog semnalul audio conține cu predilecție pasaje de vorbire iar planele audio respectă un model de tip *ABABAB*, unde *A* și *B* reprezintă două modele diferite de plane audio. Pe parcursul unei scene de acțiune, semnalul audio nu conține pasaje de vorbire iar informația vizuală are o evoluție progresivă după un model de tip *ABCDE*, unde literalele desemnează modele de plane video diferite.

Un alt exemplu este metoda propusă în [Lienhart 99b] ce localizează în secvență scenele ce au caracteristici sonore similare, scenele filmate în locuri similare și scenele de dialog. Metoda propusă presupune patru etape de analiză. În prima etapă, secvența este descompusă în plane prin detectia tranzițiilor video de tip "cut", "fade" și "dissolve". Mai departe, sunt extrase o serie de caracteristici la nivel de sunet, culoare, orientare a contururilor precum și a fețelor persoanelor prezente în imagine. Etapa următoare constă în calcularea unei măsuri de distanță între fiecare două plane video. Aceasta va fi adaptată la fiecare categorie de informații disponibile.

Distanțele astfel obținute sunt organizate sub forma unui tabel pe baza căruia planele vor fi grupate în scene în funcție de valoarea de distanță obținută. De exemplu, pentru a detecta scenele ce conțin pasaje audio similar, numite și "secvențe audio", măsura de distanță folosită este distanța minimală dintre vectorii de caracteristici audio ai planelor. Astfel, planele vecine temporal ce au valori ale distanțelor inferioare unui anumit prag (fiind astfel apropiate de distanța minimă), sunt regrupate în aceeași secvență audio. Autorii pretind că folosirea independentă a diverselor modalități ale secvenței conduce la obținerea de rezultate mai performante decât în cazul fuzionării acestora.

Din punct de vedere global, metodele de clasare automată a scenelor se bazează pe categorii predefinite, neexistând în acest moment metode ce permit localizarea oricărui tip de scenă fără a utiliza un minim de expertiză a conținutului secvenței.

Mai mult, algoritmii de detectie sunt specifici fiecărui tip de scenă în parte, nefiind aplicabili pentru a detecta și alte categorii de scene. De exemplu, un algoritm de detectie al scenelor de dialog ce folosește de regulă localizarea fețelor și detectia pasajelor de vorbire, nu poate fi folosit pentru a detecta o scenă de vânătoare.

### 2.3.2 Tehnici de descompunere în scene

Metodele de descompunere în scene nu caută să claseze conținutul scenelor ci dimpotrivă să reprezinte secvența la un nivel structural superior descompunerii în plane video (vezi Figura 2.1). În general, metodele de descompunere în scene sunt independente de aplicație sau de genul secvenței analizate.

Pentru a localiza scenele video dintr-o secvență, [Aigrain 95] definește un set de reguli globale ce trebuie să îndeplinească acestea. Regulile propuse sunt suficient de invariante la genul secvenței și se folosesc de următoarele informații:

- modalitate în care sunt **inserate tranzițiile video** graduale între tranzițiile de tip ”cut”,
- **distanța dintre planele video similare**: o schimbare de scenă este detectată atunci când se localizează în secvență un plan similar cu planul curent analizat ce se află la o distanță de numai 2 sau 3 plane. Similaritatea între plane este exprimată ca distanță între imagini cheie de luminanță ce sunt extrase din fiecare plan,
- **similaritatea planelor vecine**: continuitatea conținutului planelor video este detectată pe baza unei măsuri de similaritate, notată  $S$ . Aceasta este calculată în funcție de valoarea medie,  $m()$ , și de dispersia,  $\sigma()$ , a saturăției și respectiv a nuanței de culoare a pixelilor din anumite imagini cheie ale fiecărui plan. Mărimea  $S$  este dată de relația:

$$S_{i,j} = |m(i) - m(j)| + |\sigma(i) - \sigma(j)| \quad (2.40)$$

unde  $i$  și  $j$  reprezintă indicii planelor.

- **ritmul secvenței**: o schimbare de scenă este detectată numai dacă planul curent analizat are o durată de cel puțin trei ori mai mare decât o durată de referință,  $L$ , sau de patru ori mai mică decât  $L$ . Durata de referință  $L$  este estimată pe baza unui model autoregresiv<sup>15</sup> astfel:

$$L_n = a \cdot L_{n-1} + b \cdot L_{n-2} \quad (2.41)$$

unde  $L_n$  este valoarea lui  $L$  pentru iterată  $n$  iar parametrii  $a$  și  $b$  sunt estimati în ferestre de 10 plane, variabile temporale.

---

<sup>15</sup>un model autoregresiv, sau AR, este un predictor liniar ce încearcă să estimeze ieșirea unui sistem la un moment dat, pe baza valorilor anterioare de ieșire și respectiv intrare ale acestuia.

- **prezența pasajelor muzicale** după un pasaj lipsit de sunet,
- **similaritatea mișcării camerei video:** seriile temporale de cel puțin trei plane video ce conțin aceleași mișcări ale camerei video sunt considerate ca aparținând aceleiași scene.

Un alt exemplu este metoda propusă în [Huang 98] ce realizează o segmentare ierarhică a secvenței. Astfel, schimbările de plan și respectiv de scenă sunt detectate la diverse niveluri de detaliu. Schimbările de scenă sunt asociate schimbărilor simultane a culorilor, mișcării și a caracteristicilor semnalului sonor. Metoda propusă în [H. Sundaram 00] modelează relațiile de coerență dintre plane folosind un model de memorie cauzală de tip FIFO<sup>16</sup>.

O altă abordare constă în rezumarea conținutului fiecărui plan video cu imagini ”mozaic”<sup>17</sup> [Aner 01]. Astfel, planele sunt grupate în scene pe baza analizei similarității dintre imaginile de tip ”mozaic” asociate acestora. Metoda propusă în [Ionescu 05b] testează utilizarea histogramelor color medii clasice și respectiv ponderate, pentru a măsura similaritatea dintre planele video. Astfel, două plane sunt considerate ca aparținând aceleiași scene dacă au un conținut de culoare similar (valoare a distanței Euclidiene dintre histograme redusă) și sunt vecine temporal.

O abordare complexă inedită a problematicii descompunerii în scene este metoda hibridă propusă în [Kang 01]. Aceasta folosește pentru detecție colaborarea mai multor tehnici de detecție. Segmentarea ierarhică în scene este realizată în trei etape.

Prima etapă constă într-o *segmentare initială* în scene a secvenței. Metoda folosită este o metodă de detecție continuă. Gradul de coerență dintre diversele plane este calculat folosind modelul de memorie cauzală de tip FIFO propus în [H. Sundaram 00]. Valoarea de coerență este o valoare continuă în timp, ce este exprimată în funcție de valoarea parametrului de ”recall”. Valoarea acestuia pentru două plane  $A$  și respectiv  $B$ , notată  $recall(A, B)$ , este calculată la momentul  $S_a$  (punctul de pornire al analizei din memoria tampon) folosind ecuația:

$$recall(A, B) = (1 - dissim(A, B)) \cdot L_A \cdot L_B \cdot \left(1 - \frac{\Delta n}{N_m}\right) \cdot \left(1 - \frac{\Delta t}{T_m}\right) \quad (2.42)$$

unde  $dissim(A, B)$  este o măsură de disimilaritate a planelor  $A$  și  $B$ ,  $L_A$  și  $L_B$  reprezintă durata planelor  $A$  și respectiv  $B$  ce este ponderată de coeficienți

---

<sup>16</sup>FIFO este acronimul pentru ”First-In-First-Out” și reprezintă o modalitate de organizare și manipulare a datelor ce ține cont de timp și de prioritate. Aceasta lucrează pe principiul unei stive de date și soluționează conflictele generate de manipularea datelor folosind regula ”primul sosit este și primul servit”.

<sup>17</sup>vezi explicația de la pagina 150.

ce iau în calcul următoarele valori:  $T_m$ -dimensiunea memoriei tampon,  $N_m$ -numărul total de plane conținute în memorie,  $\Delta n$ -numărul de plane dintre planele  $A$  și  $B$ ,  $\Delta t$ -distanța temporală dintre planele  $A$  și  $B$ .

Valoarea de coerență,  $Co(S_a)$ , la momentul  $S_a$ , este dată de relația:

$$Co(S_a) = \frac{\text{recall}(A, B)}{R_{\max}(S_a)} \quad (2.43)$$

unde  $R_{\max}(S_a)$  reprezintă valoarea maximală a parametrului de "recall" ce este obținută în cazul particular în care  $dissim(A, B) = 0$ .

Astfel, schimbările de scenă sunt detectate pe baza analizei valorilor funcției de coerență într-o fereastră de decizie de dimensiune predefinită. O schimbare de scenă este detectată când un minim local al acesteia se găsește în mijlocul ferestrei considerate.

A doua etapă de analiză se folosește de această dată de o *abordare discretă*, și este folosită pentru a rafina rezultatele descompunerii în scene obținute anterior. Scopul acesteia este de a ameliora detecția prin reducerea numărului de scene rămase nedetectate. Astfel, planele conținute de scenele din prima etapă vor fi regrupate pe baza analizei unui anumit set de imagini cheie. Regruparea lor se va face cu un clasificator k-means pentru care numărul optimal de clase a fost ales pe baza unei metode de analiză a pertinenței claselor<sup>18</sup>. În urma clasificării, fiecare plan va dispune de o anumită etichetă corespunzătoare clasei din care face parte, transformând astfel scenele în serii de etichete. Planele ce au aceeași etichetă vor fi asociate între ele prin crearea unei corespondențe. Dacă în interiorul unei scene găsim două plane ce nu sunt corespondente, atunci scena va fi scindată în acel punct (vezi Figura 2.8, etichetele diverselor plane au fost figurate cu litere majuscule iar relațiile de corespondență sunt marcate cu săgeți duble).

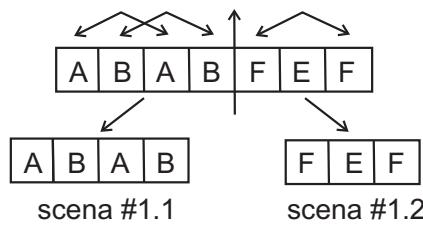


Figura 2.8: Exemplu de rafinare a detecției [Kang 01].

Ultima etapă de analiză constă în ajustarea segmentării în scene obținută până în acest punct. Aceasta, are ca scop corectarea falselor detecții.

<sup>18</sup>pentru o descriere detaliată a algoritmului k-means, vezi Secțiunea 7.1.2.

Mai întâi planele scenelor adiacente sunt regrupate pe baza unui clasificator k-means folosind același principiu ca în etapa precedentă. Falsele schimbări de scenă sunt apoi corectate pe baza examinării corespondențelor dintre planele diverselor scene.

Din punct de vedere global, metodele existente de descompunere în scene, transformă problema detecției într-o problemă de similaritate între conținutul planelor video, sau la un nivel mai general, într-o problemă de clasificare a conținutului acestora. În acest context, scenele sunt văzute ca ansambluri de plane, vecine temporal, ce au proprietăți similare ale conținutului de culoare, mișcare, etc. Definirea unei măsuri eficiente de similaritate între plane joacă în acest caz un rol decisiv în descompunerea corectă în scene a secvenței. Cu toate că o astfel de abordare este încă departe de sensul semantic pe care îl au scenele, aceasta constituie totuși un pas important spre înțelegerea automată a conținutului secvenței.

Pe lângă clasarea automată după conținut și decupajul secvenței în scene, tehniciile de analiză a scenelor video își găsesc aplicații și în problematici conexe din analiza secvențelor de imagini. În cele ce urmează vom prezenta câteva dintre acestea.

### 2.3.3 Aplicații ale analizei scenelor video

Detecția de scene, și în particular analiza similarității planelor video, permite detecția unei tehnici de filmare particulare cunoscută sub numele de **"shot-reverse-shot"** [Bimbo 99]. Această tehnică este folosită frecvent în pasajele video cu un conținut cu predilecție static, precum scenele de conversație dintre personaje, anumite pasaje ale secvențelor documentare culturale sau în scenele din sporturile de interior, precum jocul de biliard, tenisul de masă, etc.

Un **"shot-reverse-shot"** presupune folosirea alternativă a mai multor camere video. Un exemplu clasic este filmarea unui meci de biliard în care camerele video filmează alternativ jucătorul, masa de biliard și publicul (vezi Figura 2.9).



Figura 2.9: Exemplu de tehnică **"shot-reverse-shot"** dintr-un meci de biliard (pentru fiecare plan am ilustrat doar o singură imagine reprezentativă).

Tehnica "shot-reverse-shot" introduce un fel de periodicitate a conținutului secvenței. Structura de plane va avea o evoluție de tip *ABABAB*, unde *A* și *B* denotă două plane diferite din punct de vedere al conținutului. Analizând distanța, în sens de similaritate, a planului curent față de planele vecine (tehnica folosită la detectarea scenelor), putem localiza ușor semnătura temporală specifică unui "shot-reverse-shot", și anume o succesiune de valori de tip valoare redusă, valoare semnificativă, valoare redusă. Interesul asupra acestei tehnici de filmare nu se rezumă doar la detecția propriu-zisă, ci are și o aplicație în caracterizarea semantică a conținutului. Aceasta frunizează informații prețioase despre natura secvenței, despre conținutul de acțiune al acesteia precum și despre ritmului de desfășurare al secvenței.

Detectia scenelor poate fi utilă și în **cazul segmentării temporale** a secvenței, în particular pentru detectia tranzițiilor abrupte de tip "cut" (vezi Secțiunea 2.2.1). Metodele de detectie de "cut" existente se confruntă cu problema falselor detectii, ce se concretizează prin asimilarea uneia sau a mai multor schimbări vizuale drept tranzitii de tip "cut". Acestea vor produce global o supra-segmentare temporală a secvenței prin scindarea artificială a planelor video în mai multe plane. Din punct de vedere al detectiei de scene pe baza analizei similarității dintre plane, aceste situații se traduc prin scene compuse din plane succesive temporal. Astfel, dacă două plane successive sunt găsite ca aparținând aceleiași scene, atunci este foarte probabil ca acestea să fie rezultatul unei supra segmentări, fuzionarea acestora într-unul singur corectând această problemă. Detectia scenelor în acest caz, furnizează informații suplimentare metodelor clasice de detectie a tranzitiei video, datorită faptului că măsurile de disimilaritate vizuală sunt calculate de această dată la nivel de segment și nu la nivel de cadru. Totuși, detectia scenelor nu poate ține loc de detectie a schimbărilor de plan, ci trebuie văzută ca o etapă de rafinare a rezultatelor detectiei.

Descompunera în scene video își găsește aplicație și în metodele de **generare automată a rezumatelor de conținut**, permitând accesarea conținutului secvenței pe mai multe niveluri de detaliu. Rezumatul unei secvențe video este definit ca fiind o reprezentare compactă a conținutului acesteia (vezi Capitolul 5). Metodele existente de rezumare folosesc ca informație de plecare segmentarea temporală în plane a secvenței. Astfel, principiul rezumării constă în rezumarea planelor video cu un anumit număr de imagini reprezentative, numite și imagini cheie. Totuși, acest tip de abordare nu permite utilizatorului să aleagă nivelul de detaliu dorit pentru rezumat, fiind limitată în a furniza informații la nivel de plan.

Pe baza analizei similarității planelor, folosită de descompunerea în scene, se poate construi o reprezentare compactă a secvenței pe mai multe niveluri de detaliu. Principiul este ilustrat în Figura 2.10. Astfel un prim nivel de

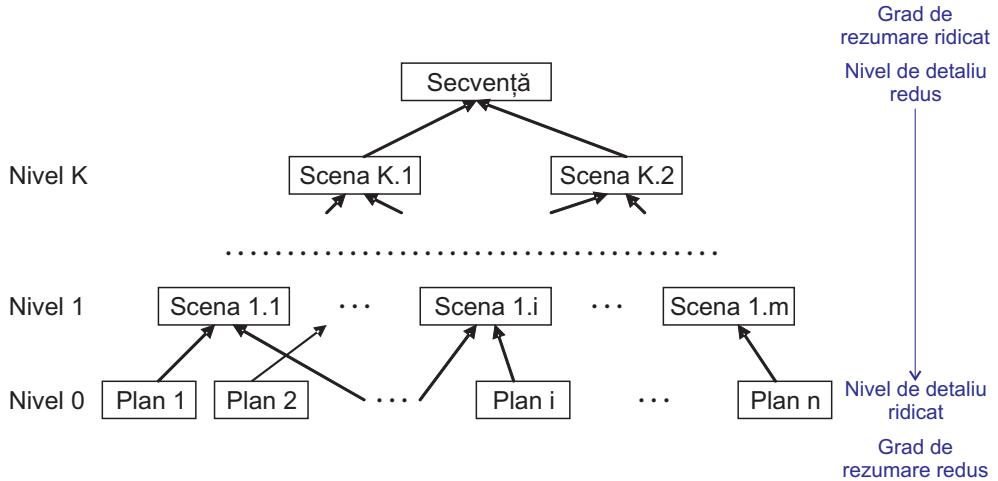


Figura 2.10: Reprezentare ierarhică a conținutului secvenței pe baza regrupării în scene.

detalii îl constituie nivelul planelor video. Acesta este urmat în mod natural de nivelul scenelor video. Nivelurile de detaliu superioare ierarhic sunt obținute pe baza regrupării scenelor similare. Pentru aceasta se poate folosi același mecanism ca pentru determinarea scenelor inițiale. Procesul poate fi repetat iterativ până se ajunge la cel mai înalt nivel ierarhic, reprezentat de secvența în totalitate.

Folosind principiul dendogramelor<sup>19</sup>, această reprezentare ierarhică a conținutului permite rezumarea secvenței în funcție de durata rezumatului dorit, cât și a cantității de informație furnizată de acesta. De exemplu, secvența poate fi rezumată cu o singură imagine cheie pentru fiecare scenă din nivelul  $k$ , obținând pentru  $k = 0, n$  imagini, pentru  $k = 1, m$  imagini, cu  $m \ll n$ , până la  $k = K$ , unde obținem doar două imagini (vezi Figura 2.10). Astfel, în funcție de cerințele aplicației, putem opta între un rezumat cu un grad de detaliu ridicat, dar de o durată mai semnificativă, sau pentru un nivel de detaliu mai scăzut, dar cu o durată mult mai redusă (în acest caz rezumatul poate conține doar câteva imagini).

Tot în cazul tehniciilor de rezumare, detectia scenelor video poate fi folosită la rafinarea conținutului rezumatelor dinamice. Spre deosebire de rezumatele ce folosesc un anumit număr de imagini cheie, rezumatele dinamice furnizează o serie de pasaje ale secvenței care sunt considerate ca fiind reprezentative pentru conținut. Problema care apare este redundanța vizuală a planelor

<sup>19</sup>O dendogramă (în grecește dendron - arbore și gramma - a desena) este o reprezentare sub formă arborescentă a claselor obținute în urma unui algoritm de clasificare.

din aceeași scenă video. Astfel, localizarea scenelor prin analiza similarității planelor, poate fi folosită pentru a localiza și elimina planele redundante, cum este cazul metodei propuse în [Laganière 08].

## 2.4 Concluzii

În acest capitol am discutat problematica segmentării temporale a secvențelor de imagini. Aceasta stă la baza majorității metodelor de analiză și prelucrare video.

Segmentarea temporală constă în descompunerea secvenței în plane video pe baza detecției tranzițiilor video. Dintre tranzиtiile existente se remarcă:

- tranzиtiile abrupte de tip ”cut”, ce reprezintă concatenarea directă a două plane succesive, fiind și tranzиtiile cel mai frecvent utilizate,
- tranzиtiile graduale, precum tranzиtiile de tip ”fade” și ”dissolve”, ce reprezintă o trecere progresivă de la un plan la altul, având un procent de apariție cu cel puțin un ordin de măsură inferior tranzиtiilor de tip ”cut”.

În ceea ce privește metodele folosite, direcția de studiu cea mai profitabilă pentru detecția tranzиtiilor video (atât abrupte cât și graduale), se dovedește a fi analiza evoluției intensității pixelilor din imagine, aceasta surclasând la nivel de performanțe celealte metode, precum metodele ce folosesc analiza de contur sau analiza mișcării. Din punct de vedere al performanțelor metodelor existente, de remarcat este faptul că cele mai precise sunt metodele de detecție a tranzиtiilor abrupte de tip ”cut”, algoritmii existenți furnizând în medie o precizia de detecție de peste 95%. Pe de altă parte, în cazul detecției tranzиtiilor graduale de tip ”fade” și ”dissolve”, precizia de detecție este inferioară, în medie situându-se undeva în jurul valorii de 75%. Acest lucru se datorează în principal modificărilor complexe ale scenei realizate de tranzиtiile graduale.

Detectia tranzиtiilor video are pe lângă rolul de a furniza descompunerea temporală în plane a secvenței și un aport în caracterizarea semantică a conținutului. Fiecare tip de tranzиcie video în parte, este folosit cu un scop precis. De exemplu, tranzиtiile de tip ”cut” sunt schimbări brusă de scenă și astfel sunt folosite pentru a face tranzиția rapidă de la un moment al scenei la altul. Tranzиtiile de tip ”dissolve” sunt tranzиti lente și sunt folosite de regulă pentru a schimba timpul acțiunii sau tranzиtiile de tip ”fade” ce introduc o pauză în desfășurarea acțiunii secvenței.

Pe de altă parte, cunoașterea structurii temporale a secvenței oferă informații și referitor la ritmul de desfășurare al acțiunii. Astfel, pasajele secvenței

bogate în schimbări de plan, reflectă un conținut bogat în acțiune și totodată un ritm alert, pe când planele video de lungă durată implică o rată de schimbare de plan redusă și astfel un ritm lent. Analiza cadenței schimbărilor de plan este utilizată cu succes de metodele existente pentru caracterizarea conținutului de acțiune al secvenței.

Dacă planele video sunt considerate ca fiind unitătile sintactice ale secvenței, fiind folosite cu predilecție pentru analiza de nivel scăzut a conținutului, scenele video sunt unitătile semantice ale secvenței ce permit o înțelegere de nivel semantic superior a conținutului. O scenă video este constituită dintr-un grup de plane video ce respectă unitatea de timp, de loc și de acțiune. Din păcate nivelul științific actual nu permite implementarea de metode automate capabile să înțeleagă în totalitate conținutul secvențelor, lucru necesar la detectarea scenelor. Din acest motiv, metodele existente se limitează cu predilecție doar la analiza similarității dintre planele video.

În concluzie, descompunerea secvenței în unități temporale, fie că este vorba de plane sau de scene, reprezintă o etapă de analiză necesară și totodată premergătoare analizei conținutului video. Aceasta oferă informații despre modalitatea în care a fost constituită secvența precum și despre modul de desfășurare al acțiunii, fiind într-un fel etapa inversă procesului de realizare a acesteia ce are loc în studio.

## CAPITOLUL 3

---

### Analiza mișcării

---

**Rezumat:** *Evoluția temporală a informației vizuale este una dintre particularitățile fundamentale a secvențelor de imagini. Din acest punct de vedere, secvențele de imagini mai sunt denumite și imagini în mișcare, fiind constituite ca o succesiune de evoluții temporale a conținutului unor imagini fixe. Pornind de la problematica estimării mișcării la nivel de pixel, în acest capitol vom face o trecere în revistă a diverselor direcții de studiu abordate de metodele de analiză și caracterizare a conținutului de mișcare din secvențele de imagini.*

Una dintre particularitățile de bază a secvențelor de imagini o constituie informația de mișcare. Raportat la imaginile statice, secvențele de imagini oferă o evoluție temporală a conținutului uneia sau a mai multor imagini fixe. În acest sens, secvențele de imagini sunt denumite și *imagini în mișcare*.

Dacă în acest moment motoarele de căutare folosite de sistemele de indexare actuale, permit accesul rapid și eficient la informațiile textuale, nu putem spune același lucru și despre accesarea conținutului multimedia. Pentru a compensa această lipsă, și anume existența unui mecanism eficient de accesare a datelor multimedia, dintre care în special a informațiilor video, grupul de dezvoltare al standardului de codare video MPEG ("Moving Picture Experts Group") lucrează la dezvoltarea și ameliorarea unui nou standard cunoscut sub numele de MPEG-7. După cum precizează creatorii aces-

tuia, ”principala ambiție a lui MPEG-7 este de a face conținutul informațiilor multimedia să fie la fel de ușor de accesat pe Internet precum informațiile textuale”. În ceea ce privește informația de mișcare, standardul MPEG-7 selectează și integrează *unele dintre cele mai performante metode existente* de analiză a mișcării. Astfel, tehniciile existente se grupează în două categorii principale [Jeannin 01]:

- pe de-o parte sunt metodele de analiză globală, bazate pe **analiza mișcării globale** a camerei video. În acest caz, analiza mișcării este realizată la nivel de segment video (pasaj al secvenței). Dintre aplicațiile analizei de mișcare globală putem enumera: recunoașterea mișcării camerei video, detectia activității de mișcare sau generarea imaginilor de tip ”mozaic”.
- pe de altă parte sunt metodele de analiză locală, ce sunt bazate pe **analiza mișcării obiectelor** din scenă. Acestea analizează mișcarea la nivel de regiuni spațiale de pixeli din imagine. De regulă, analiza locală este folosită pentru segmentarea și urmărirea temporală a obiectelor în mișcare.

Aceste două direcții de studiu sunt sintetizate în Figura 3.1. În cele ce urmează vom face o trecere în revistă a tehniciilor folosite de fiecare dintre acestea.

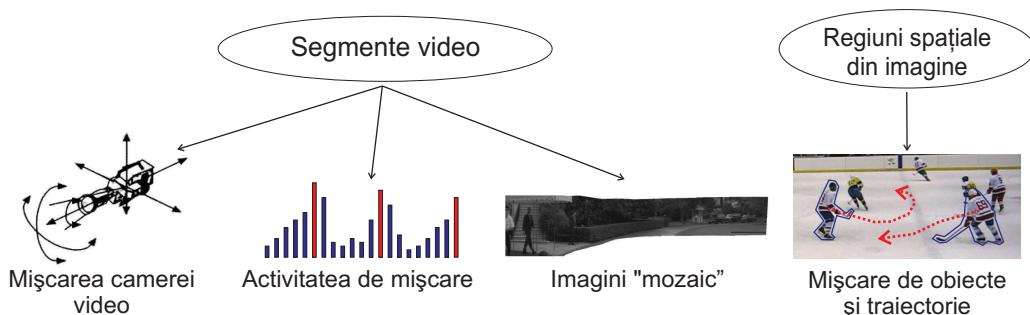


Figura 3.1: Principalele direcții de analiză a mișcării în secvențele de imagini: nivel global (segment) și nivel local (regiune) (sursă standard MPEG-7 [Jeannin 01]).

**Mișcarea globală.** Analiza mișcării globale a scenei este realizată la nivel de segment video sau de grup de imagini. O primă informație extrasă din secvență este *tipul mișcării camerei video*, ca de exemplu: mișcare translatională, mișcare de rotație, mișcare de apropiere etc. (vezi Secțiunea 3.2).

Informațiile reținute în acest caz pentru o anumită categorie de mișcare sunt de regulă amplitudinea mișcării, durata mișcării precum și localizarea acesteia în secvență. Analiza mișcării camerei video este importantă deoarece permite în anumite situații înțelegerea conținutului secvenței prin identificarea anumitor pasaje de interes din aceasta. De exemplu, focalizarea asupra unui anumit personaj se traduce printr-o mișcare a camerei de tip "zoom-in" (mărire), sau, creșterea suspansului acțiunii poate fi marcată de o mișcare de translație foarte rapidă.

O altă informație exploataată este *activitatea de mișcare*. Aceasta este o măsură a percepției vizuale pe care o avem asupra mișcării conținute în secvență. Activitatea de mișcare este determinată pe baza clasificării mișcării globale în funcție de o serie de parametri de nivel scăzut (de exemplu, dispersia amplitudinii vectorilor de mișcare). Clasificarea este realizată pe mai multe niveluri de activitate, în funcție de intensitatea acțiunii. La clasificare este luată în calcul și situația în care acțiunea este absentă, aceasta reprezentând nivelul minim de activitate. Un nivel de activitate intens corespunde evenimentelor dinamice, ca de exemplu scenele de gol din secvențele de fotbal sau scenele de urmăriră de mașini din secvențele de știri. Pe de altă parte, un nivel de activitate redus corespunde scenelor cu un conținut "sărac" în mișcare, ca de exemplu scenele de dialog dintre personaje sau scenele de interviu din secvențele de știri sau documentare.

Tot pe baza analizei mișcării globale este și *construcția imaginilor de tip "mozaic"* [Aner 01]. O imagine de tip "mozaic" este o imagine statică ce rezumă conținutul de mișcare global al unui pasaj al secvenței (de regulă un plan video). Aceasta este realizată prin regruparea și suprapunerea diverselor imagini ale segmentului, după recalarea geometrică în funcție de deplasarea globală a scenei (vezi Figura 5.2 de la pagina 150). Imaginile de tip "mozaic" sunt folosite drept rezumate compacte ale diverselor pasaje ale secvenței și de regulă au o complexitate de calcul ridicată. Totuși, aceasta poate fi redusă prin folosirea parametrilor de deformare furnizați de standardul MPEG-7.

**Mișcarea locală.** Analiza mișcării locale sau a deplasării obiectelor, este efectuată la nivel de regiuni de pixeli. Dacă pentru caracterizarea globală a mișcării, vectorii de mișcare puteau fi estimati la un nivel de detaliu mai redus (de exemplu, la nivel de blocuri de pixeli), furnizând astfel o aproximatie grosieră a fluxului optic, în cazul analizei mișcării locale a obiectelor, vectorii de mișcare sunt estimati de regulă la nivel de pixel pentru obținerea unui nivel de detaliu ridicat. Pentru analiză, metodele existente folosesc de regulă o modelare parametrică a mișcării. Aceasta permite localizarea în secvență a obiectelor cu deplasări similare, în ciuda diverselor deformări geometrice suportate de acestea. În general, rezultatul analizei mișcării obiectelor este

cuantificat prin furnizarea traectoriei acestora sub formă de evoluție temporală a unumitor puncte de interes, precum centrul de greutate sau anumite puncte de contur.

În contextul indexării după conținut a secvențelor de imagini, metodele de analiză a traectoriei obiectelor ("object tracking") sunt cu mult mai studiate decât metodele de analiză globală a mișcării camerei video. Acest lucru se datorează în principal faptului că într-o secvență mare a majoritatea evenimentelor de interes implică, și sunt legate, de mișcarea obiectelor. De exemplu, într-o secvență sportivă, va fi mult mai interesant și totodată reprezentativ pentru analiză să dispunem de traectoria unui anumit jucător care este într-o acțiune de atac, decât să caracterizăm mișcarea globală a camerei video ce urmărește jucătorul. Pentru un studiu biliografic complet al tehniciilor de analiză a mișcării obiectelor, cititorul se poate raporta la studiile [Dagtas 00], [Fablet 02] sau [Smith 04].

În concluzie, toate metodele existente de analiză a mișcării, fie că este vorba de mișcare globală sau locală, folosesc ca punct de plecare *estimarea mișcării*. Aceasta, pe baza măsurării deplasării pixelilor, sau a regiunilor de pixeli, de la un cadru la altul, furnizează un câmp vectorial de mișcare. În cele ce urmează vom face o trecere în revistă a tehniciilor de estimare a mișcării existente.

### 3.1 Estimarea mișcării

Principiul estimării de mișcare constă în determinarea deplasării unui pixel, sau a unui bloc de pixeli, între două imagini succesive ale secvenței, pe baza minimizării variației intensității acestuia, numită și DFD sau "Displaced Frame Difference". Această variație poate fi reprezentată sub forma următoare:

$$DFD(\vec{r}, \vec{d}, \Delta t) = I(\vec{r} + \vec{d}, t + \Delta t) - I(\vec{r}, t) \quad (3.1)$$

unde  $\vec{r}$  reprezintă poziția pixelului sau a blocului de pixeli în imaginea analizată,  $\vec{d}$  reprezintă vectorul de deplasare între momentele  $t$  și  $t + \Delta t$  exprimat în funcție de deplasarea pe cele două axe,  $oX$  și respectiv  $oY$ ,  $\vec{d} = (dx, dy)$ , iar  $I(t)$  reprezintă imaginea la momentul  $t$ .

Acest principiu de estimare se bazează pe ipoteza conform căreia intensitatea pixelilor nu variază semnificativ de la o imagine la alta. Un exemplu de vectori de mișcare obținuți la nivel de blocuri de pixeli este prezentat în Figura 3.2. Secvența folosită, pentru care am ilustrat câteva imagini reprezentative, conține o deplasare a camerei video către dreapta, de aici

și orientarea predominantă a vectorilor de mișcare spre stânga (imaginile se deplasează în sens invers camerei video).



Figura 3.2: Exemplu de vectori de deplasare (axa  $oX$  este axa temporală, vectorii indică direcția și amplitudinea deplasării blocurilor de pixeli, sursă imagini: filmul "The Wicker Man", Copyright 2006 Warner Bros Pictures).

Dacă considerăm imaginea ca fiind o funcție continuă, atunci putem folosi descompunerea în serie Taylor de ordinul întâi, astfel:

$$I(\vec{r} + \vec{d}, t + \Delta t) = I(\vec{r}, t) + \frac{\partial I(\vec{r}, t)}{\partial x} \cdot dx + \frac{\partial I(\vec{r}, t)}{\partial y} \cdot dy + \frac{\partial I(\vec{r}, t)}{\partial t} \cdot dt \quad (3.2)$$

și mai departe folosind ecuația 3.1, obținem:

$$DFD(\vec{r}, \vec{d}, \Delta t) = \frac{\partial I(\vec{r}, t)}{\partial x} \cdot dx + \frac{\partial I(\vec{r}, t)}{\partial y} \cdot dy + \frac{\partial I(\vec{r}, t)}{\partial t} \cdot dt \quad (3.3)$$

Minimizând funcția  $DFD$  obținem ecuația fluxului optic:

$$\frac{\partial I(\vec{r}, t)}{\partial x} \cdot u + \frac{\partial I(\vec{r}, t)}{\partial y} \cdot v + \frac{\partial I(\vec{r}, t)}{\partial t} = 0 \quad (3.4)$$

unde  $(u, v) = (\frac{dx}{dt}, \frac{dy}{dt})$  definește vectorul de deplasare în imagine. Fluxul optic astfel definit va permite estimarea mișcării doar în direcția gradientului spațial.

În general, metodele de estimare existente folosesc algoritmi de minimizare a unei funcții de cost ce este determinată pe baza funcției *DFD* de deplasare a pixelilor sau a blocurilor de pixeli. Astfel întâlnim mai multe abordări [Marichal 98]:

- **metodele diferențiale:** sunt bazate pe ecuația fluxului optic (vezi ecuația 3.4) și au ca rezultat un câmp vectorial de mișcare dens. Metodele diferențiale sunt foarte sensibile la prezența zgomotului în imagine, acesta diminuând considerabil precizia estimării. De asemenea, având o complexitate de calcul importantă, timpul necesar estimării este de regulă semnificativ.
- **metodele parametrice:** modelează deplasarea pixelilor în imagine folosind o serie de parametri. Problema estimării mișcării este astfel transformată într-o problemă de estimare a parametrilor unui anumit model, ca de exemplu modelul afin, cuadratic, etc.
- **metodele stohastice:** folosesc modele probabilistice. Explorarea spațiului parametrilor este ghidată în acest caz de procese aleatoare, precum modele Bayesiene, modele Markov sau algoritmi genetici. Metodele stohastice au o complexitate de calcul ridicată, dar rezultatele obținute corespund mult mai bine realității.
- **metodele bazate pe blocuri de pixeli:** estimează mișcarea la nivel de blocuri de pixeli. Această abordare a fost propusă pentru prima dată în [Jain 91] și s-a dovedit a fi cel mai bun compromis între complexitatea de calcul și precizia estimării obținute. Reglarea dimensiunii blocurilor de pixeli folosite la estimare permite în acest caz reglarea sensibilității și a robusteței metodei. Astfel, folosirea de blocuri de dimensiuni reduse are ca rezultat o sensibilitate ridicată a algoritmului de estimare, situație favorabilă în cazul micilor deplasări ale obiectelor din scenă. Pe de altă parte, folosirea de blocuri de dimensiuni mai mari conferă robustețe metodei la prezența zgomotului în imagine. Datorită acestor proprietăți, metodele de estimare pe blocuri de pixeli sunt folosite cu succes în marea majoritate a standardelor de codare video, precum H.263, MPEG 1, 2 și 4.

Din punct de vedere global, pentru a facilita calculele, metodele existente de estimare a mișcării adoptă o serie de *ipoteze de plecare*. În realitate, acestea nu sunt întotdeauna valabile, fapt ce duce uneori la obținerea de rezultate eronate. Dintre ipotezele cele mai importante, putem enumera următoarele:

- pixelii din blocurile de pixeli considerate, au același tip de mișcare de translație de la o imagine la alta,
- funcția  $DFD$  de deplasare a pixelilor, are o evoluție monoton crescătoare,
- intensitatea pixelilor din imagine este constantă cu mișcarea,
- mișcarea este considerată a fi constantă pentru mici volume spațio-temporale.

În funcție de nivelul de detaliu și de precizia câmpului vectorial de mișcare, întâlnim două tipuri de implementări (vezi Figura 3.3). Pe de-o parte sunt implementările *multi-rezoluție*, în care estimarea mișcării se face pentru mai multe reprezentări, de rezoluții diferite, ale aceleiași imagini analizate. Acestea sunt organizate sub formă piramidală, imaginea din vârful piramidei fiind imaginea cu rezoluția cea mai mică. Diversele rezoluții sunt obținute în urma filtrării de tip trece-jos și a subeașantionării spațiale.

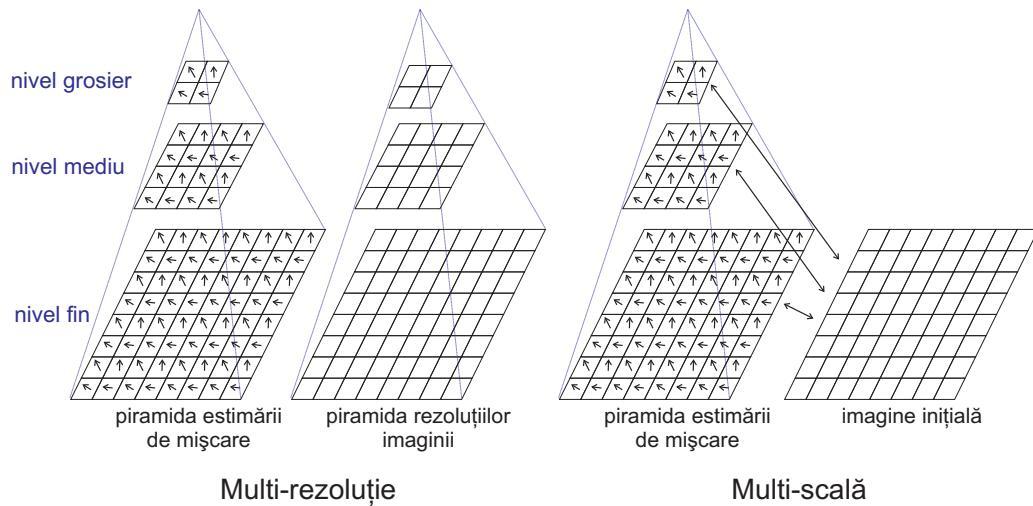


Figura 3.3: Modalități de implementare a estimării câmpului vectorial de mișcare (sursă [Marichal 98]).

Pe de altă parte, întâlnim implementarea *multi-scală*, ce folosește imaginiile în rezoluția inițială, dar furnizează mai multe niveluri de detaliu pentru vectorii de mișcare. Ca și în cazul implementării multi-rezoluție, nivelurile de detaliu sunt organizate sub formă piramidală și ordonate în funcție de densitatea câmpului vectorial furnizat.

Din punct de vedere temporal, estimarea mişcării între două imagini succesoive, poate fi realizată în două moduri:

- fie de tip *"forward"*, caz în care căutarea noii poziții a blocului de pixeli analizat din imaginea la momentul  $t$  se face în imaginea următoare la momentul  $t + 1$ ,
- sau de tip *"backward"*, când poziția blocului de pixeli la momentul  $t$  este căutată în imaginea anterioară, la momentul  $t - 1$ .

Aceste modalități de estimare își găsesc utilitate cu predilecție în metodele de codare video, precum standardele H.263 și MPEG, în care este folosită codarea bidirecțională. Astfel, un bloc de pixeli poate fi obținut atât cu predicție anterioară (*"backward"*) cât și cu predicție ulterioară (*"forward"*).

### 3.1.1 Metodele diferențiale

Metodele diferențiale de estimare a mișcării sunt metode bazate pe *calculul gradientului*. Acestea pornesc de la ipoteza conform căreia intensitatea pixelilor este constantă cu mișcarea. În această categorie putem enumera *metodele directe*, ce folosesc ca principiu anularea valorilor gradientului funcției *DFD* ce trebuie minimizată, și respectiv *metodele indirecte*, ce urmăresc convergența funcției *DFD* către o soluție în direcția gradientului. Din categoria metodelor indirecte putem menționa metodele iterative și metodele pel-recursive [Marichal 98].

#### Metodele iterative

Primele abordări ale problematicii estimării mișcării pe baza estimării fluxului optic sunt cele propuse în [Horn 81]. Pentru a rezolva ecuația 3.4 este adoptată o ipoteză suplimentară, și anume ca modulul gradientului să aibă valori mici, condiție valabilă în realitate doar pentru micile deplasări ale pixelilor din imagine. Astfel, problema estimării este transformată într-o problemă de minimizare a unei funcții de cost exprimată cu ajutorul coeficientilor Lagrange astfel:

$$\int \int [(I_x \cdot u + I_y \cdot v + I_t)^2 + \lambda \cdot (u_x^2 + u_y^2 + v_x^2 + v_y^2)^2] dx \cdot dy \quad (3.5)$$

unde  $I_i$  reprezintă derivata parțială de ordinul întâi a imaginii  $I$  în funcție de componenta  $i$ ;  $u_x, u_y, v_x$  și respectiv  $v_y$  reprezintă derivele parțiale de ordinul întâi ale celor două componente,  $(u, v)$ , ale fluxului optic, calculate în

direcția  $oX$  și respectiv  $oY$ , iar  $\lambda$  este multiplicatorul lui Lagrange ce reglează aportul erorii în ecuația de mișcare.

O soluție posibilă este dată de ecuațiile următoare:

$$\hat{u} = u_m - I_x \cdot \frac{P}{D} \quad (3.6)$$

$$\hat{v} = v_m - I_y \cdot \frac{P}{D} \quad (3.7)$$

unde  $u_m$  și  $v_m$  reprezintă valorile medii ale lui  $u$  și respectiv  $v$ , iar  $P$  și  $D$  sunt doi parametri dați de relațiile următoare:

$$P = I_x \cdot u_m + I_y \cdot v_m + I_t \quad (3.8)$$

$$D = \lambda + I_x^2 + I_y^2 \quad (3.9)$$

Calculul fluxului optic final este realizat în mod iterativ folosind pentru rafinare metoda Gauss-Seidel<sup>1</sup>. Soluția la iterația  $i$ ,  $(\hat{u}_i, \hat{v}_i)$ , este exprimată în funcție de soluția de la iterația precedentă,  $i-1$ ,  $(\hat{u}_i, \hat{v}_i) = f\{(\hat{u}_{i-1}, \hat{v}_{i-1})\}$ , soluția finală adoptată fiind soluția furnizată de iterația pentru care este îndeplinit un anumit criteriu de convergență.

Câmpul vectorial obținut în acest caz este *unul dens*. Fiecare pixel din imagine va avea asociat un vector de deplasare. Complexitatea de calcul este de asemenea semnificativă. Din această cauză, acest tip de abordare nu a fost folosită în tehniciile de codare, fiind utilizată cu predilecție în metodele de analiză a mișcării. Mai mult, datorită ipotezelor inițiale adoptate, metodele bazate pe analiza gradientului nu furnizează rezultate precise pentru deplasări importante ale pixelilor din imagine.

Ca exemple de metode de estimare a mișcării ce folosesc estimarea fluxului optic putem menționa:

- [Lim 01] ce propune îmbunătățirea calității fluxului optic obținut cu metoda Lucas-Kanade [Barron 94] pentru deplasări importante ale pixelilor din imagine,
- [Timoner 01] propune folosirea de filtre multi-dimensionale discrete pentru a ameliora precizia estimării cât și invarianța la efectul de "motion blur"<sup>2</sup>,

---

<sup>1</sup>metoda Gauss-Seidel este o tehnică iterativă folosită pentru a rezolva sisteme de ecuații liniare de tip  $Ax = b$ , unde  $A$  este o matrice de coeficienți,  $x$  reprezentă un vector de necunoscute iar  $b$  este un vector de valori. Soluțiile obținute la o anumită iterație sunt exprimate în funcție de valorile calculate anterior,  $x_i^{(k)} = \frac{b_i - \sum_{j < i} a_{i,j} x_j^{(k)} - \sum_{j > i} a_{i,j} x_j^{(k-1)}}{a_{i,i}}$ , unde  $k$  reprezintă iterația curentă. Aceasta este o versiune îmbunătățită a metodei Jacobi.

<sup>2</sup>efectul de "motion blur" corespunde urmelor de culoare ce apar în imagine în urma mișcării rapide a obiectelor sau a camerei video.

- [Elad 99] propune o soluție optimală pentru realizarea unor filtre specifice calcului gradientului, în scopul estimării mișcării.

### Metodele ”pel-recursive”

Metodele ”pel-recursive” [Netravali 79] realizează estimarea mișcării 2D a scenei la nivel de pixel și în *mod recursiv*. Dacă dispunem de o estimare inițială a mișcării pentru fiecare pixel al imaginii,  $d_i = (u_i, v_i)$ , unde  $d_i$  reprezintă vectorul de deplasare, atunci este realizată o corecție pe baza funcției *DFD* rezultate:

$$d_{i+1} = d_i + \Delta d_i \quad (3.10)$$

unde

$$\Delta d_i = (\Delta u_i, \Delta v_i) \quad (3.11)$$

reprezintă termenul de actualizare al iterăției  $i$ . Iterarea poate fi executată, fie pentru o linie de pixeli, fie pentru anumite linii de pixeli sau între imagini, operații numite și recurență orizontală, verticală și respectiv temporală.

Acest tip de abordare pornește de la ipoteza conform căreia funcția *DFD* converge local spre zero când mișcarea estimată tinde spre mișcarea reală prezentă în scenă. Prințipiu estimării constă în minimizarea recursivă a valorii cuadratică a funcției *DFD* pe baza metodei de gradient, astfel:

$$d_{i+1} = d_i + \varepsilon \cdot DFD(\vec{r}, \vec{d}, \Delta t) \cdot \nabla_{d_i}(DFD(\vec{r}, \vec{d}, \Delta t)) \quad (3.12)$$

unde  $\vec{r} = (x, y)$  reprezintă poziția inițială,  $\vec{d} = (u, v)$  reprezintă deplasarea între momentele  $t$  și  $t + \Delta t$ ,  $\nabla_{d_i}$  este operatorul de gradient calculat în funcție de  $d_i$  iar  $\varepsilon$  este în general o constantă pozitivă.

Exprimând gradientul în funcție de gradientul spațial, obținem ecuația următoare:

$$d_{i+1} = d_i + \varepsilon \cdot DFD(\vec{r}, \vec{d}, \Delta t) \cdot \nabla I(x + u, y + v, t + \tau) \quad (3.13)$$

unde  $\nabla$  este operatorul de gradient spațial iar  $I$  reprezintă imaginea. Alegerea adecvată a valorii lui  $\varepsilon$  permite obținerea convergenței estimării. Dacă valoarea utilizată este ridicată, atunci convergența este rapidă dar mai puțin precisă și vice-versa, dacă valoarea lui  $\varepsilon$  este mică, atunci convergența este sigură dar foarte lentă.

Această abordare, ca și estimarea iterativă, are ca rezultat un câmp vectorial dens. Principala problemă în acest caz este convergența metodei care este dependentă de ipoteza de plecare. Mai mult, metodele ”pel-recursive” sunt foarte sensibile la prezența zgomotului în imagine, acesta reducând drastic precizia estimării.

Ca exemple de metode de estimare a mişcării de tip "pel-recursiv" putem menŃiona:

- [Estrela 04] ce propune o metodă de estimare robustă la prezenŃa zgomotului în imagine. Vectorii de mişcare sunt estimaŃi iterativ folosind metoda EM ("Expectation-Maximization"<sup>3</sup>) și modelarea Gaussiană a datelor.
- [Hampson 00] ia în calcul influenŃa variaŃiilor de luminozitate a imaginii asupra calităŃii estimării mişcării. Pentru aceasta, estimatorul de tip "pel-recursiv" î se adaugă un coeficient multiplicativ ce modelează variaŃiile intensităŃii luminoase ale scenei.

### 3.1.2 Metodele parametrice

Metodele parametrice modelează deplasarea pixelilor din imagine pe baza unui anumit set de parametri. Problema estimării mişcării este transformată astfel într-o problemă de estimare a acestor parametrii. Numărul de parametri folosiŃi la estimare constituie caracteristica definitorie a metodei folosite. Spre deosebire de metodele din categoriile enunŃate anterior, metodele parametrice folosesc o restricŃie de natură geometrică, și anume pornesc de la ipoteza conform căreia obiectele din imagine sunt suprafeŃe plane și rigide.

Modelul cel mai simplu de parametrizare a mişcării îl constituie *modelul translational* ce definește noua poziŃie,  $(x', y')$ , a pixelului curent analizat,  $(x, y)$ , ca fiind:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} \quad (3.14)$$

unde  $t_x$  și  $t_y$  reprezintă deplasările pe cele două axe,  $oX$  și respectiv  $oY$ .

Introducerea unui factor unic de scalare, atât orizontală cât și verticală, conduce la un model cu trei parametrii, astfel:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = C \cdot \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} \quad (3.15)$$

unde  $C$  este parametrul de scală.

---

<sup>3</sup>un algoritm de tip EM - Expectation-Maximization este folosit în statistică pentru a găsi estimatorul matematic cel mai adecvat al parametrilor unui anumit model probabilistic, unde acesta din urmă depinde doar de "variabile latente" (variabile ce nu sunt direct observate, ci inferate pe baza altor variabile ce sunt observate și măsurate în mod direct).

O imbunătățire a modelului astfel obținut constă în adăugarea unui parametru suplimentar pentru a specifica factorul de scală pe fiecare axă în parte, astfel:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} C_x & 0 \\ 0 & C_y \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} \quad (3.16)$$

unde  $C_x$  și  $C_y$  sunt de această dată factorii de scală pe cele două axe,  $oX$  și respectiv  $oY$ .

O variantă similară poate fi considerat modelul în care factorul de scală este înlocuit cu un factor de rotație cu un anumit unghi  $\theta$ , astfel:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} \quad (3.17)$$

Combinând ultimele două modele, obținem un model cu cinci parametri, astfel:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix} \cdot \begin{bmatrix} C_x & 0 \\ 0 & C_y \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} \quad (3.18)$$

Separând factorii de rotație pe axa  $oX$  și respectiv  $oY$ , ajungem la modelul de mișcare cel mai cunoscut, și anume modelul dat de *transformarea afină* a conținutului imaginii, astfel:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} C_x \cdot \cos\theta_x & -C_y \cdot \sin\theta_y \\ C_x \cdot \sin\theta_x & C_y \cdot \cos\theta_y \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} \quad (3.19)$$

unde  $\theta_x$  și  $\theta_y$  reprezintă unghiurile de rotație pe cele două axe. Modelul afin este rezultatul *proiecției ortogonale* a mișcării pe o suprafață plană.

Folosind o *proiecție de perspectivă*, obținem un model al mișcării ce folosește opt parametri (*transformare de perspectivă*),  $a_i$  cu  $i = 1, \dots, 8$ , în care noile coordonate  $(x', y')$  ale pixelului curent analizat sunt date de relațiile:

$$x' = \frac{a_1 + a_2 \cdot x + a_3 \cdot y}{1 + a_7 \cdot x + a_8 \cdot y} \quad (3.20)$$

$$y' = \frac{a_4 + a_5 \cdot x + a_6 \cdot y}{1 + a_7 \cdot x + a_8 \cdot y} \quad (3.21)$$

O altă transformare folosită în mod curent este *transformarea biliniară* ce este dată de relațiile următoare:

$$x' = a_1 \cdot x + a_2 \cdot y + a_3 \cdot x \cdot y + a_4 \quad (3.22)$$

$$y' = a_5 \cdot x + a_6 \cdot y + a_7 \cdot x \cdot y + a_8 \quad (3.23)$$

Alte reprezentări iau în calcul și efectul accelerăției mișcării, ca de exemplu modelul propus în [Sanson 81]:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} a_x^x & a_x^y \\ a_y^x & a_y^y \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} b_x^{x^2} & b_x^{xy} & b_x^{y^2} \\ b_y^{x^2} & b_y^{xy} & b_y^{y^2} \end{bmatrix} \cdot \begin{bmatrix} x^2 \\ xy \\ y^2 \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} \quad (3.24)$$

unde  $a$  și  $b$  reprezintă seturile de parametri folosiți pentru modelizare.

Ca exemple de metode ce folosesc estimarea parametrică a mișcării putem menționa:

- [IRISA 05] ce propune o librărie C++ completă de funcții de calcul a estimării parametrice a mișcării ce folosesc marea parte a modelelor existente, precum modelul translational, afin și cuadratic,
- [Farnebäck 00] combină calculul tensorilor<sup>4</sup> 3D de mișcare folosind constrângerile unui model parametric. Aceștia sunt folosiți pentru a estima viteza de mișcare în imagine,
- [Wallin 01] propune o extensie a algoritmului de căutare ierarhică folosit de standardul MPEG-7 pe baza modelelor parametrice de mișcare.

### 3.1.3 Metodele stohastice

Metodele stohastice de estimare a mișcării folosesc teoria statistică. Explorarea spațiului parametrilor este ghidată în acest caz de procese aleatoare. În această categorie de metode putem enumera: *estimarea Bayesiană*, *modelele Markoviene* și *algoritmii genetici* [Marichal 98].

O soluție pentru a modela discontinuitățile câmpului vectorial de mișcare constă în folosirea câmpurilor aleatoare Markoviene sau MRF ("Markov Random Fields") [Graffigne 95] în ipotezele Bayesiene. În acest model, imaginea este reprezentată ca fiind un ansamblu de "locații", o "locație" reprezentând un pixel din imagine. Fiecare pixel este considerat astfel ca fiind o variabilă aleatoare, iar probabilitatea de apariție a acestuia cât și relațiile de vecinătate cu ceilalți pixeli sunt modelate cu ajutorul probabilităților condiționate. Câmpul aleator definit în acest fel este un câmp MRF dacă probabilitatea condiționată a fiecărui pixel depinde de un număr redus de pixeli vecini.

În [Papoulis 91] se propune modelarea discontinuităților câmpului de mișcare vectorial, pe baza introducerii noțiunii de "proces linie" sau proces de

---

<sup>4</sup>un tensor de rangul  $n$  într-un spațiu  $m$ -dimensional este un obiect matematic ce are  $n$  indici și  $m^n$  componente. Aceasta se supune unui anumit set de reguli de transformare. Fiecare indice al tensorului ia valori în numărul dimensiunilor spațiului de definiție.

discontinuitate. Acesta va permite interpretarea diferită a pixelilor vecini pixelului curent, analizat. Costul de introducere al acestui proces este modelat folosind câmpuri MRF. Astfel, elementele câmpului MRF vor avea două stări: fie sunt active, ipoteză în care a avut loc o discontinuitate în câmpul vectorial pe linia modelată, fie sunt inactive (pasive).

Estimarea mișcării va fi rezultatul minimizării unei funcții de energie ce va depinde de deplasările pe cele două axe,  $oX$  și respectiv  $oY$ , date de vectorul  $\vec{d}$ , precum și de procesul de discontinuitate ("procesul linie")  $l$ . Aceasta este modelat de patru parametri și anume:  $b$ ,  $b'$ ,  $c$  și  $c'$ . Funcția de energie astfel definită este dată de relația:

$$\begin{aligned} E(u, v|l) = & (1 - \lambda) \cdot \sum_{x,y} DFD^2(\vec{r}, \vec{d}) + \lambda^2 \cdot \left( \sum_{i,j} u_x^2(i, j)(1 - b_{i,j}) + \right. \\ & \sum_{i,j} v_x^2(i, j)(1 - b'_{i,j}) + \sum_{i,j} u_y^2(i, j)(1 - c_{i,j}) + \sum_{i,j} v_y^2(i, j)(1 - c'_{i,j}) ) + \\ & \alpha \cdot \sum_{i,j} (b_{i,j} + b'_{i,j} + c_{i,j} + c'_{i,j}) \end{aligned} \quad (3.25)$$

unde  $\vec{r} = (x, y)$  reprezintă poziția inițială a pixelului analizat,  $\vec{d} = (u, v)$  reprezintă vectorul de deplasare,  $\lambda$  este un parametru Lagrange,  $u_x$ ,  $u_y$ ,  $v_x$  și respectiv  $v_y$  reprezintă diversele deplasări exprimate în funcție de procesul de discontinuitate (vezi Figura 3.4) iar  $\alpha$  reprezintă costul de introducere a unei discontinuități, situație în care parametrii  $b$ ,  $b'$ ,  $c$  și  $c'$  iau valoarea 1.

Modelele Markoviene introduc o complexitate de calcul ridicată dar rezultatele obținute sunt apropiate de realitate [Marichal 98].

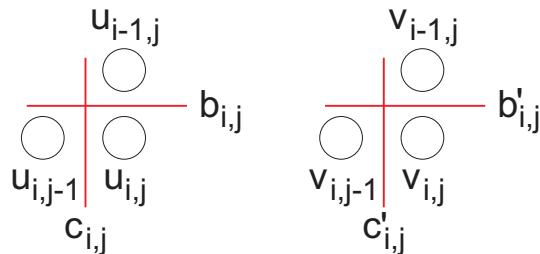


Figura 3.4: Configurațiile particulare ale procesului de discontinuitate în modelarea Markoviană (cercurile reprezintă pixelii iar liniile procesele de discontinuitate).

Un alt tip de abordare pentru estimarea mișcării constă în folosirea algoritmilor genetici. Algoritmii genetici sunt bazați pe teoria evoluției formulată de Darwin. Principiul acestora constă în evoluția artificială a unei populații

de dispozitive, proces ce este realizat pe baza unor operatori specifici, precum: selecția, creșterea sau mutația (vezi Secțiunea 7.1.2). De regulă algoritmii genetici sunt folosiți pentru probleme de optimizare a sistemelor ce implică un număr foarte mare de parametri sau obiective.

Ca exemple de metode de estimare a mișcării putem menționa metoda propusă în [Zaim 01] unde obiectele din imagine sunt parametrizate folosind vectori de caracteristici relative la forma și traectoria acestora în imagine. În acest caz, algoritmii genetici sunt folosiți pentru estimarea acestor vectori de caracteristici. Un alt exemplu este metoda propusă în [Pan 00] ce vizează estimarea mișcării 3D a oamenilor.

### 3.1.4 Metodele de estimare pe blocuri de pixeli

Tehnica de estimare a mișcării bazată pe analiza blocurilor de pixeli a fost propusă pentru prima oară în [Jain 91]. Spre deosebire de metodele prezentate în paragrafele anterioare, metode ce estimau mișcarea la nivel de pixel, estimarea pe blocuri calculează câmpul vectorial de mișcare la nivel de regiuni de pixeli, astfel furnizând un vector de deplasare pentru fiecare dintre acestea. Ameliorată de-a lungul timpului, această abordare s-a dovedit o metodă de estimare foarte eficientă ce furnizează cel mai bun compromis între complexitatea de calcul și precizia rezultatelor obținute. În funcție de dimensiunea blocurilor de pixeli folosite la estimare, întâlnim două situații posibile:

- o *sensibilitate ridicată a estimării* se obține pentru blocuri de pixeli de dimensiuni reduse, situație utilă în cazul în care estimarea mișcării trebuie să furnizeze rezultate precise, ca de exemplu pentru analiza deplasărilor fine ale obiectelor din scenă. În ciuda preciziei ridicate, în acest caz estimarea va fi mai sensibilă la prezența zgromotului în imagine.
- o *robustete ridicată a estimării* se obține pentru blocuri de dimensiuni mai mari. Folosirea mai multor pixeli pentru evaluarea funcției de cost face ca estimarea să fie mai puțin sensibilă la modificările provocate de prezența zgromotului în imagine. Totuși, având în vedere că vectorii de mișcare sunt calculați pentru fiecare bloc de pixeli, rezultatele reprezintă o aproximare mai grosieră a câmpului vectorial obținut la nivel de pixel. Această situație este utilă în cazul în care se analizează mișcarea la nivel global, ca de exemplu pentru analiza mișcării camerei video.

Datorită acestor proprietăți, după cum am mai menționat și în partea introductivă a acestui capitol, metodele bazate pe analiza blocurilor de pixeli

sunt utilizate cu predilecție de marea parte a standardelor de codare existente (MPEG-1, 2, 4, etc.).

Principiul estimării pe blocuri de pixeli este următorul. În primă fază, imaginea curentă analizată la momentul  $t$ ,  $I(t)$ , este împărțită în blocuri disjuncte de  $B \times B$  pixeli, unde  $B$  este de regulă ales ca fiind o putere a lui 2, din motive de optimizare hardware a metodei,  $B \in \{2, 4, 8, 16, \dots\}$ . Pentru fiecare bloc al imaginii  $I(t)$ , notat  $I(\vec{r}, t)$ , unde  $\vec{r} = (x', y')$  reprezintă poziția acestuia în imagine, se caută nouă sa poziție în imaginea următoare,  $I(t+l)$ , la momentul  $t+l$ , unde  $l$  reprezintă pasul de analiză (uzual  $l=1$ ).

Pentru a reduce complexitatea de calcul, căutarea acestuia nu se efectuează în toată imaginea, ci într-o *fereastră de căutare limitată*,  $S$ , de regulă de dimensiune  $(2B+1) \times (2B+1)$  pixeli. Acest lucru este posibil datorită faptului că în condițiile de continuitate a mișcării, deplasările de la un cadru la altul cadru sunt mici, fiind foarte puțin probabil ca blocul de pixeli analizat să se deplaseze în afara ferestrei  $S$ . Fereastra  $S$  este aleasă în imaginea  $I(t+l)$  ca fiind centrată pe blocul  $I(\vec{r}, t+l)$ .

Noua poziție a blocului curent analizat este dată de minimizarea unei funcții de cost,  $F_c()$ , ce estimează eroarea de aproximare a blocului curent,  $I(\vec{r}, t)$ , cu blocurile analizate din fereastra de căutare  $S$  din imaginea la momentul  $t+l$ . Astfel, vectorul de deplasare al blocului curent este dat de relația următoare:

$$\vec{d}_m = \operatorname{argmin}_{\vec{d} \in S} F_c(I(\vec{r} - \vec{d}, t + l), I(\vec{r}, t)) \quad (3.26)$$

unde  $\vec{d}_m$  reprezintă deplasarea blocului curent  $I(\vec{r}, t)$  pentru care funcția  $F_c$  este minimală. Valorile lui  $\vec{d}$  sunt toate deplasările posibile ale blocului de comparare în interiorul ferestrei de căutare. Principiul este ilustrat în Figura 3.5.

Dacă căutarea se face pentru toate valorile posibile ale vectorului de deplasare  $\vec{d}$  din fereastra de căutare  $S$ , atunci căutarea este o *căutare completă*. Căutarea completă este optimală în detrimentul vitezei de calcul. Numărul de operații estimate pentru o căutare completă este definit în [Accame 98] ca fiind:

$$\Delta_f = M \cdot N \cdot B^2 \cdot (2W + 1)^2 \quad (3.27)$$

unde  $M \cdot N$  reprezintă numărul de blocuri de pixeli din imagine,  $B^2$  este dimensiunea unui bloc de pixeli iar  $(2W+1) \times (2W+1)$  reprezintă dimensiunea ferestrei de căutare  $S$ . De exemplu, pentru o imagine de  $352 \times 288$  pixeli,  $B = 16$  și  $W = 16$ , numărul de operații necesare se ridică la  $\Delta_f = 1.1 \cdot 10^8$ , ce reprezintă un număr semnificativ, luând în calcul că o operație de filtrare cu o mască convolutivă de dimensiune  $3 \times 3$  necesită în acest caz în jur de  $\Delta_f = 1.5 \cdot 10^6$  operații.

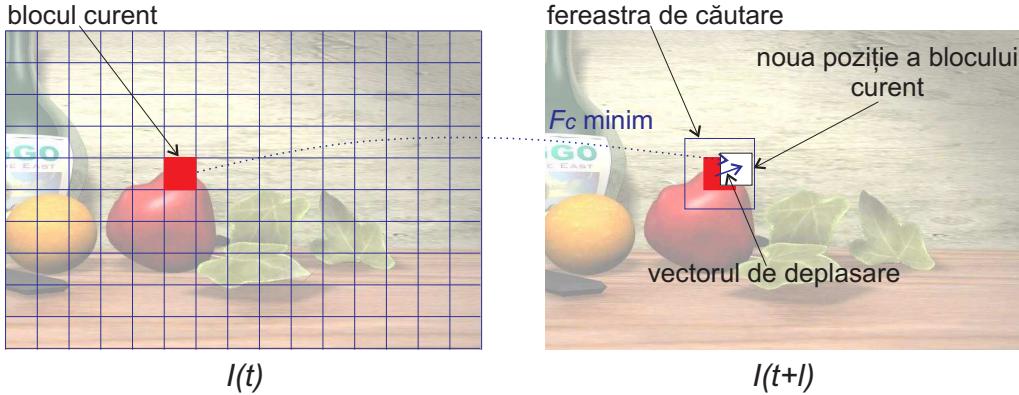


Figura 3.5: Principiul de estimare pe blocuri de pixeli a mișcării ( $I(t)$  reprezintă imaginea la momentul  $t$ ).

Dintre funcțiile de cost cel mai frecvent folosite, putem menționa următoarele ( $b_1()$  și  $b_2()$ ) reprezintă două blocuri de pixeli de dimensiuni  $M \times N$ :

- diferența medie absolută sau MAD ("Mean Absolute Difference"):

$$MAD(b_1, b_2) = \frac{1}{M \cdot N} \sum_{i=1}^M \sum_{j=1}^N |b_1(i, j) - b_2(i, j)| \quad (3.28)$$

- diferența medie pătratică sau MSD ("Mean Square Difference"):

$$MSD(b_1, b_2) = \frac{1}{M \cdot N} \sum_{i=1}^M \sum_{j=1}^N [b_1(i, j) - b_2(i, j)]^2 \quad (3.29)$$

- clasificarea distanțelor dintre pixeli sau PDC ("Pel Difference Classification"):

$$PDC(b_1, b_2) = \sum_{i=1}^M \sum_{j=1}^N ord(|b_1(i, j) - b_2(i, j)| \leq \tau) \quad (3.30)$$

unde funcția  $ord(P)$  returnează valoarea 1 dacă propoziția  $P$  este adevărată și 0 altfel, iar  $\tau$  este un prag ales arbitrar.

- proiecția integrală sau IP ("Integral Projection"):

$$IP(b_1, b_2) = \sum_{i=1}^M \left| \sum_{j=1}^N b_1(i, j) - \sum_{j=1}^N b_2(i, j) \right| + \\ \sum_{j=1}^N \left| \sum_{i=1}^M b_1(i, j) - \sum_{i=1}^M b_2(i, j) \right| \quad (3.31)$$

În general în urma minimizării funcției de cost,  $F_c()$ , pentru un anumit bloc de pixeli analizat, intervin trei situații particulare ce trebuie luate în considerare de algoritmul de estimare, astfel:

- valoarea minimală este obținută pentru blocul de pixeli din imaginea la momentul  $t + l$  ce se află în același poziție cu blocul curent analizat. În această situație vectorul de deplasare este nul, fiind un caz de "*absență a mișcării*",
- valoarea minimală a funcției de cost este foarte ridicată, fiind superioară unui anumit prag,  $\tau_{discont}$ , determinat empiric. În acest caz, mișcarea blocului analizat este catalogată ca fiind un caz de "*mișcare discontinuă*", aceasta neavând continuitate temporală în imaginea următoare,
- valoarea minimală este nenulă dar și inferioară pragului  $\tau_{discont}$ . În acest caz, deplasarea obținută corespunde *mișcării blocului analizat*.

În ceea ce privește pragul  $\tau_{discont}$ , numit și prag de discontinuitate, valoarea acestuia poate fi estimată "a priori" pe baza expertizei manuale a diverselor pasaje de discontinuitate a mișcării ce pot apărea în secvențele de imagini, ca de exemplu pasajele tranzițiilor video [Ionescu 07b].

Complexitatea de calcul a unei metode de estimare pe blocuri de pixeli este dată în primul rând de *modalitatea de căutare* a noii poziții a blocului curent analizat în fereastra  $S$ , și apoi de *dimensiunea ferestrei* de căutare precum și de funcția de cost folosită. În acest sens, căutarea completă se dovedește a fi căutarea cu gradul de complexitate de calcul cel mai ridicat. Metodele existente încearcă să găsească diverse soluții pentru a optimiza modul în care este realizată căutarea, fără a pierde însă semnificativ din precizia estimării. În cele ce urmează, vom face o trecere în revistă a diversilor algoritmi de căutare existenți [Turaga 98].

### Căutarea completă

După cum am menționat în paragrafele anterioare, căutarea completă constă în minimizarea funcției de cost pentru toate blocurile de pixeli din interiorul

ferestrei de căutare  $S$ . Astfel, căutarea completă este optimală dar are complexitatea de calcul cea mai importantă dintre metodele de căutare existente. Din această cauză, căutarea completă este folosită de regulă ca referință pentru evaluarea celorlalți algoritmi de căutare.

### Căutarea în trei etape

Principiul general constă în calcularea în avans a trei valori ale funcției de cost,  $F_c()$ , pentru trei valori diferite ale deplasării blocului curent analizat, valori ce sunt alese după un anumit criteriu. Pe parcursul căutării, aceste trei valori sunt recalculate în mod iterativ, căutarea finalizându-se în momentul în care este îndeplinită o anumită condiție de convergență [Reoxiang 94].

În această categorie se găsește algoritmul de căutare folosit de standardul de compresie video H.263+ [4i2i 06]. Pentru acesta, noua poziție a blocului curent este căutată în direcția valorii minime a funcției de cost folosind doar patru direcții de analiză, și anume: Sud, Est, Nord și Vest. Algoritmul este următorul (vezi Figura 3.6):

1. mai întâi este calculată funcția de cost între blocul curent analizat din imaginea la momentul  $t$  și cele patru blocuri vecine la distanță de un pixel pe direcțiile orizontală și respectiv verticală din imaginea următoare la momentul  $t + 1$ . Valoarea minimală a funcției de cost astfel obținută, este salvată în variabila  $D_1$  iar blocul de pixeli căruia îi corespunde această valoare devine bloc curent în imaginea la momentul  $t + 1$ .
2. folosind același principiu, se repetă calculul funcției de cost pentru vecinii noului bloc curent analizat în imaginea la momentul  $t + 1$ . Valoarea anterioară a variabilei  $D_1$ , va fi salvată într-o variabilă  $D_2$  iar noua valoare minimală a funcției de cost este salvată în variabila  $D_1$ . Blocul pentru care se obține această nouă valoare devine astfel bloc curent.
3. principiul se repetă mai departe pentru noul bloc curent. Astfel, valoarea variabilei anterioare,  $D_2$ , este salvată în noua variabilă  $D_3$ , valoarea variabilei  $D_1$  este salvată în  $D_2$  iar noua valoare minimală a funcției de cost este reținută în variabila  $D_1$ . Dacă condiția:

$$D_3 \leq D_1 \quad \cap \quad D_2 \leq D_1 \quad (3.32)$$

este îndeplinită, atunci căutarea este încheiată. În caz contrar, blocul ce a furnizat eroarea minimală devine bloc curent pentru analiză.

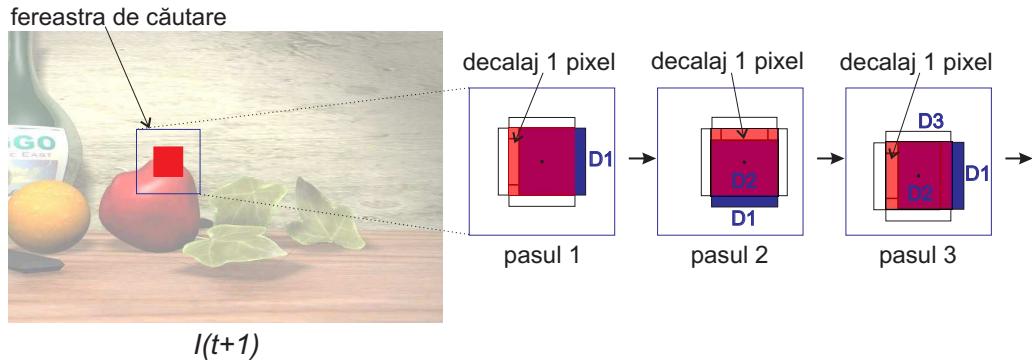


Figura 3.6: Principiul de estimare a mișcării în trei etape folosit de standardul H.273+ (blocul de pixeli curent este marcat cu roșu iar blocul ce corespunde erorii minime cu albastru).

4. se repetă în același mod calculul funcției de cost pentru noul bloc curent. Valoarea lui  $D_2$  va fi salvată în  $D_3$ , valoarea lui  $D_1$  în  $D_2$  iar noua valoare minimală a funcției de cost este reținută în  $D_1$ . Dacă valorile  $D_1$ ,  $D_2$  și  $D_3$  nu satisfac condiția 3.32, se repetă pasul 4. Dacă însă condiția 3.32 este satisfăcută, atunci căutarea se încheie iar ultimul bloc de pixeli ce a furnizat valoarea minimală a funcției de cost va determina vectorul de deplasare.

Pentru a evita poziționarea pe un minim local al funcției de cost, se poate folosi o condiție suplimentară, și anume ca valoarea minimală a funcției  $F_c()$ , obținută pentru ultimul bloc analizat din fereastra de căutare  $S$ , să fie de asemenea minimală în raport cu valoarea  $F_c()$ , obținută pentru blocul omolog blocului curent analizat din imaginea la momentul  $t$ , din imaginea la momentul  $t + 1$ , bloc pentru care se efectuează estimarea. În caz contrar, vectorul de mișcare poate fi considerat nul.

### Căutarea logaritmică

Similar cu principiul căutării în trei etape, căutarea logaritmică reduce numărul de comparații între blocuri, și astfel complexitatea de calcul, prin parcurgerea ferestrei de căutare  $S$  în direcția minimului funcției de cost  $F_c()$  [Lundmark 01].

Pentru fiecare bloc curent analizat din imaginea la momentul  $t$ , căutarea noii sale poziții în imaginea la momentul  $t + 1$  începe cu blocul omolog, pe care va fi centrată fereastra de căutare  $S$ . Se vor lua în considerare doar deplasările pe cele patru direcții fundamentale, și anume: Sud, Est, Vest și Nord. Astfel,

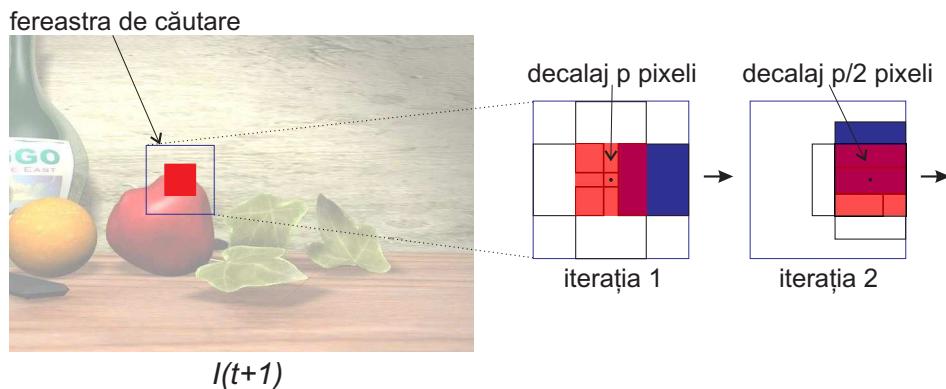


Figura 3.7: Principiul căutării logaritmice (blocul de pixeli curent este marcat cu roșu iar blocul ce corespunde eroiei minime cu albastru).

în primă etapă, funcția de cost este calculată pentru cele patru blocuri la o distanță de  $p$  pixeli de blocul curent pe cele patru direcții considerate (vezi Figura 3.7). Blocul ce corespunde valorii minimele a funcției de cost, devine bloc curent pentru etapele următoare. La fiecare etapă, deplasarea  $p$  este înjumătățită, astfel  $p \leftarrow p/2$ .

Algoritmul se repetă până în momentul în care  $p = 1$ . Ultimul bloc din imaginea la momentul  $t + 1$  pentru care s-a obținut valoarea minimală a funcției de cost, va determina vectorul de deplasare. Pentru a evita poziționarea pe un minim local al funcției de cost, se poate folosi aceeași condiție ca cea enunțată anterior pentru căutarea în trei etape, și anume ca valoarea minimală obținută pentru ultimul bloc analizat să fie minimală și în raport cu eroarea obținută pentru blocul inițial, pentru care se efectuează estimarea.

Cu această metodă de căutare, numărul de blocuri comparate este de  $2 + 7 \cdot \log_2 W$ , unde  $W$  reprezintă dimensiunea ferestrei de căutare. Dependența logaritmică de valoarea lui  $W$  a dat numele metodei de "căutare logaritmică".

### Căutarea binară

Căutarea binară este unul dintre algoritmii de căutare cei mai populari, fiind folosit și pentru estimarea mișcării în standardul MPEG [Zahariadis 96].

Principiul căutării binare constă în divizarea ferestrei de căutare într-o serie de regiuni și efectuarea unei căutări complete doar într-una dintre aceste regiuni. Algoritmul este următorul:

- mai întâi funcția de cost este estimată în imaginea la momentul  $t + 1$ , poziționându-ne pe o grilă de 9 pixeli ce sunt repartizați în fereastra de

căutare după cele 8 direcții cardinale. Folosind aceste puncte, fereastra de căutare este divizată în mai multe regiuni disjuncte. Accentul se pune pe regiunea centrală, fiind și cea mai probabilă să furnizeze noua poziție a blocului de pixeli curent analizat din imaginea la momentul  $t$  (vezi Figura 3.8).

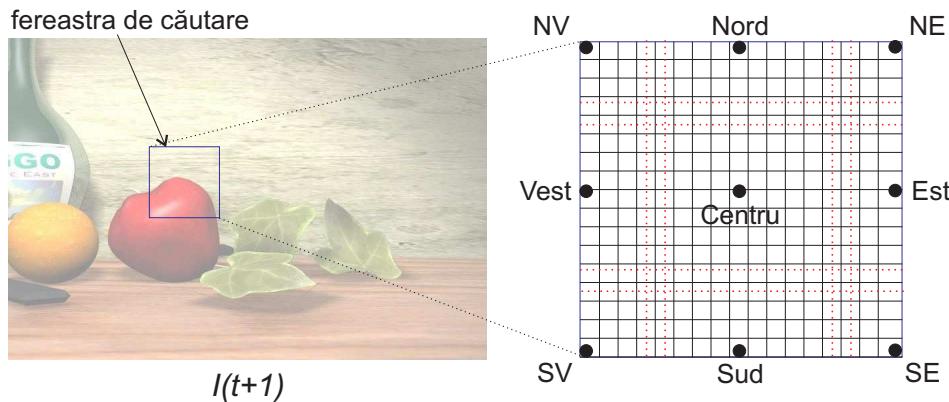


Figura 3.8: Principiul căutării binare: divizarea în regiuni a ferestrei de căutare (punctele negre reprezintă cei 9 pixeli ce formează grila de selecție iar frontieră dintre regiuni este marcată cu linia discontinuă).

2. regiunea din care face parte blocul de pixeli ce furnizează valoarea minimală a funcției de cost devine regiune curentă de analiză. În aceasta, pentru determinarea vectorului de deplasare, se va efectua o căutare exhaustivă (completă), cu mențiunea că blocurile de pixeli aflate pe frontierele dintre regiuni nu sunt luate în calcul.

În ciuda numărului redus de comparații între blocuri, căutarea binară furnizează performanțe medii datorită faptului că anumite regiuni de pixeli nu sunt luate deloc în calcul.

### Căutarea ortogonală

Căutarea ortogonală este o combinație între căutarea în trei etape și căutarea logaritmică. Aceasta implică o etapă de căutare pe verticală urmată de o etapă de căutare pe orizontală a blocului optimal. Algoritmul este următorul:

1. se alege un pas de analiză,  $p$ , de regulă ca fiind jumătate din valoarea deplasării maxime în fereastra de căutare. Se evaluatează funcția de cost,  $F_c()$ , pentru blocurile de pixeli din imaginea la momentul  $t + 1$  ce se află la distanța  $p$ , pe orizontală de centrul ferestrei de căutare. Blocul

pentru care se obține valoarea minimală a funcției de cost devine blocul curent de analiză,

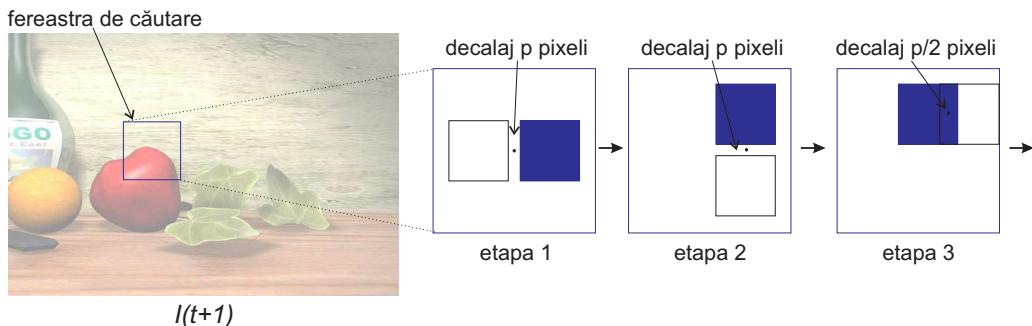


Figura 3.9: Prinzipiul căutării ortogonale (blocul de pixeli ce corespunde erorii minime este marcat cu albastru).

2. se evaluează funcția de cost pentru blocurile de pixeli ce se află de această dată la distanță  $p$ , pe verticală față de blocul curent din imaginea la momentul  $t + 1$ . Blocul ce furnizează valoarea minimală a funcției de cost devine astfel noul bloc curent de analiză (vezi Figura 3.9).
3. pasul de analiză este înjumătățit,  $p \leftarrow p/2$ . Dacă  $p > 1$ , atunci procesul se repetă. În caz contrar, căutarea se încheie, ultimul bloc curent reprezentând noua poziție căutată a blocului curent analizat din imaginea la momentul  $t$ .

### Căutarea intercalată

Căutarea intercalată se bazează de asemenea pe principiul căutării logaritmice, doar că în acest caz, blocurile de pixeli sunt alese ca formând un " $\times$ " și nu un " $+$ ". Algoritmul este următorul:

1. funcția de cost este calculată mai întâi pentru blocul de pixeli omolog, din imaginea la momentul  $t + 1$ , blocului curent analizat din imaginea la momentul  $t$ . Dacă valoarea obținută este inferioară unui anumit prag, atunci căutarea se încheie aici,
2. funcția de cost este calculată pentru 4 blocuri de pixeli din imaginea la momentul  $t + 1$ , ce formează un " $\times$ " în jurul centrului ferestrei de căutare  $S$ , la o distanță  $p$  de acesta. Blocul ce furnizează valoarea minimală a erorii devine astfel blocul curent de analiză.

3. dacă pasul de analiză,  $p$ , este mai mare ca 1, atunci acesta este înjumătățit și se repetă etapa 2. În caz contrar, se trece la etapa 4,
4. dacă blocul curent se află, fie în colțul din stânga jos al ferestrei de căutare, sau în colțul din dreapta sus, atunci se mai evaluează o dată funcția de cost pentru încă 4 blocuri, distribuite de această dată în "+" față de blocul curent și la distanța  $p$ . Dacă totuși blocul curent se află în colțul din stânga sus sau dreapta jos al ferestrei de căutare, similar cazului precedent, funcția de cost este evaluată pentru încă 4 blocuri dispuse de aceasta dată în "x" la distanța  $p$ .

Căutarea intercalată necesită un număr de aproximativ  $5 + 4 \cdot \log_2 p$  comparații, unde  $p$  reprezintă valoarea maximă a pasului de căutare. Aceasta are o complexitate de calcul redusă, dar nu este cel mai bun algoritm din punct de vedere al preciziei estimării.

### Căutarea ierarhică

Căutarea ierarhică are ca scop reducerea complexității de calcul prin furnizarea mai multor niveluri de detaliu al căutării.

Un exemplu sunt metodele piramidele în care estimarea se face pe imagini piramidele construite pe baza sub-eșantionării progresive a imaginii inițiale (vezi Figura 3.3). Estimarea mișcării este realizată astfel începând cu vârful piramidei, ce corespunde nivelului de detaliu cel mai scăzut, avansând spre baza piramidei, spre nivelul de detaliu cel mai ridicat ce finalizează cu imaginea inițială.

Pentru a reduce influența zgomotului în imaginile inferioare de detaliu, piramidele sunt construite folosind filtrări de tip "trece-jos" (FTJ). Estimarea mișcării de pe un nivel superior poate fi realizată cu una dintre metodele existente, precum căutarea în trei etape. Aceasta va fi folosită ca punct de plecare pentru nivelul imediat inferior, noii vectorii de mișcare fiind o "rafinare" a celor anteriori [Lin 98].

Avantajul căutării ierarhice constă în faptul că se poate adapta la constrângeri de timp variabile, furnizând în funcție de aplicație, nivelul de detaliu dorit. Dacă timpul de calcul este critic, estimarea se poate realiza rapid sacrificând din precizia rezultatelor.

### Căutarea hibridă

Căutarea hibridă nu este o metodă propriu-zisă, aceasta folosindu-se pentru estimare, de metodele deja existente, profitând astfel de avantajele furnizate de fiecare dintre acestea. Prințipiu constă în determinarea mai întâi a tipului

de mișcare prezent în imagine, de exemplu: mișcare lentă, mișcare rapidă, staționară, etc. Mai departe, în funcție de aceasta, se alege pentru calculul vectorilor de mișcare metoda cea mai eficientă în acest caz [Ge 02]. Astfel, precizia metodele hibride depinde de eficiența cu care se evaluează tipul de mișcare prezent în scenă.

### 3.1.5 Fluxul video MPEG

O alternativă pentru a recupera informația de mișcare constă în analiza informației din domeniul comprimat, sau în particular analiza fluxului MPEG - "Moving Picture Experts Group" [Pilu 97]. Metodele de codare video existente folosesc pentru compresie estimarea și compensarea mișcării, astfel vectorii de mișcare vor fi conținuți în fluxul video comprimat.

În mare, principiul compresiei temporale este următorul: în loc să fie stocate toate imaginile sevenței (de exemplu, 25 de imagini pe secundă), vor fi stocate integral doar anumite imagini, numite și imagini cheie, precum și vectorii de deplasare a blocurilor de pixeli din acestea în imaginile următoare. În momentul decompresiei, imaginile sevenței sunt reconstituite pornind de la o imagine de referință, ce poate fi o imagine cheie sau o imagine reconstruită anterior, pe baza vectorilor de deplasare, folosind compensarea mișcării.

Avantajul folosirii vectorilor de mișcare direct din fluxul video constă în faptul că aceștia au fost calculați în momentul codării. Astfel, imaginile folosite pentru estimare sunt calitativ net superioare imaginilor obținute în urma decompresiei, ce stau de regulă la baza estimării mișcării cu metodele prezentate anterior. Acest lucru se datorează în principal compresiei cu pierdere de informație folosită de codarea video. În funcție de constrângerile de debit impuse, aceasta poate deteriora semnificativ conținutul imaginii (necomprime). Mai mult, codarea video este o codare intensivă ce presupune o compresie atât spațială a imaginii (de regulă JPEG), cât și temporală pe baza vectorilor de mișcare.

Principiul de funcționare general al unui decodor MPEG este ilustrat în Figura 3.10. Mai întâi, datele video sunt trecute printr-un decodor Huffman<sup>5</sup> ce permite recuperarea coeficienților DCT<sup>6</sup> cuantificați cât și a vectorilor de mișcare. Mai departe, vectorii de mișcare sunt separați de fluxul de date și sunt furnizați blocului de compensare a mișcării. În același timp, coeficienții

---

<sup>5</sup>în teoria informației, codarea Huffman este un algoritm de codare bazat pe calculul entropiei folosit la compresia fără pierderi a datelor. Pentru compresie se folosește un dicționar de coduri, de dimensiune variabilă, pe baza căruia este codat fiecare simbol al sursei. Dicționarul este construit pe baza estimării probabilității de apariție a fiecarei valori posibile a simbolurilor emise de sursă.

<sup>6</sup>vezi explicația de la pagina 41.

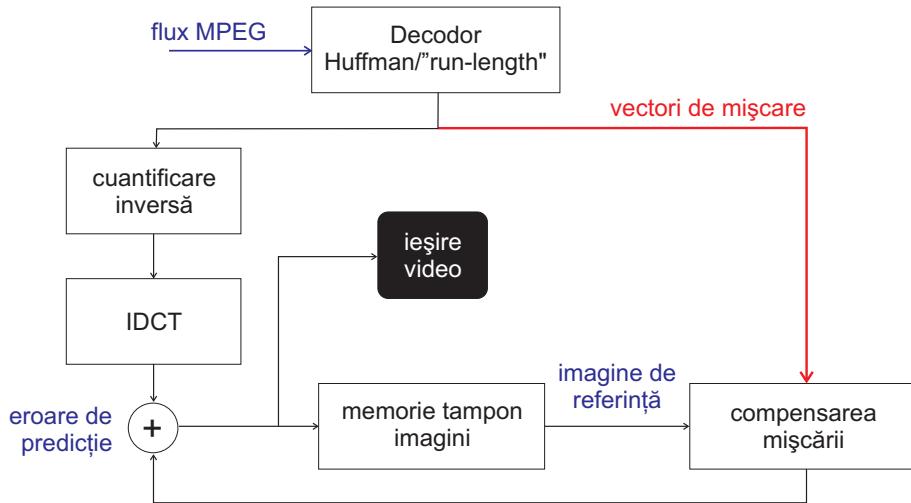


Figura 3.10: Diagrama de principiu a unui decodor MPEG-1 (sursă [Bretl 99]).

DCT sunt extrapolati și trecuți prin blocul de calcul al transformatei cosinus discretă inversă (IDCT) pentru a obține înapoi informația spațială din imagine.

În cazul imaginilor de tip P, ce sunt imagini reconstruite în funcție de imagini anterioare ("forward prediction"), sau de tip B, ce sunt imagini cu predicție bidirectională, reconstruite atât din imagini anterioare cât și următoare ("forward and backward prediction"), vectorii de mișcare sunt mutați la o anumită adresă de memorie de către blocul de predicție, pentru a accesa macro-blocul de predicție dintr-o imagine de referință stocată anterior. Sumatorul adaugă această predicție la valoarea reziduală pentru a reconstrui datele din imagine. În cazul imaginilor de tip I, ce sunt imagini fără predicție, nu există vectori de mișcare disponibili și nici imagine de referință, astfel că predicția este zero. De asemenea, pentru imaginile de tip I și P, ieșirea sumatorului este stocată ca imagine de referință pentru predicțiile ulterioare.

Astfel, vectorii de mișcare pot fi extrași direct din fluxul de date după decodarea Huffman. Pentru mai multe detalii referitor la implementarea principală cât și practică a metodelor de recuperare a vectorilor de mișcare din fluxul video MPEG, cititorul se poate raporta la lucrările [Gilvarry 99] și [Bretl 99].

În realitate, vectorii de mișcare obținuți direct din fluxul MPEG nu sunt întotdeauna coerenți. Aceștia necesită de regulă o serie de etape de corecție

pentru ameliorare. O situație frecventă ce cauzează ambiguitatea vectorilor de mișcare este estimarea mișcării pentru regiuni netexturate din imagine, unde blocurile de pixeli vecine nu prezintă suficientă variabilitate pentru a evidenția deplasarea acestora [Pilu 97]. Un exemplu de incoerență a vectorilor de mișcare extrași direct din fluxul video este prezentat în Figura 3.11, unde am ilustrat câmpul vectorial de mișcare obținut pentru o imagine dintr-o secvență ce conține deplasarea a două personaje simultan cu mișcarea globală a camerei video. Astfel, vectorii incoerenți sunt obținuți în special pentru pixelii de pe bordura imaginii, ce nu pot fi localizați în imaginea următoare datorită deplasării globale a scenei. Pentru a corecta această problemă, și astfel pentru a obține o calitate suficientă a vectorilor de mișcare, soluția cea mai frecvent adoptată este decompresia datelor până la un anumit nivel de detaliu [Pineau 05].

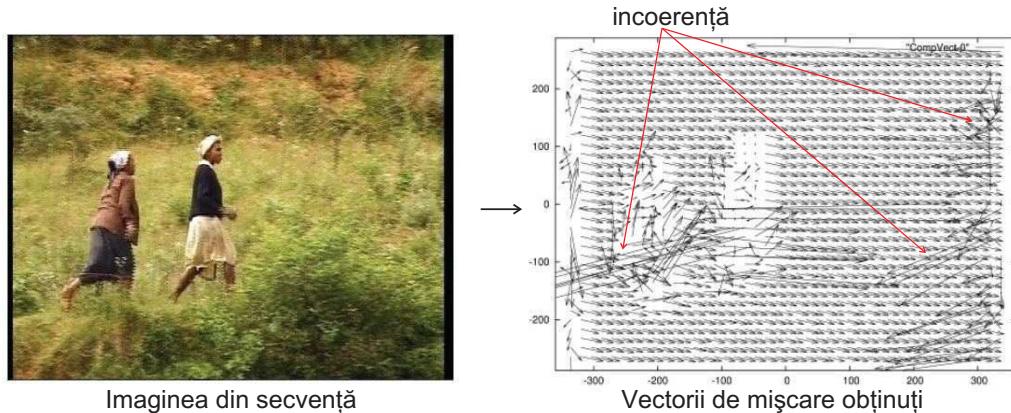


Figura 3.11: Exemplu de vectori de mișcare incoerenți extrași direct din fluxul MPEG-2 (sursă ”Projet Analyse et Indexation Vidéo” [Pineau 05]).

## 3.2 Analiza mișcării camerei video

Una dintre aplicațiile de mare interes ale estimării mișcării o constituie caracterizarea mișcării globale a scenei dată de deplasarea camerei video.

Acest tip de analiză permite obținerea unei caracterizări de ansamblu a acțiunii din secvență, anumite tipuri de mișcări ale camerei video fiind folosite voluntar pentru a marca evenimente importante din derularea secvenței. De exemplu, în anumite situații, schimbarea unui plan video este realizată folosind o mișcare globală de translație ce permite focalizarea pe un anumit punct de interes din afara scenei curente; un conținut bogat în acțiune este

deseori marcat de o mişcare rapidă a camerei video; personajele din secvenţă sunt aduse în prim plan de regulă folosind o mişcare a camerei video de tip ”zoom-in” (mărire optică a imaginii curente), etc.

Din punct de vedere global, mişcarea camerei video este o mişcare liberă în spaţiul real 3D. Totuşi, în realitate, din cauza constrângerilor tehnice şi fizice ale construcţiei camerei video, aceasta se rezumă la un număr limitat de mişcări de bază. Mişcările mai complexe, ce tind să se apropie de o mişcare liberă în spaţiul infinit 3D, sunt constituite ca o compunere de mişcări elemetare. Cele mai importante dintre acestea sunt ilustrate în Figura 3.12, astfel întâlnim:

- *mişcările translational* după cele trei axe de coordonate  $XYZ$ , precum mişcarea de tip ”dollying” (depărtare de obiectiv), ”zoom-in/ zoom-out” (mărire/micşorare), ”tracking” (translaţie spre dreapta sau stânga) sau ”booming” (translaţie în plan vertical, sus-jos). Mişcarea de tip ”zoom-in/ zoom-out” poate fi considerată ca o mişcare de translaţie doar din punct de vedere al efectului produs, acesta fiind similar cu cel al mişcării translational de tip ”dollying” (în cazul ”zoom-out”). Din punct de vedere tehnic, această mişcare se realizează fără deplasarea camerei video, prin mărirea/micşorarea optică a imaginii.

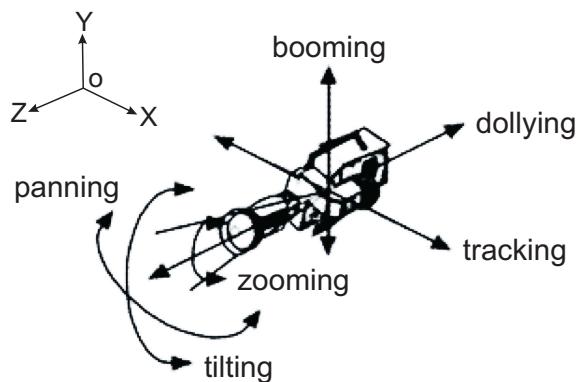


Figura 3.12: Mişcările camerei video.

- *mişcările de rotaţie* ale camerei video, precum rotaţia în sens orar şi antiorar în planul  $XoY$ , mişcarea de tip ”panning” (rotaţie în planul  $XoZ$ ) sau mişcarea de tip ”tilting” (rotaţie în planul  $YoZ$ ).

Metodele existente de analiză a mişcării camerei video pot fi regrupate în două categorii principale, astfel:

- metode ce analizează informația de mișcare direct în *domeniul comprimat* (fluxul MPEG),
- metode ce realizează analiza în *domeniul spațio-temporal* al cadrelor video ale secvenței.

Pentru un studiu complet al literaturii de specialitate al acestui domeniu, cititorul se poate raporta lucrările [Ngo 00], [Kramer 05], [Tardini 05] sau [Duan 06]. În cele ce urmează vom face o trecere în revistă a particularităților fiecărei dintre cele două categorii de metode existente.

### 3.2.1 Analiza mișcării camerei în domeniul comprimat

Din această categorie putem da ca exemplu abordarea probabilistică de detecție a mișcării de tip "zoom-in/zoom-out", propusă în [Jin 02]. Aceasta folosește algoritmul EM ("Expectation-Maximization"<sup>7</sup>) pentru estimarea probabilității de apariție a unei mișcări de tip "zoom" raportat la celelalte mișcări existente. În acest caz, informația de mișcare este extrasă direct din fluxul video MPEG-1 sau MPEG-2. Avantajul acestei abordări probabilistice constă în principal în îmbunătățirea invarianței la zgomotul ce afectează vectorii de mișcare, zgomot cauzat de regulă de erorile de cuantificare sau de prezența "artefactelor" datorate codării<sup>8</sup>.

O abordare mai generală, ce vizează detecția a șase tipuri de mișcări primare ale camerei video este propusă în [Kim 04]. Metoda propusă se bazează pe interpretarea calitativă a unui anumit număr de parametri extrași din modelele parametrice de mișcare (vezi Secțiunea 3.1.2), parametrii ce sunt estimați direct, pe baza fluxului MPEG-2. Astfel, vectorii de mișcare sunt extrași din fluxul MPEG-2 și sunt comparați cu modelul afin de mișcare:

$$\begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} a_2 & a_3 \\ a_5 & a_6 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} a_1 \\ a_4 \end{bmatrix} \quad (3.33)$$

unde  $(u, v)$  reprezintă vectorul de mișcare al blocului de pixeli centrat în punctul de coordonate  $(x, y)$  iar  $a_i$ , cu  $i = 1, \dots, 6$ , reprezintă parametrii modelului afin (vezi ecuația 3.19).

Mișcarea globală pentru o anumită imagine a secvenței, poate fi exprimată ca un vector de parametri,  $\phi = (a_1, a_2, \dots, a_6)$ , ce poate fi la rândul său

---

<sup>7</sup>vezi explicația de la pagina 83.

<sup>8</sup>"artefactele" de compresie sunt rezultatul eliminării informației utile în cazul compresiei cu pierderi. Acestea pot fi vizibile, fie ca o alterare a tranzițiilor graduale din imagine, fie ca zgomot pe contururile obiectelor, sau în cazul compresiei pe blocuri de pixeli, ca o conturare a acestora în imagine (efect de tablă de sah).

estimat pe baza câmpului vectorial de mișcare folosind metoda celor mai mici pătrate. [Kim 04] propune reprezentarea informației de mișcare cu vectori de parametri ce sunt exprimați în funcție de diversele tipuri de mișcări ale camerei video. Fiecare imagine a secvenței va fi astfel caracterizată de un vector,  $\phi_c = (pan, tilt, div, rot, hyp)$ , unde cei cinci parametri folosiți au următoarea semnificație:

- *pan* și *tilt* reprezintă mișcarea translatională în plan orizontal și respectiv vertical,
- *div* reprezintă mișcarea de tip "zoom",
- *rot* reprezintă mișcarea de rotație,
- *hyp* reprezintă ceea ce autorii numesc "flux hiperbolic", ce corespunde situațiilor de predominanță a mișcării de obiecte.

Aceștia sunt definiți în concordanță cu modelul afin de mișcare, în felul următor:

$$pan = a_1 \quad (3.34)$$

$$tilt = a_4 \quad (3.35)$$

$$div = \frac{1}{2}(a_2 + a_6) \quad (3.36)$$

$$rot = \frac{1}{2}(a_5 - a_3) \quad (3.37)$$

$$hyp = \frac{1}{4}(|a_2 - a_6| + |a_3 + a_5|) \quad (3.38)$$

Diversele tipuri de mișcări ale camerei video, respectiv "tracking", "tilting", rotație și "zoom", sunt mai departe determinate pe baza filtrării cu un anumit prag a acestor vectori de parametri.

O abordare similară ce folosește parametrizarea modelului afin de mișcare a fluxului MPEG este propusă în [Kramer 05]. Algoritmul propus are o performanță de trei, până la patru ori, mai rapidă decât timpul real<sup>9</sup>.

Metoda propusă în [Lee 02b] detectează tipul de mișcare al camerei video comparând diversele modele de mișcare prezente în secvență cu un anumit set de modele predefinite. Într-o primă etapă, câmpul de mișcare este extras la nivel de imagine din fluxul MPEG. Vectorii de mișcare astfel obținuți

---

<sup>9</sup>prin performanță în timp real a unei metode de analiză a conținutului secvențelor de imagini, înțelegem o rată de prelucrare de aproximativ 25 de imagini pe secundă (în standardul European). Metoda respectivă poate fi astfel aplicată în timp ce secvența este vizualizată fără a perturba continuitatea acesteia.

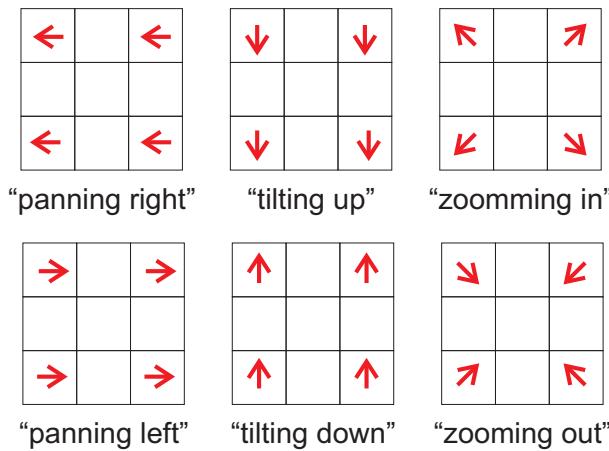


Figura 3.13: Exemple de modele predefinite de mișcare ale camerei video propuse în [Lee 02b] (săgețile indică orientarea vectorilor de mișcare).

sunt împărțiți la nivel de bloc de pixeli în două categorii: vectori de mișcare ce aparțin obiectelor din scenă, și respectiv, vectori de mișcare ce aparțin fundalului imaginii. Folosind aceste informații, mișcările camerei video sunt determinate în funcție de disponerea acestora în nouă regiuni disjuncte ale imaginii (vezi Figura 3.13). Similaritatea mișcării dintre diversele regiuni ale imaginii este exprimată ca distanță între histogramele de fază a vectorilor de mișcare conținuți în regiunea respectivă.

Principalul avantaj al metodelor de analiză a mișcării globale a camerei video ce folosesc domeniul comprimat este în primul rând timpul de calcul. Metodele din această categorie permit obținerea de performanțe superioare prelucrării în timp real. Totuși, precizia câmpului de mișcare este direct proporțională cu gradul de compresie al fluxului video folosit. Deseori, vectorii de mișcare obținuți direct din fluxul video modelează eronat mișcarea reală prezentă în sevență. Pentru a ameliora precizia, o soluție constă în decomprimarea datelor până la un anumit nivel de detaliu și re-estimarea mai precisă a mișcării folosind informația spațio-temporală.

### 3.2.2 Analiza mișcării în domeniul spațio-temporal

Metodele existente de analiză a mișcării camerei video ce folosesc informația spațio-temporală sunt foarte variate. Astfel, întâlnim o vastă diversitate de tehnici de detectie precum:

- metode ce folosesc modele predefinite de mișcare,

- metode de analiză a volumelor spațio-temporale,
- metode de clasificare cu rețele neuronale,
- metode ce folosesc descompunerea ”wavelet”<sup>10</sup>, etc.

De exemplu, [Akutsu 92] propune pentru detecția mișcării camerei video compararea repartiției vectorilor de mișcare din imagine cu modele predefinite, comparație ce este realizată de această dată în spațiul transformatei Hough<sup>11</sup>. O abordare similară ce folosește modele predefinite de mișcare este propusă în [Ionescu 07b]. Aceasta vizează detecția mișcărilor translational, de rotație și de tip ”zoom-in/zoom-out”, precum și analiza a două situații particulare de mișcare, care sunt absența mișcării și respectiv discontinuitatea mișcării. Spre deosebire de metodele existente, modelele de mișcare propuse în [Ionescu 07b] iau în calcul posibilitatea estimării eronate a vectorilor de mișcare precum și situațiile de confuzie generate de aceasta (în care mișcări diferite pot furniza modelele de mișcare similare). Un exemplu este ilustrat în Figura 3.14. Informația de mișcare este obținută în urma unei estimări pe blocuri de pixeli. Aceasta este comparată cu modelele predefinite pe baza unui set de reguli de decizie. Regulile sunt aplicate la nivel de vectori de orientare medii, ce sunt estimări în nouă regiuni disjuncte ale imaginii (vezi Figura 3.14).

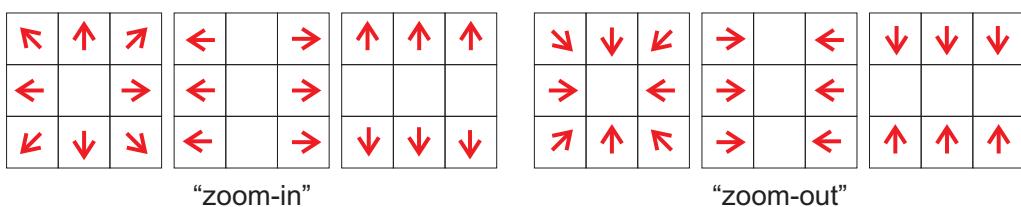


Figura 3.14: Exemplu de modele de mișcare pentru mișcarea camerei video de tip ”zoom” [Ionescu 07b] (unui anumit tip de mișcare îi pot corespunde în realitate mai multe modele de mișcare).

<sup>10</sup>O ”undisoară” (”wavelet”) reprezintă o funcție matematică ce este folosită la descompunerea unei anumite funcții sau a unui semnal continuu, în componente frecvențiale ce sunt apoi studiate la o rezoluție ce corespunde scalei acestora. Funcțiile ”wavelet” sunt copii scalate și translatate ale unei forme de undă finită sau cu atenuare rapidă, numită și funcție ”wavelet” de bază.

<sup>11</sup>Transformata Hough este o tehnică generală de identificare a orientării și a poziției anumitor tipuri de forme în imaginile digitale (de exemplu: linii, cercuri, elipse, etc.). Formele sunt determinate ca intersecții ale unor curbe sau plane în spațiul transformatei Hough, ce este determinat prin parametrizarea acestora.

O abordare inedită este propusă în [Ngo 00] unde caracterizarea mișcării camerei video și a obiectelor este realizată pe baza analizei de volume spațio-temporale extrase din imagini. O secvență de imagini poate fi văzută ca fiind un volum 3D în care primele două dimensiuni sunt date de dimensiunile spațiale  $(x, y)$  iar a treia dimensiune este timpul  $t$ . Dintr-un alt punct de vedere, acest volum poate fi reprezentat ca fiind un ansamblu de straturi temporale 2D, fie orizontale (plan  $(x, t)$ ), fie verticale (plan  $(y, t)$ ). În spațiul astfel format, diversele tipuri de mișcări ale camerei video vor fi indicate prin prezența modelelor orientate (vezi Figura 3.15). Caracterizarea acestora este realizată prin calculul a ceea ce autorii numesc histograme de tensori<sup>12</sup>, unde tensorii sunt furnizați de derivata parțială după cele trei axe,  $ot$ ,  $oX$  și  $oY$ .

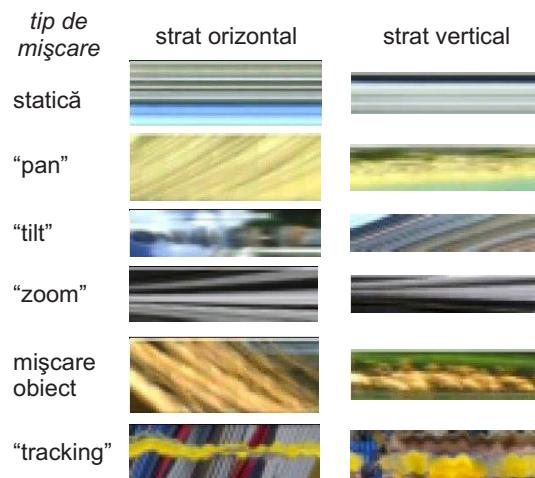


Figura 3.15: Modele orientate pentru diverse tipuri de mișcări [Ngo 00].

Raportat la metodele de analiză ce folosesc domeniul comprimat, analiza mișcării camerei video în domeniul spațio-temporal al imaginii, are o complexitate de calcul mai semnificativă, dar rezultatele obținute sunt mai precise. Domeniul spațio-temporal are avantajul de a furniza o mai mare diversitate de informații decât coeficienții MPEG.

### 3.3 Concluzii

În acest capitol am realizat o trecere în revistă a metodelor de analiză și prelucrare a uneia dintre informațiile definitorii ale unei secvențe de imagini, și anume *conținutul de mișcare*. Metodele existente folosesc ca punct de

---

<sup>12</sup>vezi explicația de la pagina 85.

plecare pentru analiză estimarea câmpului de mișcare al pixelilor din imagine. Aceasta este fie disponibil ”a priori”, cum este cazul fluxului MPEG, fie este determinat pe baza informației spațio-temporale.

În funcție de precizia dorită, metodele de estimare a mișcării se împart în două categorii. Pe de o parte sunt metodele ce propun un câmp vectorial dens, calculat la nivel de pixel, metode ce sunt în general bazate pe estimarea fluxului optic în imagine. În acest caz, complexitatea de calcul este direct proporțională cu precizia câmpului rezultat. Pe de altă parte sunt metodele ce realizează estimarea la nivel de blocuri de pixeli. De-a lungul timpului, acestea s-au dovedit a furniza cel mai bun compromis între complexitatea de calcul și precizia câmpului vectorial furnizat. Din acest motiv, metodele de estimare pe blocuri de pixeli au fost preferate de majoritatea standardelor de compresie video existente, precum MPEG, H.261, etc., unde timpul de prelucrare este critic.

Calitatea analizei conținutului de mișcare al unei secvențe de imagini, putem spune că este dependentă direct de precizia și de calitatea metodei de estimare a mișcării folosite. În acest sens, metodele de estimare a mișcării sunt alese în funcție de tipul aplicației. De exemplu, pentru segmentarea obiectelor de interes din scenă, ce sunt reprezentate în imagine cu regiuni restrâns de pixeli, este preferabil un câmp de mișcare dens, care să ofere maximum de informație. Pe de altă parte, în cazul analizei globale a mișcării, precum analiza mișcării camerei video, este preferabilă o metodă mai rapidă, cu un nivel de detaliu mai scăzut, precum estimarea la nivel de blocuri de pixeli. De asemenea, în cazul aplicațiilor în care timpul de prelucrare este limitat sau critic, se poate opta pentru recuperarea informației de mișcare direct din fluxul video, dacă acesta este disponibil.

Principalele aplicații ale estimării mișcării în secvențele de imagini sunt pe de-o parte analiza mișcării locale, ce este folosită de regulă pentru a segmenta și caracteriza proprietățile și traectoria anumitor obiecte de interes în scenă. Pe de altă parte, estimarea mișcării este folosită pentru caracterizarea mișcării globale a scenei, precum detectia diverselor mișcări ale camerei video, sau la un nivel semantic superior, a anumitor tehnici de filmare.

În concluzie, informația de mișcare este un parametru deloc neglijabil pentru analiza conținutului unei secvențe de imagini, fiind însuși motivul pentru care acestea există. Spre deosebire de alte surse de informație vizuală, precum imaginile statice, mișcarea face ca informația furnizată de secvențele de imagini să fie mai aproape de realitate.

## CAPITOLUL 4

---

### Analiza de culoare

---

**Rezumat:** *Informația de culoare joacă un rol important în percepția informației vizuale. Aceasta ne permite înțelegerea proprietăților fizice ale obiectelor ce ne înconjoară, precum și interacția cu acestea prin senzațiile de culoare ce ne sunt transmise. În ciuda faptului că informația fundamentală a unei secvențe de imagini este dată de conținutul de mișcare, din punct de vedere fiziologic, sistemul vizual uman este mult mai sensibil la schimbările de culoare. În acest capitol ne vom focaliza pe descrierea diverselor modalități de reprezentare și de caracterizare a conținutului de culoare, atât din punct de vedere sintactic cât și perceptual.*

Unul dintre cele mai importante simțuri, poate chiar cel mai important, este vederea. Simțim, explorăm și înțelegem lumea înconjurătoare folosindu-ne de percepția vizuală. Fiecărui obiect sau entitatei cu care interacționăm îi creăm mai întâi o imagine mentală a culorilor sale reprezentative, astfel: cerul este albastru, pădurea este verde, nisipul este galben, și aşa mai departe.

Acest mecanism ne facilitează recunoașterea și identificarea obiectelor ce sunt similară. Mai mult, anumite culori individuale sau aranjamente de culori ne crează senzații particulare, ca de exemplu: albastrul ne dă senzația de rece, portocaliu ne transmite o senzație de cald, negru și alb crează un contrast vizual, roșu în abundență dă o senzație de disconfort, etc. Astfel, pe

lângă informația vizuală furnizată de culoarea în sine, o importanță ridicată în percepția lumii înconjurătoare o au diversele relații ce pot exista între culori.

Inspirate de lumea reală, cercetările din domeniile vederii asistate de calculator și a prelucrării de imagini încearcă să reproducă aceste simțuri umane pentru a dezvolta sisteme capabile de a furniza o înțelegere și interpretare automată a informației vizuale. Culoarea, în particular, a fost exploatată intensiv de mai bine de trei decenii, pentru a descrie percepția vizuală a imaginilor [Smeulders 00].

În cazul sistemelor de indexare automată după conținut, informația de culoare a fost exploatată de sine stătător, aproape exclusiv în sistemele de indexare a imaginilor fixe [Bimbo 99] [Smeulders 00], aceasta fiind definitorie pentru conținutul static spațial al imaginii. În domeniul sistemelor de indexare după conținut a secvențelor de imagini, în care intervine și dimensiunea temporală, nu există multe studii care să se focalizeze pe caracterizarea în termeni de distribuție de culoare a conținutului secvenței. Acest lucru se datorează în principal faptului că doar culoarea în sine, în cele mai multe cazuri, nu este suficientă pentru înțelegerea conținutului dinamic al secvenței.

Astfel, majoritatea metodelor existente se folosesc de informația de culoare, fie în colaborare cu alte informații, precum mișcarea, distribuția de plane, textură, etc., fie pentru a caracteriza anumite proprietăți locale, la nivel de imagine, ale obiectelor de interes. Cu toate acestea, raportat la principala informație furnizată de o secvență de imagini, și anume informația de mișcare, putem spune că din punct de vedere fiziologic, sistemul uman de percepție vizuală este mai sensibil la schimbări de culoare decât la prezența mișcării.

Informația de culoare, în cazul secvențelor de imagini, este conținută la nivelul fiecărui cadru. Spre deosebire de imaginile statice, aceasta are o evoluție temporală dată de evoluția conținutului din imagine. Astfel, metodele existente de analiză a culorii în secvențele de imagini folosesc ca punct de plecare tehnici specifice analizei imaginilor statice ce sunt apoi extinse la dimensiunea temporală. De asemenea, pentru a evidenția anumite proprietăți ale culorilor, acestea sunt reprezentate în diverse spații de culoare. Spațiile de culoare au fost concepute special în funcție de necesitățile de prelucrare.

În cele ce urmează, vom face o trecere în revistă a diverselor tehnici folosite în secvențele de imagini pentru a analiza conținutul de culoare. În acest sens, mai întâi vom prezenta diversele modalități existente de reprezentare a culorilor, precum și tehniciile de caracterizare a percepției culorii la nivel de imagine.

## 4.1 Spațiile de culoare

La nivel de imagine, fiecare culoare este reprezentată ca fiind un punct într-un spațiu de referință, de regulă tridimensional, care constituie ceea ce numim *spațiul de culoare*. Ansamblul culorilor ce pot fi reprezentate într-un anumit spațiu de culoare poartă numele de *gamut de culoare*. Particularizând la nivel de imagine, echivalentul gamutului de culoare este *paleta de culoare* a imaginii, aceasta fiind definită ca totalitatea culorilor, diferite, utilizate de imagine. Astfel, fiecare imagine, în funcție de conținutul acesteia, are o anumită paletă de culoare. Paleta de culoare poate fi *general valabilă* (unică) pentru un set finit de imagini sau *adaptivă*, și astfel specifică fiecărei imagini în parte.

Conceperea diverselor spații de reprezentare a culorilor existente a fost motivată de însăși existența sistemului vizual uman, astfel fiecare culoare fiind reprezentată în general după trei componente. În funcție de natura acestora, spațiile de culoare existente se împart în următoarele categorii [Trémeau 04]:

- **sisteme de culori primare:** acestea reprezintă culorile în funcție de trei culori alese arbitrar, numite și culori primare,
- **sisteme pe bază de luminanță-crominanță:** acestea au proprietatea de a separa informația de luminanță (intensitate luminoasă) de componentele de crominanță (culoare),
- **sisteme perceptuale:** acestea sunt sisteme uniforme din punct de vedere al percepției vizuale a culorilor,
- **sisteme de axe independente:** acestea au ca scop reprezentarea culorilor într-un spațiu în care componentele de culoare sunt decorelate.

Totuși, în ciuda existenței unei diferențieri pe categorii a sistemelor de culori existente, în mod ușual anumite sisteme pot fi încadrate ca aparținând mai multor categorii deodată. Cum fiecare categorie în parte prezintă caracteristici sau proprietăți ce sunt necesare în anumite etape de prelucrare, s-a căutat conceperea de sisteme de culoare, putem spune hibride, care să reunescă aceste proprietăți pentru a răspunde simultan acestor cerințe. Se poate observa că nu există un sistem universal de culoare care să fie general valabil pentru orice aplicație, mai mult, în multe situații, nici unul dintre sistemele existente nesatisfacând cerințele de prelucrare [Trémeau 04].

În cele ce urmează vom face o trecere în revistă a particularităților sistemelor de culoare reprezentative din fiecare dintre categoriile enumerate anterior.

### 4.1.1 Sisteme de culori primare

Sistemele din această categorie folosesc pentru a reprezenta culorile principiul *trivariantei vizuale*. Astfel, fiecare culoare existentă poate fi reproducă vizual identic, în anumite condiții de observare, prin amestecul matematic, în proporții unice, a trei culori numite și *culori primare*. Aceste culori primare sunt alese arbitrar și au proprietatea că nici una dintre ele nu poate fi reproducă ca un amestec al celorlalte două [Young 02].

Astfel, fiecare culoare  $C$  poate fi exprimată pe baza principiului trivariantei vizuale ca fiind dată de ecuația:

$$C = p_1 \cdot C_1 + p_2 \cdot C_2 + p_3 \cdot C_3 \quad (4.1)$$

unde  $C_i$ ,  $i \in \{1, 2, 3\}$ , reprezintă cele trei culori primare iar coeficienții  $p_i$  reprezintă ponderea acestora la definirea culorii  $C$ .

Din această categorie de sisteme, putem menționa ca reprezentative spațiile de culoare RGB, CMY și XYZ.

#### Spațiul de culoare RGB

Spațiul de culoare RGB este un spațiu aditiv ce a fost inspirat de modalitatea fiziologică de reprezentare a culorilor în sistemul vizual uman. Fiecare culoare existentă,  $C$ , este exprimată ca un amestec aditiv a trei culori primare, și anume roșu ("Red"), verde ("Green") și respectiv albastru ("Blue"), astfel:

$$C = p_1 \cdot R + p_2 \cdot V + p_3 \cdot B \quad (4.2)$$

Totuși, această ecuație caracterizează doar nuanța culorii rezultate din sinteza aditivă și în nici un caz intensitatea sa luminoasă [Délibéré 89]. Pentru a caracteriza luminanța culorii  $C$  este necesar să cunoaștem luminanța fiecărei dintre cele trei culori primare folosite,  $R$ ,  $G$  și  $B$ . În funcție de aceasta, există mai multe sisteme bazate pe cele trei culori primare, astfel: sistemul de culori primare CIE definit luând ca referință iluminantul  $E^1$ , sistemul de culori primare NTSC definit luând ca referință iluminantul  $C$  sau sistemul de culori primare EBU ("European Broadcasting Union") ce folosește ca referință iluminantul  $D_{65}$  [Tréneau 04].

---

<sup>1</sup>sursele de lumină sunt caracterizate de repartiția spectrală a energiei sau altfel spus prin cantitatea de energie emisă pe un anumit interval de lungime de undă. Anumite surse de lumină ce corespund unor situații uzuale au fost normalizate de CIE ("Commission Internationale de l'Eclairage") sub numele de iluminanți. Astfel, de exemplu iluminantul  $E$  corespunde unei lumini de energie egală, iluminantul  $C$  corespunde unei lumini medii de zi cu o temperatură de culoare de aproximativ  $6770K$ , iluminantul  $D_{65}$  corespunde în mare unei lumini ce provine de la un cer albastru, cu  $3/5$  nori albi, măsurată în jurul orei 10 dimineață în luna septembrie.

Sistemul RGB este utilizat cu precădere în majoritatea dispozitivelor "hardware" de reproducere a culorilor ce sunt destinate marelui public, precum aparatele foto digitale, monitoarele CRT și LCD, imprimantele color, etc., și stă la baza arhitecturii sistemelor de prelucrare actuale. Din această cauză, trecerea la alte spații de culoare se face în cele mai multe cazuri pornind de la valorile RGB pe baza unei transformări matematice.

Gamutul de culoare al spațiului RGB este reprezentat de un cub în sistemul de coordonate  $R$ ,  $G$  și  $B$  (vezi Figura 4.1). În acesta, negrul se găsește la originea sistemului iar niveliurile de gri se găsesc pe diagonala principală, fiind reprezentate de valori egale ale componentelor  $R$ ,  $G$  și respectiv  $B$ .

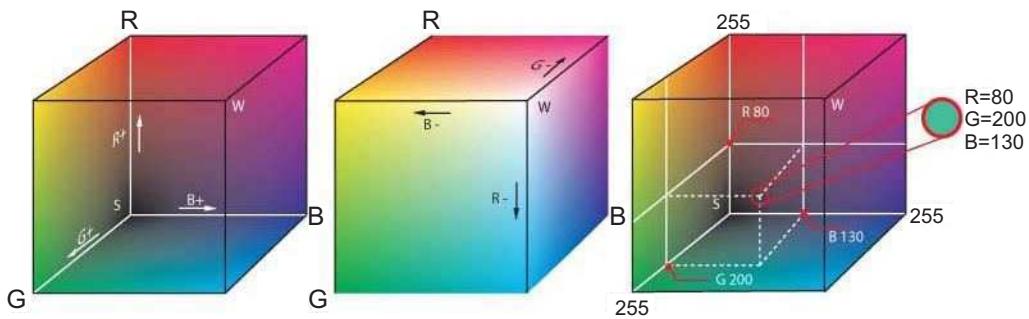


Figura 4.1: Cubul de culoare RGB (R-roșu, G-verde, B-albastru, W-alb, sursă imagine Wikipedia "[http://en.wikipedia.org/wiki/RGB\\_color\\_space](http://en.wikipedia.org/wiki/RGB_color_space)").

Spațiul de culoare astfel obținut este un spațiu cu variație neuniformă în care distanța Euclidiană dintre două culori nu corespunde distanței percepționale dintre acestea.

### Spațiul de culoare CMY

Spațiul de culoare CMY, spre deosebire de RGB, este obținut prin ceea ce se numește sinteză subtractivă a culorilor primare turcoaz (Cyan), magenta și galben (Yellow). Culorile primare folosite de cele două spații de culoare, RGB și respectiv CMY, sunt complementare. Astfel, albastrul este complementar lui galben care el însuși este obținut ca un amestec de roșu și verde, roșu este complementar lui turcoaz care este obținut din albastru și verde, iar verde este complementar lui magenta care este obținut din roșu și albastru (vezi Figura 4.2).

Trecerea de la spațiul RGB la CMY este dată de următoarea ecuație:

$$\begin{bmatrix} C \\ M \\ Y \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (4.3)$$

unde  $(C, M, Y)$  reprezintă culoarea  $(R, G, B)$  în spațiul CMY.

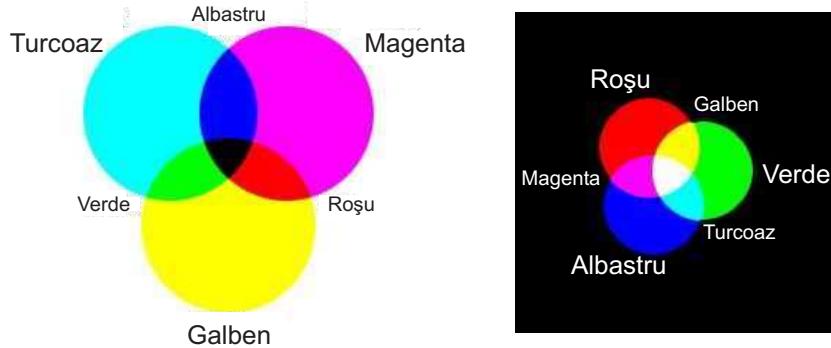


Figura 4.2: Comparație între amestecul substractiv (CMY) și aditiv (RGB), sursă imagine "[http://akarostost.com/www.whoride.com/roosto/physics/sub\\_vs\\_add.html](http://akarostost.com/www.whoride.com/roosto/physics/sub_vs_add.html)".

Spațiul de culoare CMY, în varianta CMYK, la care se adaugă separat negrul (K-”Key”), este folosit cu predilecție pentru imprimarea pe suport fizic, ca de exemplu în tipografii, deoarece modalitatea sa de constituire descrie însăși procesul de imprimare. Adăugarea separată a negrului este motivată în principal de faptul că negrul generat prin amestecul fizic al celor trei culori primare,  $C$ ,  $M$  și  $Y$ , nu oferă o calitate suficientă imprimării, cum este cazul imprimării detaliilor fine ale fonturilor de caractere.

### Spațiul de culoare XYZ

Sistemul XYZ este sistemul de referință colorimetrică definit de CIE (“Commission Internationale de l’Eclairage”). Acesta este determinat pornind de la spațiul de culoare RGB printr-o transformare liniară, transformare ce ia în calcul nu numai coordonatele tricromatice ale fiecărei culori primare,  $R$ ,  $G$  și  $B$ , precum și coordonatele tricromatice ale albului de referință  $W$ . Astfel, orice culoare  $(R, G, B)$  poate fi exprimată în coordonate  $(X, Y, Z)$  folosind următoarea ecuație:

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} 2.7690 & 1.7518 & 1.1300 \\ 1.0000 & 4.5907 & 0.0601 \\ 0.0000 & 0.0565 & 5.5943 \end{bmatrix} \times \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (4.4)$$

unde  $(2.7690, 1.0000, 0.0000)$ ,  $(1.7518, 4.5907, 0.0565)$  și  $(1.1300, 0.0601, 5.5943)$  reprezintă coordonatele XYZ ale celor trei culori primare: roșu, verde și respectiv albastru.

În funcție de iluminantul de referință folosit, trecerea la spațiul XYZ se face folosind matrice de transformări diferite, astfel:

- transformarea spațiului NTSC RGB folosind iluminantul C, care este și transformarea cea mai utilizată în domeniul prelucrării imaginilor color, este dată de:

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} 0.607 & 0.174 & 0.200 \\ 0.299 & 0.587 & 0.114 \\ 0.000 & 0.066 & 1.116 \end{bmatrix} \times \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (4.5)$$

- transformarea spațiului CIE RGB folosind iluminantul A<sup>2</sup> este dată de:

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} 0.892 & 0.330 & 0.083 \\ 0.322 & 0.863 & 0.004 \\ 0.000 & 0.011 & 0.409 \end{bmatrix} \times \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (4.6)$$

- transformarea spațiului CIE RGB folosind iluminantul C este dată de:

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} 0.166 & 0.125 & 0.093 \\ 0.060 & 0.327 & 0.005 \\ 0.000 & 0.004 & 0.460 \end{bmatrix} \times \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (4.7)$$

Spațiul XYZ este definit astfel încât să respecte o serie de constrângeri, și anume:

- orice culoare fizică monocromatică trebuie să fie caracterizată de valori tristimulus pozitive,
- componența  $Y$  trebuie să fie o măsură a intensității luminoase,
- pentru o lumină albă, valorile tristimulus trebuie să fie egale.

Matricea transformării în spațiul XYZ nu este o matrice unitară<sup>3</sup> și astfel transformarea nu reprezintă o rotație a cubului RGB cu tot cu sistemul de coordonate, gamutul de culoare XYZ fiind în acest caz un paralelipiped înclinat [Vertan 08].

---

<sup>2</sup>iluminant CIE standard de tip A corespunde unei lumini de interior tipice furnizate de o iluminare cu filament din tungsten ce are o temperatură de culoare de aproximativ 2856K.

<sup>3</sup>O transformare se numește unitară, dacă matricea transformării inverse este egală cu transpusa conjugată a matricei transformării directe.

Cum componenta  $Y$  reflectă intensitatea luminoasă, atunci informația cromatică a unei culori poate fi exprimată în funcție de doi parametri derivați, notați  $x$  și  $y$ , ce sunt date de relațiile următoare:

$$x = \frac{X}{X + Y + Z} \quad (4.8)$$

$$y = \frac{Y}{X + Y + Z} \quad (4.9)$$

Spațiul de culoare derivat, definit de cei doi parametri,  $x$  și  $y$ , precum și de componenta de intensitate luminoasă  $Y$  este cunoscut sub numele de CIE xyY și este folosit la scară largă pentru specificarea culorilor. Diagrama de cromaticitate  $xy$  este prezentată în Figura 4.3.

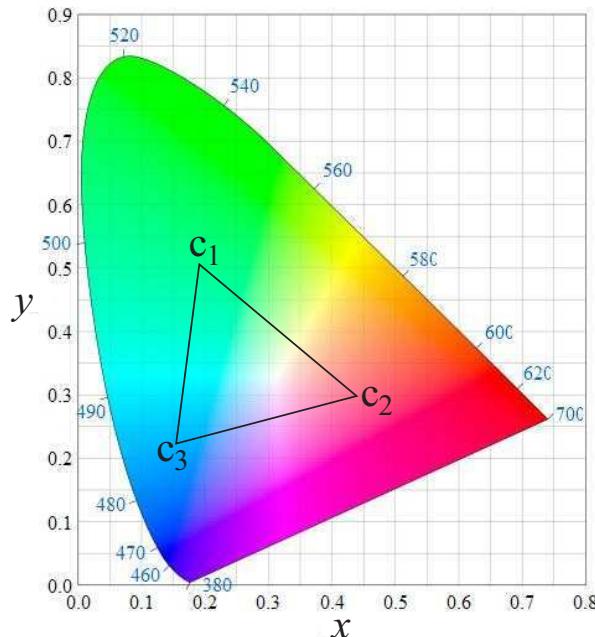


Figura 4.3: Diagrama de cromaticitate CIE (valorile de pe curba ce marchează marginea graficului reprezintă lungimi de undă exprimate în nanometrii, sursă Wikipedia "[http://en.wikipedia.org/wiki/CIE\\_1931\\_color\\_space](http://en.wikipedia.org/wiki/CIE_1931_color_space)").

Reprezentată în acest fel, diagrama de cromaticitate pune în evidență o serie de proprietăți interesante ale spațiului de culoare XYZ, astfel:

- aceasta prezintă toate culorile vizibile de către o persoană normală, sau ceea ce se numește *gamutul sistemului vizual uman*. Conturul curbat al

gamutului corespunde astfel culorilor monocromatice (vezi lungimile de undă din Figura 4.3). Linia dreaptă din partea de jos a gamutului este numită linia culorilor purpuri. Culorile mai puțin saturate se găsesc în interiorul graficului având albul în centru.

- toate culorile cromatice vizibile au valorile  $X$ ,  $Y$ , și  $Z$  nenegative.
- alegând două culori diferite arbitrar din diagrama de cromaticitate, toate culorile ce pot fi formate cu acestea se regăsesc pe linia ce unește cele două puncte. Similar, în cazul a trei puncte diferite, toate culorile ce pot fi formate cu acestea se găsesc în interiorul triunghiului format de acestea (vezi Figura 4.3), și aşa mai departe. În concluzie, gamutul de culoare reprezentat în diagrama de cromaticitate este o formă geometrică convexă.
- distanța dintre culori, în diagrama de cromaticitate  $xy$ , nu corespunde distanței perceptuale dintre culori.

#### 4.1.2 Sisteme pe bază de luminanță-crominanță

Sistemele de culori din această categorie au proprietatea de a separa componentă de intensitate luminoasă de componente cromatice. Din această categorie putem enumera sistemele de tip  $YC_bC_r$  și sistemele antagoniste [Tréneau 04].

##### Spațiile de culoare de tip $YC_bC_r$

Acest spațiu de culoare a fost creat inițial pentru a asigura compatibilitatea dintre televizoarele color și cele alb-negru. Astfel, un semnal color de tip  $YC_bC_r$  putea fi vizionat și pe un receptor alb-negru datorită separării componente de luminanță,  $Y$ , de cele cromatice,  $C_b$  și  $C_r$ . Trecerea de la spațiul RGB la spațiul  $YC_bC_r$  se face printr-o simplă transformare liniară ai cărei coeficienți diferă de la un standard de televiziune la altul.

Cele două valori de crominanță,  $C_b$  și respectiv  $C_r$ , ale unei anumite culori exprimată în coordonate RGB, pot fi calculate în funcție de patru coeficienți,  $a_1, a_2, b_1$  și  $b_2$  (specifici standardului considerat) și de informația de luminanță  $Y$ , în felul următor:

$$C_b = a_1 \cdot (R - Y) + b_1 \cdot (B - Y) \quad (4.10)$$

$$C_r = a_2 \cdot (R - Y) + b_2 \cdot (B - Y) \quad (4.11)$$

În realitate, există mai multe sisteme de tip  $YC_bC_r$ , ca de exemplu sistemul YIQ ce corespunde standardului NTSC<sup>4</sup>, sistemul YUV ce corespunde standardului PAL<sup>5</sup> și sistemul  $YD_bD_r$  ce corespunde standardului SECAM<sup>6</sup>. Dintre acestea, sistemul cel mai frecvent utilizat în domeniul imagisticii color este sistemul YIQ.

Trecerea de la spațiul de culoare RGB la spațiul YIQ este dată de transformarea următoare:

$$\begin{bmatrix} Y \\ I \\ Q \end{bmatrix} = \begin{bmatrix} 0.299 & 0.587 & 0.114 \\ 0.596 & -0.273 & -0.322 \\ 0.212 & -0.522 & 0.315 \end{bmatrix} \times \begin{bmatrix} R' \\ G' \\ B' \end{bmatrix} \quad (4.12)$$

unde  $(R', G', B')$  reprezintă coordonatele  $(R, G, B)$  după corecția de gamma<sup>7</sup> următoare:

$$\begin{bmatrix} R' \\ G' \\ B' \end{bmatrix} = \begin{bmatrix} R^{1/2.2} \\ G^{1/2.2} \\ B^{1/2.2} \end{bmatrix} \quad (4.13)$$

În cazul în care se dorește o caracterizare mai riguroasă a informației de culoare, spațiul YIQ poate fi descris și în coordonate cilindrice [A.R. Weeks 95], astfel:

$$H = \arctan \left( \frac{B' - Y}{R' - Y} \right) \quad (4.14)$$

$$S^2 = (R' - Y)^2 + (B' - Y)^2 \quad (4.15)$$

unde  $H$  reprezintă nuanța de culoare ("Hue") iar  $S$  saturăția culorii ("Saturation").

### Spațiile de culoare antagoniste

Spațiile de culoare din această categorie sunt fundamentate pe teoria culorilor opuse elaborată de Young și Hering. Aceasta afirmă că informația

---

<sup>4</sup>NTSC este prescurtarea pentru "National Television System Committee" și reprezintă un standard de codare a semnalului TV analogic color pe 525 de linii și la o cadență de 30 imagini pe secundă.

<sup>5</sup>vezi explicația de la pagina 9.

<sup>6</sup>SECAM este prescurtarea pentru "Séquentiel Couleur à Mémoire" și reprezintă un standard de codare a semnalului TV analogic color pe 625 de linii și la o cadență de 25 de imagini pe secundă.

<sup>7</sup>datorită faptului că sistemul vizual uman nu are un răspuns liniar la informația luminosă, în cazul în care luminanța este codată cu un număr redus de valori (de exemplu 256) este necesară utilizarea acestora în concordanță cu proprietățile sistemului vizual uman pentru a obține o eficiență optimală. Această operație este cunoscută sub numele de corecție de gama. Corecția de gama se poate exprima sub forma unei ecuații de tipul  $S = E^\gamma$ , unde  $E$  reprezintă semnalul căruia i se aplică corecția iar  $S$  este semnalul corectat.

de culoare recepționată de sistemul vizual uman este transmisă creierului sub forma a trei semnale, unul acromatic ce corespunde contrastului negru-alb, și două semnale cromatice, unul corespunzând contrastului verde-roșu și celălalt contrastului albastru-galben [Wyszecki 82].

Un exemplu este spațiul de culoare  $AC_1C_2$ , introdus pentru prima oară în [Faugeras 79], ce este definit de următorul sistem de ecuații:

$$A = 22.6 \cdot (0.612 \cdot \log L + 0.369 \cdot \log M + 0.019 \cdot \log S) \quad (4.16)$$

$$C_1 = 64(\log L - \log M) \quad (4.17)$$

$$C_2 = 10(\log L - \log S) \quad (4.18)$$

unde  $(L, M, S)$  reprezintă coordonatele culorii în sistemul primar LMS<sup>8</sup> iar  $(A, C_1, C_2)$  reprezintă noile coordonate în spațiul  $AC_1C_2$ .

Trecerea la spațiul LMS se poate realiza pornind de la spațiul de culoare XYZ prin transformarea următoare:

$$\begin{bmatrix} L \\ M \\ S \end{bmatrix} = \begin{bmatrix} 0.15514 & 0.54312 & -0.03286 \\ -0.15514 & 0.45684 & 0.03286 \\ 0.0 & 0.0 & 0.00801 \end{bmatrix} \times \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} \quad (4.19)$$

În ciuda proprietăților interesante ale spațiilor de culoare antagoniste, marea majoritate a acestora sunt puțin utilizate în imagistica color datorită faptului că sunt fie abordări experimentale în curs de validare, fie modele foarte simplificate [Tréneau 04].

#### 4.1.3 Sisteme perceptuale

Sistemele perceptuale de reprezentare a culorilor sunt sisteme uniforme din punct de vedere al percepției vizuale. Astfel, în acestea distanța matematică dintre două culori este proporțională cu distanța perceptuală dintre acestea. Din această categorie putem enumera ca cele mai reprezentative: spațiul de culoare  $L^*a^*b^*$ , spațiile de culoare geometrice și spațiul de culoare al lui Munsell.

##### Spațiul de culoare $L^*a^*b^*$ și $L^*u^*v^*$

Datorită popularității acestuia, spațiul de culoare  $L^*a^*b^*$  poate fi considerat ca unul dintre sistemele de referință al CIE folosit pentru evaluarea distanței

---

<sup>8</sup>sistemul de culori primare LMS este sistemul de referință în domeniul fiziolgiei sistemului vizual. Aceasta are avantajul de a folosi culori primare ce sunt în concordanță directă cu semnalul de intrare folosit de sistemul vizual uman, fiind singurul sistem primar ce reproduce activitatea reală a celulelor cu con din retina umană.

dintre culori [Tréneau 04]. Trecerea la sistemul de coordonate  $L^*a^*b^*$  se face de regulă pornind de la spațiul XYZ pe baza unei transformări, de această dată, neliniare ce ia în calcul și coordonatele tricromatice ale albului de referință,  $W$ , folosit.

Astfel, o culoare dată de tripletul  $(X, Y, Z)$  se poate exprima în coordonate  $L^*a^*b^*$  pe baza ecuațiilor următoare:

$$L^* = \begin{cases} 116 \cdot y^{1/3} - 16 & \text{dacă } y > 0.008856 \\ 903.3 \cdot y & \text{altfel} \end{cases} \quad (4.20)$$

$$a^* = 500 \cdot [f(x) - f(y)] \quad (4.21)$$

$$b^* = 200 \cdot [f(y) - f(z)] \quad (4.22)$$

unde:

$$x = \frac{X}{X_W}, \quad y = \frac{Y}{Y_W}, \quad z = \frac{Z}{Z_W} \quad (4.23)$$

$(X_W, Y_W, Z_W)$  reprezintă coordonatele XYZ ale albului de referință iar funcția  $f()$  este dată de ecuația următoare:

$$f(t) = \begin{cases} t^{1/3} & \text{dacă } t > 0.008856 \\ 7.787 \cdot t + 0.137931 & \text{altfel} \end{cases} \quad (4.24)$$

În spațiul  $L^*a^*b^*$ , fiind un spațiu perceptual, distanța perceptuală dintre două culori,  $C_1 = (L_1^*, a_1^*, b_1^*)$  și  $C_2 = (L_2^*, a_2^*, b_2^*)$ , poate fi evaluată direct pe baza distanței Euclidiene dintre coordonatele acestora, astfel:

$$\Delta E_{C_1, C_2}^2 = (\Delta L^*)^2 + (\Delta a^*)^2 + (\Delta b^*)^2 \quad (4.25)$$

unde  $\Delta E_{C_1, C_2}$  quantifică diferența dintre culorile  $C_1$  și respectiv  $C_2$  iar  $\Delta L^* = L_1^* - L_2^*$ ,  $\Delta a^* = a_1^* - a_2^*$  și  $\Delta b^* = b_1^* - b_2^*$ .

Gamutul de culoare al spațiului  $L^*a^*b^*$  este reprezentat ca o sferă de culoare (vezi Figura 4.4) în care spre polul nord sunt figurate culorile deschise, finalizându-se cu alb, iar spre polul sud culorile întunecate, finalizându-se cu negru.

Datorită construcției particulare, spațiul  $L^*a^*b^*$  poate fi încadrat ca aparținând la mai multe categorii de sisteme de reprezentare a culorilor, astfel acesta poate fi considerat ca un spațiu antagonist, deoarece componenta  $L^*$  pune în evidență contrastul negru-alb, componenta  $a^*$  contrastul verde-roșu iar componenta  $b^*$  contrastul albastru-galben (vezi Figura 4.4). De asemenea, spațiul  $L^*a^*b^*$  se încadrează și în categoria spațiilor de culoare de tip  $YC_bC_r$ , deoarece componenta de luminăță,  $L^*$ , este separată de cea de crominanță dată de  $a^*$  și respectiv  $b^*$ .

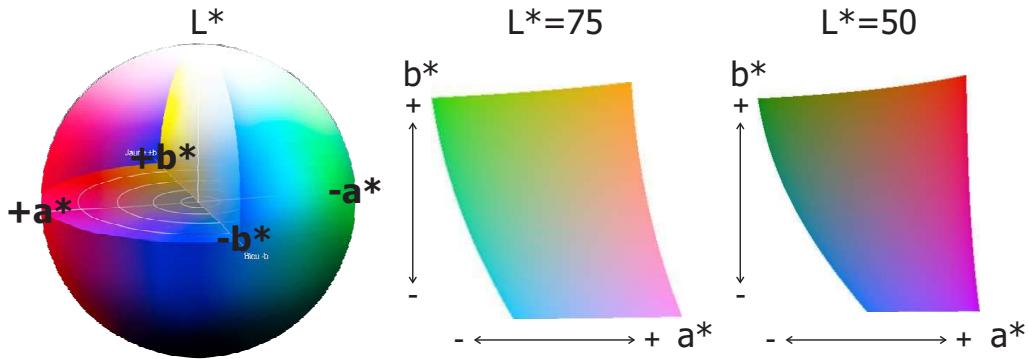


Figura 4.4: Sfera de culoare  $L^*a^*b^*$  (sursă Wikipedia ”[http://en.wikipedia.org/wiki/Lab\\_color\\_space](http://en.wikipedia.org/wiki/Lab_color_space)”).

Prin trecerea spațiului  $L^*a^*b^*$  în coordonate cilindrice obținem un nou spațiu de culoare în care culorile sunt grupate în funcție de nuanță și saturăție. Acesta este cunoscut sub numele de spațiu de culoare LCH și este dat de relațiile următoare:

$$H = \arctan\left(\frac{b^*}{a^*}\right) \quad (4.26)$$

$$C^2 = a^{*2} + b^{*2} \quad (4.27)$$

unde  $H$  reprezintă nuanță iar  $C$ , numit și ”chroma”, oferă informații despre saturăția culorii.

Tot în această categorie putem menționa un spațiu de culoare similar cu spațiul  $L^*a^*b^*$ , și anume spațiul  $L^*u^*v^*$ . Cu toate că spațiul  $L^*u^*v^*$  nu este analog spațiului  $L^*a^*b^*$ , pentru anumite intervale de culoare acesta furnizează rezultate similare [Agoston 87]. Dacă componenta  $L^*$  a spațiului  $L^*u^*v^*$  este calculată în același mod ca pentru  $L^*a^*b^*$  (vezi ecuația 4.20), componentele  $u^*$  și  $v^*$  sunt calculate în mod diferit, astfel:

$$u^* = 13 \cdot L^* \cdot (u' - u'_W) \quad (4.28)$$

$$v^* = 13 \cdot L^* \cdot (v' - v'_W) \quad (4.29)$$

unde

$$u' = \frac{4 \cdot X}{X + 15 \cdot Y + 3 \cdot Z} \quad (4.30)$$

$$v' = \frac{9 \cdot Y}{X + 15 \cdot Y + 3 \cdot Z} \quad (4.31)$$

iar  $(X_W, Y_W, Z_W)$  reprezintă coordonatele XYZ ale albului de referință folosit.

### Spațiile de culoare geometrice

Sistemele de culoare din această categorie folosesc o partitioare uniformă a culorilor în funcție de trei informații, și anume: nuanță, saturatie și luminanță. Diferența dintre acestea constă în principal în unitățile folosite la partitioare precum și în dinamica fiecărei dintre cele trei axe de culoare. Sistemele geometrice au fost concepute în scopul de a clasa culorile pe baza percepției psihovizuale a acestora, culorile fiind reperate de această dată într-un sistem de coordonate geometrice. Spațiile de culoare geometrice constituie reprezentări descriptive ale informației de culoare [Tréneau 04].

Cele mai reprezentative spații din această categorie sunt spațiile de culoare HSL (H-nuanță, S-saturație, L-luminozitate) și respectiv HSV (H-nuanță, S-saturație, V-valoare). O alternativă similară, mai puțin standardizată, a spațiilor HSL și HSV o constituie spațiile de culoare HSI și respectiv HSB, unde  $I$  reprezintă intensitatea iar  $B$  strălucirea culorii ("brightness").

Atât spațiul HSL cât și HSV reprezintă culorile din punct de vedere matematic folosind o reprezentare cilindrică (vezi Figura 4.5). Astfel, nuanța de culoare este dată de unghi, saturatia culorii este dată de distanța față de axa principală iar luminozitatea este dată de poziția de-a lungul axei principale. Nivelurile de gri se vor găsi astfel pe axa centrală, pornind de la negru și finalizând cu alb.

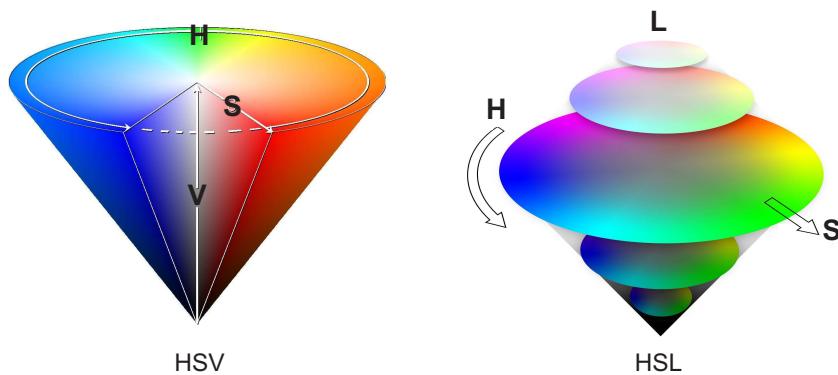


Figura 4.5: Spațiile de culoare HSV (con invers) și HSL (dublu con, sursă Wikipedia "[http://en.wikipedia.org/wiki/HSV\\_color\\_space](http://en.wikipedia.org/wiki/HSV_color_space)").

Cu toate că cele două tipuri de reprezentări sunt similare din punct de vedere al obiectivului urmărit, acestea diferă totuși din punct de vedere al abordării folosite. Gamutul de culoare al sistemului HSV reprezintă conceptual un con invers în care negrul se găsește în vîrful conului iar culorile de

saturație maximă se găsesc pe cercul bazei. Pe de altă parte, gamutul de culoare al sistemului HSL reprezintă conceptual un con dublu în care albul se găsește în vârful din Nord, negrul în vârful din Sud iar culorile de saturație maximă se găsesc pe cercul median.

Trecerea la spațiul HSV se face pornind de la reprezentarea RGB pe baza următoarelor ecuații:

$$h_{HSV} = \begin{cases} 0 & \text{dacă } max = min \\ 60 \cdot \frac{g-b}{max-min} + 0 & \text{dacă } max = r \text{ și } g \geq b \\ 60 \cdot \frac{g-b}{max-min} + 360 & \text{dacă } max = r \text{ și } g < b \\ 60 \cdot \frac{b-r}{max-min} + 120 & \text{dacă } max = g \\ 60 \cdot \frac{r-g}{max-min} + 240 & \text{dacă } max = b \end{cases} \quad (4.32)$$

$$s_{HSV} = \begin{cases} 0 & \text{dacă } max = 0 \\ \frac{max-min}{max} & \text{altfel} \end{cases} \quad (4.33)$$

$$v_{HSV} = max \quad (4.34)$$

unde  $r$ ,  $g$  și  $b$  reprezintă valorile normalize între 0 și 1 ale celor trei componente  $R$ ,  $G$  și respectiv  $B$ <sup>9</sup> iar  $max$  și  $min$  reprezintă valoarea maximă și respectiv minimă a componentelor normalize  $r$ ,  $g$ ,  $b$ . Valorile astfel obținute sunt normalize între 0 și 360 pentru componenta  $h$  și între 0 și 1 pentru componente  $s$  și  $v$ .

În ceea ce privește trecerea la spațiul HSL, componenta normalizată  $h$  este dată de aceeași ecuație ca și în cazul spațiului HSV (vezi ecuația 4.32), în schimb componentele normalize  $s$  și  $l$  sunt date de ecuațiile următoare:

$$s_{HSL} = \begin{cases} 0 & \text{dacă } max = min \\ \frac{max-min}{max+min} & \text{dacă } l \leq 0.5 \\ \frac{max-min}{2-(max+min)} & \text{dacă } l > 0.5 \end{cases} \quad (4.35)$$

$$l_{HSL} = \frac{1}{2}(max + min) \quad (4.36)$$

Principalul inconvenient al transformărilor bazate pe partitōnarea geometrică a spațiului de culoare constă în reversibilitatea limitată a acestora. Astfel, de cele mai multe ori la revenirea în spațiul inițial de culoare, de regulă RGB, culorile sunt diferite de cele originale. Acest lucru se datorează faptului că marea majoritate a metodelor de transformare sunt bazate pe calculul valorilor minime și respectiv maxime ale coordonatelor RGB (vezi ecuațiile 4.32).

---

<sup>9</sup>de regulă pentru valori RGB cuprinse între 0 și 255,  $r = \frac{R}{255}$ ,  $g = \frac{G}{255}$  și  $b = \frac{B}{255}$ .

### Spațiul de culoare al lui Munsell

Creat de profesorul Albert H. Munsell în prima decadă a secolului douăzeci, spațiul de culoare ce îi poartă numele a fost conceput în principal pentru a fi utilizat în domeniul artei. Acesta se bazează pe observații subiective ale raportului dintre culori și nu pe măsurarea directă a proprietăților de culoare [Bimbo 99]. Spațiul obținut este un spațiu perceptual, uniform, bazat pe amestecul subtractiv de culoare, în care culorile sunt reprezentate în funcție de modul în care sunt percepute de observatori umani.

Fiecare culoare este reprezentată în funcție de trei informații, și anume: luminanță ("Value"), nuanță ("Hue") și saturatie ("Chroma"), într-un sistem de reprezentare de tip cilindric. Astfel, nuanța este dată de unghiul în grade făcut de culoare în cercul perpendicular pe axa verticală a sistemului (axa nivelurilor de gri), saturarea este măsurată radial ca distanță față de axa verticală iar luminanța este măsurată vertical de-a lungul axei principale. Fiecare dintre aceste trei axe sunt discretizate într-un număr limitat de intervale, ca de exemplu: luminanța este exprimată cu valori întregi de la 0 (negru) la 10 (alb).

Repartizarea culorilor în acest mod nu este întâmplătoare ci este rezultatul măsurătorilor subiective ale percepției de culoare efectuate de Munsell. Pe fiecare dintre cele trei dimensiuni, culorile sunt alese pe cât posibil uniforme din punct de vedere al percepției, din această cauză "solidul" de culoare obținut are o formă neregulată (vezi Figura 4.6).

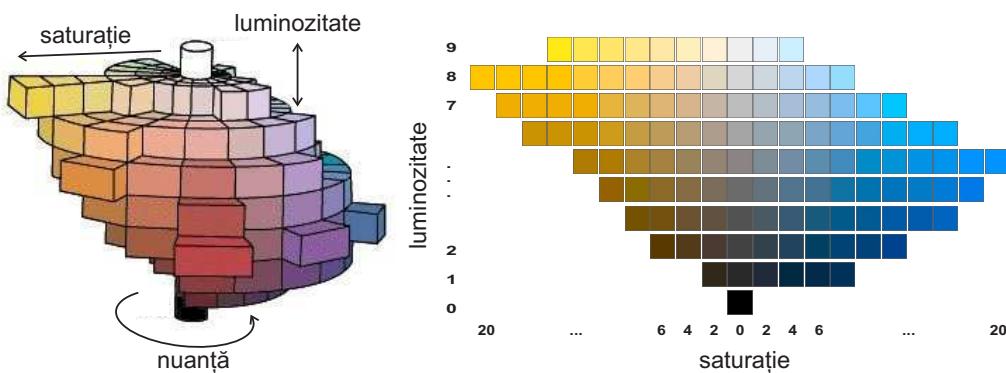


Figura 4.6: Spațiul de culoare al lui Munsell ("solidul" de culoare și o secțiune verticală a acestuia, sursă Computer Science Lab "http://www.computersciencelab.com/Direct3DTut1.htm").

Fiind un spațiu perceptual, spațiul de culoare al lui Munsell a stat la baza concepției spațiilor de culoare perceptuale actuale, precum spațiul  $L^*a^*b^*$  sau

familia de spații HSV (vezi secțiunile anterioare).

#### 4.1.4 Sisteme de axe independente

În funcție de paleta de culoare a imaginii analizate și de sistemul de reprezentare a culorilor folosit, se poate observa o corelație mai puternică sau mai slabă între diferitele componente de culoare. Astfel, pentru anumite etape de prelucrare este uneori mai "rentabil" să se folosească un sistem de reprezentare a culorilor în care componentele de culoare să fie complet decorelate<sup>10</sup>, cum este cazul sistemelor de axe independente.

Principalul inconvenient al acestor tipuri de sisteme de reprezentare a culorilor constă în dependența lor de distribuția de culoare considerată, astfel coeficienții transformării fiind diferenți de la o imagine la alta. Totuși, în urma studiilor efectuate în această direcție, s-a observat că în practică coeficienții transformării pot fi similari de la o distribuție de culoare la alta. Astfel, matricea de transformare nu este strict necesar să fie recalculată de la o imagine la alta, ci mai general, aceasta va fi specifică fiecărei aplicații în parte [Tréneau 04].

Dintre studiile semnificative realizate în această direcție putem menționa cele prezentate în [Ohta 80] unde sunt demonstate următoarele ipoteze:

- pentru majoritatea imaginilor color, axa principală obținută în urma descompunerii în componente principale (PCA - "Principal Component Analysis"<sup>11</sup>) se suprapune cu axa de luminanță,
- pentru mare parte a imaginilor color, informația de culoare, fie în totalitate, fie parțial, este dată de prima și de a doua componentă principală,
- pentru mare parte a imaginilor color, componente de culoare sunt fie strâns corelate, fie relativ decorelate. Astfel, în timp ce componente de culoare ale sistemelor RGB și XYZ sunt puternic corelate, componente sistemelor  $L^*a^*b^*$  și  $YC_bC_r$  sunt relativ decorelate,

---

<sup>10</sup>în statistică, corelația indică gradul de dependență liniară dintre două variabile aleatoare. Astfel, două variabile aleatoare sunt decorelate, sau independente, dacă nici una dintre ele nu oferă informații cu privire la valorile celeilalte, sau altfel zis, dacă acestea nu pot fi determinate una în funcție de celalaltă.

<sup>11</sup>descompunerea în componente principale, sau PCA, este definită ca fiind transformarea liniară ortogonală a datelor de intrare într-un nou sistem de coordonate în care varianța cea mai semnificativă a acestora este reprezentată pe prima coordonată, numită și prima componentă principală, a doua varianță este reprezentată pe a doua coordonată, și aşa mai departe. Astfel, transformarea PCA este teoretic optimală din punct de vedere al erorii pătratice.

- pentru marea majoritate a imaginilor color, sistemul de culoare decorrelat obținut pe baza transformatei Karhunen-Loeve<sup>12</sup> tinde să se confundă cu sistemul de axe  $I_1I_2I_3$ , ce este descris de următoarele ecuații:

$$I_1 = \frac{1}{3}(R + G + B) \quad (4.37)$$

$$I_2 = \frac{1}{2}(R - B) \quad (4.38)$$

$$I_3 = \frac{1}{4}(2 \cdot G - R - B) \quad (4.39)$$

unde  $(R, G, B)$  reprezintă coordonatele RGB.

În concluzie, dezvoltarea diverselor sisteme de reprezentare a culorilor existente a fost motivată în principal de necesitatea de a evidenția anumite proprietăți de culoare, proprietăți ce nu erau vizibile sau măsurabile în spațiile de culoare clasice folosite de dispozitivele de prelucrare a imaginilor.

De exemplu, familia de sisteme de culoare HSV permit separarea nuanței culorii de alte informații, precum saturăția sau luminozitatea, distribuția de culoare putând fi astfel interpretată pe baza culorilor primare din care aceasta a fost derivată. Un alt exemplu sunt spațiile de tip  $L^*a^*b^*$ , în care culorile sunt reprezentate în funcție de percepția vizuală a acestora, permîțând astfel evaluarea similarității dintre culori pe baza calculului distanței Euclidiene dintre acestea. Alte spațiuri de culoare, precum familia de spații  $YC_bC_r$  permit separarea informației de luminanță de informația cromatică fiind eficiente pentru metodele de prelucrare ce sunt vulnerabile la fluctuații de intensitate luminoasă în imagine.

Astfel, alegerea adecvată a spațiului de reprezentare a culorilor constituie o etapă premergătoare importantă în analiza și prelucrarea conținutului imaginilor sau a secvențelor de imagini.

## 4.2 Conținutul de culoare la nivel de imagine

La nivel de imagine, informația de culoare poate fi analizată folosind mai multe abordări. O primă metodă, și cea mai des întâlnită, o constituie analiza *parametrilor de nivel scăzut* extrași din imagine. Aceștia sunt de regulă măsuri statistice numerice ale proprietăților distribuției de culoare a imaginii globale sau a anumitor zone de interes din aceasta. Cum conceptul de culoare

---

<sup>12</sup>transformata Karhunen-Loeve discretă este cunoscută și sub numele de descompunere în componente principale sau PCA.

implică analiza senzației vizuale transmise, parametrii de nivel scăzut sunt deseori insuficienți pentru a interpreta percepția culorilor sau a tehnicielor de culoare prezente în imagine.

Un nivel semantic superior de descriere este atins atunci când sunt asociate descrierii *textuale culorilor*, precum numele culorii. Asociind un nume fiecărei culori, permite oricărei persoane să își creeze o imagine vizuală a culorii în cauză. Numele culorilor sunt astfel alese încât să fie reprezentative pentru proprietățile perceptuale cele mai semnificative ale acestora, precum nuanța, intensitatea sau saturarea.

La nivel de imagine, culorile nu sunt informații vizuale individuale, ci dimpotrivă, acestea capătă sens în relație cu alte culori sau regiuni de culoare vecine spațial. Astfel, caracterizarea conținutului de culoare implică *analiza diverselor relații perceptuale* ce pot apărea între culori. O modalitate, inspirată din domeniul artei, domeniu ce a furnizat primele studii referitoare la percepția relațiilor de culoare, constă în folosirea roțiilor de culoare ("color wheels"). Acestea constituie în esență o reprezentare grafică a unui set de culori primare ce sunt dispuse pe un cerc într-un mod perceptual uniform.

În cele ce urmează vom face o trecere în revistă a metodelor folosite de fiecare dintre direcțiile de studiu menționate.

#### 4.2.1 Analiza pe bază de histogramă

Una dintre metodele cele mai eficiente de reprezentare a conținutului de culoare o constituie *histograma*. Aceasta este o măsură statistică ce contabilizează numărul de apariții în imagine ale fiecărei culori dintr-o anumită paletă.

Din punct de vedere matematic, histograma imaginii  $I$  este dată de relația:

$$h(c) = \frac{1}{M \cdot N} \sum_{i=0}^{M-1} \sum_{j=0}^{N-1} \delta(I(i, j) - c) \quad (4.40)$$

unde  $M \times N$  reprezintă dimensiunea imaginii exprimată în pixeli,  $c$  reprezintă indicele culorii curente din paleta de culoare folosită,  $P_c$ , iar funcția delta  $\delta(x)$  este dată de ecuația următoare:

$$\delta(x) = \begin{cases} 1 & \text{dacă } x = 0 \\ 0 & \text{altfel} \end{cases} \quad (4.41)$$

Definită în acest fel, histograma imaginii  $I$  are sens de densitate de probabilitate discretă, valoarea  $h(c)$  reprezentând probabilitatea de apariție a culorii  $c$  în imagine.

O altă formă de histogramă utilizată frecvent este *histograma cumulată* [Stricker 95]. Aceasta este calculată pe baza histogramei  $h()$  în felul următor:

$$H(c) = \sum_{i=0}^c h(i) \quad (4.42)$$

unde culoarea  $c$  ia valori în paleta  $P_c$ . Definită în acest fel, histograma cumulată,  $H(c)$ , are sens de funcție de repartiție discretă, reprezentând probabilitatea ca, culoarea unui pixel din imagine să fie inferioară culorii  $c$ , unde relațiile de ordine dintre culori sunt date de ordinea de apariție a acestora în paleta de culoare  $P_c$ .

Histograma este eficientă în primul rând datorită invarianței totale sau parțiale a acesteia la anumite transformări geometrice ale imaginii, precum translatării, rotației (până la aproximativ 45 de grade), măriri ale imaginii (până la aproximativ 1.3 ori), precum și la schimbări de rezoluție (în cazul în care imaginile prezintă suficiente regiuni uniforme) sau la suprapunerile parțiale de obiecte în imagine [Bimbo 99]. Pe lângă invarianță, histograma mai are avantajul de a avea o complexitate de calcul redusă, numărul de operații fiind dat de numărul de pixeli din imagine.

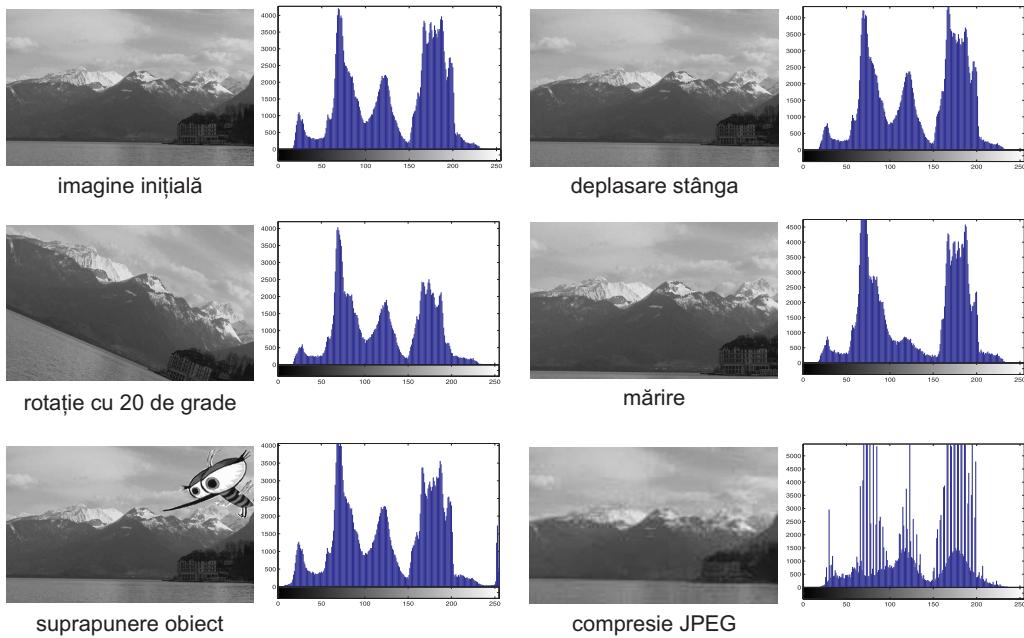


Figura 4.7: Exemple de histograme obținute pentru diverse transformări ale imaginii: translatăie, rotație, mărire, suprapunere obiect și compresie.

În Figura 4.7 am ilustrat câteva exemple de histograme obținute în urma unor transformări ale imaginii, pentru o imagine cu nivale de gri. În ciuda modificărilor, putem spune brutale, ce survin în imagine, histograma are tendința să păstreze distribuția de moduri<sup>13</sup> a imaginii originale.

În ciuda multiplelor avantaje prezentate de histogramă, această formă de reprezentare a conținutului de culoare din imagine prezintă și o serie de limitări. În primul rând, histograma de culoare este foarte sensibilă la variațiile importante de intensitate luminoasă din imagine, astfel două imagini identice din punct de vedere al conținutului, dar "achiziționate" în condiții diferite de luminozitate vor furniza histograme de culoare foarte diferite. Pentru a îmbunătății invarianța histogramei la fluctuațiile de intensitate luminoasă, soluția cea mai frecvent adoptată constă în folosirea unui alt spațiu de culoare ce permite separarea informației de luminanță de cea cromatică [Bimbo 99], cum este cazul spațiilor de culoare de tip  $YC_bC_r$  sau  $HSV$  (vezi Secțiunea 4.1). Astfel, histograma poate fi calculată doar pentru componente de crominanță, ce nu sunt afectate de variația intensității luminoase [Ionescu 08].

Un alt dezavantaj al histogramei este dat de faptul că pentru calculul acestora nu se ține cont de informația spațială din imagine, mai exact, de poziția spațială a pixelilor în imagine. Astfel, două histograme identice pot corespunde în realitate cu două imagini complet diferite din punct de vedere vizual, dar care conțin aceleași culori și cu o aceeași probabilitate de apariție (vezi Figura 4.8).

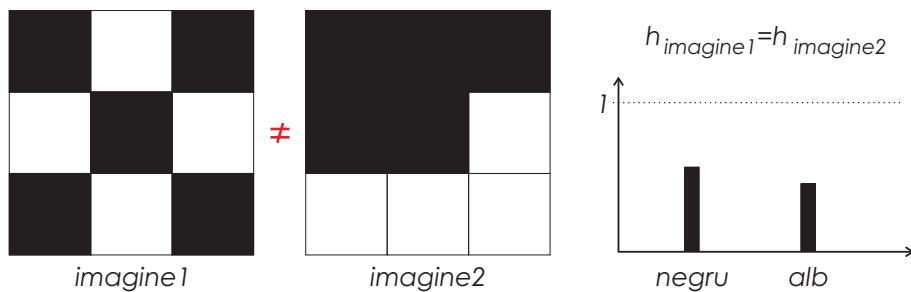


Figura 4.8: Exemplu clasic de două imagini diferențiate din punct de vedere vizual ce conduc la histograme absolut identice.

Acest lucru se datorează faptului că doi pixeli de aceeași culoare nu au

<sup>13</sup>un mod al histogramei este definit în general ca fiind plaja de valori din jurul unui maxim local, de regulă mărginită la stânga și la dreapta de minime locale, ce indică prezența în imagine a uneia sau a mai multor regiuni semnificative cu o distribuție de culoare similară.

obligatoriu și aceeași semnificație vizuală, ci dimpotrivă, aceștia în realitate pot reprezenta informații spațiale complet diferite, ca de exemplu pot fi pixeli de contur sau din regiuni uniforme de culoare.

Pentru a corecta acest lucru, *histogramele ponderate*, adaugă informații suplimentare legate de vecinătatea spațială a pixelului curent analizat. Pe baza notațiilor precedente, din ecuația 4.40, histograma ponderată este dată de ecuația următoare:

$$h_p(c) = \frac{1}{M \cdot N} \sum_{i=0}^{M-1} \sum_{j=0}^{N-1} w(i, j) \cdot \delta(I(i, j) - c) \quad (4.43)$$

unde  $w(i, j)$  reprezintă ponderea pixelului, de coordonate spațiale în imagine  $(i, j)$ , la valoarea histogramei pentru culoarea  $c$ . Funcția  $w()$  este calculată pentru o anumită vecinătate a pixelului curent, iar valorile acesteia de regulă cresc odată cu gradul de neuniformitate al regiunii.

Printre funcțiile  $w()$  cel mai frecvent folosite putem menționa operatorul Laplacian, ce reprezintă derivata de ordinul doi în spațiul bidimensional al imaginii  $I()$ :

$$\Delta I(i, j) = \nabla^2 I(i, j) = \frac{\partial^2 I(i, j)}{\partial i^2} + \frac{\partial^2 I(i, j)}{\partial j^2} \quad (4.44)$$

Acesta oferă informații referitoare la contururile prezente în imagine, conferind o pondere mai importantă pixelilor de contur, pentru care valoarea Laplacian-ului este importantă, precum și o pondere semnificativ mai mică pixelilor din regiunile uniforme, pentru care valoarea Laplacian-ului este apropiată de zero.

O altă abordare constă în definirea funcției  $w(i, j)$  ca fiind probabilitatea de reapariție a culorii pixelului de coordonate spațiale  $(i, j)$ , în vecinătatea acestuia. Dacă probabilitatea este importantă, atunci pixelul face parte dintr-o regiune uniformă, și vice-versa, dacă probabilitatea este redusă, atunci pixelul se poate afla într-o regiune de contur, poate fi un pixel izolat, etc.

O altă categorie de histograme ce țin cont de repartiția spațială a pixelilor în imagine sunt *histogramele acumulative*. Pentru calculul acestora, imaginea este mai întâi divizată în mai multe regiuni iar pentru fiecare regiune este calculată o histogramă de culoare. Histograma acumulativă este obținută prin acumularea aditivă a valorilor acestora, astfel:

$$\tilde{h}(c) = \sum_{b=0}^{N-1} f(h_b(c)) \quad (4.45)$$

unde  $b$  reprezintă indicele blocului curent,  $N$  reprezintă numărul total de blocuri din imagine iar  $f()$  este de regulă o funcție neliniară ce evidențiază

repartiția geometrică de culoare. Cu cât sunt acumulate mai multe culori, cu atât binii histogramei  $\tilde{h}()$  sunt mai diferențiabili. În ceea ce privește calitatea reprezentării conținutului de culoare, raportat la histogramele clasice, histogramele ponderate tind să furnizeze rezultate mai precise în anumite aplicații, precum indexarea după conținut a imaginilor [Ferecatu 01].

Pe lângă tipurile de histograme, putem spune, clasice, enumerate anterior, în literatura de specialitate a domeniului indexării de imagini întâlnim un caz particular de histogramă și anume *histogramele fuzzy*. Spre deosebire de histogramele clasice în care culoarea pixelilor intervine doar la calculul valorii histogramei pentru un singur bin, în histogramele fuzzy se ia în calcul similaritatea dintre culori prin distribuirea unei anumite funcții de apartenență fuzzy a fiecărui pixel către toți binii histogramei. Astfel, histograma fuzzy a imaginii  $I$  este definită ca fiind:

$$F(I) = [f_1, f_2, \dots, f_n] \quad (4.46)$$

unde

$$f_i = \sum_{j=1}^N \mu_{ij} \cdot P_j = \frac{1}{N} \sum_{j=1}^N \mu_{ij} \quad (4.47)$$

unde  $j$  reprezintă indicele pixelului curent din imagine,  $j = 1, \dots, N$  cu  $N$  numărul total de pixeli din imagine,  $P_j$  reprezintă probabilitatea de apariție a pixelului  $j$  în imagine iar  $\mu_{ij}$  reprezintă valoarea funcției de apartenență fuzzy a pixelului  $j$  pentru binul  $i$  din histogramă (pentru o descriere detaliată a conceptului de mulțime fuzzy vezi Capitolul 6).

Definită în acest fel, histograma fuzzy nu se limitează în a lua în calcul doar similaritatea dintre culorile din binii histogramei, ci și disimilaritatea culorilor asociate aceluiași bin<sup>14</sup>. În domeniul indexării după conținut a imaginilor, din punct de vedere al performanțelor, raportat la histogramele clasice, histogramele fuzzy se dovedesc a furniza rezultate mai bune cât și o robustețe mai ridicată la prezența zgomotului sau la fluctuațiile de intensitate luminoasă din imagine [Han 02a].

#### 4.2.2 Analiza pe baza denumirii culorilor

O altă modalitate de caracterizare a conținutului de culoare constă în utilizarea *denumirilor culorilor*. Asocierea de denumiri textuale culorilor existente, permite oricărui dintre noi, indiferent de cultură sau de domeniu

---

<sup>14</sup>de menționat este faptul că în general termenul de bin al histogramei se referă nu numai la o singură culoare, cum este cazul unei palete fixe de culoare (vezi ecuația 4.40), ci la un interval de culoare. În acest caz, toate culorile din imagine ce sunt cuprinse în acest interval vor modifica valoarea histogramei doar pentru binul respectiv.

de activitate, să ne creăm o imagine vizuală a proprietăților culorii despre care este vorba. Denumirile culorilor sunt de regulă alese dintr-un dicționar predefinit de nume de culori, iar asocierea între culoarea fizică și numele acesteia se realizează pe baza unui sistem de denotare a culorilor ("Color Naming System"). Problema asocierii de denumiri textuale culorilor este o problemă ce a fost îndelung cercetată în literatura de specialitate. Pentru o sinteză a acesteia, cititorul se poate raporta la studiile prezentate în [Kay 03], [Benavente 04] sau [Lay 04].

În [Berlin 91] este prezentat un studiu cu privire la definirea conceptului de culoare elementară sau de bază. Astfel, numele culorilor de bază (primare) sunt definite ca fiind numele culorilor ce respectă simultan următoarele reguli:

- folosirea acestora nu este restrictivă doar pentru o anumită categorie aparte de obiecte, de exemplu, în acest sens, culoarea numită "măsliniu" ("olive") nu este o culoare de bază,
- sensul acestora nu poate fi predictibil prin înțelegerea proprietăților anumitor obiecte, de exemplu culoarea unei frunze nu este o culoare de bază,
- sensul acestora nu este inclus în numele unei alte culori,
- acestea au o constanță a percepției, fiind percepute în același fel de persoane diferite.

Pe baza spațiului de culoare al lui Munsell (vezi Secțiunea 4.1), [Berlin 91] merge mai departe și definește ca fiind elementare, 11 culori ce sunt general valabile în cel puțin 20 de limbi diferite existente. Acestea sunt: *alb*, *negru*, *roșu*, *verde*, *galben*, *albastru*, *maro*, *roz*, *purpuriu*, *portocaliu* și respectiv *gri*. Denumirea unui set de culori elementare stă la baza definirii descrierilor textuale pentru celelalte culori derivate existente. Dacă numele culorilor elementare sunt cvasi-similare în majoritatea culturilor existente, nu se poate spune același lucru și pentru definirea frontierelor dintre diferitele categorii de culori, frontiere ce de regulă variază foarte mult de la o limbă la alta.

Sistemele existente de denumire a culorilor folosesc diverse tehnici pentru a obține o anumită universalitate a denumirilor acestora. Printre metodele folosite putem enumera:

- modelizarea fuzzy a apartenenței culorilor la anumite clase de culoare,
- denumirea culorilor în funcție de reprezentarea fizică a acestora în termeni de intervale de lungimi de undă,

- asocierea denumirii culorilor dintr-un anumit dicționar de culori, numit și "Lookup Table", ce este definit "a priori" pe baza analizei manuale a diverselor spații de reprezentare a culorilor,
- folosirea de tehnici de "clustering" pentru definirea automată a claselor de culoare pe baza similarității perceptuale a acestora.

Aproape toate metodele existente de denumire a culorilor necesită pentru definirea dicționarului de culori folosirea expertizei umane. Aceasta poate interveni, fie în totalitate, caz în care descrierile textuale sunt asociate manual culorilor, fie doar parțial prin definirea limitelor sau a unumitor parametri ce caracterizează fiecare categorie de culoare. De menționat este faptul că până în prezent nu s-a dezvoltat o metodă complet automată de denumire a culorilor dintr-un anumit spațiu de culoare [Benavente 04].

În marea majoritate a metodelor de analiză a conținutului de culoare, se preferă folosirea dicționarelor de nume de culori datorită faptului că în acest fel se elimină complexitatea de calcul specifică unei metode de denumire cvasi-automată, numele culorilor fiind disponibile direct din dicționarul de culoare. Ca exemple de astfel de dicționare de nume de culori putem menționa:

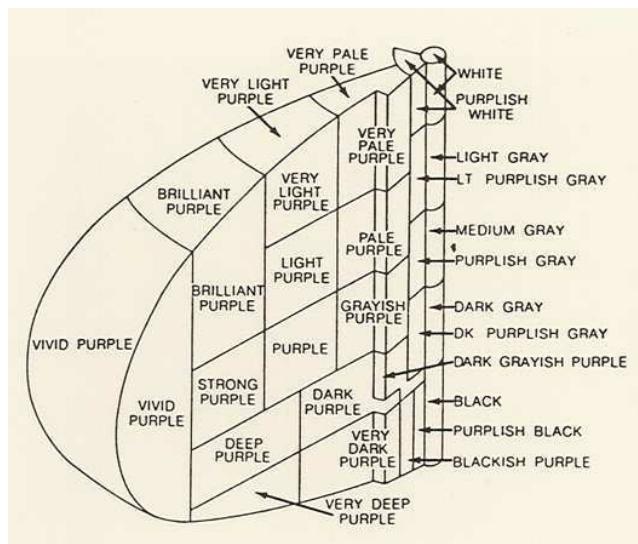


Figura 4.9: Sistemul standard de culori ISCC-NBS: exemplu de definire a culorilor purpurii prin partitarea spațiului de culoare al lui Munsell.

- sistemul standard de culoare ISCC-NBS [Kelly 76] ("Inter Society Color Council" - "National Bureau of Standards") ce este definit pe baza sferei de culoare a lui Munsell (vezi Figura 4.9),

- dicționarele de culoare ”X11 Window System Distribution”, ”Netscape Color Names”, ”HTML-4 Color Names”, ”Two4U’s Big Color Database”, ”Resene Paint Colours”, ”WebSafe” sau ”CNS Color-Naming System” [CSAIL 06].

Datorită simplității și a modului în care au fost alese culorile în dicționarul de culori disponibil, un interes aparte în descrierea conținutului de culoare îl prezintă paleta de culoare cunoscută sub numele de ”WebSafe” sau ”Webmaster” [Visibone 06] (vezi Figura 4.10).

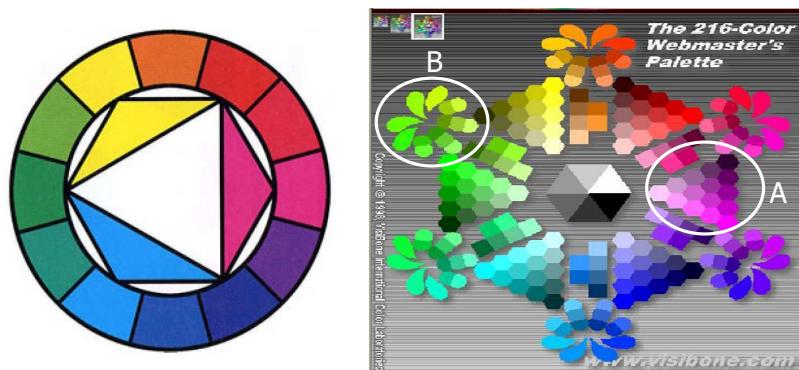


Figura 4.10: Analogia între roata de culoare a lui Itten (stânga, vezi Secțiunea 4.2.3) și paleta de culoare ”Webmaster” [Visibone 06] (dreapta, zona de tip A conține variații ale unei culori elementare iar zona de tip B conține amestecuri de culori elementare).

Aceasta prezintă o serie de avantaje ce constituie un atu pentru înțelegerea la nivel perceptual a conținutului de culoare din imagine [Ionescu 08], și anume:

- prezintă un bun compromis între numărul de culori disponibile (216) și ”bogăția” de culoare furnizată: conține 12 culori de bază și anume: *portocaliu, roșu, roz, magenta, violet, albastru, ”azure”<sup>15</sup>, turcoaz, ”teal”<sup>16</sup>, verde, ”spring”<sup>17</sup> și galben* precum și 6 niveluri de gri ce includ albul și negrul (vezi Figura 4.10 unde culorile de bază enumerate sunt figurate în sensul acelor de ceasornic începând cu portocaliu la ora 12),

<sup>15</sup>culoarea ”azure” este un albastru deschis, similar culorii cerului senin, ce se încadrează undeva între albastru și turcoaz.

<sup>16</sup>culoarea ”teal” este o combinație de verde și albastru de saturatie redusă, similară cu un turcuaz întunecat. Numele culorii vine de la denumirea unei anumite specii de rațe de apă dulce care prezintă în jurul ochilor tocmai această culoare.

<sup>17</sup>culoarea ”spring” este o combinație de verde și galben, și reprezintă culoarea mugurilor plantelor la începutul primăverii.

- prezintă o diversitate de culoare ridicată: variații ale saturăției și a intensității celor 12 culori elementare precum și amestecuri de culoare,
- fiecare culoare din paletă este denumită în funcție de nuanță, saturatie și intensitate, facilitând astfel analiza percepției de culoare (un exemplu este prezentat în Tabelul 4.1),

Culoare	$(R, G, B)$	Hexazecimal	Denumire culoare
	(255, 255, 51)	FFFF33	"Light Hard Yellow"
	(204, 0, 102)	CC0066	"Dark Hard Pink"
	(204, 204, 204)	CCCCCC	"Pale Gray"

Tabelul 4.1: Exemplu de denumiri de culori din paleta "Webmaster".

- culorile sunt reprezentate în concordanță cu roata de culoare a lui Itten (vezi Figura 4.10 și explicația de la pagina 134), facilitând astfel analiza relațiilor perceptuale dintre culori.

### 4.2.3 Analiza senzației induse de culoare

O altă modalitate de caracterizare a conținutului de culoare dintr-o imagine constă în analiza senzației transmise de către acesta observatorului uman. În acest sens, o referință în domeniu sunt cercetările lui Itten [Itten 61], care în urma experienței sale în calitate de pictor în currentul Bauhaus, definește în 1960 un set de *reguli formale* (limbaj de culoare) ce cuantificau într-un anumit număr de categorii, efectele vizuale produse la nivel perceptual de către diversele combinații de culori. Acestea sunt definite pe baza conceptelor folosite în domeniul artistic al picturii.

În domeniul artei, culoarea este deseori folosită pentru a descrie proprietățile perceptuale ale diverselor obiecte din lumea înconjurătoare, precum și ca reprezentare simbolică a culorii fizice. În acest sens, pentru a alege culorile potrivite și pentru a crea noi amestecuri de culoare ce corespund senzațiilor particulare ce vor fi transmise de opera lor, artiștii se folosesc de ceea ce numim "roți de culoare" sau "color wheels" [Birren 69].

O *roată de culoare* este în esență un spațiu de culoare particular în care relațiile dintre culori sunt inspirate de teoria contrastelor și a armoniei de culoare din domeniul picturii. Aceasta este construită pe baza unui număr arbitrar de culori elementare ce vor fi dispuse pe un cerc (roată de culoare) într-un mod perceptual liniar progresiv [Lay 04]. Definită în acest mod, o roată de culoare este o reprezentare bidimensională. Trecerea la o reprezentare tridimensională pe baza aceluiasi principiu conduce la formarea a ceea

ce numim o *"sfără de culoare"*. De-a lungul istoriei, au fost propuse mai multe astfel de sisteme de reprezentare perceptuală a culorilor, dintre acestea menționăm pe cele mai cunoscute ce poartă și numele creatorilor lor, și anume: sferă de culoare a lui Runge, roata de culoare a lui Chevreul, spațiul culorilor opuse a lui Hering, solidul de culoare a lui Munsell, roata de culoare a lui Itten, etc. (vezi Figura 4.11).

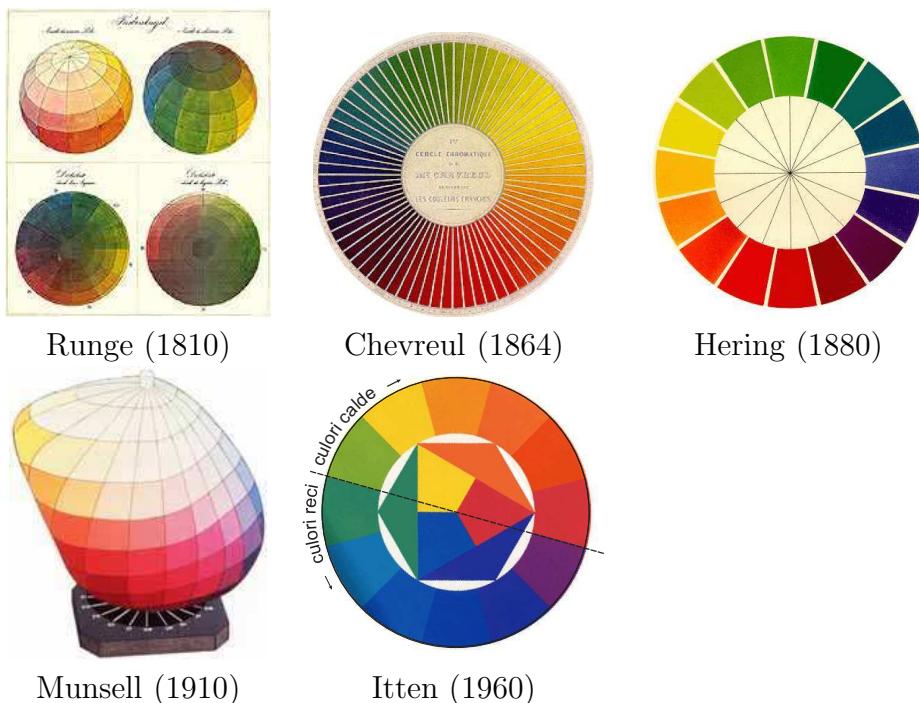


Figura 4.11: Reprezentarea culorilor sub formă perceptuală pe baza roți și a sferelor de culoare.

O astfel de reprezentare a culorilor este foarte utilă în cazul în care se dorește studierea relațiilor perceptuale dintre culori. În domeniul artei, conceptul de relație între culori, unde culorile sunt combinate pe baza unei roți de culoare, este un studiu esențial (vezi Josef Albers, Faber Birren, Johannes Itten, etc.).

Dacă luăm ca exemplu roata de culoare a lui Itten, care este și una dintre cele mai cunoscute roți de culoare, în aceasta culorile sunt aranjate cu un anumit scop, astfel: culorile considerate ca fiind calde se găsesc în prima jumătate a roții, începând cu culoarea "spring", continuând cu galben și finalizând cu magenta, în timp ce culorile considerate ca fiind reci se găsesc în cealaltă jumătate, pornind de la violet, continuând cu albastru și finalizând

cu verde (vezi Figura 4.11). Mai mult, culorile ce sunt opuse din punct de vedere al percepției se găsesc poziționate diametral opus (de exemplu albastru și galben) în timp ce culorile considerate analoage sunt culori vecine pe roata de culoare (de exemplu galben și portocaliu).

Astfel, folosind teoria de culoare dezvoltată de Itten, precum și reprezentarea perceptuală a culorilor pe baza roțiilor de culoare, putem caracteriza conținutul vizual în termeni de *contrast* și respectiv *concordanță* de culoare. Itten [Itten 61] definește percepția de culoare pe baza a *șapte contraste de culoare* ce sunt exemplificate în Figura 4.12, astfel:

- *contrastul de nuanță*: acest contrast vizual este realizat prin juxtapunerea de diverse nuanțe de culoare. Cu cât acestea sunt mai diferite din punct de vedere perceptual, cu atât mai puternic este contrastul obținut (distanța dintre culori este evaluată folosind o roată de culoare). Un exemplu este illustrat în Figura 4.12.a.

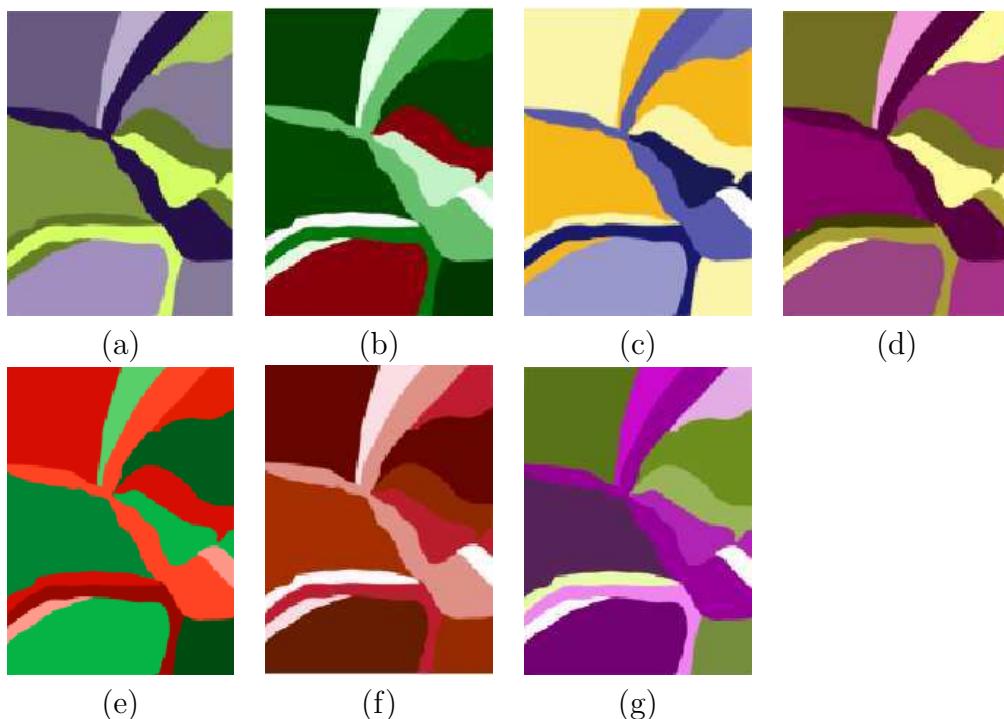


Figura 4.12: Cele șapte contraste ale lui Itten: (a) Contrastul de nuanță, (b) Contrastul închis-deschis, (c) Contrastul cald-rece, (d) Contrastul de complementaritate, (e) Contrastul de simultaneitate, (f) Contrastul de saturatie, (g) Contrastul de extensie (sursă imagini "<http://www.worqx.com/color/itten.htm>").

- *contrastul închis-deschis*: acest contrast este legat de gradul de percepție al intensității luminoase. La extreame se găsesc negrul (absența luminii) și albul (intensitatea maximă), iar între acestea sunt nivelele de gri și nuanțele cromatice. Contrastul este realizat prin juxtapunerea atât a culorilor deschise cât și închise (vezi Figura 4.12.b).
- *contrastul Cald-rece*: acest contrast corespunde senzației de căldură transmisă de anumite culori. În domeniul artei, culorile prezintă o anumită temperatură sau căldură. Astfel, galben, galben-portocaliu, portocaliu, roșu-portocaliu, roșu și roșu-violet sunt considerate ca nuanțe calde, pe când galben-verde, verde, albastru-verde, albastru, albastru-violet și violet sunt considerate ca nuanțe reci. Contrastul de culoare este realizat prin juxtapunerea atât a culorilor calde cât și reci (vezi Figura 4.12.c).
- *contrastul de complementaritate*: acest contrast corespunde relațiilor de complementaritate existente între culori. În practică, pe o roată de culoare (de exemplu roata de culoare a lui Itten, vezi Figura 4.11) perechile de culori complementare (opuse ca percepție) sunt determinate de linia dreaptă ce trece prin centrul roții și care leagă două culori diametral opuse. Contrastul de complementaritate este astfel realizat prin folosirea de culori opuse din punct de vedere al percepției, obținându-se astfel o anumită simetrie vizuală (vezi Figura 4.12.d).
- *contrastul de simultaneitate*: acest contrast se folosește de răspunsul asimetric al percepției umane la fenomenul culorilor opuse. Contrastul este realizat prin "vibrarea" percepției frontierelor dintre culori. Cu acest contrast se pot realiza o serie de iluzii optice interesante. Un exemplu este ilustrat în Figura 4.12.e.
- *contrastul de saturatie*: acest contrast este realizat prin juxtapunerea de nuanțe pure sature cu nuanțe diluate de saturatie redusă. Acest contrast se dovedește însă a fi relativ, deoarece anumite culori pot apărea ca fiind mai sature prin contrast dacă sunt alăturate unei culori mai puțin sature, și vice-versa. Un exemplu este ilustrat în Figura 4.12.f.
- *contrastul de extensie*: acest contrast este legat de proporția în care sunt folosite culorile în imagine. Percepția vizuală a unei culori este direct influențată de gradul de luminanță folosit, precum și de suprafața spațială ocupată de culoare în imagine. Contrastul de extensie este astfel realizat prin asocierea de culori regiunilor fizice din imagine ce au o suprafață proporțională cu ponderea perceptuală vizuală a culorii (vezi Figura 4.12.g).

Pe lângă cele șapte contraste definite de Itten, în domeniul artei întâlnim și anumite configurații de culoare ce au efecte vizuale particulare. Aceste ”*scheme de culoare*”, ce exprimă relațiile perceptuale dintre diferitele tipuri de contraste de culoare, sunt descrise deseori ca fiind modalitățile fundamentale de obținere a armoniei de culoare în imagine [Birren 69]. Acestea sunt următoarele:

- *schema de culoare monocromatică*: implică armonizarea unei singure nuanțe de culoare și este realizată pe baza contrastului de culoare închis-deschis. Astfel, pornind de la o anumită nuanță de culoare se obțin mai multe variații ale acesteia în modul următor: nuanța poate fi întunecată prin adăugarea de negru, se pot obține mai multe tonalități ale acesteia prin adăugarea de gri sau nuanța poate fi deschisă prin adăugarea de alb. Mai mult, pentru a se evita o posibilă monotonie a culorilor, se poate folosi și contrastul de saturatie prin varierea saturatiei nuanței folosite.
- *schema de culori adiacente*: implică armonizarea nuanțelor de culoare similară sau adiacente. Aceasta este realizată de regulă prin amestecul a maxim trei culori adiacente, ce sunt alese pe baza unei roți de reprezentare a culorilor (ca de exemplu roata lui Itten). Diferența între acestea constă în faptul că una dintre cele trei culori va fi pusă în evidență prin folosirea sa în imagine într-o proporție semnificativă față de celelalte nuanțe.
- *schema de culori complementare*: implică armonizarea echilibrului dintre culorile opuse ca percepție sau complementare. Modalitatea cea mai simplă în care poate fi realizată este pe baza a două culori ce sunt complementare pe o roată de reprezentare a culorilor. O altă modalitate de realizare implică folosirea a două culori adiacente ce sunt contrastate de perechea de culori opuse acestora. Această alegere a culorilor corespunde unei alegeri în  $X$  pe o roată de culoare și este cunoscută și sub numele de *dublă complementaritate*. O *complementaritate divizată* este obținută prin asocierea unei nuanțe de culoare cu două culori vecine nuanței complementare, formându-se astfel o structură de tip  $Y$ .

Ca exemple de sisteme ce folosesc descrierea conținutului de culoare la nivel perceptual pe baza conceptelor enumerate mai sus, putem menționa sistemul QBIC al ”State Hermitage Museum” [QBIC 03] și sistemul PICASSO [Corridoni 99]. În sistemul QBIC, picturile sunt căutate într-o bază de date folosind culorile predominante, precum și relațiile existente între regiunile de culoare din pictură. Acestea din urmă sunt specificate de utilizator prin

desenarea interactivă a unei anumite scheme de culoare. Un alt exemplu este sistemul SoloArt propus în [Lay 04] unde concepțele artistice de culoare sunt extrase tot pentru indexarea după conținut a picturilor. Relațiile dintre culori sunt analizate de această dată într-un spațiu perceptual, și anume spațiul de culoare LCH (vezi Secțiunea 4.1), pe baza căruia sunt determinate în mod automat tehniciile de culoare folosite: contrastele de culoare, celeșapte contraste ale lui Itten precum și schemele de armonie a culorilor.

### 4.3 Conținutul de culoare în secvențele de imagini

După cum am menționat deja în partea introductivă a acestui capitol, majoritatea tehniciilor de caracterizare a conținutului de culoare se limitează în mod natural la analiza acestuia doar la nivel de imagine, fără a lua în considerație dimensiunea temporală. Astfel, majoritatea metodelor existente sunt destinate analizei imaginilor statice. Totuși, putem menționa o serie de abordări bazate mai mult sau mai puțin integral pe informația de culoare, abordări ce încearcă să descrie conținutul temporal de culoare din secvențele de imagini în diverse domenii de aplicație. Dintre acestea, ne vom limita la detalierea metodelor propuse în [Colombo 99] (secvențe publicitare), [Detyniecki 03] (secvențe de știri) și [Ionescu 08] (secvențe de animație).

Sistemul propus în [Colombo 99] folosește informația de culoare în relație cu alți parametri, specifici secvențelor de imagini, pentru a descrie din punct de vedere semantic conținutul secvențelor publicitare. Aceasta propune două niveluri perceptuale de analiză, și anume: "nivelul expresiv" și "nivelul emoțional".

La *nivel expresiv*, secvențele publicitare sunt clasificate în patru categorii perceptuale, și anume: "tipul practic", "animat", "utopic" și "critic". Clasificarea automată a acestora este realizată pe baza a mai mulți parametrii de nivel scăzut, precum prezența tranzițiilor de tip "cut" și "dissolve" (vezi Capitolul 2), apariția și dispariția repetitivă a anumitor culori din imagine, prezența liniilor verticale și orizontale în imagine precum și prezența culorilor saturate sau nesaturate. La *nivel emoțional*, secvențele publicitare sunt caracterizate din punct de vedere al acțiunii, suspansului, senzației de calm, senzației de relaxare, de veselie precum și al entuziasmului transmis. Informația de culoare este utilizată în acest caz pentru a atinge anumite niveluri emoționale. De exemplu, nivelul de acțiune al unei secvențe de publicitate poate fi crescut prin folosirea culorilor roșu și purpuriu. În mod similar, senzația de calm poate fi amplificată prin folosirea culorilor: albas-

tru, portocaliu, verde sau alb, sau poate fi diminuată prin prezența culorilor purpurii și a negrului.

Fiecare categorie de descrieri semantice propusă este caracterizată de o anumită configurație a valorilor unui set de parametri de nivel scăzut, ce vor fi asociați secvențelor publicitare pe baza reprezentării fuzzy. De exemplu, secvențele publicitare din categoria "practică" sunt caracterizate de absența culorilor saturate,  $\phi_{saturated} = 0$ , unde  $\phi$  reprezintă funcția fuzzy de apartenență la simbolul ce modelizează conceptul de prezență a culorilor saturate în secvență (vezi Capitolul 6); de prezența liniilor orizontale și verticale,  $\phi_{hor/vert} = 1$ , precum și de prezența tranzițiilor de tip "dissolve",  $\phi_{dissolves} = 1$ .

Sistemul propus în [Detyniecki 03] folosește arbori de decizie fuzzy pentru a extrage informații din secvențele de știri, procedeu cunoscut și sub numele de "data mining"<sup>18</sup>. Arborii de decizie sunt folosiți în acest caz la extragerea în mod automat a unui set de reguli ce vor fi folosite pentru clasificarea conținutului anumitor segmente tematice din secvențele de știri. Informația folosită în acest caz este exclusiv informația de culoare.

În primă fază de analiză, fiecare plan video al secvenței este rezumat cu un anumit număr de imagini reprezentative pentru conținutul acestuia (imagini cheie). Culorile fiecărei imagini reținute sunt proiectate pe o paletă de culoare redusă (64 sau 256 culori) prin cuantificarea liniară a spațiului RGB. Conținutul de culoare la nivel de imagine este reprezentat mai departe folosind histograme color ce formează vectori de caracteristici de culoare în spațiul determinat de paleta de culoare folosită. Mai departe, vectorii de caracteristici astfel obținuti sunt folosiți pentru a construi arborele de decizie fuzzy pe baza căruia se va realiza detectarea a două evenimente semantice din secvențele de știri, și anume: prezența de text încrustat în imagine ("inlays") și prezența prezentatorului sau a jurnalistului.

Un sistem particular de analiză a conținutului de culoare, ce ia în calcul dimensiunea temporală, este propus în [Ionescu 08]. Acesta folosește exclusiv informația de culoare pentru a detecta și analiza tehniciile de culoare prezente în filmele artistice de animație. Metoda propusă se bazează pe ipoteza că aproape fiecare film de animație conține o paletă de culoare particulară, paletă ce a fost aleasă în mod intenționat de autor în momentul conceperii filmului pentru a transmite anumite sentimente sau pentru a exprima anumite concepte artistice.

Sistemul propus reduce în primă fază redundanța temporală a secvenței

---

<sup>18</sup>"data mining" reprezintă procesul de căutare în conținutul unor colecții foarte vaste de date, de regulă după criterii semantice, a anumitor informații relevante sau de noi "cunoașteri" ("knowledge") în cazul în care natura conținutul datelor precum și relațiile dintre acestea nu sunt cunoscute.

prin reținerea a doar un procent,  $p\%$ , din imaginile fiecărui plan video. Mai departe, folosind un algoritm de reducere de culoare de tip "dithering"<sup>19</sup>, metodă ce furnizează o fidelitate vizuală ridicată a imaginii rezultate, conținutul temporal de culoare al secvenței este reprezentat prin calcularea histogramei globale ponderate,  $h_{GW}$ , propusă în [Ionescu 05c]. Aceasta este dată de relația:

$$h_{GW}(c) = \sum_{i=0}^M \left[ \frac{1}{N_i} \sum_{j=0}^{N_i} h_{shot_i}^j(c) \right] \cdot w_i \quad (4.48)$$

unde  $M$  reprezintă numărul total de plane video al secvenței,  $N_i$  reprezintă numărul total de imagini reținute pentru planul de indice  $i$  (reprezentând un procent  $p\%$  din numărul total de imagini ale acestuia),  $h_{shot_i}^j()$  reprezintă histograma color a imagini de indice  $j$  din planul video  $i$ ,  $c$  reprezintă o culoare din paleta de culori folosită iar  $w_i$  reprezintă ponderea planului video  $i$  la calculul histogramei ponderate și este dată de relația următoare:

$$w_i = \frac{N_{shot_i}}{N_{total}} \quad (4.49)$$

unde  $N_{shot_i}$  reprezintă numărul total de imagini ale planului video  $i$  iar  $N_{total}$  reprezintă numărul total cumulat de imagini ale tuturor planelor video din secvență. Astfel, cu cât planul curent analizat este mai lung, cu atât contribuția acestuia la distribuția globală de culoare a secvenței este mai importantă.

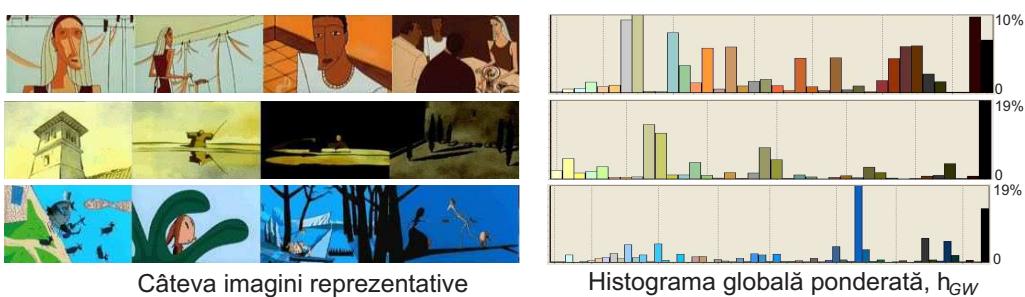


Figura 4.13: Exemple de histograme globale ponderate obținute pentru  $p = 15\%$  (sursă imagini [Folimage 06]).

<sup>19</sup>procesul de "dithering" constă în adăugarea intenționată a unei anumite forme de zgomot necorelat pentru a distribui eroarea de cuantizare în cazul procesului de cuantificare a valorilor unei funcții. În cazul reducerii de culoare, eroarea obținută prin substituirea culorilor este distribuită progresiv în vecinătatea acestora, evitând astfel formarea de artefacte în imagine.

Definită în acest fel, histograma ponderată redă procentul de apariție temporală în secvență a fiecărei culori, fiind o măsură cantitativă a distribuției de culoare globală a secvenței (vezi un exemplu în Figura 4.13).

Mai departe, conținutul de culoare este caracterizat la nivel sintactic cu o serie de parametri statistici de nivel scăzut, precum procentul de culori închise și deschise, procentul de culori calde și reci, diversitatea de culoare, procentul de culori adiacente, etc., parametrii ce sunt extrași pe baza histogramei ponderate folosind dicționarului de nume de culori furnizat de paleta "Webmaster" (vezi Figura 4.10). Nivelul semantic de descriere este obținut pe baza unui set de reguli de reprezentare fuzzy a acestor parametrii. Astfel, secvențelor le sunt asociate descrieri simbolice, exprimate textual în felul următor: "culorile predominante sunt calde" sau "secvența prezintă un contrast de saturație", descrieri al căror grad de veridicitate este modelizat de gradul de apartenență fuzzy.

## 4.4 Concluzii

În acest capitol am discutat diversele modalități de analiză și caracterizare a conținutului de culoare, atât la nivel spațial, în imagine, cât și la nivel temporal, în secvențe de imagini.

Un rol important în analiza de culoare îl are spațiul de reprezentare al culorilor folosit. Aceasta este de regulă ales în aşa fel încât să pună în evidență anumite proprietăți de culoare ce vor facilita implementarea algoritmului de prelucrare dorit. Pornind de la spații de culoare, putem spune clasice, precum spațiul RGB folosit în aproape toate dispozitivele hardware de reprezentare a culorilor (plăci grafice, ecrane, imprimante, etc.), spațiile de culoare evoluează spre spații perceptuale sau inspirate de modalitatea fizică de percepție a culorii din sistemul vizual uman, precum spațiul LCH. Dezvoltarea acestora este motivată în principal de necesitatea de caracterizare a informației de culoare la un nivel perceptual, apropiat de cel uman, ținta finală fiind atingerea unui nivel semantic de descriere a conținutului de culoare din imagine, direcție de studiu de foarte mare actualitate în acest moment.

Pentru a înțelege automat conținutul de culoare, în absența unei analize vizuale a acestuia, o modalitate foarte eficientă constă în folosirea numelor atribuite culorilor. Prin asocierea unei descrieri textuale fiecărei culori, ne putem face o imagine mentală asupra proprietăților fizice ale acesteia. Problema care apare, având în vedere diversitatea putem spune infinită de culoare, este modalitatea de denumire a culorilor existente în natură. Soluția adoptată este realizarea pe baza expertizei manuale a unumitor spații de culoare a unor dicționare de nume de culori. Acestea, în funcție de aplicație, pot

conține o paletă mai vastă sau mai redusă de culoare precum și un nivel de detaliu al descrierii textuale variabil. Dicționarele de culoare sunt preferate metodelor de denumire automată a culorilor, pe de-o parte datorită preciziei mult mai bune a acestora (fiind constituite de operatori umani) dar și datorită faptului că nu necesită operații de calcul, numele culorilor fiind disponibile în mod direct prin accesarea dicționarului.

Cunoașterea percepției conținutului de culoare la nivel de pixel nu este suficientă pentru înțelegerea acestuia. În transmiterea informației vizuale, datorită modului în care este construit sistemul vizual uman, un rol definitiv îl au relațiile existente între diversele regiuni de culoare din imagine. De exemplu, percepția unei regiuni dintr-o imagine în care culorile pixelilor au o distribuție similară unei table de săh, unde în loc de alb se află de exemplu o culoare arbitrară deschisă (efect ce poate fi obținut în urma unei reducerii de culoare de tip "dithering"), este de regiune întunecată, în ciuda faptului că analizate individual, procentul culorilor deschise este egal cu cel al culorilor închise (negru în acest caz). Acest lucru se datorează faptului că la nivel vizual culorile sunt interpretate în corelație cu regiunile vecine de culoare. O modalitate eficientă de analiză a relației de culoare este pe baza metodelor folosite în domeniul artei, domeniu ce a dat naștere primelor studii de percepție de culoare. În artă, relațiile dintre culori sunt analizate fie prin proiecția acestora pe roți de reprezentare a culorilor, fie folosind teoria contrastelor și a armoniei de culoare.

Ca tendință generală a sistemelor existente, putem menționa că informația de culoare este folosită cu predilecție pentru caracterizarea conținutului static al imaginilor fixe, cum este cazul indexării după conținut a acestora. În domeniul indexării după conținut a secvențelor de imagini, metodele bazate exclusiv pe culoare sunt foarte puține. Acest lucru se datorează în principal faptului că informația de culoare, analizată individual, nu are suficientă putere discriminatorie pentru a oferi o caracterizare globală a conținutului. Din această cauză, în marea majoritate a metodelor de analiză existente, culoarea este folosită alături de alte informații definitorii pentru secvențele de imagini, precum mișcarea și structura temporală.

## CAPITOLUL 5

---

### Rezumarea automată de conținut

---

**Rezumat:** Accesarea informației utile dintr-o colecție vastă de secvențe de imagini este de regulă o operație dificilă ce necesită un timp important. Acest lucru se datorează în principal volumului foarte mare de date ce trebuie vizualizate. O soluție constă în reducerea redundanței informaționale prin generarea automată de reprezentări compacte de conținut, numite și rezumate. În funcție de rezultatul obținut, există două categorii distincte de rezumat: rezumatele statice ce furnizează o colecție de imagini reprezentative ale secvenței și rezumatele dinamice ce furnizează o colecție de pasaje ale secvenței ce sunt asamblate tot sub forma unei secvențe. În acest capitol vom discuta diversele metode și tehnici de generare automată a rezumatelor din cele două categorii, precum și aplicațiile acestora. De asemenea, vom discuta un punct critic al metodelor de rezumare automată ce îl constituie problematica evaluării calității și pertinenței acestora.

În general, bazele sau colecțiile de secvențe de imagini conțin la ordinul a mii de secvențe. Cum un sigur minut video, la o cadență de 25 de imagini pe secundă (standardul PAL<sup>1</sup>), este echivalentul a 1500 de imagini statice, putem doar să ne imaginăm numărul gigantesc de imagini conținute într-o astfel de colecție de date. Pentru a avea acces la informație, utilizatorul

---

<sup>1</sup>pentru o descriere a standardelor de televiziune vezi explicațiile de la pagina 116.

în cele mai multe cazuri, are nevoie să vizualizeze conținutul secvențelor. Problema nu este una dificilă în cazul în care ar fi vorba de doar câteva secvențe, dar vizualizarea a mii de secvențe este un lucru aproape imposibil de realizat. Una dintre soluțiile adoptate constă în folosirea *rezumatelor de conținut*.

Un rezumat al unei secvențe de imagini poate fi definit în linii mari ca fiind *o reprezentare compactă a conținutului acesteia* [Li 01]. Mai riguros, rezumatul unei secvențe de imagini reprezintă *o colecție, de dimensiuni reduse, de imagini fixe (colecție de imagini) sau în mișcare (colecție de segmente), ce redă conținutul secvenței în așa fel încât partea esențială a acestuia să fie conservată iar utilizatorul să fie informat rapid și concis* [Pfeiffer 96].

Interesul în a dispune de o reprezentare compactă a secvenței nu se rezumă doar la reducerea timpului necesar căutării și navigării în conținutul bazei de date. Rezumatul de conținut poate fi folosit și pentru a reduce timpul de calcul în anumite metode de analiză și prelucrare a secvențelor de imagini, prin reducerea volumului de date ce trebuie prelucrate. Mare majoritatea a metodelor existente nu folosesc integral conținutul secvenței, ci numai informația furnizată de un anumit număr sau grup de imagini reprezentative ("imagini cheie"). Acestea sunt alese astfel încât, pentru prelucrarea vizată, informația necesară din secvență să nu fie alterată. De exemplu, pentru a calcula și analiza distribuția globală de culoare a unei secvențe de imagini, folosirea unei singure imagini pentru fiecare plan video oferă o precizie similară folosirii tuturor imaginilor din secvență [Ionescu 05c].

Din punct de vedere al procesului de generare, un rezumat poate fi construit *manual, semi-automat* (intervenția umană este parțială, de regulă folosită ca validare) sau în mod *automat*. Având în vedere volumul mare de date conținute chiar la nivelul unei singure secvențe, metodele de generare manuală, și chiar semi-automată, sunt evitate datorită implicării unui număr important de resurse umane în procesul de selecție a conținutului reprezentativ al secvenței. Tendința actuală este de automatizare completă a procesului de rezumare pentru a putea fi astfel folosit în timp real la indexarea conținutului bazelor de secvențe de imagini.

După cum reiese din însăși definiția conceptului de rezumat, este posibilă generarea a două categorii distincte de rezumat, și anume:

- pe de-o parte sunt **rezumatele în imagini** sau *rezumatele statice*: acestea reprezintă un fel de "storyboard"<sup>2</sup> simplificat al secvenței și sunt

---

<sup>2</sup>un "storyboard" reprezintă o modalitate de organizare grafică sub forma unei serii de ilustrații sau imagini, prezentate similar unei benzi desenate ("comics"), a momentelor importante din conținutului unui document video, film, etc. Aceasta are ca scop previzualizarea conținutului și precede procesul de creare propriu-zisă.

la bază o colecție de imagini reprezentative pentru conținutul secvenței. În literatura de specialitate acestea sunt cunoscute sub numele de ”video summaries”.

- pe de altă parte sunt **rezumatele în mișcare** sau *rezumatele dinamice*: acestea reprezintă o colecție de segmente ale secvenței, fiind ele însese niște secvențe de imagini, dar de o durată mult inferioară secvenței inițiale. Uzual, dacă este vorba de rezumarea unui document video, în rezumatul dinamic este prezentă și informația audio. Rezumatele dinamice sunt cunoscute în literatura de specialitate sub numele de ”video skims”.

Acstea două modalități de rezumare a conținutului unei secvențe de imagini prezintă fiecare o serie de avantaje și dezavantaje. Astfel, în ceea ce privește rezumatele în imagini, principalele avantaje ale acestora pot fi sintetizate cu următoarele:

- pot fi *generate rapid* deoarece nu iau în calcul decât informația vizuală (sunetul și textul nu sunt prezente),
- pot fi *vizualizate foarte ușor*, fiind doar un ansamblu de imagini ce nu necesită sincronizarea sau temporizarea datelor (de exemplu, sincronizarea sunetului și a imaginii),
- pot facilita construirea *imaginilor de tip ”mosaic”*<sup>3</sup> [Aner 01],
- pot fi ușor de imprimat pe un suport fizic pentru a ține loc, de exemplu, de ”*Storyboard*” al secvenței,
- permit *reducerea complexității de calcul* pentru anumite metode de analiză ce pot fi aplicate direct acestora.

Pe de altă parte, rezumatele dinamice prezintă și ele o serie de avantaje fundamentale, astfel:

- acestea *au mai mult sens* decât rezumatele statice deoarece conțin informație temporală de mișcare, informație ce este pierdută în rezumatele statice,
- un rezumat dinamic este *mai bogat în informație*, acesta putând conține și alte informații precum sunetul,

---

<sup>3</sup>vezi exlicația de la pagina 150.

- vizualizarea rezumatelor dinamice este mult mai *naturală și atractivă*: este mult mai interesant pentru utilizator să vizualizeze, de exemplu, reclama unui nou film în curs de apariție ("trailer"), decât să vizualizeze o succesiune sacadată de imagini statice din acesta [Li 01].

În ciuda faptului că rezumatele dinamice oferă un conținut informațional mult mai bogat decât rezumatele statice, acestea implică în cele mai multe situații o complexitate de calcul mult mai importantă, precum și un proces de genereare mult mai laborios (sincronizare imagine și sunet, respectarea continuității și a coerenței vizuale, etc.).

În acest punct, analizând avantajele și dezavantajele celor două tipuri de rezumate, am fi tentați să alegem doar unul dintre ele ca fiind soluția optimală a problemei de rezumare de conținut. În practică, *ambele tipuri de rezumate sunt necesare* într-o aplicație de indexare după conținut, deoarece fiecare dintre acestea este adaptat unei cerințe diferite. Astfel, rezumatul în imagini permite o reprezentare rapidă și concisă, în doar câteva imagini, a conținutului vizual, ideală în cazul în care utilizatorul dorește doar să "răsfoiască" conținutul bazei de date, în timp ce rezumatul dinamic permite o reprezentare rapidă și concisă a conținutului dinamic al secvenței, permitând utilizatorului ca într-un timp relativ scurt să-și facă o idee asupra acțiunii secvenței.

Mai mult, cu toate că nu este cea mai bună strategie, fiecare tip de rezumat poate fi generat pornind de la celălalt. De exemplu, un rezumat dinamic poate fi construit pe baza unui rezumat static prin concatenarea unor segmente de o anumită durată ce conțin imaginile rezumatului static. Similar, un rezumat static poate fi generat dintr-un rezumat dinamic prin păstrarea doar a unor imagini reprezentative sau chiar prin sub-eșantionarea temporală a acestuia.

Pentru un studiu bibliografic complet al literaturii de specialitate din această direcție de studiu, cititorul se poate raporta la lucrările [Li 01] sau [Truong 07]. În cele ce urmează vom face o trecere în revistă a particularităților metodelor folosite de fiecare dintre cele două categorii de rezumat.

## 5.1 Construcția rezumatelor statice

După cum am menționat în partea introductivă a acestui capitol, un rezumat static este la bază o colecție de imagini fixe ce sunt considerate ca fiind reprezentative pentru conținutul secvenței. Acestea sunt numite și "imagini cheie". Din punct de vedere formal, rezumatul static al secvenței  $S$ , notat

$R_{img}(S)$ , poate fi definit în felul următor:

$$R_{img}(S) = \{image_1, image_2, \dots, image_N\} \quad (5.1)$$

unde  $image_i$  reprezintă imaginea cheie de indice  $i$ , cu  $i = 1, \dots, N$  iar  $N$  reprezintă numărul total de imagini din rezumat sau dimensiunea rezumatului.

Valoarea parametrului  $N$  are un rol important asupra calității rezumatului. Dacă  $N$  este cunoscut ”*a priori*”, atunci dimensiunea rezumatului va fi folosită drept constrângere inițială a algoritmului de extragere a imaginilor cheie. În acest caz, dimensiunea rezumatului va fi aceeași pentru toate secvențele analizate. Pe de altă parte, dacă  $N$  nu este fixat inițial, atunci valoarea acestuia va fi determinată automat de algoritmul de calcul. În acest caz dimensiunea rezumatului va fi adaptată la conținutul fiecărei secvențe în parte (de exemplu, un conținut bogat în acțiune va fi reprezentat cu mai multe imagini decât un conținut static).

### 5.1.1 Clasificarea metodelor existente

Din punct de vedere al modului în care sunt extrase imaginile cheie, metodele existente de calcul al rezumatelor statice pot fi clasificate în următoarele categorii [Li 01]:

- *extragerea imaginilor prin eșantionare,*
- *extragerea imaginilor la nivel de plan,*
- *extragerea imaginilor la nivel de segment video,*
- *alte abordări.*

#### Extragerea imaginilor cheie prin eșantionare

Extragerea prin eșantionare constă în selectarea imaginilor cheie, fie în mod aleator, fie ca fiind uniform distribuite temporal în secvență [Taniguchi 95]. Avantajul unei astfel de abordări constă în primul rând în simplitatea și rapiditatea acesteia, complexitatea de calcul fiind în acest caz total neglijabilă.

Totuși, rezumatul obținut nu este neapărat reprezentativ pentru conținutul secvenței deoarece acesta nu este luat în calcul la momentul generării. De exemplu, folosind această strategie este foarte probabil ca anumite plane video de scurtă durată, dar reprezentative pentru acțiunea secvenței, să nu fie reprezentate în rezumat, în timp ce anumite plane de lungă durată, cu un conținut informational redundant, să fie reprezentate în rezumat cu mai multe imagini cheie similare.

### Extragerea imaginilor cheie la nivel de plan

O metodă mai elaborată constă în extragerea imaginilor cheie la nivelul planelor video.

Într-o variantă simplificată, selectarea imaginilor cheie se poate realiza prin alegerea în mod arbitrar a unei singure imagini pentru fiecare plan video, de exemplu: prima imagine, imaginea de mijloc, o imagine aleatoare, ultima imagine, etc. Această strategie dă rezultate vizuale satisfăcătoare în cazul în care conținutul secvenței prezintă o variabilitate medie sau redusă, oferind un bun compromis între calitate și complexitatea de calcul. Pe de altă parte, complexitatea de calcul în acest caz este dată în principal doar de segmentarea în plane. Luând în calcul faptul că de regulă orice metodă de analiză a conținutului video necesită segmentarea temporală a secvenței, se poate presupune că în momentul calculării rezumatului, aceasta este deja disponibilă. În aceste condiții, putem considera complexitatea de calcul a rezumatului ca fiind neglijabilă.

Totuși, în cazul unui conținut bogat în acțiune, o singură imagine nu este întotdeauna suficientă pentru a reprezenta conținutul vizual la nivel de plan, existând riscul ca aceasta să fie o imagine de tranziție a unui efect vizual sau a unei mișcări rapide. Faptul că imaginile cheie sunt alese la nivel de plan video, va asigura totuși că acestea nu sunt imagini de tranziție ale tranzițiilor video prezente în secvență (de exemplu: "fades", "dissolves", etc., vezi Secțiunea 2.1).

Pentru a ilustra cele menționate anterior, în Figura 5.1 am comparat rezultatele obținute cu strategiile simplificate de extragere a imaginilor cheie la nivel de plan cu o metodă de rezumare ce ia în calcul conținutul vizual al planului, și anume alegerea imaginii mediane<sup>4</sup> ca imagine cheie. Se poate observa că strategiile de alegere a imaginilor cheie ce nu iau în calcul analiza conținutul planului furnizează rezultate satisfăcătoare pentru un conținut cvazi-constant al planului (vezi primul plan din Figura 5.1), în timp ce pentru un conținut complex (vezi al doilea plan din Figura 5.1) imaginea mediană este mai pertinentă pentru acțiunea globală a planului. De notat este faptul că anumite studii au demonstrat că din punct de vedere statistic, uneori o strategie de alegere aleatoare poate furniza rezultate similare sau chiar superioare unei strategii adaptate conținutului. Această ipoteză are sens în virtutea faptului că prin definiție un plan video prezintă o omogenitate a conținutului (temporală, spațială și de acțiune), astfel că teoretic oricare dintre imaginile acestuia este reprezentativă pentru conținut.

O altă metodă, ce ține cont de această dată de conținutul vizual al

<sup>4</sup>Imaginea mediană a unui anumit plan este definită ca fiind imaginea cea mai apropiată, în sensul unei măsuri de distanță, de ansamblul celorlalte imagini din plan [Chanussot 98].

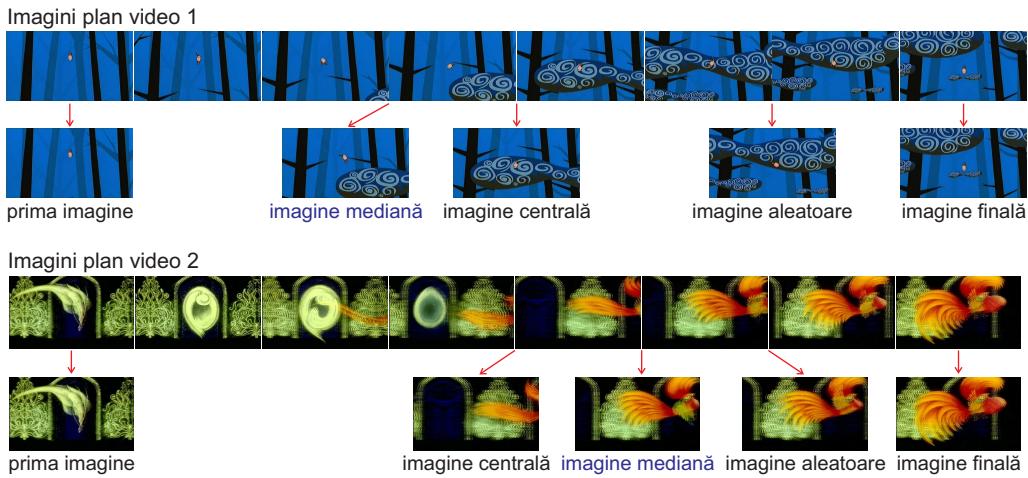


Figura 5.1: Exemplu de rezumare cu o singură imagine pentru o mișcare de obiecte (primul plan), precum și pentru o schimbare complexă a scenei (al doilea plan, axa orizontală este axa temporală, sursă imagini [Folimage 06]).

planelor video, constă în extragerea imaginilor cheie pe baza *anализei статистиче-ской и дистрибуции цвета*. Această strategie a fost motivată în principal de eficiența și robustețea histogramelor de culoare în reprezentarea globală a conținutului de culoare la nivel de imagine (vezi Secțiunea 4.2.1). Un exemplu este metoda propusă în [Zhang 97] în care diferențele dintre histogramele color ale imaginilor succesive sunt filtrate cu un anumit prag pentru a extrage imaginile cele mai diferite din punct de vedere vizual. O altă abordare este propusă în [Zhuang 98] unde imaginile cheie sunt extrase folosind un algoritm de clasificare nesupervizată în care măsura de disimilaritate dintre imagini este dată de distanța dintre histogramele de culoare.

O abordare interesantă inspirată de modalitatea de calcul a medianului vectorial este propusă în [Ott 07]. Metoda propusă rezumă conținutul fiecărui plan video cu un număr de imagini ce este adaptat variabilității conținutului vizual al acestuia. Alegerea acestora este realizată în funcție de numărul de moduri<sup>5</sup> ale unei histograme de distanțe cumulate, unde distanța cumulată a unei imagini este definită ca fiind suma distanțelor dintre histograma acesteia și histogramele tuturor celorlalte imagini ale planului.

O limitare a metodelor din această categorie este dată de faptul că nu iau în calcul conținutul de mișcare. Astfel, imaginile selectate, chiar dacă sunt diferite între ele, este posibil să aparțină unei aceleasi mișcări, ca de exemplu a unei mișcări globale a camerei video. Mai mult, aceste metode sunt

<sup>5</sup>vezi explicația de la pagina 127.

dependente de alegerea unui anumit număr de praguri folosite la definirea similarității dintre imagini.

Abordările bazate pe *analiza mișcării* sunt concepute pentru a adapta numărul de imagini cheie extrase la dinamica temporală a scenei. Un exemplu este metoda propusă în [Wolf 96] unde imaginile cheie sunt extrase ca fiind minime locale ale unei anumite funcții temporale ce măsoară activitatea de mișcare a secvenței.

Acest gen de abordare prezintă și el o serie de dezavantaje. În primul rând este dificil să determinăm importanța unui segment al secvenței folosind doar criterii de mișcare. Mai mult, ca și în cazul precedent, este posibil cu această strategie să obținem imagini de tranziție, de regulă încețoșate, specifice unei mișcări foarte rapide.

Un alt tip de abordare, care se folosește de asemenea de informația de mișcare, este construcția *imaginilor de tip "mozaic"*. Acestea sunt imagini panoramice ce reprezintă întregul conținut al unui plan video sau chiar al unei porțiuni a secvenței cu o singură imagine [Irani 95] (vezi Figura 5.2). Avantajul imaginilor "mozaic" constă în faptul că acestea înglobează întreaga dinamică a scenei în doar o singură imagine. Pe de altă parte, acestea nu au sens decât pentru pasaje ale secvenței ce conțin o mișcare globală a camerei video, neputând fi astfel generate pentru pasaje statice sau pasaje ce conțin mișcări complexe ale scenei. Pentru a compensa acest neajuns, de regulă este adoptată o strategie hibridă: imaginile "mozaic" sunt determinate doar în situațiile în care acest lucru este posibil, rămânând ca în celelalte cazuri conținutul planelor video să fie reprezentat în mod clasic cu imagini cheie.



Figura 5.2: Exemplu de imagine de tip "mozaic" obținută pentru o deplasare la dreapta a camerei video (sursă "[http://scien.stanford.edu/2002projects/ee392j/Roman\\_Gilat/](http://scien.stanford.edu/2002projects/ee392j/Roman_Gilat/)").

Pe lângă metodele enumerate anterior, putem menționa și o serie de abordări mixte ce folosesc colaborarea mai multor surse de informație din secvență, precum culoarea și mișcarea, profitând astfel de avantajele fiecăreia dintre acestea. Un exemplu în acest sens este metoda din [Doulamis 00a],

unde imaginile cheie sunt extrase ca fiind puncte similare ale curbei formate de vectorii de caracteristici ai fiecărei imagini. Aceştia sunt obținuți pe baza unui proces de segmentare aplicat, atât la nivel de culoare cât și de mișcare.

În general metodele mixte, cu toate că sunt mai dificil de implementat datorită complexității de calcul mai ridicate, oferă rezultate mai bune decât celelalte metode datorită diversității surselor de informații folosite.

### **Extragerea imaginilor cheie la nivel de segment**

Una dintre constrângerile strategiei de extragere a imaginilor cheie la nivel de plan video este dată de lipsa de scalabilitate a unei astfel de abordări. În cazul secvențelor de imagini cu o durată semnificativă, numărul de imagini obținut cu o astfel de abordare este ridicat, ajungând ușor la ordinul sutelor, ceea ce face nerentabilă vizualizarea integrală de către utilizator a acestora.

Soluția la această problemă constă în folosirea unităților structurale de nivel ierarhic superior planelor video (vezi Secțiunea 2.1), numite generic și *segmente video*: scene, episoade, anumite evenimente sau chiar secvență în totalitate. Astfel, în funcție de aplicație, se poate opta pentru un nivel de granularitate, de la ridicat (număr mai mare de imagini, de exemplu obținut prin extragerea la nivel de plane sau scene), până la un nivel redus (număr redus de imagini, de exemplu prin extragerea la nivel de episoade sau evenimente).

O astfel de metodă ce folosește un nivel structural superior planelor video este propusă în [Sun 00]. Într-o primă etapă, întreaga secvență este segmentată în mod uniform, în ceea ce autorii numesc "unități de lungă durată" (segmente). Pentru fiecare segment astfel obținut este calculată o măsura de schimbare prin evaluarea distanței dintre histogramele primei și respectiv ultimei imagini a segmentului. Segmentele din întreaga secvență sunt ordonate și clasificate pe baza valorilor acestei măsură, în două clase, și anume: clasa *A* ce va conține valorile de schimbare redusă și clasa *B* ce va conține valorile de schimbare importantă. Mai departe, imaginile cheie sunt extrase în felul următor: pentru segmentele din clasa *A*, prima și ultima imagine a segmentului vor fi selectate ca imagini cheie, iar pentru segmentele din clasa *B*, toate imaginile segmentului vor fi selectate ca imagini cheie. Dacă rezumatul obținut satisface numărul de imagini dorit, atunci algoritmul se încheie. În caz contrar, rezumatul obținut va ține loc de secvență, iar procesul de selecție este repetat. Limitarea acestei metode este dată de măsura de similaritate folosită, diferența dintre prima și ultima imagine a segmentului nefiind suficientă pentru a evalua variabilitatea conținutului acestuia. Astfel, metoda va neglija conținutul anumitor segmente importante pentru care imaginea de început și de sfârșit sunt similare, și respectiv va reprezenta

cu un număr important de imagini redundante segmentele cu un conținut omogen, dar care încep și se termină cu imagini relativ diferite.

Metoda propusă în [Doulamis 00b] introduce o strategie fuzzy pentru extragerea imaginilor cheie. În primă fază, fiecare imagine a secvenței este segmentată folosind algoritmul RSST ("Recursive Shortest Spanning Tree", vezi [Kwok 04]). Atributele de culoare și de mișcare obținute în urma segmentării sunt clasificate într-un număr predefinit de clase folosind o clasificare fuzzy (vezi Secțiunea 7.1.2). Imaginile cheie sunt extrase mai departe prin minimizarea unei funcții de cost, determinată cu un algoritm genetic pe baza unei măsuri a intercorelației dintre date. În acest caz, principalul inconvenient al acestei metode este dat de complexitatea de calcul ridicată.

Un alt exemplu este metoda propusă în [Gong 00] ce folosește descompunerea în valori singulare sau SVD ("Singular Value Decomposition"<sup>6</sup>). Pentru întreaga secvență este creată mai întâi o matrice  $A$  de vectori de caracteristici. Aceștia sunt obținuți pentru un anumit set, suficient de vast, de imagini din secvență. Algoritmul SVD este aplicat matricei  $A$ . Pe baza acestuia se obține pe de-o parte rafinarea spațiului de caracteristici ce va facilita o mai bună clasificare a imaginilor similare din punct de vedere vizual, dar și determinarea unei metrici ce va permite măsurarea cantitativă a schimbărilor vizuale din fiecare clasă de imagini. Mai departe, în spațiul de caracteristici obținut după aplicarea SVD, clasa cu conținutul cel mai puțin variabil este luată ca referință, iar valoarea de schimbare a conținutului vizual al acesteia în corelație cu distanța dintre imagini sunt folosite ca praguri pentru clasificarea restului de imagini. În cele din urmă, imaginile cheie sunt alese ca fiind imaginile cu vectorii de caracteristici cei mai apropiati de centroizii claselor. Această abordare asigură o minimă redundanță informațională a rezumatului, precum și posibilitatea ca numărul de imagini din rezumat să fie reglat de utilizator. Totuși, o constrângere este dată de faptul că datorită clasificării, în rezumatul obținut nu se va mai respecta ordinea temporală inițială a secvenței.

## Alte abordări

Alte abordări de extragere a imaginilor cheie se folosesc de alte spații de reprezentare a informației vizuale, precum spațiul transformatei "wavelet"<sup>7</sup> [Campisi 99] sau se bazează pe localizarea în secvență a anumitor pasaje de

---

<sup>6</sup>în algebra liniară, descompunerea în valori singulare sau SVD, reprezintă descompunerea ("factorizarea") unei matrice pătratice de numere reale sau complexe la forma canonică (forma standard). Ca exemple de aplicații ce folosesc SVD putem menționa: calcularea matricei pseudo-inverse, aproximarea unei matrice sau determinarea rangului.

<sup>7</sup>vezi explicația de la pagina 104.

interes, ca de exemplu prezența fețelor umane sau a culorii pielii [Dufaux 00]. Metoda propusă în [Dufaux 00] se folosește de analiza de mișcare precum și de analiza activității spațiale, ce integrează detecția fețelor umane și a prezenței culorii pielii în imagine, pentru a rezuma conținutul secvenței cu o singură imagine. În acest caz, sunt considerate ca reprezentative imaginile ce conțin personaje sau portrete.

O altă abordare ce se folosește de numărul de obiecte prezente în scenă pentru a localiza pasajele importante ale secvenței, este propusă în [Kim 00a]. O imagine cheie este extrasă de fiecare dată când numărul de obiecte prezente în imaginea curentă este diferit de numărul de obiecte din imaginea ulterioară temporal.

Principalul inconvenient al metodelor bazate pe localizarea anumitor pasaje de interes din secvență este dat de faptul că acestea sunt dependente și adaptate fiecărui domeniu de aplicație, sau în unele situații, chiar tipului de secvență folosit, neavând astfel un caracter generic.

### 5.1.2 Mecanismul de extragere a imaginilor cheie

Metodele existente de construcție a rezumatelor statice folosesc o serie de mecanisme de selectare a imaginilor cheie din secvență. Acestea, în funcție de strategia folosită, se împart în următoarele categorii [Truong 07]:

- selectarea imaginilor pe baza *schimbării suficiente* a conținutului,
- selectarea imaginilor pe baza *egalizării varianței temporale*,
- selectarea imaginilor astfel încât acestea să asigure o *acoperire maximală* a conținutului,
- selectarea imaginilor cu metode de *clasificare*,
- selectarea imaginilor pe baza *corelației minimale* dintre acestea,
- selectarea imaginilor astfel încât *eroarea de reconstrucție* a secvenței să fie minimală,
- selectarea imaginilor prin *simplificarea curbei de caracteristici*,
- selectarea imaginilor pe baza *localizării evenimentelor* importante din secvență.

În cele ce urmează, vom prezenta particularitățile fiecărei dintre aceste abordări aşa cum au fost prezentate în [Truong 07].

### Schimbarea suficientă de conținut

Această abordare analizează imaginile secvenței în mod secvențial. O imagine analizată este selectată ca imagine cheie dacă conținutul vizual al acesteia prezintă o schimbare importantă, relativ la imaginile cheie extrase anterior.

În general, metodele ce folosesc criteriul schimbării suficiente de conținut aleg o nouă imagine cheie, de indice  $r_{i+1}$  în rezumat, în funcție de imaginea cheie extrasă anterior, de indice  $r_i$ , pe baza relației următoare:

$$r_{i+1} = \operatorname{argmin}|_t \{C(I_t, I_{r_i}) > \varepsilon, i < t \leq N\} \quad (5.2)$$

unde  $I_t$  reprezintă imaginea de indice temporal  $t$  în pasajul analizat al secvenței (de exemplu, acesta poate fi un plan video, o scenă, etc.),  $N$  reprezintă numărul total de imagini conținute în acesta,  $C()$  este o funcție ce măsoară schimbarea conținutului vizual iar  $\varepsilon$  reprezintă pragul ce definește o schimbare de conținut ca fiind semnificativă. De regulă, prima imagine a rezumatului, de indice  $r_1$ , este aleasă ca fiind prima imagine a pasajului analizat. Dintre cele mai frecvent folosite funcții  $C()$  putem menționa: diferența dintre histogramele de culoare, funcția de energie cumulativă sau schimbarea proprietăților geometrice ale obiectelor, etc.

Unul dintre principalele avantajele ale acestei abordări este dat de faptul că rezumatul static obținut va avea un număr variabil de imagini, fiind adaptat la conținutul secvenței și reflectând evoluția dinamică a acesteia, indiferent de durata sau de tipul conținutului de acțiune. Pe de altă parte, principalul inconvenient este dat de faptul că imaginile cheie sunt extrase fără a ține cont de conținutul secvenței ce precede imaginea cheie curentă, ceea ce face ca rezumatul obținut să nu ofere o acoperire maximală a conținutului secvenței [Truong 07].

În [Xiong 97] este propusă o extensie a acestei abordări către metoda numită "Seek and Spread", ce nu restricționează alegerea noii imagini cheie în funcție de imaginea cheie curentă din rezumat. Dacă  $b_i$  reprezintă indicele imaginii de referință curente, cu  $b_1 = 1$ , atunci  $r_i$ , indicele imaginii cheie, și  $b_{i+1}$ , indicele noii imagini de referință, sunt calculate în felul următor:

$$r_i = \operatorname{argmin}|_t \{C(I_t, I_{b_i}) > \varepsilon, t > b_i\} \quad (5.3)$$

$$b_{i+1} = \operatorname{argmin}|_t \{C(I_t, I_{r_i}) > \varepsilon, t > r_i\} \quad (5.4)$$

unde imaginea cheie  $I_{r_i}$ , din rezumat, este imaginea considerată ca reprezentativă pentru segmentul din secvență delimitat de indicii  $[b_i, b_{i+1}-1]$ . Măsura  $C()$  este calculată în acest caz ca fiind distanța Euclidiană dintre o serie de vectori de caracteristici obținuți pe baza descompunerii "wavelet". Totuși, în ciuda îmbunătățirii aduse, această metodă nu este scalabilă iar rezumatul obținut este asimetric.

Un alt exemplu este metoda propusă în [Rasheed 03]. Rezumatul obținut cu aceasta este mai concis și se bazează pe extragerea unei imagini cheie de fiecare dată când imaginea curentă analizată diferă suficient față de toate imaginile cheie extrase anterior, astfel fiind luat în calcul, chiar dacă la un nivel de detaliu redus, și conținutul precedent imaginii curente din rezumat.

### Egalizarea varianței temporale

În această abordare, numărul de imagini cheie al rezumatului static este cunoscut ”a priori”. Principiul metodei de egalizare a varianței temporale constă în alegerea imaginilor cheie (indice  $r_i$  în rezumat) ca fiind reprezentative pentru anumite segmente ale secvenței (intervale  $[b_i, b_{i+1} - 1]$ ) ce prezintă o varianță temporală similară.

Valorile indicilor  $r_i$  și  $b_i$  vor fi determinate independent, în doi pași. Din punct de vedere teoretic, frontierele  $b_i$  ale segmentelor secvenței sunt selecționate astfel încât să satisfacă egalitatea următoare:

$$\nu(b_1, b_2) = \nu(b_2, b_3) = \dots = \nu(b_k, b_{k+1}) \quad (5.5)$$

unde  $\nu(b_i, b_{i+1})$  reprezintă varianța temporală a segmentului  $[b_i, b_{i+1} - 1]$ , segment ce este reprezentat în rezumat cu imaginea cheie  $I_{r_i}$ .

Cum în realitate egalitatea exactă a valorilor varianței este aproape imposibil de obținut, [Sun 00] propune transformarea acesteia într-o problemă de optimizare, astfel:

$$\{b_1, b_2, \dots, b_{k+1}\} = \operatorname{argmin}_{b_i} \sum_{i=1}^k \sum_{j=1}^k |\nu(b_i, b_{i+1}) - \nu(b_j, b_{j+1})| \quad (5.6)$$

Funcția de varianță temporală,  $\nu()$ , poate fi aproximată cu o măsură a schimbării cumulative a conținutului dintre imagini sau chiar cu diferența dintre prima și ultima imagine a fiecărui segment. Această ipoteză permite soluționarea ecuației 5.6 într-un mod recursiv.

După localizarea segmentelor ce îndeplinesc criteriul egalității varianței temporale, imaginile cheie pot fi extrase folosind mai multe strategii. De exemplu, în [Fauvet 04] imaginile cheie sunt imaginile de început și de mijloc ale fiecărui segment. O abordare mai complexă este propusă în [Lee 02a], unde imaginea  $I_{r_i}$  este selectată ca imagine cheie reprezentativă pentru conținutul segmentului  $[b_i, b_{i+1} - 1]$  prin minimizarea diferenței cumulative următoare:

$$r_i = \operatorname{argmin}_{b_i \leq t < b_{i+1}} \sum_{j=b_i}^{b_{i+1}-1} D(I_j, I_t) \quad (5.7)$$

unde  $D()$  reprezintă un anumit operator de distanță între imagini.

În general, metodele din această categorie prezintă o complexitate de calcul superioară metodelor ce folosesc schimbarea suficientă a conținutului vizual, drept criteriu de selecție al imaginilor cheie, dar, în revanșă, rezumatul obținut este de această dată optimal din punct de vedere global.

### Acoperirea maximală a conținutului

Această strategie este cunoscută și sub numele de metoda de selecție a imaginilor cheie pe baza criteriului de fidelitate. Principiul metodei constă în a determina pentru fiecare imagine din secvență, sau dintr-un anumit pasaj al acesteia, a unei liste de imagini pe care aceasta le poate substitui, listă ce poartă numele de acoperire a imaginii.

Dacă notăm cu  $C_i(\varepsilon)$  lista imaginilor din pasajul  $V$  al secvenței ce pot fi substituite cu imaginea cheie  $I_i$  pe baza unei valori de toleranță  $\varepsilon$ , atunci, setul optimal al imaginilor cheie din  $V$ , determinat fără constrângeri de debit (număr de imagini), este dat de ecuația următoare:

$$\{r_1, r_2, \dots, r_k\} = \operatorname{argmin}_{|r_i|} \{k / C_{r_1}(\varepsilon) \cup C_{r_2}(\varepsilon) \cup \dots \cup C_{r_k}(\varepsilon) = V\} \quad (5.8)$$

unde  $r_i$  reprezintă indicele imaginii cheie din rezumat.

Dacă alegem drept constrângere numărul de imagini din rezumat, atunci problema selecției imaginilor cheie poate fi văzută din două perspective. Pe de-o parte, putem căuta valoarea minimală a lui  $\varepsilon$  pentru care toate imaginile pot fi reprezentate cu cel puțin o singură imagine cheie, astfel:

$$\{r_1, r_2, \dots, r_k\} = \operatorname{argmin}_{|r_i|} \{\varepsilon / C_{r_1}(\varepsilon) \cup C_{r_2}(\varepsilon) \cup \dots \cup C_{r_k}(\varepsilon) = V\} \quad (5.9)$$

sau pe de altă parte, putem căuta un grup de imagini ce pot fi reprezentative pentru cea mai mare parte a imaginilor disponibile, astfel:

$$\{r_1, r_2, \dots, r_k\} = \operatorname{argmin}_{|r_i|} \{|C_{r_1}(\varepsilon) \cup C_{r_2}(\varepsilon) \cup \dots \cup C_{r_k}(\varepsilon)|\} \quad (5.10)$$

Astfel, în procesul de construcție al rezuma lui apar două etape distincte, și anume: determinarea *acoperirii imaginilor* în secvență și respectiv *optimizarea selecției* făcute.

În [Yahiaoui 01] acoperirea unei imagini este dată de numărul de extrase din secvență ("excerpts"), de dimensiune fixă, ce conțin cel puțin o imagine similară cu imaginea curent analizată. Imaginile rezumatului sunt selectate mai departe folosind procedeul programării dinamice. Un alt exemplu este metoda propusă în [Rong 04] ce face analogie între procesul de selecție al imaginilor cheie din secvență și procesul de extragere a cuvintelor pe baza

metodei TF-IDF ("Term-Frequency Inverse Document Frequency"<sup>8</sup>) folosit în sistemele de căutare a textului. Astfel, o imagine a secvenței va fi selectată ca fiind imagine cheie reprezentativă pentru un anumit plan video, dacă aceasta are o acoperire redusă pentru restul conținutului secvenței.

Unul dintre principalele avantaje ale metodelor din această categorie este dat de faptul că nu se impune ca acoperirea unei imagini să fie un segment continuu din punct de vedere temporal, cum este cazul metodelor ce folosesc ca criteriu schimbarea suficientă de conținut sau egalitatea varianței temporale. Acest lucru, va permite obținerea unui rezumat mai concis. Pe de altă parte, complexitatea de calcul crește datorită calculului pentru fiecare pereche de imagini, a unei măsuri de similaritate.

### **Clasificarea imaginilor**

Metodele de selecție a imaginilor cheie din această categorie se folosesc de metode de clasificare automată a datelor pentru a regrupa imaginile în funcție de similaritatea conținutului acestora (pentru un studiu detaliat al metodelor de clasificare existente cititorul se poate raporta la Capitolul 7).

Principiul clasificării constă în reprezentarea fiecărei imagini cu un vector de atrbute (caracteristici) ce sunt extrase la nivel de imagine, ca de exemplu: histograme de culoare, parametri de textură, etc. Acestea formează în spațiul de caracteristici un nor de puncte. Pe baza evaluării unei anumite măsuri de similaritate între vectorii de atrbute, imaginile (punctele din spațiul de caracteristici) vor fi mai departe alocate unei anumite clase. Imaginile cheie ale rezumatului static vor fi alese din clasele considerate ca reprezentative pentru conținutul secvenței.

Procesul de constituire a claselor de imagini implică în general patru etape de analiză, și anume:

- *o etapă de pre-analiză*: aceasta are ca scop să reducă redundanța datelor de intrare și astfel să îmbunătățească eficiența algoritmului de clasificare. Printre metodele cele mai cunoscute de decorelare a datelor, putem menționa analiza în componente principale sau PCA<sup>9</sup> [Gibson 02] ce se folosește de descompunerea în vectori primi. O altă strategie constă în a reduce, înaintea clasificării, numărul de imagini

---

<sup>8</sup>metoda TF-IDF presupune calculul unei măsuri statistice folosită la evaluarea importanței unui anumit cuvânt într-un document textual dintr-o anumită colecție sau corpus. Importanța acestuia va crește proporțional cu numărul de apariții în document. Variații ale metodei TF-IDF sunt folosite frecvent în motoarele de căutare pentru a calcula relevanța și rangul unui document relativ la cererea de căutare formulată de utilizator.

<sup>9</sup>vezi explicația de la pagina 123.

din secvență folosind de exemplu imaginile cheie obținute cu una dintre metodele de rezumare prezentate anterior.

- *o etapă de clasificare propriu-zisă:* aceasta efectuează repartiția imaginilor în clase. Printre metodele cel mai frecvent folosite putem enumera următoarele: clasificarea secvențială, metoda "Hierarchical Complete Link", clasificarea Fuzzy C-Means [Yu 04] sau folosirea modelelor Gaussiene mixte (GMM - "Gaussian Mixture Models") [Gibson 02].
- *o etapă de filtrare a claselor obținute:* în urma procesului de împărțire în clase este foarte probabil ca anumite clase, de exemplu, afectate de prezența zgomotului, să nu conțină informații relevante. Pentru aceasta, este necesară o etapă de triere a claselor. Astfel, putem considera că relevante doar clasele cu o dimensiune (număr de imagini) superioară dimensiunii medii, sau, pentru a evita prezența "artefactelor", putem selecta doar clasele ce conțin cel puțin o secvență continuuă de imagini, cu o durată superioară unui anumit prag (de exemplu 9 secunde [Truong 07]).
- *o etapă de selecție a imaginilor cheie:* soluția cea mai intuitivă și cea mai frecvent folosită constă în alegerea ca imagini cheie a imaginilor ai căror vectori de caracteristici se află cel mai aproape de centroidul clasei [Yu 04]. În acest fel, avem siguranță că imaginile selectate sunt cele mai reprezentative pentru conținut, deoarece, centroidul clasei are sens de vector mediu în spațiul n-dimensional de caracteristici.

Una dintre principalele limitări ale selecției imaginilor cheie pe baza clasificării imaginilor secvenței este dată de faptul că în general, informația semantică vizuală a secvenței implică o varianță intra-clasă, în același timp importantă cât și redusă, ceea ce face dificilă separarea în clase. Mai mult, în acest caz rezumatul obținut în urma clasificării nu va respecta evoluția temporală a secvenței.

### **Corelația minimală a imaginilor cheie**

Metodele din această categorie consideră că imaginile rezumatului static trebuie să fie cât mai puțin corelate între ele. Totuși, în practică, deseori se ia în calcul doar corelația dintre imagini vecine. Prințipiu de selecție a imaginilor cheie pe baza corelației minime poate fi formulat în felul următor:

$$\{r_1, r_2, \dots, r_k\} = \operatorname{argmin}_{|r_i|} \{\operatorname{Corr}(I_{r_1}, I_{r_2}, \dots, I_{r_k})\} \quad (5.11)$$

unde  $r_i$  reprezintă indicele imaginii cheie în rezumat iar  $\operatorname{Corr}()$  reprezintă operatorul de corelație.

În [Doulamis 98] valoarea corelației pentru ansamblul imaginilor cheie este redefinită pentru a ține cont de contribuția corelației dintre fiecare pereche de imagini, astfel:

$$\text{Corr}^2(I_{r_1}, I_{r_2}, \dots, I_{r_k}) = \sum_{i=1}^{k-1} \sum_{j=i+1}^k \text{Corr}^2(I_{r_i}, I_{r_j}) \quad (5.12)$$

unde  $\text{Corr}(I_i, I_j)$  reprezintă funcția de corelație dintre vectorii de caracteristici ai imaginilor de indice  $i$  și  $j$ . Pentru ca soluția ecuației 5.11 să fie una optimală, se pot folosi o serie de metode, ca de exemplu căutarea logaritmică, căutarea stochastică sau algoritmii genetici [Doulamis 00b].

O soluție aproximativă este propusă în [Liu 02c]. Complexitatea de calcul este redusă la ordinul  $O(n \cdot \log(n))$ , unde  $n$  reprezintă numărul de imagini analizate, pe baza algoritmului "Greedy"<sup>10</sup>. La început se pornește cu toate cele  $n$  imagini ale secvenței drept imagini cheie, ceea ce constituie nivelul  $n$  de analiză. Mai departe, se trece iterativ de la un nivel de analiză la altul. La trecerea de la nivelul de analiză curent,  $t$ , la un nivel inferior,  $t - 1$ , imaginile ce prezintă o corelație minimală în raport cu imaginile vecine sunt eliminate din rezumat.

Criteriul corelației minime asigură un nivel redus de redundanță a imaginilor din rezumat. Totuși, acest criteriu este sensibil la prezența variațiilor importante din secvență (numite în literatura de specialitate și "outliers"<sup>11</sup>). Pentru a diminua acest efect, soluția propusă în [Porter 03] constă în a reprezenta secvența sub forma unui graf<sup>12</sup>. În acesta, fiecare imagine corespunde unui nod iar corelația dintre imagini, calculată în acest caz în funcție de informația de mișcare, reprezintă legăturile dintre noduri. Setul optimal de imagini cheie va fi dat de drumul cel mai scurt ce unește prima și ultima imagine a grafului.

### Eroarea minimală de reconstrucție

Metodele din această categorie se folosesc pentru construcția rezumatului de o anumită metrică, numită și SRE (eroarea de reconstrucție a secvenței),

---

<sup>10</sup>algoritmul "Greedy" este un algoritm general de soluționare a unei probleme de calcul. Aceasta este abordată în mai multe faze. Pentru fiecare fază este adoptată o soluție optimală local, fără a se ține cont de implicațiile ulterioare ale acesteia. La finalizarea algoritmului, este de așteptat ca optimul local să conveargă spre optimul global. În caz contrar, soluția obținută este considerată ca fiind suboptimală.

<sup>11</sup>în statistică matematică, un "outlier" reprezintă o observație a cărei reprezentare numerică este foarte diferită de restul valorilor datelor considerate.

<sup>12</sup>în terminologia matematică, un graf este definit ca fiind o colecție de puncte și linii ce interconectează anumite subseturi ale acesteia (posibil vide). Punctele unui graf sunt numite și noduri în timp ce liniile ce le conectează se numesc muchii.

ce măsoară capacitatea rezumatului de a reproduce conținutul original al secvenței. Folosirea măsurii SRE se dovedește a fi eficientă în două situații, și anume: dacă numărul de imagini al rezumatului este cunoscut ”a priori” sau dacă se dorește ca imaginile din rezumat să păstreze evoluția temporală din secvență.

Dacă considerăm o funcție de interpolare a imaginilor,  $F_{int}(t, R)$ , ce permite reconstrucția unei imagini la momentul  $t$  din unitatea  $V$  a secvenței (de exemplu: un plan video, un segment, etc.) pornind de la imaginile rezumatului  $R$ , atunci mărimea SRE, notată  $\varepsilon(V, R)$ , este dată de relația următoare:

$$\varepsilon(V, R) = \sum_{i=1}^n D(I_i, F_{int}(i, R)) \quad (5.13)$$

unde  $n$  reprezintă numărul total de imagini din unitatea  $V$  a secvenței iar  $D()$  este un operator de distanță între imagini.

Imaginiile cheie ale rezumatului sunt selectate folosind drept criteriu ipoteza că un rezumat optimal trebuie să aibă un scor SRE minim, astfel:

$$\{r_1, r_2, \dots, r_k\} = argmin|_{r_i} \{\varepsilon(V, R), 1 \leq r_i \leq n\} \quad (5.14)$$

unde  $r_i$  reprezintă indicele imaginii cheie din rezumat iar  $k$  reprezintă numărul de imagini al rezumatului ce este fixat ”a priori”.

Dintre metodele existente ce folosesc măsura SRE putem exemplifica metodele propuse în [Liu 02b] sau [Li 04].

În [Liu 02b] rezumatul static este generat folosind metoda ”Heap-Based Greedy”<sup>13</sup>. Algoritmul de selecție a imaginilor cheie pornește cu toate cele  $n$  imagini ale secvenței ca imagini cheie. Aceasta constituie nivelul  $n$  de analiză. Mai departe, imaginile cheie sunt eliminate iterativ la trecerea de la un nivel de analiză la altul. Astfel, imaginile cheie din nivelul de analiză  $t - 1$  sunt determinate pe baza imaginilor cheie din nivelul superior  $t$  prin eliminarea imaginilor ce minimizează scorul SRE.

Metoda propusă în [Li 04] încearcă să determine un rezumat optimal pe baza metodei de programare dinamică. Aceasta folosește drept constrângere egalitatea:

$$b_i = r_i \quad (5.15)$$

unde  $r_i$  reprezintă indicele imaginii cheie din rezumat ce este reprezentativă pentru segmentul secvenței definit de indicii  $[b_i, b_{i+1} - 1]$ . Cu alte cuvinte, fiecare imagine cheie este determinată în funcție de imaginile ce o succed.

---

<sup>13</sup>algoritmul ”Heap-Based Greedy” este un algoritm ”Greedy” (vezi explicația anterioară de la pagina 159) ce se folosește de acumularea valorilor într-o coadă, reducând astfel complexitatea de calcul.

Rezumatul optimal este obținut cu o căutare prin înjumătățire ("bisection search"<sup>14</sup>).

### Simplificarea curbei de caracteristici

Similar metodelor de clasificare, metodele din această categorie reprezintă secvența sub forma unei curbe de caracteristici în spațiul multidimensional format de vectorii de atribute ai imaginii. Una dintre dimensiunile acestui spațiu va fi dată de axa temporală a secvenței. Prințipiu simplificării curbei de caracteristici constă în izolarea unui set de puncte considerate ca fiind reprezentative pentru forma curbei și care păstrează o bună aproximare a acesteia. Aceste puncte, nu este obligatoriu să fie distribuite uniform după axa temporală.

Ca exemplu de astfel de metode putem menționa abordările propuse în [Latecki 01] și [Calic 02a], unde selecția imaginilor cheie este realizată cu ajutorul algoritmul "Discrete Contour Evolution" de analiză a evoluției conurilor din imagine. Acesta este aplicat curbei de evoluție temporală a diferențelor dintre histogramele imaginilor din secvență. Prințipiu folosit constă în înlocuirea iterativă a fiecărei perechi de segmente consecutive ale curbei de caracteristici ce prezintă un scor de relevanță minim, cu segmentul ce reunește primul nod al primului segment cu ultimul nod al celui de-al doilea segment din pereche. Procesul se oprește în momentul în care numărul dorit de imagini al rezumatului este atins. În [Calic 02a] imaginile cheie astfel obținute sunt rafinate mai departe prin reținerea doar a imaginilor ce corespund punctelor de minim local ale curbei de caracteristici simplificate.

Rezumatul obținut cu această strategie are proprietatea de a respecta evoluția temporală a secvenței.

### Localizarea evenimentelor importante

Metodele din această categorie selectează imaginile cheie ca fiind imagini reprezentative ale anumitor pasaje ale secvenței ce prezintă un interes semantic (de exemplu: anumite evenimente, pasaje de acțiune, prezența personajelor, etc.).

Metoda propusă în [Dufaux 00] rezumă întreaga secvență cu o singură imagine cheie. Algoritmul propus implică două etape. În prima etapă sunt selectate din secvență doar planele video considerate ca importante. Acestea

<sup>14</sup>căutarea de tip "bisection search" este un algoritm simplu de căutare ce se bazează pe localizarea soluției prin înjumătățirea iterativă a intervalului de căutare. La fiecare iterație, intervalul curent de căutare devine jumătatea intervalului din iterarea anterioară în care este cel mai probabil să se afle soluția.

sunt planele ce îndeplinesc următoarele cerințe: durata lor este importantă, au un conținut de mișcare semnificativ, o activitate spațială importantă iar în acestea apar personaje umane. Mai departe, în a doua etapă, imaginea cheie este aleasă folosind criterii similare cu cele menționate anterior, cu excepția faptului că pentru a evita ca imaginea selectată să fie o imagine încețoșată (datorită efectului de "motion blur"<sup>15</sup>) sau să conțină "artefacte"<sup>16</sup>, aceasta nu va fi aleasă dintr-un pasaj de mișcare, ci mai degrabă dintr-un pasaj cu o activitate spațială importantă.

Un alt exemplu este metoda propusă în [Liu 03] unde extragerea imaginilor cheie este realizată pe baza conținutului de mișcare. Pentru aceasta, planele video sunt mai întâi descompuse în segmente ce conțin o mișcare continuă din punct de vedere al accelerării și respectiv decelerării mișcării. Imaginele cheie sunt alese, mai departe, ca fiind imaginile ce fac joncțiunea între un segment în care mișcarea este accelerată și un segment cu mișcare decelerată. Dacă planul video analizat nu prezintă un conținut de mișcare, atunci este aleasă doar o singură imagine cheie și anume imaginea de început a planului.

Metoda propusă în [Calic 04] selectează imaginile la nivel de plan video prin localizarea anumitor evenimente de interes din interiorul planului. Acestea sunt marcate de comportamentul particular al anumitor regiuni din imagine. Astfel, un eveniment este detectat dacă o regiune din imagine dispăr din scenă, dacă aceasta se combină cu o alta sau dacă aceasta devine plan principal al scenei.

Una dintre limitările unei astfel de abordări este dată de faptul că localizarea evenimentelor considerate ca importante pentru conținut este guvernată de reguli eurisitice. Acestea sunt de regulă definite pe baza observării proprietăților unui set redus de date, fiind astfel mai mult sau mai puțin adaptate unui anumit domeniu de aplicație.

## 5.2 Construcția rezumatelor dinamice

Un rezumat dinamic, sau "video skim", reprezintă o colecție de segmente sau pasaje ale secvenței. Acesta constituie el însuși o secvență de imagini, dar, de o durată mult inferioară secvenței inițiale.

Din punct de vedere formal, rezumatul dinamic  $R_{seq}(S)$  al secvenței  $S$  poate fi exprimat astfel:

$$R_{seq}(S) = seg_1 \cup seg_2 \cup \dots \cup seg_M \quad (5.16)$$

---

<sup>15</sup>vezi explicația de la pagina 81.

<sup>16</sup>vezi explicația de la pagina 101.

unde  $seg_i$  reprezintă segmentul video de indice  $i$  în rezumat, cu  $i = 1, \dots, M$ , iar  $M$  reprezintă numărul total de segmente prezente în rezumat.

Spre deosebire de mecanismele de rezumare în imagini, generarea unui rezumat dinamic este o operație mult mai laborioasă, având o complexitate de calcul ridicată datorată în principal necesității de analiză și înțelegere a conținutului secvenței la un nivel semantic superior. Dacă în cazul rezumării statice, unitatea de analiză de bază o constituie imaginea, în cazul rezumării dinamice, aceasta este segmentul video (sub-secvență de imagini).

După cum am menționat și în partea introductivă a acestui capitol, rezumatele dinamice pot fi generate pornind de la rezumatele statice. O metodă imediată constă în înlocuirea fiecărei imagini cheie din rezumatul static cu un interval de imagini din secvență, ca de exemplu intervalul centrat pe imaginea cheie. Rezumatul astfel obținut este o reprezentare compactă a conținutului secvenței, dar calitatea acestuia este dependentă de calitatea rezumatului static. Datorită limitărilor metodelor de rezumare în imagini (vezi secțiunea anterioară) este foarte probabil ca rezumatul dinamic obținut astfel să nu fie reprezentativ pentru conținutul de mișcare al secvenței. Din această cauză, în cele mai multe situații este de preferat ca rezumatul să fie extras direct din secvență cu metode specifice.

O altă strategie rapidă de generare a unui rezumat dinamic constă în utilizarea segmentării în plane a secvenței. Astfel, un posibil rezumat poate fi constituit prin reprezentarea fiecărui plan video cu un anumit pasaj al acestuia. Rezumatul astfel obținut este o foarte bună sinteză a conținutului global al secvenței. Pe de altă parte, acesta va conține atât pasajele importante cât și cele mai puțin relevante din secvență, ceea ce face ca durata rezumatului să fie în cele mai multe situații foarte ridicată (de regula, în condițiile asigurării unei continuități vizuale a rezumatului se obține în medie o compresie de 1/3 a duratei secvenței inițiale).

Un rezumat dinamic eficient din punct de vedere al duratei precum și al calității informației furnizate, necesită în mod ideal o analiză de nivel semantic a conținutului secvenței, similară modului de percepție uman. Datorită limitărilor metodelor de analiză semantică existente, ce nu au ajuns la un nivel de dezvoltare suficient pentru a facilita o astfel de analiză, metodele existente de rezumare dinamică tind să simplifice problema rezumării prin adoptarea de condiții particulare. Din această cauză, metodele existente sunt în mare parte specifice domeniului de aplicație. Astfel întâlnim metode de rezumare distincte pentru: secvențele sportive [Coldefy 04], secvențele documentare [Yu 03], materialele video personale [Zhao 03], secvențele de animație [Ionescu 06b], etc. Dificultatea analizei semantice este simplificată în acest caz prin folosirea expertizei domeniului respectiv.

Din punct de vedere al popularității metodelor de rezumare, literatura

de specialitate este mult mai bogată în tehnici de rezumare statică, decât în tehnici de rezumare dinamică. Acest lucru se datorează în principal complexității de calcul ridicate a acestora din urmă [Li 01]. Pentru un studiu complet al literaturii de specialitate, cititorul se poate raporta la lucrările [Lecce 99], [Truong 07] sau [Benoit 07].

### 5.2.1 Informația conservată de rezumat

O primă constrângere la generarea unui rezumat dinamic o constituie felul în care informația din secvență va fi reprezentată în rezumat. Aceasta depinde în general de aplicația vizată, căreia îi este destinat rezumatul, și va determina modalitatea de generare a acestuia. În funcție de modul de reprezentare al conținutului secvenței, metodele existente de rezumare dinamică se împart în trei categorii [Truong 07]:

- metode ce generează rezumate ce propun o acoperire *integrală* a conținutului secvenței,
- metode ce generează rezumate ce reproduc doar *anumite pasaje* ale unor evenimente importante pentru conținutul secvenței,
- metode ce produc rezumate *personalizate*, generate în funcție de cerințele utilizatorului.

#### Acoperirea integrală a conținutului secvenței

Metodele ce rezumă secvența astfel încât conținutul acesteia să fie acoperit în totalitate de rezumat, au ca scop de a informa utilizatorul cu privire la conținutul global al secvenței [Sundaram 02] [Gong 03]. În acest caz, înțelegerea conținutului inițial al secvenței nu este aproape deloc alterată de rezumat.

Acest tip de rezumat este util în cazul în care utilizatorul dorește informații mai detaliate despre conținutul dinamic al secvenței, dar totuși nu dispune de timpul necesar pentru a vizualiza secvența în totalitate. Putem spune că un astfel de rezumat este un compromis realizat asupra duratei rezumatului în favoarea completitudinii informației furnizate de acesta.

#### Redarea evenimentelor importante ale secvenței

Metodele de construcție a rezumatelor dinamice cel mai frecvent întâlnite sunt cele care reproduc anumite evenimente importante ale secvenței, fiind numite și ”video highlights”. Acestea sunt de regulă specifice fiecărui domeniu de aplicație.

În funcție de tipul evenimentelor redate de rezumat, putem menționa următoarele abordări [Truong 07]:

- [Xiong 03]: rezumatul este constituit pe baza detectiei evenimentelor ce antrenează reacții particulare din partea publicului, ca de exemplu aplauze sau aclamații,
- [Coldefy 04]: rezumatul este construit pe baza pasajelor secvenței ce provoacă o reacție de entuziasm a naratorului,
- [Pan 01]: rezumatul este construit folosind pasajele din secvență ce sunt evidențiate de producător prin tehnici de montaj specifice, ca de exemplu o frecvență ridicată a tranzițiilor de tip "cut", prezența textului sau reluarea anumitor scene ale secvenței,
- [Radhakrishnan 04]: rezumatul este constituit pe baza detectiei evenimentelor ce corespund unor anumite modele predefinite,
- [Yu 03]: rezumatul este construit pe baza pasajelor secvenței ce sunt preferate de utilizator (de exemplu, vizualizate de mai multe ori).

O categorie aparte de "video highlights" o constituie rezumatele de tip "movie trailer". Acestea sunt rezumate publicitare, destinate de regulă promovării unui nou film, ce redau un colaj al pasajelor celor mai captivante și mai reprezentative ale secvenței. Localizarea automată a acestora în secvență este o operație dificilă, din această cauză se folosește expertiza domeniului de aplicație. De exemplu, este evident că pentru a construi rezumatul de tip "trailer" al unui meci de fotbal, evenimentele căutate vor fi momentele de gol precum și aclamațiile publicului.

Ca exemplu de metodă de generare a unui astfel de rezumat, putem menționa metoda propusă în [Ionescu 06b] ce se folosește de analiza distribuției temporale a tranzițiilor video pentru a localiza pasajele de acțiune în filmele artistice de animație. În primă fază, activitatea temporală a secvenței este reprezentată grafic sub forma unui semnal continuu temporal, de amplitudine arbitrară 1, ce trece prin valoarea 0 la apariția unei tranziții video (vezi Figura 5.3). Aceasta constituie nivelul de analiză macro, la nivel de segment.

Un pasaj de acțiune este definit mai departe ca fiind un segment al secvenței pentru care numărul de schimbări de plan survenite în ferestre temporale de durată  $T = 5s$ , este superior vitezei medii de schimbare de plan,  $\bar{v}_T$ . Valoarea lui  $\bar{v}_T$  este estimată pentru întreaga secvență folosind

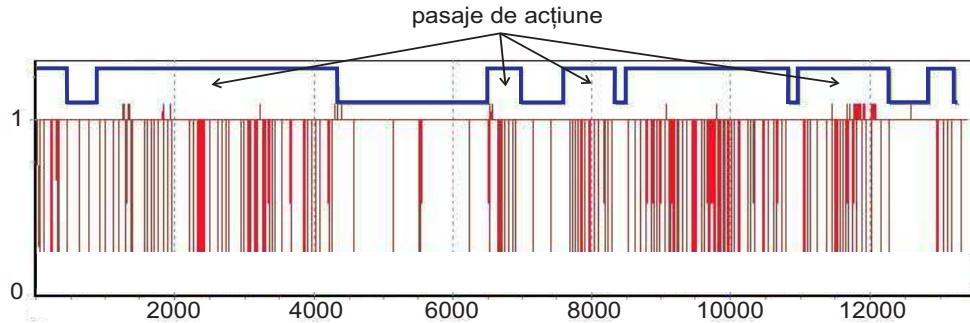


Figura 5.3: Exemplu de localizare a pasajelor de acțiune pe baza analizei frecvenței schimbărilor de plan (axa oX este axa temporală, liniile roșii verticale reprezintă tranzițiile video iar linia albastră reprezintă un semnal binar ce ia valoarea 1 pentru un pasaj de acțiune și valoarea 0 în rest, secvență "François le Vaillant" [Folimage 06]).

relația următoare:

$$\bar{v}_T = E\{\zeta_T\} = \sum_{t=1}^{T \cdot 25} t \cdot f_{\zeta_T}(t) \quad (5.17)$$

unde  $f_{\zeta_T}()$  reprezintă densitatea de probabilitate a variabilei aleatoare discrete  $\zeta_T()$  ce reprezintă numărul de schimbări de plan survenite în fereastra temporală de durată  $T$ , și este dată de ecuația:

$$f_{\zeta_T}(t) = \frac{1}{N} \sum_{i=1}^N \delta(\zeta_T(i) - t) \quad (5.18)$$

unde  $N$  reprezintă numărul total de ferestre de analiză de durată  $T$  secunde,  $\delta(t) = 1$  pentru  $t = 0$  și 0 altfel, iar  $i$  reprezintă indicele temporal al ferestrei de analiză curente pentru care se calculează valoarea lui  $\zeta_T()$  (vezi Figura 5.3).

Rezumatul de tip "movie trailer" este constituit mai departe pe baza segmentelor de acțiune folosind o analiză intra-plan. Fiecare plan video dintr-un segment de acțiune este rezumat cu o sub-secvență a cărei durată este proporțională cu nivelul de activitate vizuală al imaginilor acestuia. Nivelul de activitate este evaluat pe baza calculului unei histograme de distanțe cumulate [Ionescu 06b].

### Personalizarea conținutului rezumatului

O altă categorie de rezumate dinamice sunt cele care folosesc personalizarea conținutului. Acestea sunt generate în funcție de specificațiile utilizatorului

cu privire la conținutul secvenței ce urmează să fie prezentat în rezumat. Pentru aceasta, utilizatorul își alege, fie un model predefinit de conținut, fie își exprimă criterile de selecție sub forma unei cereri de tip ”query”. Ca exemple de astfel de metode de rezumare putem menționa metoda propusă în [Lu 03], destinată rezumării secvențelor de știri. În aceasta, utilizatorul poate opta pentru un rezumat care să prezinte pasajele secvenței ce conțin fețe umane, pasajele de dialog, pasajele ce conțin o mișcare a camerei video de tip ”zoom-in/zoom-out” sau pasajele ce conțin text încrustat în imagine. Un alt exemplu este metoda propusă în [Li 03] unde evenimentele vizate pentru a face parte din rezumat sunt scenele de dialog dintre două sau mai multe personaje, precum și scenele hibride.

În cazul rezumatelor personalizate, procesul de selectare a segmentelor secvenței ce vor face parte din rezumat este simplificat. Rezumatul obținut nu va conține decât informația care corespunde cerințelor utilizatorului. Acest tip de rezumat poate fi văzut ca un rezumat semi-automat, deoarece în procesul de generare este necesară intervenția utilizatorului. Rezumatul obținut cu o astfel de strategie are un caracter generic redus, fiind dificil de aplicat într-un context necunoscut.

### 5.2.2 Procesul de generare a rezumatului dinamic

În general, procesul de constituire a unui rezumat dinamic implică o serie de etape de prelucrare. Acestea pot fi sintetizate în următoarele [Truong 07]:

- *descompunerea secvenței în segmente (plane video, scene, etc.),*
- *selecția segmentelor ce vor fi folosite la constituirea rezumatului,*
- *reducerea dimensiunii segmentelor prin reducerea redundanței informaționale sau pe baza anumitor constrângeri de conținut,*
- *integrarea multimodală a segmentelor,*
- *construcția propriu-zisă a rezumatului.*

În funcție de aplicație, este posibil ca anumite metode de rezumare să folosească doar unele dintre cele cinci etape enumerate mai sus, sau ca anumite etape să fie efectuate simultan.

#### Descompunerea în segmente

Prima etapă, indispensabilă procesului de constituire a rezumatului dinamic, constă în *descompunerea secvenței în segmente*. Un segment poate fi un plan

video, o scenă sau un pasaj al secvenței ce conține un eveniment considerat drept important pentru conținut. Noțiunea de segment nu se limitează doar la informația vizuală, descompunerea în segmente fiind efectuată, după caz, și pentru celelalte modalități ale unui document video, precum sunetul (descompunere în segmente audio) și textul (descompunere în cuvinte sau propoziții).

### Selectia segmentelor

Etapa următoare descompunerii în segmente o reprezintă *selecția segmentelor* ce vor fi folosite pentru constituirea rezumatului. Felul în care aceasta este realizată va influența atât *coerența* rezumatului cât și *gradul de acoperire* al conținutului secvenței, și va determina *contextul aplicației* căreia îi este destinat rezumatul.

Ca exemple de tehnici de selecție a segmentelor, putem menționa metodele propuse în [Gong 03] și [Ngo 03]. În [Gong 03] proprietățile temporale și spațiale ale secvenței sunt caracterizate pe baza descompunerii în valori singulare (SVD<sup>17</sup>) a unei matrice de caracteristici de culoare a imaginilor. Pe baza acesteia, planele video sunt selectate în funcție de gradul de schimbare vizuală, de uniformitatea distribuției de culoare și în funcție de similaritate.

Metoda propusă în [Ngo 03] folosește un algoritm generalizat de detectie a tranzițiilor de tip "cut" pentru a determina planele video, sau mai general, clasele de imagini din secvență. Acestea sunt analizate mai departe folosind un model al perceptiei umane de tip "Motion Attention Model". Valorile obținute pentru gradul de atenție sunt structurate sub forma unui graf și sunt prelucrate folosind metode specifice analizei lanțurilor Markov<sup>18</sup>. Graful obținut este folosit pentru a regrupa clasele de imagini similare în scene, iar gradul de atenție este folosit drept criteriu de selecție a segmentelor ce vor face parte din rezumat.

### Reducerea dimensiunii segmentelor

În general, segmentele obținute în urma etapei anterioare sunt redundante din punct de vedere al conținutului, iar durata acestora tinde să fie prea ridicată pentru rezumat.

Etapa de *reducere a dimensiunii segmentelor* va asigura optimizarea acestora prin păstrarea doar a conținutului esențial, și va conduce astfel la

---

<sup>17</sup>vezi explicația de la pagina 152.

<sup>18</sup>un lanț Markov, numit astfel după numele matematicianului rus Andrey Markov, reprezintă o colecție de variabile aleatoare,  $\{X_t\}$ , unde  $t = 0, 1, \dots$ , ce au proprietatea că pentru o anumită stare prezentă a acestora, valorile viitoare sunt condițional independente de cele anterioare.

obținerea unui rezumat mai concis. Pentru aceasta, s-au propus diferite soluții, ca de exemplu metoda propusă în [Cooper 02] ce folosește o matrice de autosimilaritate. Din fiecare segment este reținut doar acel pasaj al căruia conținut este cel mai reprezentativ pentru conținutul segmentului întreg (din punct de vedere al similarității). Alte abordări folosesc metode de compresie a segmentelor prin eliminarea anumitor imagini sau pasaje redundante. Eliminarea acestora se face ținând cont și de informația audio, dacă aceasta este prezentă. Rezumatul obținut trebuie să respecte limitele de coerență vizuală cât și sonoră [Li 03].

Totuși, datorită procesului de eliminare a informației, această etapă poate introduce momente de discontinuitate vizuală în rezumat.

### **Integrarea multimodală**

În general, segmentele extrase din secvență sunt *unimodale*. Acestea sunt, fie segmente vizuale, fie segmente audio sau segmente textuale. Etapa de *integrare multimodală* are ca scop tocmai fuzionarea și sincronizarea tuturor acestor informații astfel încât rezumatul final să respecte evoluția temporală a evenimentelor secvenței. În acest scop, *alinierea segmentelor* joacă un rol important. Astfel, frontierele dintre segmente sunt ajustate astfel încât să fie conservat fluxul vizual precum și coerența secvenței. Mai mult, continuitatea sunetului este asigurată prin evitarea intreruperilor ce pot surveni, de exemplu, în mijlocul propozițiilor.

În funcție de modul în care este realizată integrarea audio-vizuală, putem distinge două categorii de rezumate dinamice, și anume: *rezumatele sincrone* și *rezumatele asincrone* [Truong 07].

În rezumatele sincrone, informația vizuală este sincronizată cu cea audio folosind ca sistem de referință axa temporală a secvenței. Acest tip de rezumat este mai adaptat filmelor, deoarece în momentul vizualizării, sunetul este în corespondență directă cu imaginile. În cazul în care segmentele vizuale și audio au fost generate separat din cele două modalități ale secvenței, integrarea acestora poate fi efectuată pe baza unui set de reguli de fuziune guverنate de operatorii clasici de conjuncție ("și") și disjuncție ("sau") [Erol 03]. De exemplu, putem selecționa pasajele importante ale secvenței ce sunt conținute într-un segment audio "sau" într-unul vizual.

Pe de altă parte, rezumatele asincrone sunt mai adaptate în cazul rezumării secvențelor documentare și de știri, deoarece acestea tend să maximizeze acoperirea conținutului secvenței. Acest gen de rezumat este de regulă generat în primă fază folosind doar una dintre modalitățile secvenței, urmând ca celelalte să fie adăugate ulterior. Un exemplu este metoda propusă în [Smith 98] unde rezumatul este extras folosind informația audio. Elementele

vizuale sunt adăugate ulterior pe baza unui set de reguli euristice determinate în funcție de mișcarea obiectelor sau a camerei video, de prezența fețelor umane sau a textului încrustat în imagine.

### Construcția propriu-zisă a rezumatului

În ceea ce privește *construcția rezumatului final*, strategia cel mai frecvent folosită constă în agregarea tuturor segmentelor obținute în etapele anterioare folosind ca referință axa temporală. Un caz particular îl constituie rezumatele de tip ”movie trailer”, ce prezintă succint doar pasajele dinamice și de acțiune ale secvenței, în scopul de a capta atenția utilizatorului. În acest caz, ordinea temporală nu este întotdeauna respectată.

## 5.3 Metodele de evaluare a rezumatelor

Pe lângă dificultățile de analiză implicate de procesul de generare a rezumatelor de conținut, o problemă la fel de importantă o constituie *evaluarea calității acestora*.

Întrebările care se pun sunt: ”*Este rezumatul obținut pertinent pentru conținutul secvenței? Este acesta coerent pentru persoana care îl vizualizează?*”. În realitate, pentru o aceeași secvență se pot genera o multitudine de rezumate care să fie corecte din acest punct de vedere. Faptul că rezumatul a fost generat în conformitate cu anumite criterii obiective, inspirate de percepția umană, nu asigură întotdeauna o coerență vizuală a acestuia.

În momentul de față, în literatura de specialitate, *nu există o metodologie validată* și general valabilă de evaluare a calității rezumatelor automate de conținut, cum sunt de exemplu erorile de tip ”precision” și ”recall” pentru evaluarea detectiei de caracteristici (vezi Secțiunea 2.2.4). Tendința este de a testa metodele propuse folosind strategii de evaluare mai mult sau mai puțin particulare, lucru ce nu încurajează un studiu comparativ și competitiv al rezultatelor obținute cu diverse metode.

Acest lucru se datorează în principal *dificultății* și a *subiectivității* construirii de rezumate de referință, generate manual pe baza expertizei umane, care să constituie ceea ce numim ”realitate de teren” sau ”groundtruth”<sup>19</sup> pentru procesul de evaluare. Dispunând de aceste date de referință, rezu-

<sup>19</sup>termenul de ”groundtruth” își are originea în domeniul cartografiei și implică procesul de colectare de informații despre un anumit fenomen, prin observarea practică pe teren a acestuia. Datele obținute constituie ”realitatea de teren” folosită pentru calibrarea, validarea și interpretarea observațiilor sau a măsurătorilor de la distanță a fenomenului în cauză sau a altor fenomene similare.

matele obținute urmează să fie comparate cu acestea pentru a judeca calitatea lor. Chiar și în cazul în care constituirea unui "groundtruth" ar fi posibilă, apare ca problemă subiectivitatea comparării a două rezumate, deoarece în multe situații chiar și pentru o persoană este dificil să decidă dacă un anumit rezumat este mai bun decât un altul.

Totuși în ciuda acestor dificultăți de evaluare, în literatura de specialitate putem evidenția o serie de strategii folosite pentru evaluarea atât subiectivă cât și obiectivă a rezumatelor automate. Acestea sunt [Truong 07]:

- *analiza descriptivă* a rezultatelor obținute,
- folosirea unei anumite *măsuri matematice* pentru evaluarea obiectivă,
- *teste de evaluare* efectuate de utilizatori.

### 5.3.1 Analiza descriptivă a rezultatului

Metoda cea mai simplă de evaluare, ce nu necesită compararea cu alte metode, o constituie analiza descriptivă a rezultatului. Rezumatele obținute cu o anumită metodă sunt analizate manual pentru anumite secvențe reprezentative, fiind justificate astfel avantajele și performanțele metodei propuse.

O astfel de evaluare este propusă în [Yu 04] unde rezultatele obținute sunt exemplificate și analizate manual din punct de vedere al pertinenței acestora. Alte abordări constau în explicarea și ilustrarea, într-o manieră descriptivă, a avantajelor și a superiorității metodei propuse în raport cu alte strategii de bază, cum ar fi alegerea aleatoare a imaginilor cheie sau eșantionarea uniformă a secvenței [Vermaak 02] (vezi Figura 5.4, se observă eficiența rezumatului propus în raport cu rezumatul de comparație prin prezența unui număr redus de imagini redundante).

Din punct de vedere global, o astfel de strategie de evaluare este foarte subiectivă, deoarece nu există dovezi că metoda propusă este eficientă și aplicată altor secvențe decât cele testate și exemplificate. Evaluarea experimentală propusă în acest caz este insuficientă pentru o evaluare globală. Mai mult, acest tip de abordare este dificil de utilizat la evaluarea rezumatelor dinamice, deoarece volumul de date în acest caz este mult prea ridicat pentru a permite o analiză descriptivă.

### 5.3.2 Utilizarea unei măsuri matematice

Pentru evaluarea rezumatelor statice, măsura matematică folosită este de regulă o funcție de fidelitate ce este estimată atât pentru imaginile cheie din rezumat cât și pentru imaginile din secvență.

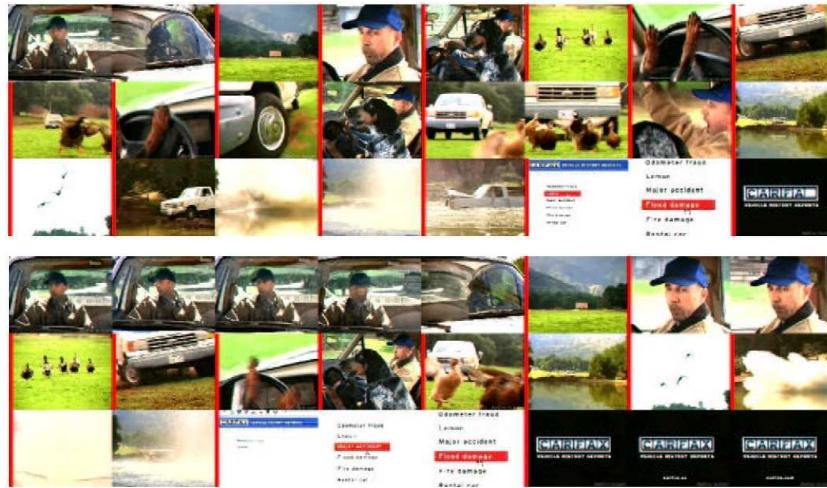


Figura 5.4: Exemplu de analiză descriptivă și comparare a rezumatului obținut cu criteriul "Bayes Information Criterion" propus în [Vermaak 02] (prima imagine), cu rezumatul obținut prin eșantionarea uniformă a secvenței (a două imagine, liniile roșii marchează schimbările de plan).

Această strategie permite compararea rezultatelor obținute cu mai multe metode. Pe de altă parte, similar evaluării prin descrierea rezultatului, această măsură nu oferă certitudinea că este o măsură pertinentă din punct de vedere al modului de evaluare uman. Ca exemplu, putem menționa metoda propusă în [Liu 04] ce folosește ca măsură de evaluare distanța Hausdorff<sup>20</sup> aplicată între erorile de reconstrucție SRE<sup>21</sup> obținute cu metodele de rezumare evaluate. Un alt exemplu este metoda propusă în [Liu 02c] ce folosește pentru evaluare conceptul de "imagină cheie bine distribuită". Aceasta este definită de autori ca fiind o imagine cheie ce nu este redundantă și care nu aparține unei tranzitii video. Evaluarea calității rezumatelor este realizată prin reprezentarea grafică a numărului de imagini "bene distribuite" în funcție de numărul total al imaginilor cheie din rezumat.

În cazul particular al rezumatelor dinamice de tip "video highlight", ce sunt determinate pe baza anumitor pasaje ale secvenței ce conțin evenimente de interes, evaluarea calității este relativ mai ușor de realizat deoarece se poate construi echivalentul unei "realități de teren". Aceasta este constituită prin

<sup>20</sup>numită după matematicianul Felix Hausdorff, distanța ce îi poartă numele este definită între două mulțimi  $A$  și  $B$  ca fiind distanța maximă a mulțimi  $A$  față de cel mai apropiat punct din mulțimea  $B$ , astfel:  $d_H(A, B) = \max_{a \in A} \{ \min_{b \in B} \{ d(a, b) \} \}$ , unde  $d()$  reprezintă o anumită metrică.

<sup>21</sup>vezi explicația de la pagina 159.

localizarea manuală în secvență a evenimentelor de interes, iar evaluarea constă în verificarea prezenței acestora în rezumat [Xiong 03] [Ariki 03]. Ca măsură cantitativă a evaluării se pot folosi în acest caz erorile de tip "precision" și "recall" (vezi Secțiunea 2.2.4) folosite la scară largă în toate metodele de detecție.

Alte abordări ale procesului de evaluare încearcă să îmbunătățească obiectivitatea acestuia prin folosirea unei "realități de teren". Aceasta este constituită de regulă de o serie de rezumate de referință, ce au fost create manual de specialiști pentru un set redus de secvențe ce vor fi folosite la evaluare. De exemplu, în [He 99] evaluarea metodei propuse de rezumare a secvențelor de conferințe este realizată prin compararea rezultatelor cu o serie de rezumate create manual de autori, sau în [Miura 03] unde rezumatele propuse ale programelor televizate dedicate gastronomiei sunt evaluate folosind ca referință comentarile furnizate de producător. Totuși, cu toate că putem vorbi de existența unei "realități de teren", aceasta nu poate fi unică determinată, fapt ce poate conduce la o evaluare total eronată. Mai mult, evaluarea devine subiectivă în momentul comparării rezumatului obținut, cu rezumatul sau rezumatele de referință. După cum am menționat și în partea introductivă, compararea a două rezumate este dificilă chiar și pentru operatorul uman.

Din punct de vedere global, metodele de evaluare ce folosesc măsuri matematice sunt mai obiective decât metodele ce descriu rezultatele obținute. Acestea încearcă constituirea unei "realități de teren" prin expertiza secvențelor folosite, lucru ce permite o evaluare competitivă raportată la alte metode de rezumare. Evaluarea calității este realizată prin compararea pe baza unei anumite măsuri de distanță a rezumatului obținut, cu rezumatele de referință din "realitatea de teren". De notat este faptul că în această abordare, analiza perceptiei umane nu intervine în procesul propriu-zis de evaluare.

### 5.3.3 Testele de evaluare

În cazul *testelor de evaluare*, sarcina de a decide asupra calității rezumatului îi revine operatorului uman.

Un test de evaluare implică de regulă un grup de persoane, specialiste sau non specialiste în domeniu, ce sunt desemnate pentru a vizualiza rezumatele obținute cu o anumită metodă. Opinia personală a acestora asupra conținutului și a calității rezumatelor este de regulă analizată pe baza răspunsurilor la un anumit chestionar. Chestionarul în cauză este constituit în funcție de cerințele și de constrângerile inițiale impuse tipului de rezumat evaluat. În ciuda subiectivității aparente a acestei strategii de evaluare, teste de evaluare se dovedesc a fi cele mai apropiate de realitate, deoarece

rece evaluarea este realizată de însuși ”consumatorul produsului”, și anume utilizatorul.

Din punct de vedere practic, organizarea unei astfel de campanii de evaluare este dificilă în primul rând ca logistică (pregătirea filmelor, pregătirea protocolului de evaluare, strângerea unui număr suficient de evaluatori, etc.), dar și din punct de vedere al timpului necesar vizionării de către evaluatori a rezumatelor propuse, cât și a secvențelor originale. Vizualizarea secvențelor originale este foarte importantă pentru procesul de evaluare deoarece doar pe baza acestora evaluatorul poate să-și facă o idee asupra conținutului real al secvenței.

Ca exemple de astfel de strategii de evaluare, putem menționa în cazul rezumatelor statice metoda propusă în [Dufaux 00] unde fiecare imagine cheie a rezumatului este clasată de utilizator folosind trei niveluri de apreciere, și anume: ”bună”, ”corectă” sau ”slabă”. Pe baza acestor aprecieri este determinat un scor global pentru întregul rezumat. O abordare similară, dar care implică un test mai complex, este propusă în [Liu 03]. Imaginele cheie sunt analizate în acest caz în contextul planelor video din care provin. Pe baza acestora, utilizatorii atribuie un scor de satisfacție cuantificat pe trei niveluri de apreciere: ”bine”, ”acceptabil” sau ”nesatisfăcător”.

În cazul evaluării rezumatelor dinamice, strategia cel mai frecvent folosită constă în a cere utilizatorului direct opinia acestuia cu privire la calitatea globală a rezumatului propus. În [Ionescu 06b], o astfel de abordare este folosită pentru evaluarea rezumatelor de tip ”movie trailer” în cazul filmelor artistice de animație. Evaluarea este efectuată în acest caz pe baza scorurilor atribuite răspunsurilor la întrebările unui chestionar inspirat de teoria sondajelor. O primă întrebare este relativă la calitatea rezumatului, și anume: ”Considerați că rezumatul ”movie trailer” propus conține părțile cele mai atractive ale filmului?”. Răspunsului îi este atribuit un scor numeric cu valori de la 1 la 10, având semnificația: 1, 2= ”deloc”, 3, 4= ”foarte puține”, 5, 6= ”unele dintre acestea”, 7, 8= ”aproape toate” și respectiv 9, 10= ”toate”. O a doua întrebare este relativă la durata rezumatului, și anume: ”Cum vi se pare durata rezumatului propus?”. Scorul atribuit răspunsurilor la această întrebare variază de la 0 la 4 și are semnificația următoare: 0= ”foarte scurtă”, 1= ”scurtă”, 2= ”corectă”, 3= ”lungă” și respectiv 4= ”foarte lungă”. Evaluarea calității rezumatului se face pe baza scorului mediu și a varianței acestuia obținute pentru ansamblul evaluatorilor. Alte abordări mai elaborate încearcă să aprecieze în ce măsură rezumatele propuse pot fi utile pentru o serie de aplicații practice, precum navigarea și căutarea într-o bază de filme [Ngo 03]. De asemenea, o altă strategie constă în încercarea de a evalua calitatea unui rezumat pe baza capacitatei utilizatorului de a identifica conținutul original al secvenței, pornind de la rezumatul propus [Erol 03].

Pe lângă avantajele prezentate din punct de vedere al relevanței procesului de evaluare, testele de evaluare prezintă și o serie de limitări [Ionescu 07a]. Acestea pot fi sintetizate cu următoarele:

- *timpul necesar* evaluării este ridicat, mai ales în cazul rezumatelor dinamice,
- *dificultatea realizării unui test comparativ* pentru mai multe metode de rezumare, datorată subiectivității evaluatorului pus în situația de a decide între două rezumate diferite,
- *subiectivitatea* modului de percepție a fiecărui evaluator, ce este strict dependentă de formarea profesională a acestuia (de exemplu, un artist va avea o percepție diferită față de un inginer),
- *dificultatea de vizualizare* a rezumatelor în imagini: modul de prezentare al acestora influențează calitatea evaluării. De exemplu, prezentarea imaginilor cheie sub forma unui "slideshow" are ca efect crearea artificială a senzației de secvență sacadată, pe când prezentarea acestora sub forma unei planșe are ca rezultat tendința de a neglija evoluția temporală a secvenței.

## 5.4 Concluzii

În acest capitol am discutat problematica rezumării automate a conținutului secvențelor de imagini în contextul general al indexării după conținut. În funcție de tipul rezumatului furnizat, tehniciile de rezumare de conținut existente se împart în două mari categorii, și anume: metode de *rezumare în imagini* și metode de *rezumare în mișcare*.

Metodele de rezumare în imagini generează ceea ce numim rezumate statice. Un rezumat static reprezintă o colecție de imagini statice, numite și imagini cheie, ce sunt considerate ca fiind reprezentative pentru conținutul secvenței. Tehnicile existente variază de la metode simple ce aleg imaginile cheie, uniform, la nivel de plan video, până la metode adaptive ce selectează imaginile cheie în funcție de nivelul de activitate al secvenței. Cu toate că acest tip de reprezentare este foarte aproximativă, se dovedește totuși a fi foarte eficientă în multe situații. Rezumatul în imagini are în general o complexitate redusă de calcul și furnizează utilizatorului o imagine de ansamblu asupra conținutului vizual al secvenței. Rezumatul static este util în cazul navigării în baza de secvențe, când utilizatorul poate fi informat instantaneu, prin simpla vizualizare a cătorva imagini cheie, asupra conținutului secvenței.

Mai mult, rezumatul în imagini permite reducerea redundanței vizuale, necesară în multe metode de prelucrare specifice analizei video.

Pe de altă parte, metodele de rezumare în mișcare generează rezumate dinamice. Acestea sunt la bază o colecție de pasaje ale secvenței, ce sunt asamblate și sincronizate pentru a forma ele însăși o secvență. Alegerea pasajelor este realizată în acest caz în funcție de tipul informației ce va fi redată de rezumat. De exemplu, un rezumat de tip "movie trailer" va prezenta succint doar momentele de acțiune ale secvenței, în timp ce un rezumat dinamic de tip "movie highlight" va prezenta doar anumite evenimente de interes din secvență. Spre deosebire de rezumatele statice, rezumatele dinamice au mai mult sens din punct de vedere al vizualizării, deoarece redau informația de mișcare a secvenței. Mai mult, în funcție de caz, acestea pot prezenta și informația audio. Din această cauză, un dezavantaj al rezumatelor dinamice este dat de complexitatea de calcul superioară metodelor de extragere a rezumatelor statice.

De notat este faptul că cele două categorii de rezumate *nu sunt concurențiale*, sau cu alte cuvinte, metodele dezvoltate nu au ca scop dovedirea superiorității uneia dintre categorii. Ambele tipuri de rezumate sunt necesare indexării după conținut a secvențelor de imagini. Rezumatele în imagini sunt ușor de generat și vizualizat în detrimentul conținutului de mișcare, în timp ce rezumatele dinamice sunt mai greu de generat și necesită un timp mai ridicat pentru vizualizare în avantajul furnizării informației dinamice și audio a secvenței.

O problemă delicată care apare este *evaluarea calității* unui rezumat. Aceasta este în cele mai multe cazuri subiectivă și depinde în principal de modul de percepție al fiecărui dintre noi, dar și de tipul conținutului redat de rezumat. De exemplu, evaluarea unui rezumat ce prezintă conținutul global al secvenței nu poate fi realizată în același mod ca pentru un rezumat ce redă doar anumite evenimente importante ale acesteia, sau, evaluarea unui rezumat static nu poate utiliza aceleași criterii ca evaluarea unui rezumat dinamic, deoarece prezintă informații diferite. În acest moment, nu se poate spune că există o metodologie standardizată de evaluare a rezumatelor, existând mai multe metodologii mai mult sau mai puțin relevante de evaluare. Dintre acestea, cele mai pertinente se dovedesc a fi teste de evaluare, deoarece acestea implică în procesul de evaluare chiar utilizatorul căruia îi sunt destinate.

# CAPITOLUL 6

---

## Formalizarea fuzzy

---

**Rezumat:** Descrierea matematică, cu valori numerice, a proprietăților de interes a datelor multimedia analizate, nu reprezintă o modalitate de descriere accesibilă publicului larg, ci mai degrabă specialiștilor din domeniu. Mai mult, aceste descrieri nu au nici un sens din punct de vedere al percepției umane. Tendința actuală a metodelor de analiză existente este tocmai de a dezvolta metode capabile să înțeleagă conținutul datelor multimedia într-un mod similar modului în care percepem lumea înconjurătoare. Una dintre soluțiile larg acceptate o reprezintă formalizarea fuzzy pe baza conceptului de incertitudine. În acest capitol vom prezenta modul în care informațiile de nivel scăzut, numerice, pot fi convertite în concepte semantice folosind inferența fuzzy, precum și impactul acesteia asupra metodelor existente de analiză și prelucrare a secvențelor de imagini.

După cum am menționat în primul capitol al acestei lucrări, direcția de studiu ce face obiectul numeroaselor cercetări actuale din domeniul sistemelor de indexare după conținut, o constituie *analiza semantică* a percepției datelor.

Această direcție, putem spune, aflată încă la începuturi, are ca obiectiv dezvoltarea de metode și algoritmi de modelare a sistemului de percepție uman, în încercarea de înțelegere automată a sensului datelor. Furnizarea unui nivel de descriere semantic, conferă sistemelor de prelucrare și analiză,

în primul rând, mult mai mult sens, acestea devenind mai ușor accesibile publicului larg. De exemplu, descrierea distribuției de culoare a unei imagini cu parametri de nivel scăzut, precum  $I = \{P_{(255,0,0)} = 50\%, P_{(0,0,255)} = 50\%\}$  (unde  $P_{(c)}$  reprezintă procentul de apariție al unei culori  $c$ ), va fi accesibilă doar unei persoane avizate în domeniu, în timp ce descrierea la un nivel semantic perceptual, precum "imaginea contrastează culoarea roșie și albăstră", va permite oricărui dintre noi să-și creeze o imagine mentală a conținutul de culoare al imaginii în cauză.

În sistemele de indexare după conținut a secvențelor de imagini, descrierile semantice își găsesc un spectru de aplicabilitate mai larg [Ionescu 07a], astfel:

- permit *simplificarea navigării în baza de date*: acestea vin cu informații suplimentare despre conținutul secvențelor, facilitând astfel înțelegerea rapidă a acestuia de către utilizator. De exemplu, pe lângă rezumatele automate de conținut, o secvență poate fi acompaniată și de informații textuale, cum ar fi: genul acesteia (acțiune, dramă, documentar), tipul conținutului (natură, oraș, studio), etc.
- permit *simplificarea căutării în baza*: fiind exprimate sub formă textuală, pot fi folosite ca indici de căutare în baza de date. Astfel, utilizatorul își poate formula cererea de căutare într-un limbaj natural, apropiat de limbajul uman, de exemplu "caută secvențele de gol din meciul echipelor X și Y" sau "caută filmele de ficțiune".
- constituie *un ajutor pentru specialiști*: descrierile semantice de conținut pot însăși descrierile de nivel scăzut pentru a furniza un pachet de informații complet cu privire la tehniciile folosite în secvență: structura temporală, conținut de mișcare, conținut de culoare, etc.

După cum se poate observa din cele enunțate anterior, modalitatea cea mai expresivă de exprimare a sensului semantic al datelor constă în *reprzentarea textuală a acestuia*. Dintre metodele de asociere de descrieri textuale de conținut datelor numerice, o pondere importantă o au metodele ce se folosesc de conceptul de incertitudine. Una dintre acestea o reprezintă *formalizarea fuzzy* a conceptelor semantice textuale.

## 6.1 Introducerea conceptului de incertitudine

Printre schimbările spectaculoase ale paradigmelor existente în diversele domenii științifice, una dintre cele mai importante o constituie fundamentarea și dezvoltarea *conceptului de incertitudine* a datelor.

Această schimbare majoră a modului de percepție științific s-a materializat prin tranzitia modului tradițional de gândire, care insista asupra faptului că noțiunea de incertitudine este o proprietate a datelor ce nu se dorește să apară și care trebuia evitată pe cât posibil, spre o viziune alternativă ce toleră tocmai această noțiune, fiind considerată de această dată inevitabilă și foarte utilă pentru analiză.

În general, în momentul proiectării unui anumit sistem de analiză, întrebarea care se pune este: *"Cum trebuie gerat sistemul și problemele asociate acestuia în cazul în care complexitatea proceselor ce trebuesc modelizate depășește cu mult posibilitățile noastre de prelucrare?"*. Cu alte cuvinte, volumul informațional disponibil este foarte ridicat pentru a putea fi controlat în totalitate iar înțelegerea proceselor este limitată. Soluția în acest caz constă tocmai în introducerea noțiunii de incertitudine pentru situațiile în care soluția ce trebuie adoptată nu este deloc evidentă, ci mai degrabă incertă.

În momentul construcției unui anumit model, se încearcă întotdeauna să se maximizeze utilitatea acestuia. Acest obiectiv este strâns legat de relațiile ce pot exista între cele trei categorii de caracteristici cheie ale unui model, și anume: *complexitatea, credibilitatea și incertitudinea* acestuia. Aceste relații nu sunt întotdeauna înțelese în totalitate. Știm doar că incertitudinea, fie că este predictivă, prescriptivă, etc., joacă rolul esențial pentru efortul de maximizare a utilității sistemului.

Totuși în cele mai multe situații, dar nu întotdeauna, incertitudinea nu reprezintă un punct forte dacă aceasta este considerată independent de alți parametri. Incertitudinea devine o informație prețioasă a sistemului dacă este analizată în raport cu alte caracteristici ale acestuia. În general, cu toate că aparent este paradoxal, cu cât se adaugă mai multă incertitudine în modelarea sistemului, cu atât complexitatea acestuia este redusă și în efect credibilitatea modelului crește [Klir 95].

În concluzie, conceptul de incertitudine este un instrument important pentru modelarea unui anumit sistem sau pentru soluționarea unei anumite probleme. Aceasta permite obținerea de caracteristici "avantajoase" pentru modelul vizat, caracteristici ce vor conduce ulterior la maximizarea utilității acestuia relativ la scopul pentru care a fost creat.

Conceptul de incertitudine a fost materializat pentru prima dată în lucrările publicate de Lotfi A. Zadeh [Zadeh 65] (anticipat de filozoful Max Black în 1937). Aceasta propunea o nouă teorie bazată pe reprezentarea datelor cu mulțimi fuzzy. Mulțimile fuzzy sunt mulțimi pentru care frontierele dintre date nu sunt exakte. Apartenența datelor, în acest caz, la o astfel de mulțime, nu mai este o problemă de confirmare sau negare, ci o problemă de *grad de apartenență*.

Dacă teoria probabilităților este fondată pe definirea a două valori logice de adevăr, și anume *Adevărat* (1) și *Fals* (0), în logica fuzzy, gradul de adevăr este formulat în felul următor:

*dacă A reprezintă o mulțime fuzzy iar x este un obiect de interes, atunci propoziția "x este inclus în A" nu este obligatoriu să fie Adevărată sau Falsă, lucru impus de logica booleană, ci aceasta poate fi adevărată într-un anumit grad.*

Acest grad de adevăr este exprimat de regulă ca o valoare cuprinsă în intervalul  $[0; 1]$ , unde limitele acestuia reprezintă negația totală (limita inferioară, valoare de adevăr 0) și respectiv, afirmația totală (limita superioară, valoare de adevăr 1).

Capacitatea mulțimilor fuzzy de a exprima tranziția graduală între apartenența totală și non apartenență, și vice-versa, își găsește o vastă utilitate în marea majoritate a domeniilor existente. Mulțimile fuzzy nu numai că propun o reprezentare discriminantă și plină de sens a conceptului de incertitudine, ci și o reprezentare pertinentă a conceptelor vagi ce sunt exprimate într-un *limbaj natural*.

Pentru a înțelege avantajul folosirii mulțimilor fuzzy la descrierea proprietăților anumitor procese, vom considera exemplul următor [Klir 95]: în loc să descriem prognoza meteo a zilei curente specificând procentajul exact de acoperire al cerului cu nori,  $P_{nori}$ , putem adopta o soluție mai eficientă spunând că ziua va fi, fie "însorită", fie "cu un cer acoperit".

Această descriere este o descriere vagă și puțin exactă, dar în cele mai multe cazuri este mult mai utilă decât prima modalitate de descriere. Sensul termenului de "însorit" nu este în totalitate arbitrar. O acoperire cu nori în procent de  $P_{nori} = 100\%$  indică faptul că ziua nu este însorită, dar același lucru este valabil și pentru  $P_{nori} = 80\%$ . Astfel, pentru a desemna o "zi însorită" vom considera o serie de valori intermediare pentru  $P_{nori}$ , de exemplu:  $P_{nori} \in [10\%; 20\%]$ . Problema care apare este cum alegem aceste frontiere? Dacă considerăm că un procent de acoperire de mai puțin de 25% corespunde unei "zi însorite", atunci o acoperire de 26% corespunde sau nu aceluiași caz? Este evident inacceptabil în această modalitate de descriere binară ca o singură valoare a lui  $P_{nori}$  să facă diferența între două concepte opuse: "zi însorită" și respectiv "zi cu cer acoperit".

Pentru a soluționa acest conflict, termenul de "zi însorită" necesită un anumit grad de incertitudine ce va fi obținut prin introducerea unei tranziții graduale între valorile lui  $P_{nori}$ , folosite pe de-o parte pentru a desemna noțiunea de "zi însorită" și respectiv pentru conceptul opus. Aceasta constituie exact principiul de bază al logicii fuzzy ce reprezintă o generalizare a logicii booleene.

## 6.2 Logica booleană și logica fuzzy

În teoria clasice a mulțimilor, ce este bazată pe *logica booleană*, apartenența unui obiect la o anumită mulțime este exprimată în felul următor: dacă  $U$  reprezintă universul de discuție iar  $A$  reprezintă o submulțime a acestuia,  $A \subset U$ , atunci pentru fiecare obiect  $x$  din  $U$ , apartenența acestuia la mulțimea  $A$  este definită de funcția specifică  $\alpha_A$ :

$$\alpha_A(x) = \begin{cases} 1 & \text{dacă } x \in A \\ 0 & \text{dacă } x \notin A \end{cases} \quad (6.1)$$

Astfel, în acest caz sunt posibile doar două situații, și anume: fie  $x$  aparține lui  $A$ , fie acesta nu aparține.

Pe de altă parte, în *logica fuzzy* conceptul de apartenență se bazează pe noțiunea de incertitudine. Dacă presupunem că  $\alpha_A$  reprezintă funcția de apartenență fuzzy a unui obiect din universul de discuție la mulțimea fuzzy  $A$ , atunci  $\alpha_A$  este dat de relația următoare:

$$\alpha_A(x) \in [0; 1] \quad (6.2)$$

pentru toate valorile lui  $x \in U$ . Funcția  $\alpha_A$  exprimă în acest caz *gradul de apartenență*, sau cu alte cuvinte, valoarea de adevăr, unde valoarea 1 reprezintă apartenență sigură iar valoarea 0 non-apartență sigură.

Pentru a ilustra diferența dintre cele două concepte, vom considera exemplul următor: dacă presupunem că un pacient bolnav de gripă prezintă o temperatură  $T$  a corpului ridicată, atunci conceptul de "temperatură ridicată" poate fi asociat valorilor măsurate ale lui  $T$  folosind, fie modelul clasic, fie modelul fuzzy (vezi Figura 6.1).

În primul caz al logicii clasice, decizia de apartenență este una netă (transplantă), astfel, fie pacientul are o "temperatură ridicată" dacă  $T > 39$  de grade (valoare de adevăr  $\alpha(T) = 1$ ), fie nu are o "temperatură ridicată" în caz contrar ( $\alpha(T) = 0$ ). În logica fuzzy, pacientul are o "temperatură ridicată" cu o valoare de adevăr 1 doar dacă  $T > 41$  de grade. Pentru  $T \in [37; 41]$  funcția de apartenență este o funcție liniară ce ia valori reale între 0 și 1. Astfel, în acest caz, valoarea de adevăr a apartenenței la conceptul de "temperatură ridicată" va fi una graduală. De exemplu, pentru  $T = 39$  de grade putem afirma că pacientul are o "temperatură ridicată" cu o valoare de adevăr de numai 0.5 (vezi Figura 6.1).

Decizia netă nu este eficientă în acest caz deoarece pacientul va fi diagnosticat ca având o temperatură ridicată doar în cazul în care valoarea temperaturii  $T$  este deja foarte ridicată (de exemplu,  $T > 39$ ). Mai mult, este cu totul artificial să spunem că pacientul are o temperatură ridicată

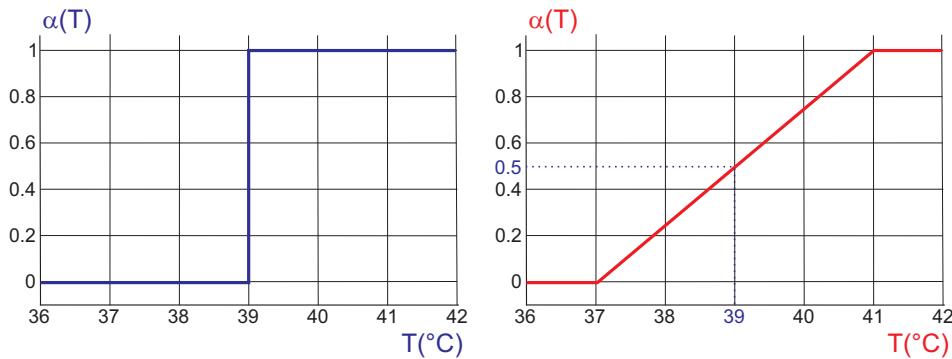


Figura 6.1: Reprezentarea funcției de apartenență  $\alpha(T)$  a conceptului de "temperatură ridicată" în logica booleană (prima imagine) și respectiv în logica fuzzy (a doua imagine).

dacă  $T > 39$  și în cazul în care  $T = 38.8$  acestă afirmație să nu fie valabilă. Avantajul reprezentării fuzzy este evident în cazul în care asociem conceptul de "febră ridicată" cu diagnosticul de gripă. Astfel, în logica clasică deducem că dacă pacientul nu are o "temperatură ridicată" atunci pacientul nu este bolnav de gripă, afirmație care nu este întotdeauna adevărată. Pe de altă parte, în logica fuzzy vom deduce că dacă pacientul prezintă o "temperatură ridicată" cu o valoare de adevăr  $x \in [0; 1]$ , atunci pacientul este bolnav de gripă cu o valoare de adevăr  $\tilde{x} \in [0; 1]$ , ce va fi calculată în funcție de  $x$ .

Modul în care funcția de apartenență fuzzy  $\alpha()$  ia valori în intervalul  $[0; 1]$  este dependent de aplicația vizată. Acest lucru reprezintă un avantaj important al acestui tip de abordare și anume de a fi adaptabil la contexte diferite. De exemplu, definirea conceptului de "temperatură ridicată", nu poate fi realizată în același fel pentru două contexte diferite, precum prognoza vremii și măsurarea temperaturii unui reactor nuclear. Dacă o temperatură de 40 de grade poate fi considerată ca fiind ridicată pentru condițiile meteo, nu putem spune același lucru și pentru reactorul nuclear. Pentru aceasta, cu toate că conceptul va fi formulat textual în același fel, funcțiile de apartenență  $\alpha()$  vor avea variații diferite, acestea fiind determinate pe baza expertizei fiecărui context în parte.

În Figura 6.2 am ilustrat câteva dintre funcțiile de apartenență fuzzy cel mai frecvent folosite. În ciuda diferențelor importante dintre acestea, putem identifica o serie de proprietăți constructive comune:

- în general valoarea maximală este limitată la valoarea 1, în timp ce valoarea minimală este 0, ceea ce corespunde celor două grade de adevăr folosite în logica booleană,

- funcția de apartenență nu trebuie să atingă de mai multe ori, în mod alternativ discontinuu, valoarea maximală,

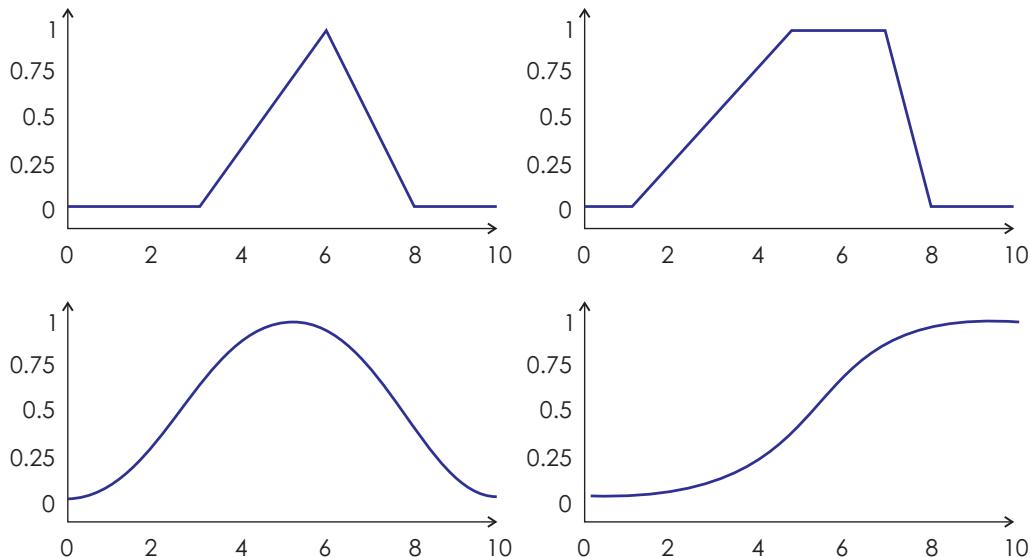


Figura 6.2: Exemple de funcții fuzzy de apartenență: triangulară, trapezoidală, Gausiană și crescătoare (ordine de la stânga la dreapta, și de sus în jos, axa  $oY$  corespunde valorii de adevăr, în timp ce axa  $oX$  reprezintă universul de discurs).

- o funcție de apartenență trebuie să conțină cel puțin o tranziție între valoarea maximală și cea minimală (între 1 și 0).
- în ceea ce privește simetria funcției, întâlnim două situații: fie funcția are o anumită simetrie față de valoarea maximală, fie aceasta este monoton crescătoare sau descrescătoare.

Proprietățile enumerate mai sus nu sunt întâmplătoare, acestea având rolul de a conserva veridicitatea conceptului reprezentat. Astfel, nu orice funcție matematică poate fi folosită ca funcție de apartenență fuzzy.

După cum am menționat anterior, forma și tipul de variație al funcției de apartenență, și astfel a valorilor de adevăr asociate valorilor parametrului analizat, sunt dependente de aplicație. Modul de variație al parametrului vizat influențează forma funcției de apartenență. De exemplu, dacă măsurăm o anumită temperatură, funcția de apartenență la conceptul de "temperatură medie" va fi cel mai probabil o funcție simetrică față de valoarea maximală (vezi funcția Gausiană din Figura 6.2). Acest lucru se datorează faptului că

acest concept implică ca parametrul măsurat, și anume temperatura  $T$ , să ia valori la mijlocul gamei, când temperatura va fi considerată ca fiind medie cu o valoare de adevăr de 1. Dacă în acest caz s-ar folosi, de exemplu o funcție monotonă crescătoare, atunci toate temperaturile ridicate vor fi catalogate ca fiind medii cu o valoare de adevăr de 1, afirmație ce este complet eronată. De asemenea, folosirea unei funcții nesimetrice nu își are sensul în acest caz deoarece felul în care valoarea de adevăr descrește pentru valorile de temperatură scăzute și respectiv ridicate nu are nici un motiv să fie diferit, cele două variații fiind, putem spune, echiprobabile.

În cele mai multe situații, funcția de apartenență este definită *experimental* (empiric) pe baza expertizei manuale a modului de variație al parametrilor vizați. De exemplu, pentru a defini funcția de apartenență la conceptul de "temperatură ridicată" în contextul unui reactor nuclear, avem nevoie de expertiza specialiștilor din domeniu. În funcție de experiența acestora în domeniul fizicii nucleare, funcția de apartenență va fi determinată în conformitate cu modul în care temperatura variază în funcție de procesele fizice din reactor. Generarea automată a funcțiilor de apartenență, cu excepția unor situații simplificate, va duce în mod sigur la obținerea de rezultate eronate. Astfel, formalizarea unui concept folosind logica fuzzy este dependentă de intervenția expertizei umane.

### 6.3 Formalizarea pe baza regulilor fuzzy

Procesul de formalizare pe baza regulilor de decizie fuzzy implică crearea de noi concepte prin combinarea a o serie de concepte fuzzy inițiale. Aceasta implică de regulă două etape, și anume:

- în prima etapă, pornind de la parametrii numerici de nivel scăzut ai sistemului, se definesc o serie de *variabile fuzzy* (*mulțimi fuzzy*) ce vor parametriza sistemul folosind concepte semantice,
- etapa a doua constă în folosirea de operatori logici specifici, ca de exemplu, operatorii de conjuncție și disjuncție logică, ce sunt aplicații variabilelor fuzzy definite anterior pentru a lua decizii cu privire la semnificația și la relațiile existente între acestea. Deciziile sunt luate folosind o bază de reguli fuzzy. Regulile sunt formulate sub forma unor propoziții de tip "dacă ... atunci" ("if ... then") și vor fi definite pentru totalitatea variabilelor fuzzy. Această a doua etapă poartă numele de *inferență fuzzy*.

În cele ce urmează, vom detalia fiecare dintre aceste două etape specifice formalizării fuzzy.

### 6.3.1 Variabilele fuzzy

Înaintea caracterizării semantice fuzzy propriu-zise a unui sistem, sunt definite o serie de variabile fuzzy ce vor caracteriza anumite proprietăți importante ale acestuia. Astfel, fiecărui parametru numeric de nivel scăzut al sistemului îi va fi asociată o anumită variabilă fuzzy. Aceasta este de regulă exprimată sub forma unui concept lingvistic ce poate lua anumite valori textuale.

Logica fuzzy este bazată pe aceste variabile fuzzy, numite și *variabile lingvistice*, ce au o valoare lingvistică în universul de discurs  $U$ . Fiecare valoare lingvistică constituie o submulțime fuzzy a universului de discurs. De exemplu, dacă  $U$  este dat de gama de temperaturi de la 0 la 200 de grade Celsius, atunci o posibilă variabilă lingvistică este conceptul de "temperatură", concept ce poate lua valorile lingvistice: "foarte rece", "rece", "temperat", "cald" și respectiv "foarte cald".

Corespondența între valorile numerice ale parametrilor de nivel scăzut măsurăți și conceptele lingvistice asociate este realizată pe baza *funcțiilor de apartenență fuzzy*, ce au fost prezentate în secțiunea anterioară (vezi Figura 6.2). Funcțiile de apartenență sunt alese în aşa fel încât să redea cât mai bine relația dintre variația parametrului modelat și simbolurile textuale asociate.

Pentru a înțelege mai bine mecanismul de definire a variabilelor fuzzy, în cele ce urmează vom ilustra un exemplu concret (sursă [Lescieux 06]). Să presupunem că vrem să evaluăm cât de înaltă este o anumită persoană în funcție de înălțimea acesteia. Parametrul de nivel scăzut măsurat va fi în acest caz înălțimea  $H$  exprimată în metri. Conceptul lingvistic "înălțimea unei persoane" va fi astfel asociat parametrului  $H$  pe baza unei reprezentări fuzzy. Acest concept va avea, de exemplu, trei valori lingvistice posibile, și anume: "înălțime mică", "înălțime medie" și "înălțime mare", valori ce constituie submulțimile fuzzy ale parametrului  $H$ . O modalitate de asociere a valorilor parametrului  $H$  celor trei simboluri poate fi realizată folosind funcțiile de apartenență ilustrate în Figura 6.3.

Pe baza acestei definiții, putem spune că înălțimea unei persoane este considerată ca fiind: "mică" cu valoarea de adevăr 1 dacă  $H < 1.6$  metri, "medie" cu valoarea de adevăr 1 dacă  $H = 1.7$  metri și respectiv "mare" cu o valoare de adevăr 1 dacă  $H > 1.8$  metri. Partiția fuzzy a universului de discurs este astfel determinată de ansamblul funcțiilor de apartenență la simboluri, în cazul nostru de cele trei funcții ilustrate în Figura 6.3 (vezi ultimul grafic).

Pe baza acestei reprezentări putem transforma valorile numerice ale parametrului  $H$  din intervalul de interes  $[1.5; 1.9]$ , într-o reprezentare simbolică, ce este similară modului de percepție uman. De exemplu, dacă o anumită

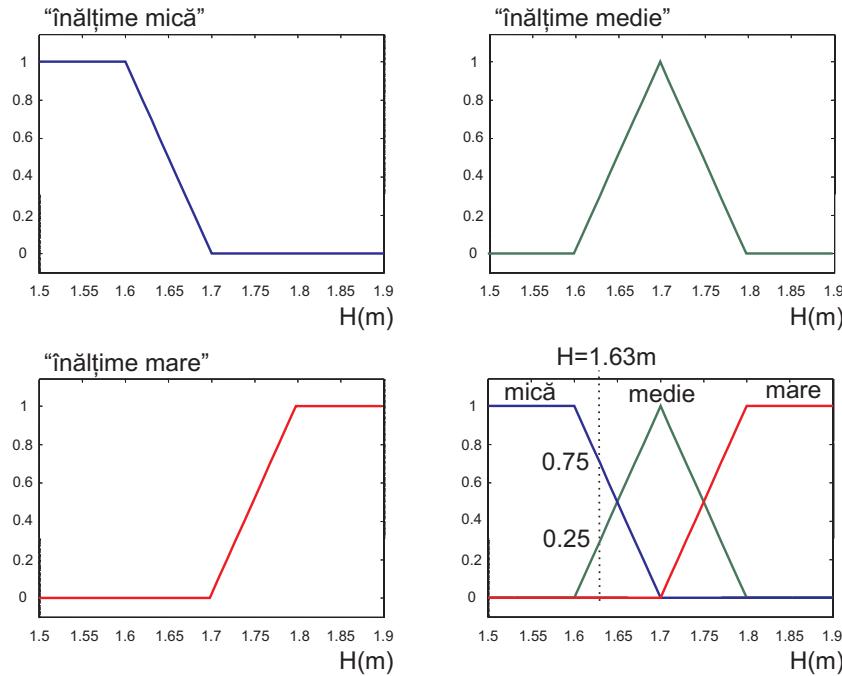


Figura 6.3: Funcțiile de apartenență la valorile lingvistice ale conceptului fuzzy de ”înălțimea unei persoane” și partiția fuzzy a universului de discurs (ultimul grafic). Axa  $oY$  corespunde valorii de adevăr. Exemplu de apartenență la fiecare simbol pentru  $H = 1.63$  metri.

persoană are înălțimea de 1.63 metri, folosind partiția fuzzy definită anterior, putem conchide că:

- propoziția ”persoana are o înălțime mică” are valoarea de adevăr 0.75 (foarte probabil),
- propoziția ”persoana are o înălțime medie” are valoarea de adevăr 0.25 (puțin probabil),
- propoziția ”persoana are o înălțime mare” are valoarea de adevăr 0 (imposibil).

Procedând în același fel pentru toate variabilele numerice din sistem, vom defini totalitatea mulțimilor fuzzy ale acestuia. Pe baza acestora, vom putea analiza mai departe relațiile vizibile, precum și ascunse, care există între variabilele fuzzy, ceea ce ne va permite să luăm decizii semantice cu privire la comportamentul sistemului.

### 6.3.2 Principiul inferenței fuzzy

Inferența fuzzy este definită ca fiind procesul de *luare de decizii* pornind de la o anumită *bază de reguli* de tip ”dacă ... atunci”. Această bază de reguli este construită, fie pe baza expertizei manuale a relațiilor existente între diversele mulțimi fuzzy folosite pentru a reprezenta universul de discurs  $U$ , sau folosind informații ”a priori” despre domeniul de aplicație. Dacă variabilele fuzzy caracterizează proprietățile sistemului, inferența fuzzy caracterizează din punct de vedere semantic relațiile care există între diversele variabile ale sistemului.

#### Proprietățile unei baze de reguli fuzzy

O bază de reguli este caracterizată de o serie de proprietăți. În general, sistemele de inferență fuzzy pot fi considerate ca fiind *aproximatori universali* [Wang 92], ceea ce a condus la tendința ca metodele de inducere a bazei de reguli să fie în principal ghidate și evaluate în raport cu acest unic criteriu de performanță numerică [Guillaume 01].

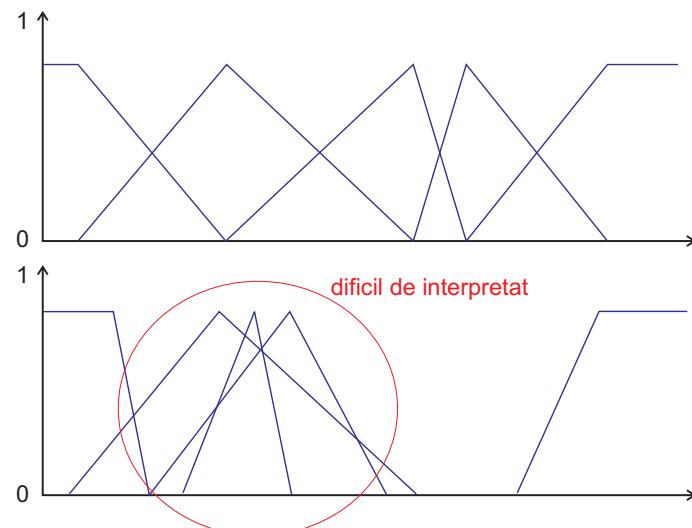


Figura 6.4: Exemple de partiții fuzzy (axa  $oX$  corespunde universului de discurs iar axa  $oY$  valorii de adevăr).

Pe de altă parte, obiectivul de ”extragere a cunoașterii” din universul de discurs al parametrilor măsurăți ai sistemului, impune bazei de reguli o proprietate suplimentară, și anume ca aceasta să respecte *sensul semantic*. Submulțimile fuzzy trebuie să poată fi interpretate în termeni lingvistici,

lucru ce nu este garantat doar de utilizarea unui formalism fuzzy. În Figura 6.4 am ilustrat tocmai acest lucru. Modul de suprapunere a submulțimilor fuzzy din prima imagine a Figurii 6.4 permite ordonarea acestora și astfel interpretarea lor cu termeni lingvistici, ca de exemplu, de la ”foarte mic”, la, ”foarte mare”. În a doua imagine a Figurii 6.4, partitia fuzzy este imposibil de a fi interpretată deoarece nu este posibilă etichetarea distinctă a celor trei submulțimi fuzzy centrale (vezi regiunea marcată), acestea fiind mult prea suprapuse.

Pe langă cele două proprietăți enumerate anterior, o bază de reguli fuzzy mai este caracterizată de trei proprietăți fundamentale, și anume: *coerență*, *continuitate* și respectiv *completitudine*, astfel:

- *coerența bazei de reguli* este una dintre proprietățile esențiale ale acesteia, indiferent de domeniul de aplicație al sistemului. O bază de reguli este coerentă, dacă concluziile regulilor ce sunt simultan activate<sup>1</sup> de un anumit exemplu nu sunt contradictorii. O anumită regulă prezintă incoerențe, dacă distribuția valorilor de ieșire a exemplelor activate de aceasta este foarte heterogenă. Aceste incoerențe traduc faptul că regula nu este suficient de specifică pentru a ține cont de datele de antrenare<sup>2</sup>. Aproximatorii universali divizează spațiul de intrare prin adăugarea de reguli tocmai pentru a elimina aceste incoerențe.
- proprietatea de *continuitate a unei baze de reguli* presupune faptul că mici variații ale intrării sistemului nu produc variații importante ale ieșirii acestuia.
- proprietatea de *completitudine a bazei de reguli* asigură faptul că fiecare dintre valorile posibile ale intrării sistemului activează cel puțin o regulă, astfel încât, în orice situație va fi inferată o valoare de ieșire. În cazul aplicațiilor ce presupun definirea bazei de reguli pe baza expertizei domeniului respectiv, definiția completitudinii este mai puțin restrictivă. În acest caz, completitudinea nu va fi măsurată raportat la totalitatea valorilor posibile ale intrării sistemului, ci doar în raport cu cele conținute în setul de exemple de antrenare furnizate de expert.

În funcție de tipul aplicației vizate, o bază de reguli poate prezenta, fie toate proprietățile enumerate anterior, fie doar o parte dintre acestea. În Tabelul 6.1 am sintetizat cerințele pentru o bază de reguli în cazul câtorva dintre cele mai frecvente aplicații ale acestora.

<sup>1</sup>prin definiție, un exemplu activează o regulă, sau o regulă activează un exemplu, dacă valoarea de adevăr a regulii pentru exemplul în cauză, nu este nulă.

<sup>2</sup>antrenarea supervizată constă în inducerea de relații între intrarea și ieșirea unui sistem pe baza unui set de exemple cunoscute ”a priori” (vezi și Secțiunea 7.2).

Aplicație	Semantică	Coerență	Continuitate	Completit.
Comandă	+	+++	+++	+++
Clasificare	+	++	+	+++
Decizie	+++	+++	+++	++
Caracterizare	+++	+++	++	+
Diagnostic	+++	+++	++	+
Simulare	+++	+++	+++	+

Tabelul 6.1: Proprietățile necesare unei baze de reguli fuzzy pentru diverse tipuri de aplicații (fiecărei proprietăți i-am asociat un scor empiric ce este proporțional cu importanța acesteia pentru aplicația respectivă, sursă [Guillaume 01]).

Astfel, în contextul indexării semantice după conținut, pentru a ajunge la o caracterizare semantică a conținutului datelor, baza de reguli trebuie în primul rând să respecte sensul semantic și să fie coerentă. Aceasta, în mod obligatoriu, trebuie să poată fi capabilă de a furniza rezultate necontradictorii ce pot fi interpretate lingvistic (textual). În ceea ce privește proprietatea de completitudine, aceasta nu este strict necesară deoarece regulile vor fi determinate în acest caz folosind expertiza domeniului respectiv.

### Relațiile dintre submulțimile fuzzy

După cum am menționat anterior, construcția unei baze de reguli fuzzy implică analiza relațiilor existente între diversele mulțimi fuzzy (variabile fuzzy) ce modeleză sistemul. Relațiile dintre acestea sunt exprimate folosind *operatori specifici logicii fuzzy*. Dintre aceștia, operatorii cei mai frecvenți folosiți sunt:

- **operatorul de intersecție:** acesta este un operator de conjuncție logică, reprezentat cu  $\cap$  sau  $\wedge$ , și care este denotat de cuvântul cheie *SI*,
- **operatorul de reuniune:** acesta este un operator de disjuncție logică, reprezentat cu  $\cup$  sau  $\vee$ , și care este denotat de cuvântul cheie *SAU*,
- **operatorul de complementaritate:** acesta este un operator de negație logică și este denotat de cuvântul cheie *NOT*.

Folosirea acestor operatori pentru a caracteriza relațiile dintre variabilele fuzzy ale unui anumit sistem implică ”tradicerea” acestora cu o serie de ope-

ratori matematici specifici. Acești operatori sunt definiți pe baza noțiunilor de *t-conormă* și respectiv *t-normă*.

Astfel, prin definiție, o *t-normă* este o aplicație  $T(x, y)$  între variabilele  $x$  și  $y$  ce satisface următoarele proprietăți:

- valoarea 1 este element neutru:

$$\forall x \in [0; 1] \quad T(x, 1) = T(1, x) = x \quad (6.3)$$

- este comutativă:

$$T(x, y) = T(y, x) \quad (6.4)$$

- este asociativă:

$$T(x, T(y, z)) = T(T(x, y), z) \quad (6.5)$$

- este monotonă:

$$\text{daca } x \leq z, y \leq w \text{ atunci } T(x, y) \leq T(z, w) \quad (6.6)$$

În mod similar, o *t-conormă* este o aplicație  $S(x, y)$  ce satisface aceleași proprietăți ca *t-normă*, cu excepția faptului că în acest caz elementul neutru este valoarea 0, astfel:

$$\forall x \in [0; 1] \quad S(x, 0) = x \quad (6.7)$$

Orice t-normă poate fi folosită pentru a defini *operația de intersecție fuzzy*, iar orice t-conormă poate fi folosită pentru a defini *operația de reuniune fuzzy*.

În ecuațiile următoare am prezentat câteva dintre t-normele și respectiv t-conormele cel mai frecvent folosite, astfel:

- operatorii lui Zadeh:

$$T(x, y) = \text{Min}(x, y) \quad (6.8)$$

$$S(x, y) = \text{Max}(x, y) \quad (6.9)$$

- operatorii probabiliști:

$$T(x, y) = x \cdot y \quad (6.10)$$

$$S(x, y) = x + y + x \cdot y \quad (6.11)$$

- operatorii Lukasiewicz:

$$T(x, y) = \text{Max}(0, x + y - 1) \quad (6.12)$$

$$S(x, y) = \text{Min}(1, x + y) \quad (6.13)$$

- operatorii Zadeh condiționali:

$$T(x, y) = \begin{cases} \text{Min}(x, y) & \text{dacă } x = 1, y = 1 \\ 0 & \text{altfel} \end{cases} \quad (6.14)$$

$$S(x, y) = \begin{cases} \text{Max}(x, y) & \text{dacă } x \cdot y = 0 \\ 1 & \text{altfel} \end{cases} \quad (6.15)$$

- operatorii Yager ( $p > 0$ ):

$$T(x, y) = 1 - \text{Min}([(1 - x)^p + (1 - y)^p]^{1/p}, 1) \quad (6.16)$$

$$S(x, y) = \text{Min}((x^p + y^p)^{1/p}, 1) \quad (6.17)$$

- operatorii Weber ( $\lambda > -1$ ):

$$T(x, y) = \text{Max}\left(0, \frac{x + y - 1 + \lambda \cdot x \cdot y}{1 + \lambda}\right) \quad (6.18)$$

$$S(x, y) = \text{Min}(x + y + \lambda \cdot x \cdot y, 1) \quad (6.19)$$

- operatorii Hamacher ( $\gamma > 0$ ):

$$T(x, y) = \frac{x \cdot y}{\gamma + (1 - \gamma) \cdot (x + y - x \cdot y)} \quad (6.20)$$

$$S(x, y) = \frac{x + y - x \cdot y - (1 - \gamma) \cdot x \cdot y}{1 - (1 - \gamma) \cdot x \cdot y} \quad (6.21)$$

În cele ce urmează, folosind exemplul prezentat la pagina 185, unde conceptul fuzzy de ”înălțimea unei persoane” era asociat parametrului măsurat  $H$ , vom ilustra modul în care operatorii fuzzy sunt folosiți pentru a crea noi variabile fuzzy în cadrul bazei de reguli (sursă [Lescieux 06]).

Dacă  $A$  reprezintă submulțimea fuzzy ”înălțime mică”, dată de funcția de apartenență fuzzy  $\alpha_A()$ , iar  $B$  este submulțimea fuzzy ”înălțime medie”, dată de funcția de apartenență  $\alpha_B()$  (vezi Figura 6.5.a), atunci noua mulțime fuzzy ”înălțime mică” SAU ”înălțime medie” va avea funcția de apartenență,  $\alpha_{A \cup B}()$ , dată de ecuația următoare:

$$\alpha_{A \cup B}(x) = \alpha_A(x) \cup \alpha_B(x), \quad \forall x \in U \quad (6.22)$$

unde  $x$  ia valori în universul de discurs  $U$ , cu  $A, B \subset U$ .

Dacă operatorul de reuniune este exprimat, de exemplu, cu t-conorma Zadeh, atunci funcția de apartenență  $\alpha_{A \cup B}()$  poate fi exprimată în felul următor:

$$\alpha_{A \cup B}(x) = \text{Max}(\alpha_A(x), \alpha_B(x)), \quad \forall x \in U \quad (6.23)$$

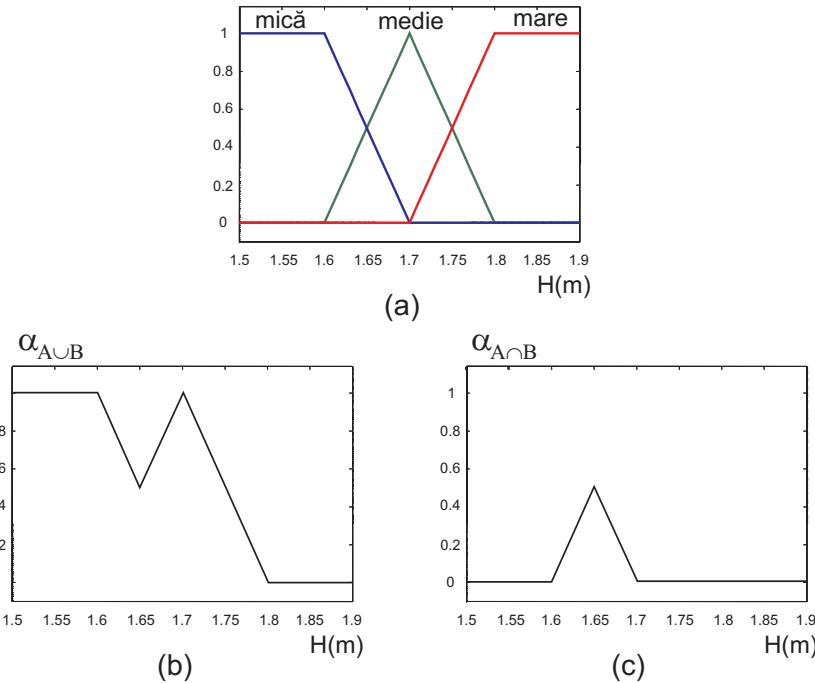


Figura 6.5: Operatori fuzzy: (a) partiția universului de discurs, (b) funcția de apartenență  $\alpha_{A \cup B}$ , (c) funcția de apartenență  $\alpha_{A \cap B}$  ( $A, B$  reprezintă submulțimile fuzzy "înălțime mică" și respectiv "înălțime medie").

Similar, mulțimea fuzzy "înălțime mică" și "înălțime medie" va avea funcția de apartenență dată de relația:

$$\alpha_{A \cap B}(x) = \alpha_A(x) \cap \alpha_B(x) = \text{Min}(\alpha_A(x), \alpha_B(x)), \quad \forall x \in U \quad (6.24)$$

unde operatorul de intersecție a fost definit folosind t-norma Zadeh. Noile funcții de apartenență astfel obținute pe baza operatorilor Zadeh sunt ilustrate în Figura 6.5.

Baza de reguli fuzzy ce va modela comportamentul dinamic al sistemului este definită în funcție de context pe baza diferitelor combinații relevante dintre mulțimile fuzzy ce descriu caracteristicile sistemului.

### Generarea regulilor fuzzy

În general, o regulă fuzzy este exprimată folosind următoarea structură:

$$\begin{aligned} \textbf{DACĂ } &(X_1 \text{ este } A_1) \text{ op } \dots \text{ op } (X_n \text{ este } A_n) \text{ ATUNCI} \\ &(Y_1 \text{ este } B_1) \text{ op } \dots \text{ op } (Y_m \text{ este } B_m) \end{aligned} \quad (6.25)$$

unde  $X_i$  și  $Y_j$  reprezintă variabilele lingvistice de intrare și respectiv de ieșire ale sistemului,  $A_i$  și  $B_j$  reprezintă submulțimile fuzzy ce definesc partitioarea spațiului de intrare și respectiv de ieșire al sistemului,  $i = 1, \dots, n$ ,  $j = 1, \dots, m$ , o propoziție de tip ” $X$  este  $A$ ” este cuantificată de gradul de apartenență al variabilei lingvistice  $X$  la submulțimea fuzzy  $A$  iar  $op$  reprezintă un operator fuzzy (de exemplu, intersecție sau reuniune). Această modalitate de formulare este cunoscută și sub numele de Mamdani, având particularitatea că valoarea returnată este tot o mulțime fuzzy, similară intrării.

Astfel, în formularea unei reguli fuzzy întâlnim trei părți distincte, și anume: *premisele* regulii, *implicația* acestora precum și *concluzia*:

- **premisele** reprezintă ipotezele de plecare ce sunt exprimate sub forma unei succesiuni de afirmații, interconectate de operatori fuzzy. Acestea sunt încadrate în regula fuzzy de cuvintele cheie ”DACĂ” și ”ATUNCI”,
- **implicația** premiselor este anunțată de cuvântul cheie ”ATUNCI”,
- **concluzia** este formulată de regulă în mod similar cu premisele, și anume folosind o succesiune de afirmații. Aceasta se găsește după cuvântul cheie ”ATUNCI” și reprezintă rezultatul premiselor inițiale.

Dacă luăm ca exemplu următoarea regulă simplă:

**DACĂ** (*Vremea este ”frumoasă”*) **și** (*Momentul este ”dimineață”*)

**ATUNCI** (*Moralul este ”ridicat”*)

atunci premisele acesteia vor fi date de variabilele *Vremea* și *Momentul* ce iau valorile ”frumoasă” și respectiv ”dimineață”, premise ce duc la implicația concluziei că variabila *Moralul* va lua valoarea ”ridicat”.

Similar variabilelor fuzzy, regulile fuzzy au și ele o valoare de adevăr. Fiind dată o regulă,  $r$ , valoarea de adevăr a acesteia pentru un anumit exemplu, numită și pondere a regulii, notată cu  $\omega_r$ , rezultă dintr-o operație logică între elementele premisei, astfel:

$$\omega_r = \alpha_{A_1}(x_1) op \alpha_{A_2}(x_2) op \dots op \alpha_{A_n}(x_n) \quad (6.26)$$

unde  $\alpha_{A_i}(x_i)$  reprezintă gradul de apartenență a valorii  $x_i$  din universul de discurs la submulțimea fuzzy  $A_i$ , cu  $i = 1, \dots, n$  iar  $op$  este un operator fuzzy.

O altă modalitate de definire a unei reguli este folosind modelul Takagi-Sugeno. Diferența față de modelul Mamdani constă în faptul că în acest caz concluzia regulii va fi una netă. Astfel, valoarea de ieșire  $j$  a sistemului pentru regula  $r$  este calculată ca fiind o combinație liniară a valorilor de intrare, astfel:

$$y_j^r = b_{j,0} + b_{j,1} \cdot x_1 + b_{j,2} \cdot x_2 + \dots + b_{j,n} \cdot x_n \quad (6.27)$$

unde  $b_{j,i}$ , cu  $i = 1, \dots, n$  și  $j = 1, \dots, m$ , reprezintă coeficienții de pondere ce sunt determinați în funcție de importanța fiecărei valori de intrare,  $x_i$ , pentru valoarea ieșirii sistemului.

Procesul de inducere a bazei de reguli fuzzy a unui sistem constă în formularea tuturor mulțimilor de reguli ce vor face legătura între ieșirile sistemului și intrările acestuia. Pentru un studiu detaliat a literaturii de specialitate din teoria mulțimilor fuzzy, cititorul se poate raporta la lucrările [Guillaume 01] și [Klir 95].

## 6.4 Avantajele reprezentării fuzzy

Avantajele furnizate de logica fuzzy la formalizarea conceptuală a informațiilor numerice, cât și simbolice, de nivel scăzut, nu sunt deloc neglijabile. Acestea pot fi sintetizate cu următoarele:

- sistemele de inferență fuzzy sunt **aproximatori universali**. Universul de discurs, care de regulă este foarte vast sau chiar infinit, este convertit cu ajutorul formalizării fuzzy într-un număr limitat de concepte. Cercetările din domeniu au demonstrat de-a lungul timpului că sistemele de inferență fuzzy sunt cel puțin la fel de performante ca alte tehnici consacrate de aproximare a datelor [Wang 92].
- formalizarea fuzzy **reduce complexitatea** sistemului modelat. În cazul în care cantitatea de informație ce trebuie prelucrată este mult prea mare pentru a putea fi controlată în totalitate, iar înțelegerea proceselor existente în sistem este limitată, formalizarea fuzzy permite reducerea complexității prin introducerea conceptului de incertitudine. În general, s-a dovedit faptul că, cu cât este tolerată mai multă incertitudine în modelarea sistemului, cu atât complexitatea acestuia tinde să scadă, iar credibilitatea modelului obținut tinde să crească [Klir 95].
- datorită conceptului de incertitudine, variabilele fuzzy **reprazintă realitatea**, care este în general una incertă, mult mai fidel decât o fac variabilele clasice nete. Această proprietate a fost enunțată chiar de Albert Einstein în 1921: *pe cât legile matematice se raportează la realitate, pe atât acestea nu sunt veridice. Și pe cât sunt veridice, acestea nu se raportează la realitate.*
- în logica fuzzy, concepțile vagi sunt reprezentate într-un **limbaj natural**. Această proprietate este una dintre cele mai importante. Formalizarea fuzzy transformă mărimi numerice în concepte lingvistice ce

sunt exprimate într-un mod similar modului de percepție uman. Dacă universul de discurs este unul numeric, formalizarea fuzzy este realizată cu variabile lingvistice ce iau valori textuale [Lescieux 06].

- modul de funcționare al logicii fuzzy este foarte similar cu modul de percepție uman. Însuși creierul uman funcționează după principiile logicii fuzzy. Creierul uman apreciază valorile stimulilor de intrare (valori continue) într-un mod aproximativ. De exemplu, apreciem un obiect ca fiind apropiat sau depărtat și nu ca fiind la distanță punctuală  $d$ . Astfel, formalizarea fuzzy tinde să fie **în concordanță cu percepția semantică**, facilitând astfel "extragerea de cunoaștere" ("knowledge") din date numerice de nivel scăzut.
- faptul că formalizarea fuzzy este construită în general pe baza expertizei și a cunoașterii "a priori" a domeniul de aplicație, conferă modelului obținut mult mai multă **coerență** decât unui model clasic.
- formalizarea fuzzy **normalizează valorile datelor** de intrare. Conceptelor lingvistice le sunt asociate valori de adevăr ce sunt valori reale, de regulă normalizează între 0 și 1, unde 0 reprezintă negația sigură iar 1 afirmația sigură. În acest fel, comparația parametrilor cu plaje de valori diferite este simplificată.
- reprezentarea clasă a datelor pe baza logicii booleene este un caz particular al logicii fuzzy. Astfel, **formalizarea fuzzy include formalizarea netă**.

În cele ce urmează, vom prezenta câteva exemple de aplicații practice ale formalizării fuzzy, atât în domeniul indexării după conținut a secvențelor de imagini, cât și în domenii conexe.

## 6.5 Aplicabilitatea sistemelor fuzzy

Fiind foarte apropiată de modul în care percepem realitatea înconjurătoare, logica fuzzy a câștigat rapid teren în aproape toate domeniile de activitate existente, înlocuind parțial sau chiar în totalitate logica clasă booleană. Acest lucru se datorează în mare parte faptului că aceasta modelează însuși modul de funcționare al creierului uman, care este unul fuzzy și nu net.

Creierul uman apreciază diversele variabile din mediul înconjurător într-un mod aproximativ, ca de exemplu folosind atrbute precum redus/ridicat, aproape/departe, etc. Același lucru este valabil și pentru răspunsul la anumiți stimuli exteriori, de exemplu când conducem o mașină frânăm mai puțin sau

mai tare în funcție de indicația semaforului și de distanța până la acesta. Astfel, pe baza acestor date aproximative, creierul uman își dezvoltă o serie de reguli, concepute într-un mod similar modului de inferare a regulilor fuzzy, reguli ce vor determina răspunsul nostru în funcție de diversele valori ale variabilelor de intrare.

Sistemele fuzzy își găsesc aplicație practic în toate domeniile existente ce se bazează pe luarea de decizii, ca de exemplu: asistare la diagnostic și luarea de decizii (domeniul medical, orientare profesională, etc.), baze de date (obiecte fuzzy și formularea cererilor de căutare cu multimi fuzzy), recunoaștere automată a formelor, vedere asistată de calculator, prelucrare de imagini, sisteme de agregare multicriteriu, sisteme de optimizare, comandă fuzzy a sistemelor, etc. Standardizarea anumitor metode de analiză fuzzy a dus chiar la implementarea hardware sau software a acestora în diverse produse finite destinate publicului larg sau specialiștilor din domeniu, ca de exemplu: aparate electrocasnice, sisteme audio-vizuale, procesoare dedicate, interfețe de dezvoltare specifice precum procesorul Motorola 68HC12 sau Thomson WARP [Lescieux 06].

În domeniul prelucrării de imagini și al vederii asistate de calculator, aporțul logicii fuzzy este incontestabil. Aceasta constituie puntea de legătură între datele numerice de nivel scăzut obținute în urma diverselor măsurători ale proprietăților fizice ale fenomenelor analizate, sau în urma anumitor etape de prelucrare, și modul de percepție al acestora. Conceptele fuzzy sunt exprimate într-un limbaj natural ceea ce face ca "cifrele să capete sens".

Datorită diversității foarte mare de aplicații posibile ale logicii fuzzy, este aproape imposibil să realizăm o clasificare a metodelor existente. Astfel că, în cele ce urmează ne vom limita la prezentarea a doar câteva exemple semnificative de aplicații ale logicii fuzzy în domeniul analizei și prelucrării secvențelor de imagini.

Un exemplu de folosire a logicii fuzzy pentru a simplifica și a da sens comparației datelor numerice cu plaje de valori diferite, este metoda propusă în [Doulamis 00c]. În aceasta, variabilele fuzzy sunt folosite pentru a caracteriza proprietățile de culoare și de adâncime a secvențelor de imagini stereoscopice<sup>3</sup> în vederea rezumării automate a conținutului acestora. Pentru aceasta, conținutul vizual al fiecărei imagini este mai întâi reprezentat cu o serie de parametri extrași la nivel de segment (obiect), formând astfel vectorul de caracteristici al imaginii. Cum numărul de segmente variază de la o imagine la alta, dimensiunea vectorului de caracteristici va fi de asemenea

---

<sup>3</sup>imaginile stereoscopice sau imaginile 3D sunt imagini ce sunt capabile să redea informația spațială 3D prin crearea iluziei de adâncime. În cazul imaginilor sau a secvențelor de imagini, efectul de adâncime este obținut prin înregistrarea a două imagini din perspective un pic diferite, ce vor fi proiectate independent fiecărui ochi.

variabilă, fapt ce duce la imposibilitatea de a evalua similaritatea imaginilor pe baza comparației vectorilor de caracteristici.

Pentru a soluționa această problemă sunt folosite variabile fuzzy. Segmentele sunt mai întâi clasificate într-un anumit număr de clase predefinite, formând astfel o histogramă multidimensională. Fiecare element al vectorului de caracteristici va corespunde în acest fel unei anumite clase sau cu alte cuvinte, unui anumit bin al histogramei. Pentru a evita clasificarea a două segmente similare în clase diferite, clasificare ce poate duce la erori de comparație semnificative, fiecare clasă va fi reprezentată de un grad de apartenență fuzzy. Astfel, fiecărui segment i se va permite să aparțină la mai multe clase, dar cu grade diferite de apartenență.

Din punct de vedere matematic, fiecare segment  $S_i$ , cu  $i = 1, \dots, K$ , unde  $K$  reprezintă numărul de segmente extrase din imagine, va fi reprezentat de vectorul de caracteristici  $s_i$  următor:  $s_i = [c_{S_i}^T d_{S_i} l_{S_i}^T a_{S_i}]$ , unde matricea  $c_{S_i}$  conține componente medii de culoare ale segmentului  $S_i$ ,  $d_{S_i}$  și  $a_{S_i}$  reprezintă adâncimea și respectiv dimensiunea acestuia iar  $l_{S_i}$  este o matrice ce conține coordonatele orizontale și verticale ale centrului segmentului ( $T$  denotă operația de transpunere matricială).

Domeniul de valori al elementului  $j$  din vectorul de caracteristici  $s_i$  va fi partionat în  $Q$  regiuni folosind funcțiile de apartenență fuzzy  $\mu_{n_j}(s_{i,j})$  (funcții triunghiulare cu suprapunere de 50%), unde  $n_j = 1, \dots, Q$ . Astfel,  $\mu_{n_j}(s_{i,j})$  va reprezenta gradul de apartenență al lui  $s_{i,j}$  la clasa de indice  $n_j$ . Ansamblul claselor  $n_j$  pentru toate valorile lui  $j$ , cu  $j = 1, \dots, L$ , va fi dat de clasa L-dimensională  $n = [n_1 n_2 \dots n_L]^T$ . În acest fel, gradul de apartenență al vectorului de caracteristici  $s_i$  la clasa  $n$ ,  $\mu_n(s_i)$ , poate fi definit ca un produs de funcții de apartenență fuzzy, astfel:

$$\mu_n(s_i) = \prod_{j=1}^L \mu_{n_j}(s_{i,j}) \quad (6.28)$$

Mai departe, histograma fuzzy de caracteristici,  $H_f()$ , este construită ca fiind:

$$H_f(n) = \frac{1}{K} \sum_{i=1}^K \mu_n(s_i) \quad (6.29)$$

unde valoarea  $H_f(n)$  va reprezenta gradul de apartenență al întregii imagini la clasa  $n$ . Vectorul de caracteristici al imaginii va fi dat astfel de valorile lui  $H_f()$  pentru toate clasele posibile  $n$ . Aceasta va fi folosit mai departe pentru a evalua similaritatea imaginilor în vederea selectării imaginilor cheie ale rezumatului.

În [Dorado 04] variabilele fuzzy sunt folosite pentru a ”copia” modul de percepție uman în încercarea de adnotare automată a conținutului secvențelor de imagini. Astfel, datele de nivel scăzut extrase din secvență vor fi convertite într-o serie de concepte semantice. Aceste concepte sunt definite de un anumit lexicon format, atât din simboluri lingvistice, cât și din simboluri grafice considerate ca fiind semnificative pentru descrierea conținutului secvenței. Etichetarea secvenței cu simbolurile din lexicon va fi realizată folosind o bază de reguli. Aceasta este obținută într-o etapă de învățare, pe baza expertizei domeniului de aplicație.

Principiul de extragere a ”cunoașterii” (“knowledge”) din parametrii numerici de nivel scăzut este următorul: dacă  $W^f$  reprezintă setul de cuvinte asociate parametrilor de nivel scăzut considerați,  $W^f \subset L$ , unde  $L$  reprezintă lexiconul semantic, atunci procesul de asociere a conceptelor semantice poate fi definit ca fiind o aplicație  $\tilde{M}(w, \tilde{A})$  ce face trecerea de la un set de cuvinte din  $L$  la un set de interpretări ale acestora definite de universul  $Y$ . Fiecare cuvânt,  $w \in L$ , corespunde unei mulțimi fuzzy,  $\tilde{A} \in Y$ , mulțime ce reprezintă interpretarea semantică a lui  $w$ . Gradul de adevăr al acestei apartenențe va fi furnizat de funcțiile de apartenență fuzzy:  $\mu_{\tilde{M}}(w, y) = \mu_{\tilde{A}}(y)$ .

De exemplu, în cazul secvențelor de știri, folosind ca parametri de nivel scăzut decupajul în plane video și culorile predominante ale anumitor imagini cheie ale secvenței, se poate defini un lexicon care să conțină două concepte, și anume: conceptul de ”anchorpersoană”<sup>4</sup> ce survine în cazul în care culorile predominante prezintă variații slabe pe parcursul secvenței, și respectiv conceptul de ”raportaj” în cazul în care culorile predominante au o variație semnificativă. În acest caz, reprezentarea semantică este dată de două mulțimi fuzzy ce sunt legate de valoarea medie a schimbării de culoare.

Un alt exemplu este sistemul propus în [Ionescu 08] unde mulțimile fuzzy sunt folosite pentru a descrie percepția conținutului de culoare din secvențele artistice de animație. Metoda folosită este una clasică de inferență fuzzy. Conținutul de culoare al secvenței este mai întâi caracterizat cu o serie de parametri statistici de nivel scăzut, precum distribuția globală de culoare a secvenței, distribuția de culori elementare, gradul de variabilitate al culorilor, procentul de culori închise sau deschise, etc.

Trecerea la un nivel semantic al descrierii se realizează în două etape. Mai întâi, un nivel simbolic de descriere este obținut prin asocierea de mulțimi fuzzy fiecărui parametru de nivel scăzut. Funcțiile de apartenență sunt definite pe baza expertizei domeniului particular al filmului de animație. De

---

<sup>4</sup>”anchorpersoană” este un termen folosit în domeniul televiziunii pentru a desemna reporterul ce coordonează o anumită transmisie TV în care intervin mai mulți corespondenți.

exemplu, variabila lingvistică ”light color content” este asociată parametrului  $P_{light}$  ce reprezintă procentul de culori deschise prezente în secvență. Conceptul este descris mai departe cu trei simboluri, și anume: ”low-light color content”, ”mean-light color content” și respectiv ”high-light color content”. Astfel, secvența are o distribuție de culoare săracă în culori deschise (valoare de adevăr 1) dacă  $100 \cdot P_{light} < 33\%$ , o distribuție cu un conținut mediu de culori deschise (valoare de adevăr 1) dacă  $100 \cdot P_{light} > 50\%$  și  $100 \cdot P_{light} < 60\%$ , și respectiv o distribuție de culoare bogată în culori deschise (valoare de adevăr 1) dacă  $100 \cdot P_{light} > 66\%$ .

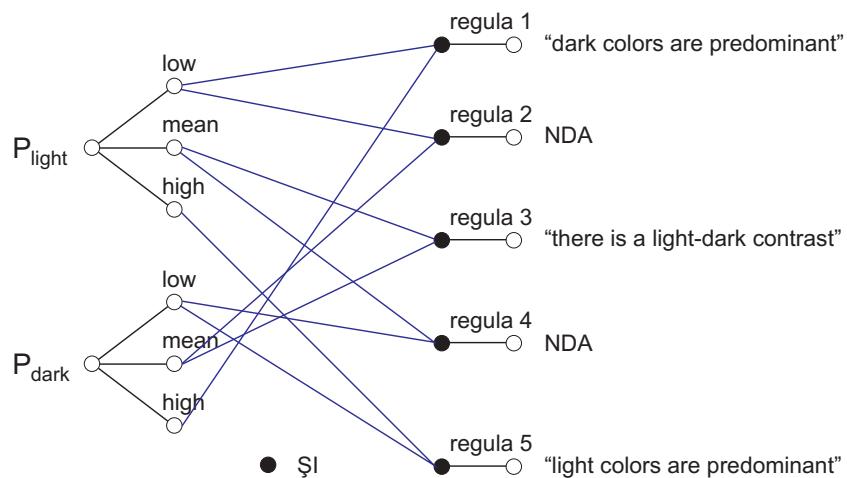


Figura 6.6: Exemplu de bază de reguli fuzzy folosită la descrierea intensității de culoare în filmele artistice de animație (NDA = ”fără descriere”, sursă [Ionescu 08]).

Nivelul de descriere semantic propriu-zis este obținut prin introducerea de reguli fuzzy. Descrierile vizate sunt tehniciile artistice de culoare prezente în filmele de animație, precum celește contraste ale lui Itten sau schemele de armonie a culorilor (vezi Secțiunea 4.2.3). Un exemplu este prezentat în Figura 6.6. Descrierile obținute, atât simbolice cât și semantice sunt folosite ulterior pentru indexarea automată după conținut a unei baze de secvențe de imagini, precum și ca informații de conținut pentru facilitarea navigării în baza de date.

## 6.6 Concluzii

În acest capitol am prezentat modul în care formalizarea fuzzy poate fi folosită la descrierea semantică, pe bază de concepte lingvistice, a conținutului

datelor. Logica fuzzy, spre deosebire de logica clasică booleană, prin introducerea conceptului de incertitudine are avantajul de a fi apropiată modului de percepție uman, care este unul aproximativ. În logica fuzzy, percepția fenomenelor înconjurătoare nu este o problemă de validare sau invalidare, ci mai degrabă o problemă de grad de adevăr.

Formalizarea fuzzy implică de regulă două etape distincte. Prima etapă este etapa de "fuzzyficare" și constă în asocierea de concepte lingvistice parametrilor numerici ce modeleză sistemul. Acestea sunt definite de o serie de valori lingvistice ce constituie submultimile fuzzy ale universului de discurs. Etapa a doua este etapa de inferență fuzzy și constă în inducerea unei baze de reguli fuzzy pentru a modeliza, atât relațiile vizibile, cât și cele ascunse ce pot exista între diversii parametri ai sistemului.

Unul dintre principalele avantaje ale formalizării fuzzy, relativ la utilizarea acesteia la descrierea anumitor procese de interes, constă în trecerea de la o reprezentare numerică de nivel scăzut a proprietăților sistemului modelat la o reprezentare lingvistică de nivel semantic superior. Astfel, formalizarea fuzzy *dă sens* valorilor numerice obținute în diverse etape de prelucrare. Dacă universul de discurs este unul numeric, universul conceptelor fuzzy este unul textual.

Unul dintre principalele dezavantaje ale formalizării fuzzy este dat de faptul că în cele mai multe cazuri, exceptând sistemele cu un grad de complexitate foarte redusă, aceasta nu este un proces automat. Formalizarea fuzzy necesită intervenția unui operator sau a unui expert. Acesta, pe baza expertizei domeniului de aplicație va determina mecanismul de "fuzzyficare", și eventual pe cel de inferență fuzzy, astfel încât sistemul modelat să răspundă cerințelor de prelucrare. Pe de altă parte, tocmai implicarea raționamentului uman în procesul de proiectare face ca formalizarea fuzzy să fie în conformitate cu realitatea.

Datorită caracterului semantic, evaluarea unui model fuzzy este un proces subiectiv. Totuși, implicarea umană în procesul de proiectare al mecanismului fuzzy implică într-o anumită măsură și veridicitatea rezultatelor obținute.

## CAPITOLUL 7

---

### Clasificarea după conținut a datelor

---

**Rezumat:** Căutarea informației utile într-o colecție mare de date este un lucru relativ dificil de realizat. Aceasta se datorează imposibilității noastre de a procesa și tria volume mari de date. Gruparea automată a datelor după anumite criterii de similaritate, transferă problema căutării globale într-o problemă de localizare a datelor într-o colecție restrânsă de date de același tip. Metodele de clasificare automată, vizează tocmai partaționarea automată a datelor în clase omogene din punct de vedere al conținutului. În acest capitol vom prezenta tehniciile existente de clasificare, atât supervizată cât și nesupervizată, punând accentul pe utilitatea acestora în procesul de indexare după conținut.

Tehnicile de clasificare a datelor permit, pe baza definirii unor criterii de similaritate, regruparea automată în funcție de conținut a volumelor vaste de date. Aceste date, accesate în mod direct, sunt de regulă greu accesibile utilizatorului.

Din punct de vedere matematic, problematica abordată de metodele de clasificare existente este următoarea: având la dispoziție o mulțime de  $N$  obiecte ale căror proprietăți fizice sunt sintetizate prin intermediul unor vectori multidimensionali de caracteristici, problema care se pune este definirea unei anumite partații, cât mai relevante, a acestor obiecte într-un anumit număr de clase [Hinneburg 00]. În acest scop se definește conceptul de simi-

laritate între vectorii de caracteristici asociați obiectelor, concept pe baza căruia obiectele aşa zise similare vor fi atribuite aceleiași clase. Principiul clasificării este sintetizat în Figura 7.1.

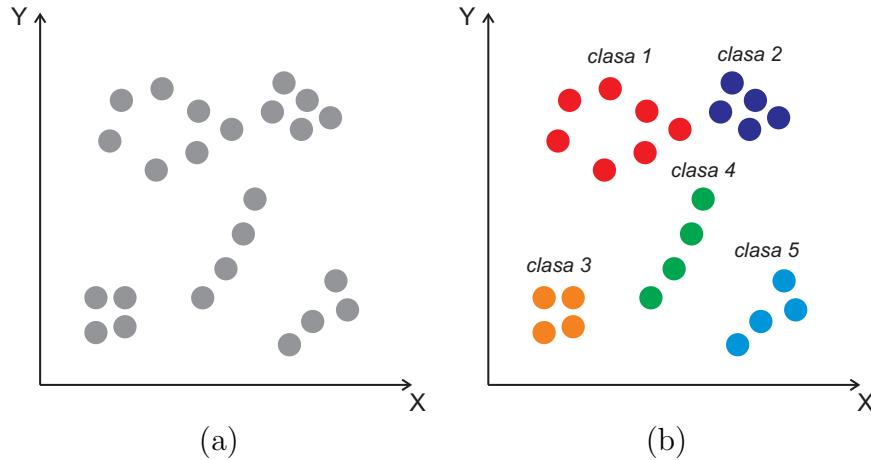


Figura 7.1: Principiul clasificării datelor: (a) datele de intrare reprezentate în spațiul de caracteristici, (b) repartitia în clase obținută în urma clasificării (obiectele din aceeași clasă sunt reprezentate cu aceeași culoare).

Tehnicile de clasificare existente sunt utilizate într-o gamă foarte largă de aplicații. Acestea deservesc diverse obiective, dintre care putem menționa drept cele mai importante următoarele:

- **reducerea volumului informațional:** tehniciile de clasificare a datelor permit regruparea unui ansamblu de date în grupuri omogene, lucrul ce facilitează reducerea volumului informațional disponibil. De exemplu, fiecare grup (clasă) de date poate fi reprezentat, în etapele de prelucrare ulterioare, doar de informația cea mai reprezentativă a grupului. La un alt nivel, clasificarea datelor permite eliminarea redundanței informaționale prin reducerea spațiului de caracteristici.
- **punerea în evidență a "cunoașterii":** tehniciile de clasificare permit localizarea într-un volum mare de date a unor grupuri de informații ce prezintă anumite caracteristici de interes. Localizarea acestora furnizează utilizatorului o cunoaștere nouă a relațiilor existente între date, cunoaștere ce nu era disponibilă anterior căutării. Acest proces este cunoscut în literatura de specialitate și sub numele de "data mining".

- **punerea în evidență a relevanței claselor:** tehnici de clasificare permit localizarea anumitor grupuri de date ce sunt reprezentative pentru ansamblul datelor analizate,
- **punerea în evidență a datelor atipice:** tehnici de clasificare permit de asemenea localizarea datelor ce nu corespund niciunui criteriu de similaritate, date ce sunt considerate ca fiind atipice pentru criteriile considerate. Acestea sunt importante, deoarece sunt un caz particular și trebuie analizate separat. Un astfel de exemplu sunt datele ce se găsesc pe frontiera dintre două clase diferite, date ce pot fi considerate ca aparținând ambelor clase cât și ca o clasă independentă.

Din punct de vedere al problematicii indexării datelor, subiect ce face obiectul acestei cărți, tehnici de clasificare sunt indispensabile unui sistem de indexare după conținut. Clasificarea datelor intervine în general în însuși procesul de căutare al informației. Utilizatorul, prin formularea cererii de căutare va defini spațiul de caracteristici ce va fi folosit pentru localizarea datelor dorite. Pe baza acestuia, datele din baza de date pot fi grupate în funcție de similaritate sau cu alte cuvinte în funcție de asemănarea dintre vectorii de caracteristici asociați. Astfel, grupul sau grupurile de date ce sunt suficient de similare vectorului de caracteristici asociat cererii de căutare vor fi furnizate utilizatorului drept rezultat.

Metodele de clasificare existente se împart în două mari categorii. Prima categorie de metode o constituie *metodele probabilistice* sau de *clasificare supervizată*. Clasificarea supervizată implică clasarea datelor pe baza unor modele predefinite de clase sau "date de antrenament". Acestea reprezintă de regulă o clasificare de referință ce corespunde realității (similară unei "realități de teren" sau "groundtruth"<sup>1</sup>), folosită inițial pentru antrenarea sistemului înaintea clasificării propriu-zise a datelor. În literatura de specialitate, termenul asociat metodelor din această categorie este de "metode de clasificare" sau "classification methods".

O a doua categorie de metode de clasificare sunt *metodele de clasificare nesupervizată sau automată*, desemnate în literatura de specialitate prin termenul de "clustering"<sup>2</sup>. Clasificarea nesupervizată, spre deosebire de clasificarea supervizată, propune o partiție optimă a spațiului de caracteristici din punct de vedere al unui anumit criteriu matematic, fără a folosi informații "a priori" (de exemplu, o partiție de referință). Avantajul acestor metode este dat de faptul că sunt complet automate (nu necesită intervenția utilizatorului) și pot fi folosite pentru clasarea datelor despre care nu dispunem de

---

<sup>1</sup>vezi explicația de la pagina 170.

<sup>2</sup>de notat este faptul că în limba română, termenul de clasificare este folosit generic pentru a desemna, în funcție de context, atât o clasificare supervizată cât și nesupervizată.

informații relative la conținutul acestora (număr de clase, prototipul clasei, etc.). Pe de altă parte, fiind un proces automat, relevanța claselor tinde să fie mai redusă decât în cazul clasificării supervizate, aceasta fiind dependentă de metoda folosită cât și de puterea discriminatorie a spațiului de caracteristici folosit.

În cele ce urmează, vom face o trecere în revistă a tehniciilor de clasificare supervizată și nesupervizată existente punând în evidență avantajele cât și dezavantajele fiecărei abordări.

## 7.1 Clasificarea nesupervizată a datelor

O catalogare interesantă a metodelor de clasificare nesupervizată existente în funcție de proprietățile contrastante ale acestora este propusă în [Jain 99], astfel, metodele de clasificare nesupervizată sunt:

- **acumulative sau partitioane**: această proprietate este legată de modul de structurare al algoritmului folosit. Metodele acumulative pornesc clasificarea de la o anumită partitură în clase, clase care pe parcursul algoritmului sunt fuzionate iterativ până când este satisfăcut un anumit criteriu de convergență. Pe de altă parte, metodele partitioane pornesc de la o singură clasă ce este divizată iterativ până când criteriul de convergență considerat este satisfăcut.
- **monotetice sau politetice<sup>3</sup>**: aceste proprietăți sunt legate de modul de utilizare a vectorilor de caracteristici în procesul de clasificare, care poate fi secvențial sau simultan. Mare parte a metodelor existente sunt politetice, astfel că pentru estimarea distanței dintre obiecte sunt folosiți toți parametrii disponibili ("features"). De asemenea, decizile de clasare sunt luate pe baza acestei măsuri de distanță. Pe de altă parte, metodele monotetice folosesc parametrii în mod secvențial pentru a constitui progresiv clasele, de exemplu, parametrul  $x_1$  este folosit pentru a diviza datele în două clase, mai departe, parametrul  $x_2$  este folosit pentru divizarea claselor anterioare, și aşa mai departe.
- **nete sau fuzzy**. O clasificare netă presupune alocarea fiecărui obiect unei singure clase, astfel apartenența fiind sigură. Pe de altă parte, o clasificare fuzzy asociază fiecărui obiect un grad de apartenență la una sau mai multe clase, apartenența la clase fiind de această dată incertă.

---

<sup>3</sup>termenul de monotetic desemnază o anumită clasă ai cărei membrii sunt identici din punct de vedere al tuturor caracteristicilor acestora, în timp ce termenul de politetic desemnează o anumită clasă ai cărei membrii sunt similari dar nu identici.

Clasificarea fuzzy permite trecerea la o clasificare netă prin asocierea obiectelor claselor cu gradul de apartenență cel mai semnificativ.

- **deterministe sau stohastice:** aceste proprietăți sunt relevante mai ales pentru abordările partiționale ce se bazează pe optimizarea unei erori pătratice. Această optimizare poate fi realizată, fie în mod clasic (determinist), fie pe baza unei căutări aleatoare în spațiul format de toate clasificările posibile (stochastic).
- **incrementale sau non-incrementale:** aceste proprietăți intervin atunci când volumul de date este foarte mare și constrângerile impuse (temp de execuție, putere de calcul, etc.) tind să influențeze arhitectura algoritmului. Astfel, în general, metodele incrementale încearcă să minimizeze numărul de citiri al datelor, să reducă numărul de repartiții în clase analizate pe parcursul algoritmului sau să reducă dimensiunea structurilor de date folosite de algoritm.

În general, diversitatea metodelor de clasificare existente lasă suficientă flexibilitate clasificării astfel încât aceasta să poată fi implementată optimal pentru aplicația dorită.

### 7.1.1 Etapele clasificării nesupervizate

În ciuda diversității abordărilor existente, metodele de clasificare nesupervizată implică o serie de etape comune [Jain 99], astfel:

- **reprezentarea datelor:** această primă etapă constă în alegerea parametrilor metodei de clasificare, cât și a datelor folosite, de exemplu: numărul de clase anticipate, numărul de modele disponibile, tipul și dimensiunea vectorilor de caracteristici ce vor defini spațiul în care se va efectua clasificarea. De asemenea, tot în această etapă se poate realiza *selecția de caracteristici* sau/și *extragerea de caracteristici*. Selectia caracteristicilor reprezintă procesul de identificare a unei submulțimi optimale a setului inițial de caracteristici. Această etapă permite astfel reducerea dimensiunii datelor de intrare (reducerea dimensiunii vectorilor de caracteristici) și astfel a complexității de calcul. Pe de altă parte, extragerea de caracteristici folosește o anumită transformare pentru a projecța vectorii de caracteristici inițiali într-un alt spațiu de reprezentare ce oferă o topologie mai eficientă pentru clasificare.
- **alegerea unei măsuri de distanță:** similaritatea vectorilor de caracteristici este cuantificată pe baza unei anumite măsuri de distanță sau funcții de similaritate. Distanța Euclidiană este una dintre măsurile

cele mai des întâlnite, datorită simplității și eficienței acesteia. Alte măsuri uzuale sunt distanța Mahalanobis<sup>4</sup> sau distanța Hausdorff<sup>5</sup>. Această etapă permite definirea noțiuni de "proximitate" între date, astfel că alegerea măsurii de distanță este definitorie pentru performanțele algoritmului de clasificare. Pentru un studiu bibliografic al măsurilor de distanță folosite la clasificarea nesupervizată a datelor, cititorul se poate raporta la lucrarea [Jain 99].

- **clasificarea propriu-zisă:** această etapă permite regruparea datelor în clase omogene în funcție de vectorii de caracteristici asociați acestora. Abordările folosite au fost discutate în paragraful anterior (vezi Secțiunea 7.1).
- **reducerea vectorilor de caracteristici:** această etapă este realizată doar dacă este necesară și constă în reducerea dimensiunii vectorilor de caracteristici, de regulă  $n$ -dimensionali cu  $n > 3$ , la o dimensiune ce poate fi reprezentată grafic, de regulă  $n = 3$ . Procesul de reducere a informației trebuie să conserve pe cât posibil informația utilă conținută în clase. Un exemplu clasic de astfel de metodă este analiza în componente principale sau PCA ("Principal Components Analysis"<sup>6</sup>).
- **evaluarea rezultatelor:** evaluarea pertinenței repartiției în clase este de asemenea o etapă optională. Metodele folosite sunt, fie *matematice* precum evaluarea tendinței de fragmentare a claselor sau analiza validității claselor, fie *empirice* precum analiza manuală a repartiției în clase pe baza vizualizării grafice a acestora.

### 7.1.2 Metodele existente de clasificare nesupervizată

Tehnicile de clasificare în general, au fost studiate intensiv de mai bine de 30 de ani. Diversitatea metodelor existente face dificilă găsirea metodei care să corespundă cel mai bine necesităților de prelucrare ale aplicației vizate.

În cele ce urmează vom realiza o trecere în revistă a metodelor de clasificare nesupervizată cel mai frecvent întâlnite în literatura de specialitate, evidențiind particularitățile și avantajele fiecărei dintre acestea. Astfel, vom prezenta:

---

<sup>4</sup>distanța Mahalanobis a fost introdusă pentru prima oară de statisticianul P.C. Mahalanobis în anul 1936. Aceasta se bazează pe analiza corelației dintre date și este folosită pentru a măsura gradul de separare a două grupuri de date. De exemplu, distanța Mahalanobis dintre doi vectori aleatori de aceeași distribuție,  $x$  și  $y$ , este dată de relația următoare:  $d_M^2(x, y) = (x - y)^T \Sigma^{-1} (x - y)$ , unde  $\Sigma$  reprezintă matricea de covarianță.

<sup>5</sup>vezi explicația de la pagina 172.

<sup>6</sup>vezi explicația de la pagina 123.

- principiul *clasificării ierarhice*,
- algoritmul ”*k-means*”,
- clasificarea pe bază de *grafuri*,
- clasificarea *pe bază de modele*,
- *clasificarea fuzzy*,
- abordările bazate pe *teoria evoluției*.

### Clasificarea ierarhică

Având la dispoziție un set de  $N$  obiecte și o matrice de similaritate (distanță) între acestea, de dimensiune  $N \times N$ , principiul clasificării ierarhice al acestora constă în următorul algoritm [Johnson 67]:

1. în primă fază, fiecare obiect este atribuit unei clase diferite, astfel inițial vom dispune de  $N$  clase. Similaritatea dintre clase va fi dată în acest punct de similaritatea dintre obiectele claselor,
2. se caută cele mai similare două clase ce vor fi fuzionate într-o singură, astfel numărul de clase disponibile devenind  $N - 1$ ,
3. se recalculează similaritatea dintre noua clasă și fiecare dintre celelalte clase existente,
4. se repetă pașii 2 și 3 până când toate obiectele se vor regăsi într-o clasă unică de dimensiune  $N$ .

În funcție de metoda folosită, pasul 3 poate fi realizat în mai multe moduri. Astfel întâlnim: *înlănțuire simplă* (“simple-linkage”), *înlănțuire completă* (“complete-linkage”) și respectiv *înlănțuire medie* (“average-linkage”).

În *înlănțuirea simplă* se consideră ca distanță între două clase distanța minimă dintre membrii unei clase la membrii celeilalte clase. Acest lucru se traduce prin aprecierea similarității între două clase ca fiind dată de cele mai similare două obiecte din aceste două clase.

Opus *înlănțuirii simple*, *înlănțuirea completă*, consideră ca distanță între două clase, distanța cea mai semnificativă dintre membrii unei clase la membrii celeilalte clase.

În cele din urmă, *înlănțuirea medie* consideră distanța dintre două clase ca fiind media distanțelor dintre toți membrii unei clase la toți membrii celeilalte clase. O altă versiune a *înlănțuirii medii* constă în alegerea distanței ca fiind

valoarea mediană a valorilor de distanță obținute. Aceasta se dovedește în cele mai multe situații a fi mai eficientă decât valoarea medie, reducând numărul de clase atipice.

Clasificarea ierarhică este o metodă de tip acumulativ, deoarece clasele sunt fuzionate progresiv. Totuși, clasificarea ierarhică poate fi implementată și ca metodă partitională prin inversarea algoritmului de clasificare. Astfel, se pornește cu o singură clasă ce conține toate obiectele, clasă ce este divizată progresiv în clase mai mici în funcție de disimilaritatea obiectelor.

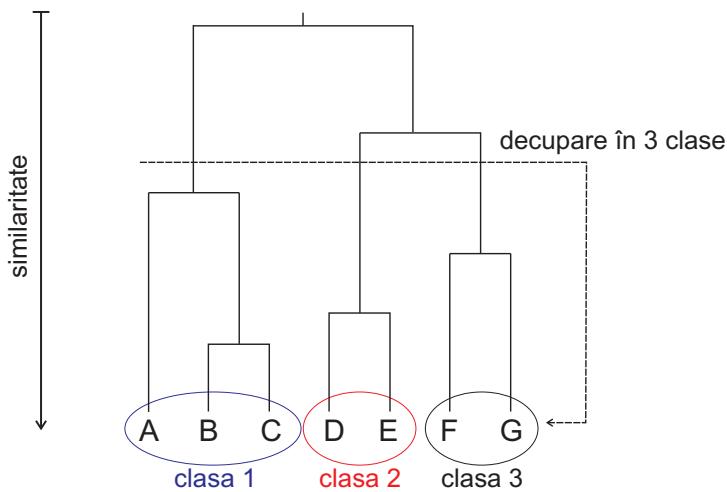


Figura 7.2: Exemplu de clasificare ierarhică, obiectele clasificate sunt notate cu litere de la A la G.

În raport cu celelalte metode existente, principalul avantaj al metodelor de clasificare ierarhică este dat de adaptabilitatea numărului de clase. În urma clasificării ierarhice, clasele obținute pot fi reprezentate sub forma unei dendograme<sup>7</sup>. Un exemplu este prezentat în Figura 7.2. Astfel, în funcție de cerințele aplicației, prin alegerea unui anumit prag de similaritate între clase, utilizatorul poate opta pentru o partiție în clase ce variază de la o singură clasă, până la un număr de clase egal cu numărul de obiecte  $N$ . De notat este faptul că această repartiție variabilă este disponibilă imediat fără a mai fi necesară relansarea algoritmului de clasificare.

### Clasificarea K-means

Clasificarea k-means este una dintre cele mai populare metode de clasificare nesupervizată a datelor [MacQueen 67]. Acest lucru se datorează în mare

<sup>7</sup>vezi explicația de la pagina 70.

parte simplității și eficienței acesteia. Clasificarea k-means este folosită pentru a partaja un anumit număr de obiecte într-un număr de clase,  $k$ , ce este fixat ”a priori”.

Pentru aceasta, în primă fază se vor defini centrele celor  $k$  clase, numite și centroizi. Modul în care sunt aleși centroizii va influența rezultatul clasificării. Cea mai eficientă soluție constă în alegerea acestora cât mai depărtați unii de alții în spațiul de caracteristici.

Următorul pas al clasificării constă în asocierea fiecărui obiect clasei cu centroidul cel mai apropiat de acesta. După ce fiecare obiect a fost asociat unei clase, pozițiile centroizilor vor fi recalculate în funcție de membrii clasei, iar obiectele vor fi realocate în funcție de distanța față de noii centroizi. Pe măsură ce procesul este repetat, centroizii claselor își vor schimba succesiv pozițiile, până în momentul în care nu se va mai produce nici o schimbare. Repartiția în clase astfel obținută corespunde rezultatului final al clasificării.

Având în vedere că rezultatul clasificării depinde de modul în care au fost aleși centroizii inițiali, astfel că seturi diferite de centrozi produc clase diferite, în cele mai multe cazuri la finalul clasificării se calculează o anumită funcție de cost,  $J$ . Clasificarea poate fi repetată de mai multe ori folosind seturi diferite de centroizi, optându-se în final pentru soluția care minimizează valoarea lui  $J$ .

În cazul k-means, funcția  $J$  este aleasă în sensul erorii pătratice, și este definită astfel:

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2 \quad (7.1)$$

unde  $k$  reprezintă numărul de clase,  $n$  reprezintă numărul de obiecte ce trebuie clasificate și  $\|x_i^{(j)} - c_j\|$  reprezintă o măsură de distanță între obiectul  $x_i^{(j)}$  al clasei  $j$  și centroidul  $c_j$  al acesteia. Definită în acest fel, funcția  $J$  este o măsură a distanței globale a obiectelor față de centrul claselor. Cu cât valoarea lui  $J$  este mai importantă, cu atât clasele sunt mai neomogene.

Printre diversele variante ale clasificării k-means putem menționa algoritmul ISODATA [Jain 99]. Algoritmul ISODATA încearcă să rezolve problema numărului fix de clase ce este impus în algoritmul k-means, prin reducerea claselor aşa zise redundante. Astfel că, atunci când unui centroid nu i se asociază un număr suficient de obiecte, acesta este vizat pentru a fi eliminat. În acest fel, numărul de clase obținute va fi mai mult sau mai puțin optimă. Pentru ca selecția să fie eficientă, numărul inițial de clase  $k$  trebuie ales suficient de mare.

Raportat la celelalte metode de clasificare nesupervizată existente, clasificarea k-means este una dintre cele mai eficiente din punct de vedere al complexității de calcul și al calității rezultatelor obținute. De exemplu, pen-

tru o clasificare a 60 de obiecte în 5 clase, marea parte a metodelor enunțate anterior sunt de 500 de ori, până la 2500 de ori, mai lente decât algoritmul k-means [Jain 99]. Din acest motiv, clasificarea k-means și corespondentul acesteia din categoria de metode ce folosesc rețele neuronale (rețelele neuronale Kohonen) au fost folosite cu succes la clasificarea volumelor foarte vaste de date. Celealte tehnici de clasificare mai complexe sunt de regulă folosite pentru volume de date mai restrânse din cauza timpului de execuție ridicat.

### Clasificarea pe bază de grafuri

Din categoria metodelor de clasificare ce folosesc teoria grafurilor<sup>8</sup>, algoritmul cel mai cunoscut este bazat pe construcția arborelui minimal, numit și *"Minimum Spanning Tree"* sau MST.

Având la dispoziție un graf nedirecțional, definim ca arbore al acestuia, subgraful ce conectează toate muchiile<sup>9</sup> grafului. Un anumit graf poate avea astfel mai mulți arbori. De asemenea, fiecarei muchii a grafului i se asociază o anumită pondere ce este aleasă ca fiind proporțională cu costul parcurgerii acesteia în graf. În aceste condiții, un arbore minimal este definit ca fiind arboarele cu ponderea totală cea mai mică sau cel mult egală cu ponderea unui alt arbore posibil, unde ponderea unui arbore este dată de suma ponderilor muchiilor acestuia.

Reprezentând obiectele ce trebuie clasificate sub forma unui graf în spațiul de caracteristici, unde ponderea muchiilor poate fi dată de o anumită măsură de distanță între noduri, atunci putem construi arboarele minimal MST al acestuia. O strategie de clasificare pe baza MST constă mai departe în eliminarea progresivă a segmentelor arborelui cu lungimea cea mai semnificativă. Un exemplu este prezentat în Figura 7.3 [Jain 99], în care am ilustrat arboarele minimal obținut pentru 9 obiecte (notate cu litere de la *A* la *I*), ce sunt reprezentate într-un spațiu de caracteristici bidimensional. Astfel, prin ruperea legăturii cu ponderea cea mai importantă, și anume segmentul *CD* (pondere 6) vom obține două clase, clasa *c*<sub>1</sub> ce conține obiectele {*A, B, C*} și respectiv clasa *c*<sub>2</sub> ce conține obiectele {*D, E, F, G, H, I*}. Mai departe, prin ruperea legăturii *EF* (pondere 4.5), clasa *c*<sub>2</sub> poate fi divizată în două subclase, și anume: *c*<sub>21</sub> = {*D, E*} și *c*<sub>22</sub> = {*F, G, H, I*}. Procesul poate continua astfel până la obținerea granularității (numărului de clase) dorite.

Principalul avantaj al metodelor de clasificare pe bază de grafuri este dat de faptul că o clasificare complexă multidimensională este redusă la o

---

<sup>8</sup>vezi explicația de la pagina 159.

<sup>9</sup>într-un graf, o muchie este dată de o pereche de noduri. Dacă graful este orientat, atunci aceasta se numește arc.

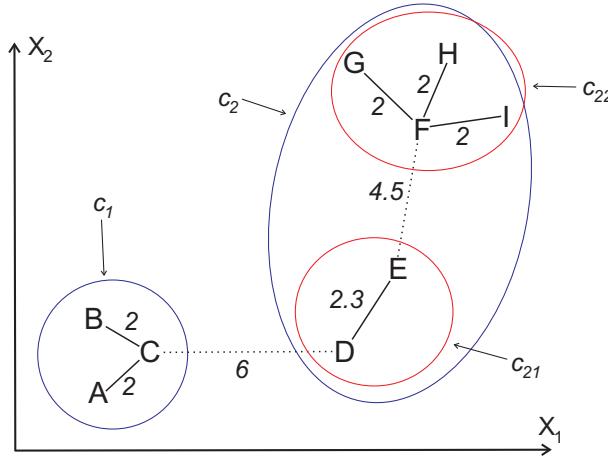


Figura 7.3: Exemplu de clasificare pe baza construcției arborelui minimal MST ( $X_1oX_2$  reprezintă spațiul de caracteristici iar numerele asociate muchiilor reprezintă ponderea acestora).

problemă de partitiorare a unui arbore bidimensional, fară a pierde însă din informațiile esențiale pentru clasificare. Un alt avantaj important este dat de faptul că clasificarea pe bază de grafuri nu depinde de forma geometrică a claselor. În ultimă instanță, aceasta poate fi folosită ca etapă intermediară la implementarea altor tehnici de clasificare mai riguroase, dar care necesită un timp de calcul important.

### Clasificarea pe bază de modele statistice

O altă abordare a problematicii clasificării datelor constă în modelarea statistică a claselor și încercarea de optimizare a modelului obținut, relativ la datele reale. Astfel, clasificarea pe bază de modele folosește reprezentarea matematică a fiecărei clase cu o anumită distribuție parametrică, ca de exemplu, distribuția Gaussiană<sup>10</sup> (continuă) sau Poisson<sup>11</sup> (discretă). În acest fel, întregul set de date va fi modelat de un amestec de astfel de distribuții ce

<sup>10</sup>distribuția Gaussiană sau normală a unei variabile aleatoare  $x$  este dată de funcția de densitate de probabilitate continuă:  $\varphi_{\mu,\sigma^2}(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-(x-\mu)^2/2\sigma^2}$ , unde  $\mu$  reprezintă valoarea medie a lui  $x$  iar  $\sigma^2$  reprezintă varianța.

<sup>11</sup>distribuția Poisson este o distribuție de probabilitate discretă ce exprimă probabilitatea ca o serie de evenimente să aibă loc într-o anumită perioadă de timp, dacă aceste evenimente au o rată de apariție medie cunoscută și apar independent de timpul scurs de la ultimul eveniment. Probabilitatea ca un eveniment să apară de  $k$  ori într-un interval de timp pentru care valoarea medie de apariții este  $\lambda$ , este dată de:  $f(k, \lambda) = \frac{\lambda^k \cdot e^{-\lambda}}{k!}$ .

constituie și modelul datelor.

Principiul matematic este următorul [Fraley 96]: având la dispoziție o populație de date ce conține un număr  $G$  de sub-populații diferite (potențiale clase), atunci densitatea de probabilitate a observației  $p$ -dimensionale  $x$ , din sub-populația de indice  $k$ , este dată de funcția  $f_k(x; \theta)$ , unde  $\theta$  reprezintă un anumit vector de parametri necunoscut. Având la dispoziție observațiile  $x = \{x_1, \dots, x_n\}$  ce reprezintă datele ce trebuie clasificate, se definește vectorul de etichete al claselor ca fiind:

$$\gamma = (\gamma_1, \dots, \gamma_n)^T \quad (7.2)$$

unde  $\gamma_i = k$  dacă  $x_i$  aparține sub-populației de indice  $k$ .

Pentru clasificare, parametrii  $\theta$  și etichetele  $\gamma$  sunt astfel aleși încât să maximizeze o anumită funcție de plauzibilitate ("likelihood function"), și anume:

$$L(x; \theta, \gamma) = \prod_{i=1}^G f_{\gamma_i}(x_i; \theta) \quad (7.3)$$

De cele mai multe ori, ca distribuție este aleasă distribuția normală sau Gaussiană multidimensională, parametrizată de vectorul de medii  $\mu_k$  și matricea de varianță  $\Sigma_k$ . Astfel, vectorul de parametri  $\theta$  va fi dat de parametrii distribuțiilor pentru fiecare sub-populație,  $\theta = (\mu_1, \dots, \mu_G, \Sigma_1, \dots, \Sigma_G)$ .

Înlocuind în relația anterioară, obținem funcția de probabilitate următoare:

$$L(x; \mu_1, \dots, \mu_G, \Sigma_1, \dots, \Sigma_G, \gamma) = \prod_{k=1}^G \prod_{i \in I_k} (2\pi)^{-\frac{p}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right] \quad (7.4)$$

unde  $I_k$  reprezintă setul de indici ce corespund observațiilor din grupul de indice  $k$ ,  $I_k = \{i / \gamma_i = k\}$ .

Simplificări ale funcției de plauzibilitate  $L()$  pot fi obținute prin înlocuirea parametrilor distribuției cu anumiți estimatori, ca de exemplu:

$$\mu_k = \hat{\mu}_k = \frac{\sum_{i \in I_k} x_i}{n_k} \quad (7.5)$$

unde  $n_k$  reprezintă numărul de elemente al lui  $I_k$ , substituție ce conduce la o expresie logaritmică, astfel:

$$L(x; \hat{\mu}_1, \dots, \hat{\mu}_G, \Sigma_1, \dots, \Sigma_G, \gamma) = -\frac{pn \log(2\pi)}{2} - \frac{1}{2} \sum_{k=1}^G [tr(W_k \Sigma_k^{-1}) + n_k \log |\Sigma_k|] \quad (7.6)$$

unde  $W_k = \sum_{i \in I_k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$  iar  $tr(A)$  ("trace") reprezintă operatorul ce returnează suma elementelor de pe diagonala principală a matricei  $A$ .

Totuși, chiar și în această situație, maximizarea funcției  $L()$  poate fi dificilă, având o complexitate de calcul ridicată. Din acest motiv, în practică, se preferă o metodă derivată și anume algoritmul EM - "Expectation-Maximization" [Dempster 77].

Datorită faptului că în momentul clasificării datelor, nu se cunosc exact, nici distribuția acestora și nici parametrii modelului mixt, similar cu algoritmul k-means (vezi pagina 208), algoritmul EM se folosește de iterare. Astfel, clasificarea pornește cu un set de parametri aleși arbitrar. Valorile acestora sunt folosite pentru estimarea probabilităților claselor, iar mai departe, aceste probabilități sunt folosite pentru re-estimarea parametrilor, procesul repetându-se în mod iterativ.

Prima etapă de calcul a probabilităților se numește "expectation", deoarece acestea reprezintă valorile *așteptate* ale claselor. A doua etapă ce constă în estimarea parametrilor modelului poartă numele de "maximization" și constă în *maximarea* funcției de plauzibilitate.

Spre deosebire de algoritmul k-means, în care iterarea se oprește atunci când clasele nu se mai modifică, în cazul algoritmului EM convergența este mai dificilă. Teoretic, algoritmul converge către un punct fix dar practic nu ajunge niciodată acolo. Totuși, se poate estima cât de aproape suntem de convergență prin calcularea unei funcții globale de plauzibilitate pentru ansamblul datelor folosite, prin înmulțirea probabilităților pentru fiecare iterație.

În practică, se preferă calculul logaritmului acestei funcții. Decizia de oprire a iterăției se ia atunci când pe parcursul a  $n$  iterății diferența dintre valorile funcției de plauzibilitate este inferioară valorii  $10^{-d}$ , un exemplu de astfel de valori fiind  $n = 10$  și  $d = 10$  [Witten 05].

Global, metodele de clasificare din această categorie prezintă o serie de avantaje. În primul rând, folosind drept criteriu de convergență maximizarea funcției de plauzibilitate, acesta duce la obținerea de clase compacte ce vor îngloba tipurile de date dominante (furnizate de distribuția componentelor modelului mixt). Un alt avantaj al folosirii modelării statistice a claselor constă în baza teoretică foarte avansată și bine instrumentată disponibilă în acest domeniu, și anume al inferenței statistice. Acest lucru oferă de asemenea și o flexibilitate în alegerea distribuției statistice adaptate datelor ce trebuie clasificate.

Pe de altă parte, convergența unei astfel de abordări este complexă, atât din punct de vedere al modelării matematice, cât și al complexității de calcul.

### Clasificarea fuzzy

Dacă metodele de clasificare tradiționale furnizează o partitioare netă a datelor în clase, astfel, fiecare obiect aparținând unei singure clase, iar clasele fiind disjuncte, *clasificarea fuzzy* se folosește de teoria mulțimilor fuzzy pentru a introduce conceptul de incertitudine (vezi Capitolul 6). În acest fel, problema apartenenței la clase nu mai este o problemă de afirmare sau negare, ci o problemă de grad de apartenență. Fiecare obiect va apartine astfel unei clase într-un anumit grad, ce este furnizat de coeficientul fuzzy.

Principiul este ilustrat în Figura 7.4, în care datele ce trebuie clasificate au fost reprezentate cu cifre de la 1 la 9. Astfel, în cazul unei clasificări nete, clasele obținute sunt de tipul claselor disjuncte  $N_1$  și  $N_2$ , în timp ce clasele obținute cu o clasificare fuzzy,  $F_1$  și  $F_2$ , se întrepătrund.

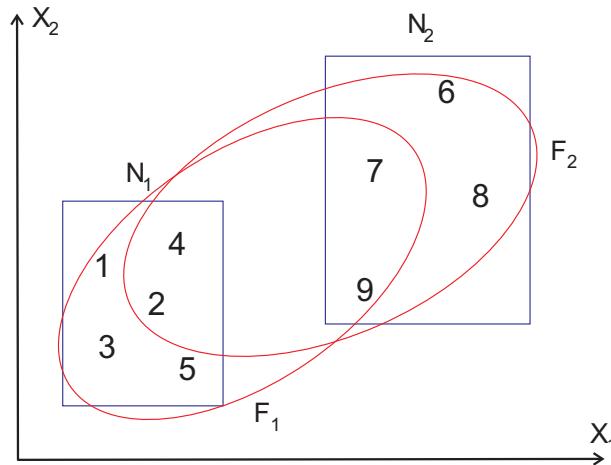


Figura 7.4: Principiul clasificării fuzzy ( $X_1 \circ X_2$  reprezintă spațiul de caracte-

ristici, dreptunghiurile reprezintă clasele nete, în timp ce elipsele marchează clasele fuzzy).

Dintre algoritmii de clasificare fuzzy cei mai populari, putem menționa algoritmul *Fuzzy C-Means* sau FCM [Dunn 73]. Aceasta poate fi considerat ca fiind echivalentul în logică fuzzy al algoritmului k-means (vezi pagina 208). FCM se bazează pe minimizarea unei funcții obiective,  $J_m$ , ce este definită în felul următor:

$$J_m = \sum_{i=1}^N \sum_{j=1}^C u_{ij}^m \|x_i - c_j\|^2 \quad (7.7)$$

unde  $m$  este un număr real,  $m \geq 1$ ,  $u_{ij}$  reprezintă gradul de apartenență al lui  $x_i$  la clasa de indice  $j$ ,  $u_{ij} \in [0; 1]$  (valoarea 1 indică apartenență sigură, iar

valoarea 0 negația sigură),  $x_i$  reprezintă data de indice  $i$  din setul de  $N$  date ce trebuie clasificate,  $c_j$  reprezintă centrul clasei  $j$ ,  $C$  reprezintă numărul de clase iar  $\|\cdot\|$  este o anumită normă ce exprimă similaritatea dintre  $x_i$  și  $c_j$ .

Partiționarea fuzzy este realizată prin optimizarea iterativă a funcției  $J_m$ , timp în care funcțiile de apartenență,  $u_{ij}$ , și centrele claselor,  $c_j$ , sunt recalculate progresiv, astfel:

$$u_{ij} = \left[ \sum_{l=1}^C \left( \frac{\|x_i - c_l\|}{\|x_i - c_j\|} \right)^{\frac{2}{m-1}} \right]^{-1} \quad (7.8)$$

$$c_j = \frac{\sum_{i=1}^N u_{ij}^m \cdot x_i}{\sum_{i=1}^N u_{ij}^m} \quad (7.9)$$

Iterarea se va opri atunci când condiția următoare este satisfăcută:

$$\max_{i,j} \{|u_{ij}^{(k+1)} - u_{ij}^{(k)}|\} < \epsilon \quad (7.10)$$

unde  $k$  reprezintă indicele iterației iar  $\epsilon$  este un anumit prag ales între 0 și 1. În acest fel, algoritmul va converge către un minim sau o valoare fixă a funcției  $J_m$ .

Astfel, algoritmul de clasificare FCM poate fi sintetizat în felul următor:

1. în primă fază este inițializată matricea  $U = U^{(0)}$ , matrice ce conține gradele de apartenență fuzzy la clase,  $U = [u_{ij}]$ , iar indicele iterației devine  $k = 0$ ,
2. pe baza matricei  $U^{(k)}$ , se calculează centrele claselor,  $C^{(k)} = [c_j]$ , folosind relația 7.9,
3. mai departe, matricea de la iterația  $k + 1$ ,  $U^{(k+1)}$ , este calculată prin reactualizarea valorilor matricei  $U^{(k)}$  pe baza relației 7.8,
4. se testează condiția de oprire și anume dacă  $\|U^{(k+1)} - U^{(k)}\| < \epsilon$ , caz în care clasificarea se oprește. În caz contrar, algoritmul trece la iterația următoare,  $k \leftarrow k + 1$ , și acesta se reia de la punctul 2.

Raportat la algoritmul echivalent în logică booleană, și anume algoritmul k-means, algoritmul FCM converge rapid către un minim local. Principalul dezavantaj al clasificării fuzzy constă totuși în dificultatea de alegere adecvată a funcțiilor de apartenență fuzzy, alegere care de cele mai multe ori este determinantă pentru calitatea clasificării.

### Clasificarea bazată pe teoria evoluției

Abordările bazate pe teoria evoluției sunt inspirate de *principiul evoluției naturale a lumii*<sup>vii</sup>. Principiul care stă la baza acestora, constă în aplicarea, asupra unei anumite populații de soluții, a o serie de operatori genetici specifici. În acest fel, populația tinde să evolueze spre o soluție global optimală de partizionare a datelor [Jain 99].

Potențialele soluții ale problemei de clasificare sunt codate folosind cromozomi<sup>12</sup>. Dintre operatorii genetici cei mai frecvent utilizați, putem menționa: *selecția, recombinarea și mutația*. Fiecare dintre aceștia transformă unul sau mai mulți cromozomi de intrare în unul sau mai mulți cromozomi de ieșire. De asemenea, pentru fiecare cromozom se evaluează o anumită funcție obiectivă ce determină probabilitatea de supraviețuire a acestuia în următoarea generație.

Astfel, un algoritm generic de clasificare pe baza teoriei evoluției poate fi enunțat în felul următor [Jain 99]:

1. în primă fază se alege aleator o populație de soluții, fiecare dintre soluții corespunzând unei anumite partiții valide a datelor. Fiecarei soluții i se asociază o funcție obiectivă, de regulă aleasă în aşa fel încât să fie invers proporțională cu eroarea pătratică. Astfel, o soluție cu o eroare pătratică mică, va corespunde unei valori ridicate a funcției obiective,
2. folosind operatorii de selecție, recombinare și mutație, se generează o nouă populație de soluții. Se re-evaluează funcția obiectivă a acestei noi populații,
3. se repetă pasul 2 până când este îndeplinit un anumit criteriu de convergență.

Dintre tehniciile bazate pe teoria evoluției cele mai cunoscute, putem enumera *algoritmii genetici* (GAs), *strategiile evoluționiste* (ESs) și *programarea evoluționistă* (EP) [Jain 99]. Dintre acestea, algoritmii genetici sunt folosiți cel mai frecvent în probleme de clasificare. În GAs, soluțiile sunt exprimate binar iar propagarea de la o generație la alta este realizată pe baza funcției obiective prin intermediul operatorului de selecție. Selecția se folosește de o strategie probabilistică pentru a desemna soluțiile cu o valoare semnificativă a funcției obiective ca soluții cu o probabilitate importantă de reproducere. Operatorul de mutație este folosit în GAs doar pentru a se asigura că întreg

---

<sup>12</sup>cromozomul, în limba greacă chromo-nuanță și soma-obiect, reprezintă macromolecule de ADN care conțin mai multe gene și secvențe nucleotide cu rol în păstrarea informației ereditare a celulei.

spațiul soluțiilor este explorat. Spațiul soluțiilor este explorat folosind operatorul de recombinare ("crossover"). Pe de altă parte, ESs și EP, diferă de GAs, prin modul de reprezentare al soluțiilor și prin operatorii de mutație folosiți. De exemplu, în EP nu se folosește operatorul de recombinare, ci doar cei de mutație și selecție.

Global, abordările bazate pe teoria evoluției sunt tehnici de căutare a unei soluții globale a problemei clasificării, ceea ce le diferențiază de marea parte a celorlalte tehnici existente ce propun de regulă o căutare locală a soluției. Mai mult, abordările bazate pe teoria evoluției sunt capabile să găsească soluția optimală chiar și în cazul în care funcția asociată criteriori de selecție este discontinuă.

## 7.2 Clasificarea supervizată

După cum am menționat și în partea introductivă a acestui capitol, diferența dintre clasificarea nesupervizată și cea supervizată constă în faptul că aceasta din urmă se folosește de o serie de exemple de clasificări (date de antrenare), de regulă obținute în urma expertizei domeniului de aplicație<sup>13</sup>, pentru a identifica clasele ce prezintă o densitate de probabilitate importantă relativ la o singură clasă de referință. Mai mult, clasificarea supervizată, tinde să păstreze un număr de clase relativ scăzut [Eick 04].

Din punct de vedere al modului în care obiectele sunt asociate claselor, în clasificarea supervizată se folosește noțiunea de "apropiere" în sensul unei anumite măsuri de distanță, spre deosebire de clasificarea nesupervizată în care de regulă se calculează o anumită eroare de clasificare ce minimizează distanța dintre obiectele din interiorul unei clase.

Pentru a înțelege mai bine diferența dintre cele două tipuri de abordări, în Figura 7.5 am ilustrat exemplul următor [Eick 04]: să presupunem că obiectele ce trebuie clasificate reprezintă subspecii ale plantei din familia Iris, și anume Setosa (culoare neagră) și Virginica (culoare albă). În urma unei clasificări nesupervizate, este foarte probabil ca obiectele să fie împărțite în patru clase, aşa cum am ilustrat în Figura 7.5.a. Acest lucru se datorează faptului că, clasificarea nesupervizată ține cont în primul rând de minimizarea distanței obiectelor în spațiul de caracteristici. Astfel, clasificarea obținută nu este foarte interesantă din punct de vedere obiectiv deoarece clasele obținute, fie combină cele două subspecii, fie separă aceeași specie în clase diferite. De

---

<sup>13</sup>datele de antrenare pot fi văzute în sensul unei realități de teren ("groundtruth") folosite de regulă pentru validarea rezultatelor anumitor algoritmi de detecție (vezi și explicația de la pagina 170).

exemplu, clasa A conține atât subspecia Virginica cât și Setosa, iar clasele B și C conțin aceeași subspecie Virginica.

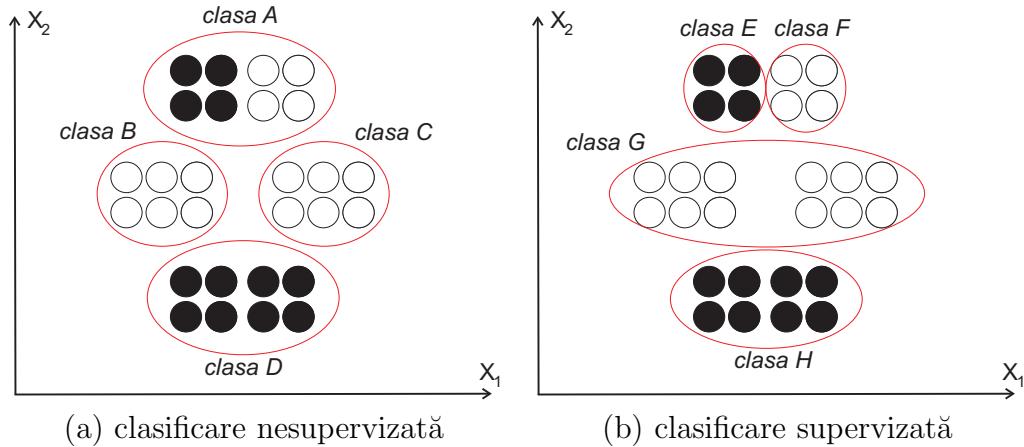


Figura 7.5: Diferența dintre clasificarea nesupervizată și clasificarea supervizată ( $X_2 \times X_1$  reprezintă spațiul de caracteristici, obiectele de același tip sunt figurate cu aceeași culoare) [Eick 04].

Pe de altă parte, o metodă de clasificare supervizată ce are ca scop maximizarea "puritatei" claselor, va diviza clasa A obținută anterior în două clase, și anume: clasa E și respectiv F (vezi Figura 7.5.b). De asemenea, în scopul reducerii numărului de clase fără însă a compromite omogenitatea acestora, clasele B și C obținute anterior vor fuziona într-o singură clasă G.

Global, clasele obținute vor fi pe cât posibil omogene, conținând obiecte de același tip, lucru ce se datorează în principal antrenării prealabile a clasificatorului cu exemple din fiecare clasă.

### 7.2.1 Etapele clasificării supervizate

În general, un clasificator supervizat implică introducerea noțiunii de învățare sau "machine learning". Aceasta presupune adaptarea algoritmului la domeniul de aplicație și implică mai multe etape de prelucrare. În Figura 7.6 am sintetizat diagrama de funcționare generală a unui clasificator supervizat [Kotsiantis 07].

Astfel, prima etapă constă în identificarea setului de date ce trebuie clasificat. Tot în această etapă sunt alese și attributele ce vor fi folosite pentru definirea spațiului de caracteristici în care se va efectua clasificarea propriu-zisă. Acestea trebuie alese în aşa fel încât să fie cât mai reprezentative pentru criteriile de selecție vizate.

A doua etapă constă în pre-prelucrarea acestor date. De regulă operațiile efectuate constau în reducerea zgomotului existent în date cât și simplificarea procesului de manevrare a seturilor voluminoase de date de antrenare ("instance selection"), sau reducerea dimensiunii setului de date prin eliminarea datelor ce sunt irelevante sau redundante pentru clasificare ("feature subset selection").

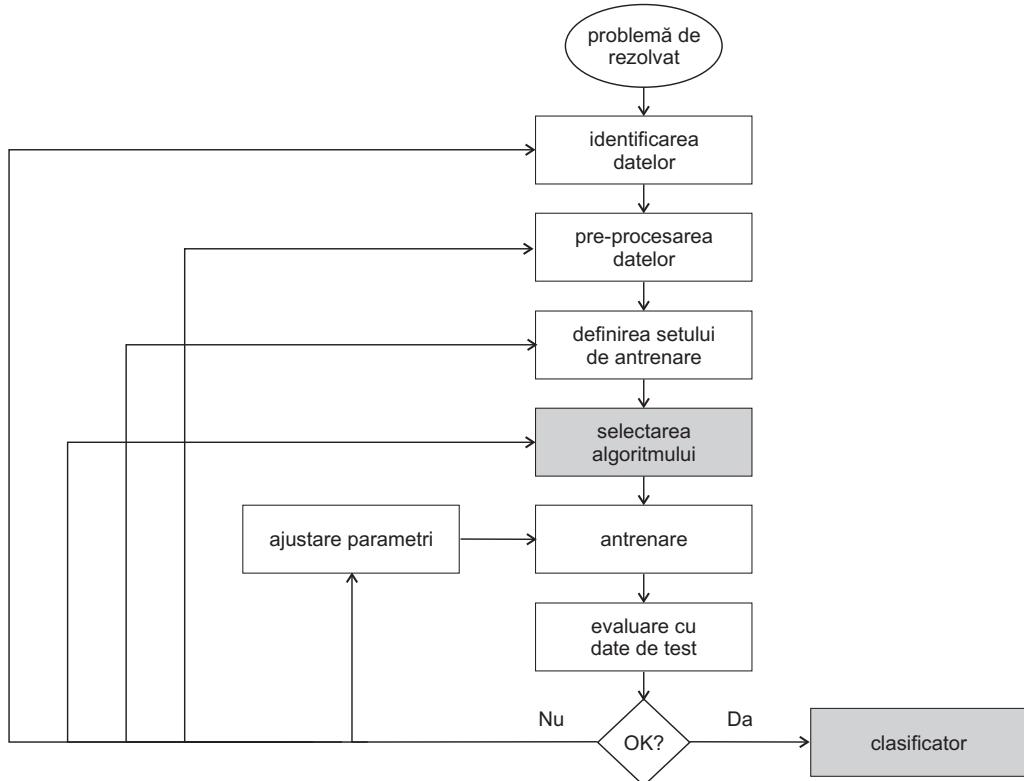


Figura 7.6: Procesul de clasificare supervizată [Kotsiantis 07].

După ce datele de antrenare au fost definite, acestea fiind de regulă disponibile "a priori" clasificării, se trece la alegerea algoritmului de clasificare supervizată ce corespunde cel mai bine aplicației vizate. Acest pas este foarte important deoarece influențează decisiv rezultatele clasificării, algoritmi diferenții furnizând rezultate uneori foarte diferite. Din această cauză se preferă testarea rezultatelor clasificării folosind setul de antrenare. În momentul în care clasificatorul testat este considerat ca fiind satisfăcător pentru aplicație, de abia atunci se poate trece la folosirea acestuia cu datele de intrare ne-etichetate (datele propriu-zise).

Evaluarea clasificatorului constă de regulă în estimarea preciziei predictiei,

aceasta fiind definită ca raportul dintre numărul de predicții corecte și numărul total de predicții<sup>14</sup>.

Dintre metodele existente, menționăm trei dintre cele mai frecvent folosite. O primă abordare constă în divizarea setului de date de antrenare în trei părți, dintre care două sunt folosite pentru antrenare iar una dintre acestea pentru estimarea preciziei predicției clasificatorului. O altă abordare, cunoscută și sub numele de "validare încrucișată" ("cross-validation"), constă în împărțirea setului de antrenare în subseturi disjuncte și egale. Pentru fiecare subset, clasificatorul va fi antrenat folosind setul format prin reuniunea tuturor celorlalte subseturi existente, iar evaluarea se va realiza pe acesta. Eroarea clasificatorului va fi dată de media tuturor erorilor obținute pentru fiecare subset în parte. Un caz particular al validării încrucișate este algoritmul "leave-one-out". În acest caz, fiecare subset de test va fi constituit dintr-o singură instanță. Această metodă este cea mai complexă din punct de vedere al calculelor, dar pe de altă parte și cea mai precisă din punct de vedere al estimării erorii de predicție a clasificatorului.

În cazul în care eroarea de predicție nu este suficient de mică, atunci etapele anterioare trebuie revizuate. Mai mulți factori pot fi responsabili pentru aceasta, ca de exemplu: este posibil să nu fi ales attributele cele mai relevante, setul de antrenare să fie prea redus, complexitatea problemei să fie prea ridicată, clasificatorul ales să nu fie cel mai adaptat problemei vizate, parametrii acestuia să fie reglați incorrect, și așa mai departe.

### 7.2.2 Metodele existente de clasificare supervizată

Clasificarea supervizată este realizată cu predilecție de ceea ce numim Sisteme Inteligente. Dintre tehniciile existente, ne vom focaliza în continuare pe prezentarea celor mai semnificative dintre acestea, și anume:

- *algoritmi logici* (arbori de decizie, seturi de reguli),
- *algoritmi bazați pe perceptroni* (uni-strat, rețele neuronale),
- *algoritmi statistici* (clasificatori Bayes, rețele Bayes),
- *algoritmi bazați pe instanțe* (kNN),
- algoritmi de tip "*Support Vector Machine*".

---

<sup>14</sup>termenul de predicție este folosit în acest caz pentru a desemna modul în care clasificatorul etichetează datele de antrenare. Dacă etichetele asociate de clasificator corespund etichetelor deja existente, atunci predicția este corectă. În caz contrar, aceasta este eronată.

### Algoritmi logici

În această categorie de metode de clasificare supervizată putem menționa construcția *arborilor de decizie* și a *regulilor de decizie*.

Un **arbore de decizie** ("decision tree") [Yuan 95] este un arbore ce clasifică anumite instanțe ale datelor de intrare pe baza trierii lor în funcție de valorile atributelor acestora. Fiecare nod al arborelui reprezintă un atribut al unei anumite instanțe ce trebuie clasificată. Fiecare ramură a arborelui reprezintă o anumită valoare pe care nodul și-o poate atribui. Instanțele datelor de intrare sunt clasificate pornind de la rădăcina arborelui, fiind sortate progresiv în funcție de valorile atributelor. În Figura 7.7 am exemplificat arboarele de decizie asociat datelor de antrenare din Tabelul 7.1 [Kotsiantis 07].

Atributul 1	Atributul 2	Atributul 3	Atributul 4	Clasa nr.
<i>A</i> 1	<i>A</i> 2	<i>A</i> 3	<i>A</i> 4	1
<i>A</i> 1	<i>A</i> 2	<i>A</i> 3	<i>B</i> 4	1
<i>A</i> 1	<i>B</i> 2	<i>A</i> 3	<i>A</i> 4	2
<i>A</i> 1	<i>B</i> 2	<i>B</i> 3	<i>B</i> 4	3
<i>A</i> 1	<i>C</i> 2	<i>A</i> 3	<i>A</i> 4	4
<i>A</i> 1	<i>C</i> 2	<i>A</i> 3	<i>B</i> 4	5
<i>B</i> 1	<i>B</i> 2	<i>B</i> 3	<i>B</i> 4	2
<i>C</i> 1	<i>B</i> 2	<i>B</i> 3	<i>B</i> 4	6

Tabelul 7.1: Exemplu de date de antrenare în 6 clase (valorile atributelor sunt reprezentate cu litere majuscule de la A la C).

Astfel, instanța  $\{Atributul\ 1 = A1, Atributul\ 2 = B2, Atributul\ 3 = A3, Atributul\ 4 = B4\}$  va parcurge arboarele prin nodurile *Atributul 1*, *Atributul 2* și respectiv *Atributul 3*, ceea ce duce la clasificarea instanței ca apartinând clasei 2.

Una dintre problemele unei astfel de abordări de clasificare este dată de adaptarea arborelui strict la datele de antrenare, ceea ce va duce la o posibilă clasificare eronată a datelor reale. Acest fenomen poartă numele de "overfitting". Fiind dat un spațiu de ipoteze  $H$ , o ipoteză  $h \in H$  (de exemplu, un arbore de decizie) este considerată că se "supra-adaptează" datelor de antrenare dacă există o altă ipoteză alternativă  $h' \in H$ , iar  $h$  obține o eroare inferioară lui  $h'$  pe datele de antrenare, dar  $h'$  duce la o eroare mai mică decât  $h$  pentru întreaga distribuție de instanțe de intrare [Venkataraman 99].

De regulă metodele folosite pentru combaterea acestui fenomen constau, fie în *reducerea algoritmului de antrenare* prin folosirea doar a unui subset de date de antrenare, sau folosirea în totalitate a acestora și "rafinarea"

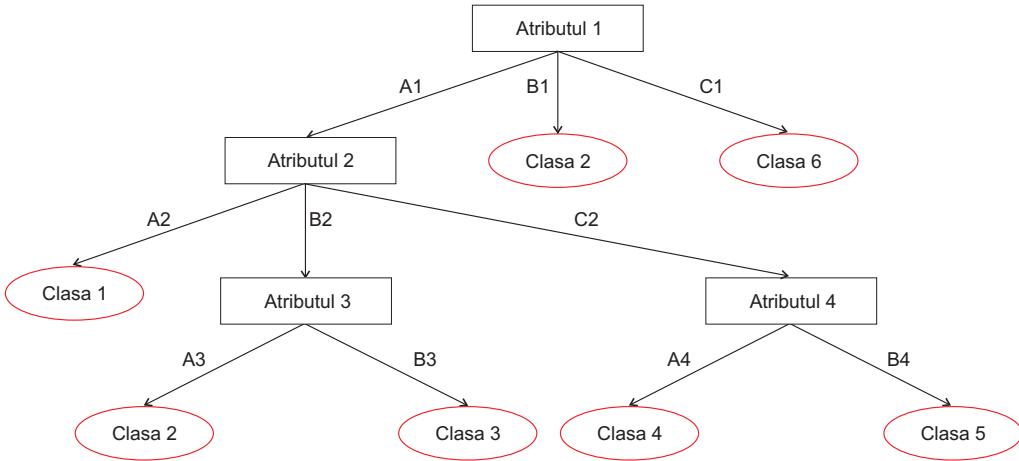


Figura 7.7: Arborele de decizie asociat datelor de antrenare din Tabelul 7.1.

*ulterioră a arborelui* de decizie obținut prin eliminarea unui ramur sau frunze<sup>15</sup> ("pruning"). De regulă, dacă doi arbori de decizie diferă, testați pe aceleași date, obțin aceeași precizie a predicției, atunci se preferă arborele ce conține cele mai puține frunze.

Spre deosebire de alte metode, din punct de vedere al complexității calculelor, arborii de decizie au avantajul de a fi în general uni-variabilă, deoarece divizarea ramurilor este realizată prin analiza unui singur atribut. Totuși, din această cauză, arborii de decizie se limitează la a furniza o partitioare a spațiului de intrare ortogonală pe una dintre axele spațiului de caracteristici și paralelă cu toate celelalte. Astfel, clasele rezultate în urma partitioarei sunt hiperdreptunghiuri<sup>16</sup>.

O altă metodă de clasificare supervizată din această categorie sunt **regulile de decizie**. Acestea pot fi obținute indirect pe baza arborilor de decizie, prin asocierea unei reguli pentru fiecare cale de parcurs a arborelui de la rădăcină până la o frunză, sau pot fi induse direct pe baza datelor de antrenare (vezi studiul prezentat în [Furnkraz 99]).

### Algoritmi bazați pe perceptroni

Perceptronul [Rosenblatt 62] poate fi considerat ca fiind cel mai simplu exemplu de rețea neuronală. Aceasta este un clasificator liniar. În versiunea simplificată, uni-strat, perceptronul nu are decât o singură ieșire la care sunt

<sup>15</sup>termenul de "frunză" a unui arbore este folosit pentru a desemna un nod terminal.

<sup>16</sup>în geometrie, un hiperdreptunghi reprezintă generalizarea noțiunii de dreptunghi pentru dimensiuni de ordin superior lui 2, fiind definit ca un produs Cartezian de intervale.

conectate toate intrările (vezi Figura 7.8).

Principiul de funcționare este următorul: dacă  $x_i$ , cu  $i = 1, \dots, n$ , reprezintă valorile atributelor de intrare, iar  $w_i$  sunt niște ponderi asociate conexiunilor de intrare (de regulă alese ca fiind numere reale în intervalul  $[-1; 1]$ , numite și vector de predicție), atunci perceptronul va calcula suma ponderată a intrărilor:

$$y = \sum_{i=1}^n x_i \cdot w_i \quad (7.11)$$

iar ieșirea acestuia va fi dată de valoarea binară ce rezultă din filtrarea cu un anumit prag  $\tau$  a valorii  $y$  obținute. Astfel, dacă  $y > \tau$  valoarea de ieșire este 1, iar în caz contrar se returnează valoarea 0.

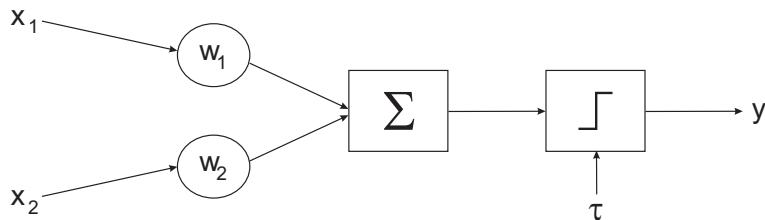


Figura 7.8: Exemplu de perceptron uni-strat cu două intrări ( $x$  reprezintă intrările,  $w$  ponderile acestora,  $\tau$  pragul folosit iar  $y$  ieșirea perceptronului).

Modalitatea cea mai simplă de antrenare a perceptronului pentru un anumit set de date de antrenare, constă în rularea repetitivă a algoritmului prezentat anterior pentru acest set de date, până când se obține un vector de predicție care este corect pentru toate datele de antrenare.

Un alt algoritm de antrenare mai elaborat este metoda WINNOW. Aceasta modifică ponderile conexiunilor în felul următor: dacă valoarea predicției la un moment dat este  $y' = 0$  iar valoarea actuală este  $y = 1$ , atunci înseamnă că ponderile sunt valori prea mici. În acest caz, ponderile vor fi mărite folosind următoarea relație:

$$w_i = \alpha \cdot w_i \quad (7.12)$$

unde  $\alpha > 1$  reprezintă parametrul de creștere a ponderii ("promotion parameter").

Dacă valoarea predicției este  $y' = 1$  iar valoarea actuală este  $y = 0$ , atunci înseamnă că ponderile au valori prea ridicate. În acest caz, se recurge la diminuarea acestora, astfel:

$$w_i = \beta \cdot w_i \quad (7.13)$$

unde  $0 < \beta < 1$  reprezintă parametrul de diminuare a ponderii ("demotion parameter").

Perceptroni sunt clasificatori binari ce pot fi folosiți pentru a clasifica doar seturi liniar separabile de date de intrare. Astfel, aceștia oferă doar o partajare liniară a spațiului de intrare. În cazul în care datele de intrare nu sunt liniar separabile, se poate opta pentru perceptroni multi-strat sau ceea ce numim rețele neuronale artificiale (ANN - Artificial Neural Networks) [Rumelhart 86].

O rețea neuronală este constituită dintr-un număr mare de unități de prelucrare, numite și neuroni, ce sunt interconectate pe straturi pentru a forma o anumită structură complexă de prelucrare. Un exemplu este ilustrat în Figura 7.9. De regulă, neuroni unei rețele neuronale sunt de trei tipuri:

- *neuroni de intrare* ce primesc informația ce trebuie prelucrată,
- *neuroni ascunși* ce efectuează anumite etape intermediare de prelucrare și înregistrează datele,
- *neuroni de ieșire* ce furnizează rezultatele obținute.

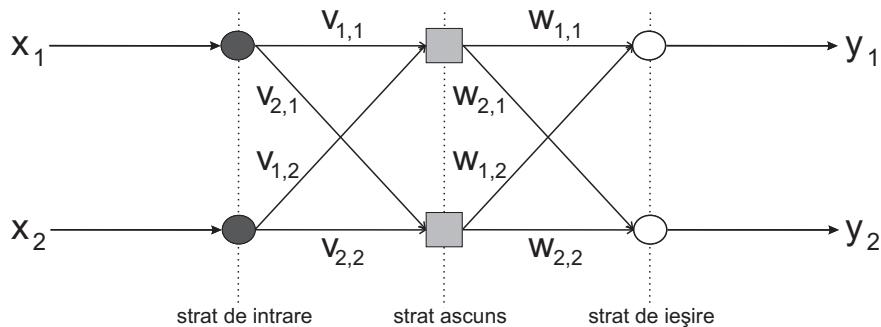


Figura 7.9: Exemplu de rețea neuronală cu două intrări și două ieșiri ( $x$  reprezintă intrările,  $v$  și  $w$  reprezintă ponderile conexiunilor iar  $y$  reprezintă ieșirile).

Un neuron primește o serie de valori de intrare, fie de la o sursă externă, fie de la un alt neuron. Ca și în cazul perceptronului, fiecare valoare de intrare are asociată o pondere care poate fi modificată astfel încât acesta să se adapteze modelului de învățare. Valorile de ieșire ale neuronului pot să devină valori de intrare pentru alți neuroni.

Dacă considerăm că  $x_j$ , cu  $j = 1, \dots, m$ , reprezintă valorile de intrare ale neuronului  $i$ , iar  $w_{ij}$  reprezintă ponderea legăturii de la neuronul  $j$  la neuronul

$i$ , atunci ieșirea neuronului  $i$  este dată de ecuația următoare:

$$y_i = f\left(\sum_{j=1}^m w_{ij} \cdot x_j\right) \quad (7.14)$$

unde  $f()$  reprezintă funcția de activare a neuronului iar suma  $\sum_{j=1}^m w_{ij} \cdot x_j$  poartă numele de intrare a rețelei și este de regulă notată cu  $net_i$ , astfel:

$$y_i = f(net_i) \quad (7.15)$$

Antrenarea unei rețele neuronale se realizează simplu prin folosirea datelor de antrenare pentru a modifica ponderile rețelei. Datele de antrenare sunt de regulă perechi de intrări și ieșiri ale rețelei. Ponderile obținute sunt stocate și folosite apoi pentru clasificarea datelor reale.

De menționat este faptul că rețelele neuronale pot exista și în varianta nesupervizată, cum ar fi de exemplu rețelele cu auto-organizare de tip Kohonen. Dintre rețelele neuronale existente, menționăm pe cele mai importante, și anume: rețelele RBF (Radial Basis Function) în care funcțiile de activare sunt funcții radiale, rețelele recurente RNN (Recurrent Neural Networks) în care conexiunile dintre neuroni sunt circulare, rețelele neurale stohastice ce introduc variații aleatoare în rețea sau rețelele neuronale modulare.

Global, ca metode de clasificare, rețelele neuronale folosesc doar date de intrare numerice, ceea ce restricționează aplicabilitatea acestora doar la clasificarea obiectelor ce pot fi reprezentate numeric. Mai mult, varietatea importantă de rețele existente, precum și numărul important de parametri ce trebuie reglați (de exemplu: numărul de straturi, numărul de neuroni, funcțiile de activare, etc.), fac dificilă alegerea metodei care să fie cea mai adaptată problemei de clasificare vizată. Totuși, un avantaj important al acestora îl constituie arhitectura pur paralelă ce permite procesarea simultană a unui volum important de date (de exemplu, pentru clasificarea conținutului unei imagini, se poate asocia un neuron fiecărui pixel al acesteia, calculele realizându-se astfel simultan).

### Algoritmi statistici

Spre deosebire de rețelele neuronale, abordările statistice se evidențiază prin modelarea repartiției în clase pe baza unui model statistic. Astfel, datelor de intrare li se asociază o anumită probabilitate de apartenență la fiecare dintre clase, și nu doar unei singure clase ca în cazul metodelor prezentate anterior.

Dintre metodele statistice cele mai cunoscute putem menționa *rețelele Bayes*. În varianta cea mai simplă sunt clasificatorii NB sau "Naive Bayes Classifiers". Un clasificator NB este un clasificator probabilistic simplu ce se

bazează în mod ”naiv” pe ipoteza de independentă a datelor. Cu alte cuvinte, un clasificator NB consideră că prezența (sau absența) unui anumit atribut al unei clase este total independentă de prezența (sau absența) oricărui alt atribut.

Modelul probabilistic folosit de clasificatorii Bayes este un model condițional. Acesta poate fi exprimat pe baza probabilității ”a posteriori”,  $P(C|F_1, \dots, F_n)$ , unde  $C$  reprezintă variabila clasă iar  $F_i$ , cu  $i = 1, \dots, n$ , reprezintă attributele claselor.

Pe baza teoremei Bayes, această probabilitate poate fi rescrisă în felul următor:

$$P(C|F_1, \dots, F_n) = \frac{P(C)P(F_1, \dots, F_n|C)}{P(F_1, \dots, F_n)} \quad (7.16)$$

unde  $P(C)$  reprezintă probabilitatea ”a priori”,  $P(F_1, \dots, F_n|C)$  reprezintă verosimilitatea (”likelihood”) iar  $P(F_1, \dots, F_n)$  reprezintă evidența.

Mai departe, expresia  $P(C)P(F_1, \dots, F_n|C)$  poate fi rescrisă în felul următor:

$$\begin{aligned} P(C)P(F_1, \dots, F_n|C) &= P(C)P(F_1|C)P(F_2|C, F_1)P(F_3|C, F_1, F_2) \cdot \\ &\quad \dots \cdot P(F_n|C, F_1, \dots, F_{n-1}) \end{aligned} \quad (7.17)$$

Folosind ipoteza de independentă a atributelor  $F_i$  obținem:

$$P(F_i|C, F_j) = P(F_i|C) \quad (7.18)$$

și mai departe:

$$P(C|F_1, \dots, F_n) = \frac{1}{Z}P(C) \prod_{i=1}^n P(F_i|C) \quad (7.19)$$

unde  $Z = P(F_1, \dots, F_n)$  reprezintă un factor de scală ce depinde doar de attributele  $F_i$ , fiind constant dacă valorile acestora sunt cunoscute.

Pentru a transforma modelul probabilistic prezentat în ecuația 7.19 într-un clasificator, se vor adăuga anumite reguli de decizie. Una dintre regulile cel mai frecvent folosite este dată de alegerea ipotezei cea mai probabilă, abordare cunoscută și sub numele de ”maximum a posteriori” sau MAP.

Clasificatorul obținut în acest fel este descris de funcția următoare:

$$Clasif(f_1, \dots, f_n) = argmax_c P(C = c) \prod_{i=1}^n P(F_i = f_i|C = c) \quad (7.20)$$

unde  $Clasif()$  reprezintă funcția de clasificare a datelor,  $f_i$ , cu  $i = 1, \dots, n$ , reprezintă valorile atributelor  $F_i$  iar  $c$  valorile variabilei clasă  $C$ . Astfel, datele

de intrare reprezentate prin valorile  $f_i$  ale atributelor, vor fi clasate în clasa  $c$ , pentru care probabilitatea dată de expresia 7.19 este maximă.

O abordare mai complexă o constituie *rețelele Bayes* sau BN ("Bayesian Networks"). O rețea Bayes este un model grafic de reprezentare a relațiilor probabilistice ce pot exista între valorile unui set de atrbute. O astfel de rețea este structurată sub forma unui graf orientat aciclic sau DAG ("Directed Acyclic Graph"<sup>17</sup>) în care nodurile corespund atrbutei  $X$ , pe baza cărora se realizează clasificarea. Un exemplu este ilustrat în Figura 7.10.

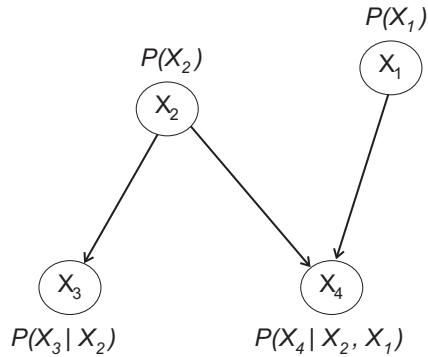


Figura 7.10: Exemplu de rețea Bayes cu patru noduri ( $X$  reprezintă atrbutele iar  $P(X_i|X_j)$  probabilitățile condiționale).

Muchiile grafului corespund relațiilor cauzale ce pot exista între atrbute, în timp ce lipsa muchiilor indică independență statistică dintre atrbute. Dacă există o conexiune de la nodul  $A$  la nodul  $B$ , atunci nodul  $A$  este considerat părintele nodului  $B$  iar acesta din urmă copilul nodului  $A$ . Astfel, un graf orientat aciclic este o rețea Bayes dacă probabilitatea nodurilor poate fi exprimată ca un produs de probabilități locale ale fiecărui nod și a părintilor acestuia, astfel:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Părinti}(X_i)) \quad (7.21)$$

unde funcția  $\text{Părinti}(X_i)$  returnează părintii nodului  $X_i$  [Heckerman 96].

Pentru a răspunde unei problemele de clasificare supervizată, o rețea Bayes trebuie mai întâi antrenată. În cazul cel mai simplu, un expert în domeniu va specifica parametrii rețelei astfel încât aceasta să fie adaptată

<sup>17</sup>un graf orientat aciclic este un graf orientat (sensul de parcurgere este luat în calcul) în care nu există bucle. Cu alte cuvinte, pentru orice nod  $v$  nu există o cale nevidă prin care se pornește din  $v$  și se ajunge tot în  $v$ .

cerințelor de clasificare dorite. În cazul în care acest lucru nu este posibil sau dacă rețeaua este mult prea complexă pentru a fi concepută manual, parametrii acesteia precum și distribuțiile locale trebuie extrase automat din datele de antrenare. Problema care se pune este de a găsi rețea Bayes care modeleză cel mai bine datele de antrenare, și anume corelațiile de referință dintre atributele de intrare și respectiv de ieșire. Pentru o descriere detaliată a metodelor de învățare folosite pentru rețelele Bayes, cititorul se poate raporta la lucrarea [Niculescu 06].

Clasificarea unui set de date de intrare reprezentate prin mulțimea de atrbute  $F_1, \dots, F_n$  se traduce într-o rețea Bayes prin clasarea acestora în clasa de etichetă  $c$  care maximizează probabilitatea "a posteriori", și anume  $P(c|F_1, \dots, F_n)$  (vezi ecuația 7.16).

Principalul avantaj al rețelelor Bayes, relativ la celelalte metode de clasificare (de exemplu: rețelele neuronale sau arborii de decizie), constă în posibilitatea acestora de a se modela sau adapta problemei de clasificare prin luarea în calcul a tuturor informațiilor disponibile "a priori" despre aceasta. Astfel, conexiunile dintre noduri pot fi modificate pentru a lua în calcul eventualele ipoteze adiționale. Pe de altă parte, în cazul în care setul de atrbute este complex și conține un număr considerabil de atrbute, folosirea unui clasificator pe bază de rețele Bayes nu este eficientă datorită volumului rețelei, cât și a timpului de calcul. De asemenea, trebuie luat în calcul faptul că într-o rețea Bayes de cele mai multe ori valorile atrbutelor de intrare trebuie eșantionate [Kotsiantis 07].

### **Algoritmi bazați pe instanțe**

Această categorie de metode de clasificare supervizată intră tot sub incidența metodelor de clasificare statistică. Totuși, metodele bazate pe instanțe sunt algoritmi de învățare "leneș" ("lazy-learning algorithms") întrucât aceștia întârzie procesul de inducție și de generalizare până în momentul în care se realizează clasificarea. Din această cauză, etapa de învățare a unui algoritm "leneș" are de regulă o complexitate de calcul mai redusă decât în cazul algoritmilor "intensivi" ("eager-learning algorithms", precum arborii de decizie sau rețelele Bayes). Pe de altă parte, procesul de clasificare se dovedește a fi mai complex.

Dintre metodele cele mai populare, putem menționa metoda de clasificare de tip "cel mai apropiat vecin" sau k-NN ("k-Nearest Neighbor") [Fix 51]. Aceasta se bazează pe principiul că instanțele dintr-un anumit set de date (valorile atrbutelor) se vor găsi în general în proximitatea altor instanțe cu proprietăți similare. Astfel, dacă o instanță este clasată ca aparținând unei anumite clase, atunci o altă instanță, ce nu a fost încă clasificată, poate fi

atribuită unei clase doar prin observarea clasei din care face parte cel mai apropiat vecin al acesteia<sup>18</sup>. Mai general, metoda k-NN, analizează cele mai apropiate  $k$  instanțe, iar noua instanță va fi în principiu alocată clasei predominante. Clasele vor fi diferențiate între ele pe baza analizei similarității dintre datele conținute de acestea.

Etapa de antrenare a metodei k-NN constă doar în stocarea datelor de antrenare, și anume a perechilor: valori atribuite - etichetă clasă. Clasificarea propriu-zisă a unei noi instanțe se realizează prin calcularea în primă fază a unei măsuri de distanță între aceasta și instanțele deja stocate. Dintre măsurile de distanță folosite, putem enumera următoarele (unde  $x$  și  $y$  reprezintă doi vectori  $m$ -dimensionali de caracteristici):

- distanța Minkowski:

$$D(x, y) = \left( \sum_{i=1}^m |x_i - y_i|^r \right)^{1/r} \quad (7.22)$$

unde  $r$  este o valoare întreagă.

- distanța Manhattan:

$$D(x, y) = \sum_{i=1}^m |x_i - y_i| \quad (7.23)$$

- distanța Chebyshev:

$$D(x, y) = \max_{i=1}^m |x_i - y_i| \quad (7.24)$$

- distanța Euclidiană:

$$D(x, y) = \left( \sum_{i=1}^m |x_i - y_i|^2 \right)^{1/2} \quad (7.25)$$

- distanța Canberra:

$$D(x, y) = \sum_{i=1}^m \frac{|x_i - y_i|}{|x_i + y_i|} \quad (7.26)$$

---

<sup>18</sup>noțiunea de apropiere între instanțe se referă la proximitatea acestora în spațiul de caracteristici considerat.

- corelația Kendall:

$$D(x, y) = 1 - \frac{2}{m(m-1)} \sum_{i=j}^m \sum_{j=1}^{i-1} sign(x_i - x_j) sign(y_i - y_j) \quad (7.27)$$

unde funcția  $sign(E)$  returnează semnul expresiei  $E$ .

Calculând  $D$  pentru toate perechile de instanțe, se va obține o matrice de distanțe ce este folosită mai departe pentru a analiza gradul de vecinătate dintre acestea. Fiecare instanță  $i$ , va fi însotită de eticheta clasei din care face parte,  $c_i \in \{c_1, \dots, c_n\}$ , unde  $n$  reprezintă numărul total al instanțelor.

Mai departe, pentru instanța curentă analizată, folosind matricea de distanțe, sunt localizate cele mai apropiate  $k$  instanțe. Eticheta clasei predominante va fi atribuită instanței analizate. În cazul în care una sau mai multe etichete de clase sunt prezente în număr egal, atunci testul este relansat de această dată pentru  $k-1$  instanțe. Dacă egalitatea se păstrează, numărul de obiecte este diminuat recursiv până când se ajunge la  $k=1$ , moment în care eticheta clasei rămase va fi atribuită implicit instanței analizate.

În ciuda simplității sale, metoda k-NN se dovedește a fi o metodă de clasificare eficientă pentru multe domenii de aplicație. Totuși, aceasta prezintă și o serie de dezavantaje, printre care cele mai importante sunt: dependența acestoria de noțiunea de proximitate (concept vag, dependent de măsura de distanță folosită) precum și dificultatea de alegere adecvată a valorii parametrului  $k$  (metodele existente propuse având o complexitate de calcul ridicată).

## Support Vector Machines

”Support Vector Machines” sau SVM realizează clasificarea datelor prin construcția hiperplanului<sup>19</sup> ce separă în mod optimal datele de intrare în două categorii [Welling 05].

Ca și în cazul celorlalte metode de clasificare, datele sunt reprezentate ca fiind vectori  $n$ -dimensionali în spațiul de caracteristici, iar SVM încearcă să determine dacă aceste puncte pot fi separate de un hiperplan ( $n-1$ ) dimensional. Aceasta este o problemă de clasificare liniară. Având în vedere că există o multitudine de hiperplane ce pot separa datele, SVM restricționează căutarea la acele hiperplane ce permit o separare maximă între cele două clase (maximizarea ”marginii” dintre date). Cu alte cuvinte, se caută hiperplanul cu proprietatea ca acesta să maximizeze distanța față de cel mai apropiat

---

<sup>19</sup>un hiperplan este un concept folosit în domeniul algebrei liniare pentru a generaliza noțiunea de linie, folosită în geometria Euclidiană a planului, sau de plan, folosită în geometria Euclidiană tridimensională, pentru cazul  $n$ -dimensional, cu  $n > 3$ .

punct din spațiul de caracteristici. Acesta este cunoscut și sub numele de "hiperplanul marginii maximale" ("maximum-margin hyperplane").

Formalizarea problemei de clasificare abordată de SVM este următoarea: având la dispoziție un set de date de antrenare,  $D$ , constituit ca fiind:

$$D = \{(X_i, c_i) | X_i \in R^n, c_i \in \{-1, 1\}\} \quad (7.28)$$

unde  $X_i$  este un vector  $n$ -dimensional (vector de caracteristici),  $c_i$  indică clasa din care face parte vectorul  $X_i$  (valori etichete  $-1$  și  $1$ ),  $i$  reprezintă indicele vectorului curent, cu  $i = 1, \dots, p$ , iar  $p$  reprezintă numărul de vectori considerați; se caută hiperplanul marginii maximale ce permite separarea punctelor din clasa  $c_i = 1$  de cele din clasa  $c_i = -1$  (vezi Figura 7.11).

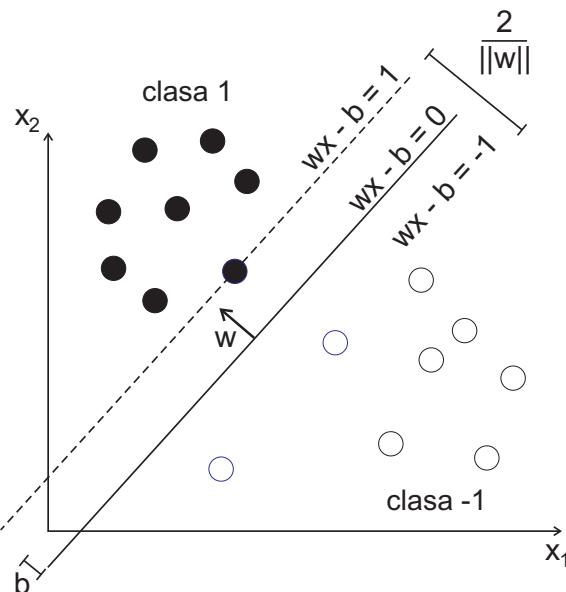


Figura 7.11: Hiperplanul marginii maximale în cazul a două clase (cercurile reprezintă vectorii de caracteristici,  $X_1$  și  $X_2$  formează spațiul de caracteristici, sursă Wikipedia "[http://en.wikipedia.org/wiki/Support-vector\\_machine](http://en.wikipedia.org/wiki/Support-vector_machine)").

Un hiperplan oarecare poate fi definit ca un set de puncte  $X$  ce satisfac următoarea relație:

$$W \cdot X - b = 0 \quad (7.29)$$

unde  $W$  reprezintă un vector normal (perpendicular pe hiperplan) iar parametrul  $\frac{b}{\|W\|}$  va defini decalajul hiperplanului față de originea axei de coordinate, de-a lungul vectorului  $W$  (vezi Figura 7.11).

În scopul definirii marginii maximale, căutăm valorile lui  $W$  și  $b$  astfel încât acestea să maximizeze distanța dintre hiperplanele paralele, cele mai depărtate, dar care încă separă datele. Acestea sunt date de ecuațiile:

$$W \cdot X - b = 1 \quad (7.30)$$

$$W \cdot X - b = -1 \quad (7.31)$$

Distanța dintre acestea este  $\frac{2}{\|W\|}$ , astfel că problema maximizării se transformă într-o problemă de minimizare a valorii  $\|W\|$ .

De asemenea, pentru a preveni ca punctele să se găsească pe margini, se folosesc o serie de constrângeri suplimentare, astfel marginea maximală este determinată de condițiile următoare:

$$W \cdot X_i - b \geq 1, \quad X_i \in c_1 \quad (7.32)$$

$$W \cdot X_i - b \leq -1, \quad X_i \in c_{-1} \quad (7.33)$$

sau

$$c_i \cdot (W \cdot X_i - b) \geq 1 \quad (7.34)$$

pentru oricare  $i \in [1; p]$ .

Transformată într-o problemă de optimizare, clasificarea SVM poate fi enunțată astfel: alege parametrii  $W$  și  $b$  astfel încât să minimizeze valoarea  $\|W\|$  cu constrângerea că  $c_i \cdot (W \cdot X_i - b) \geq 1$ , pentru oricare  $i$ . Această clasificare este valabilă totuși doar în cazul în care datele sunt liniar separabile.

Pentru a crea un clasificator SVM neliniar, la maximizarea marginii dintre clase se folosesc ceea ce numim funcții nucleu sau "kernel functions". Algoritmul rezultat este unul similar din punct de vedere al principiului, doar că operațiile de înmulțire sunt înlocuite acum de nuclee de funcții neliniare. În acest fel, hiperplanul marginii maximale va fi potrivit datelor într-un spațiu de caracteristici transformat neliniar. Dintre nucleele,  $k()$ , cel mai frecvent folosite putem menționa următoarele:

- nucleu polinomial omogen:

$$k(X, X') = (X \cdot X')^d \quad (7.35)$$

unde  $d$  este un număr întreg,

- nucleu polinomial neomogen:

$$k(X, X') = (X \cdot X' + 1)^d \quad (7.36)$$

- funcție radială:

$$k(X, X') = \exp(-\gamma \|X - X'\|^2) \quad (7.37)$$

unde  $\gamma > 0$ ,

- funcție radială Gaussiană:

$$k(X, X') = \exp\left(-\frac{\|X - X'\|^2}{2\sigma^2}\right) \quad (7.38)$$

unde  $\sigma^2$  reprezintă varianța statistică,

- funcție sigmoidă:

$$k(X, X') = \tanh(\kappa \cdot X \cdot X' + c) \quad (7.39)$$

unde  $\tanh$  reprezintă tangenta hiperbolică,  $\kappa > 0$  iar  $c < 0$ .

Din punct de vedere al performanțelor de clasificare, SVM sunt capabile să clasifice eficient date multidimensionale precum și vectori de caracteristici continui. Totuși, pentru a obține o precizie a predicției maximală, este de regulă necesar un set de date de antrenare suficient de vast. În ceea ce privește limitările SVM, una dintre cele mai importante este faptul că SVM este un clasificator binar, datele fiind clasificate în doar două clase. Pentru a realiza o clasificare în mai mult de două clase, aceasta trebuie descompusă într-un set de clasificări binare [Kotsiantis 07].

### 7.3 Concluzii

În acest capitol am prezentat problematica clasificării automate după conținut a datelor. Metodele de clasificare existente au ca obiectiv general regruparea datelor unei colecții mari de date în categorii cât mai *omogene*. Gradul de omogenitate este controlat de utilizator prin introducerea noțiunii de similaritate între date. Aceasta este definită de regulă într-un spațiu  $n$ -dimensional, dat de o serie de atrbute numerice reprezentative a datelor, spațiu numit și *spațiu de caracteristici*.

În ciuda diversității ridicate de metode existente, acestea se împart în două mari categorii, și anume:

- vorbim de *clasificare nesupervizată* în cazul în care partiția optimală a spațiului de caracteristici, din punct de vedere al unui anumit criteriu matematic, este realizată direct, fără a folosi o analiză prealabilă a

conținutului datelor sau alte informații referitoare la acestea. Avantajul acestei abordări este dat de automatizarea procesului de clasificare, lucru ce tinde să diminueze relevanța claselor obținute,

- pe de altă parte, vorbim de *clasificare supervizată* în cazul în care clasarea datelor se face pornind de la o serie de modele predefinite de clase sau clasificări de referință. Clasificarea supervizată implică astfel o etapă prealabilă clasificării în care sistemul învață din exemple. Avantajul acestui tip de abordare este dat în principal de relevanța în general ridicată a conținutului claselor obținute. Totuși, o limitare a acestor metode este dată de faptul că datele de antrenare nu sunt întotdeauna disponibile.

Din punct de vedere al problematicii indexării datelor, tehniciile de clasificare sunt indispensabile unui sistem de indexare după conținut. Clasificarea datelor intervine în general în însuși procesul de căutare al informației. Utilizatorul, prin formularea cererii de căutare va defini spațiul de caracteristici ce va fi folosit pentru localizarea datelor dorite. Pe baza acestuia, datele din baza de date pot fi grupate în funcție de similaritate sau cu alte cuvinte în funcție de asemănarea dintre vectorii de caracteristici asociați. Astfel, grupul sau grupurile de date ce sunt suficient de similare vectorului de caracteristici asociat cererii de căutare vor fi furnizate utilizatorului drept rezultat.

La indexarea bazelor de date generice se caută folosirea unei metode de clasificare nesupervizată, deoarece în acest caz, nu dispunem de informații referitoare la conținutul acesteia sau la o posibilă repartiție în clase. Astfel, se caută un clasificator care să poată pune în evidență noi relații dintre date, sau ceea ce numim "cunoaștere". Clasificarea supervizată este folosită atunci când domeniul de aplicație este cunoscut "a priori" și o expertiză a acestuia este disponibilă. Astfel, algoritmul de clasificare poate fi antrenat pentru a răspunde la anumite cerințe de selecție.

Diversitatea metodelor de clasificare disponibile face dificilă alegerea metodei adecvate aplicației dorite. Și pentru ca lucrurile să fie și mai complicate, găsirea metodei optimale nu garantează clasificarea corectă. Aceasta mai depinde, în principal, și de alegerea eficientă a atributelor (cât mai reprezentative), de măsura de similaritate folosită cât și de validarea corectă a rezultatelor (de regulă dificil de reprezentat).

---

## Bibliografie

---

- [4i2i 06] 4i2i. *H.263 Video Coding Tutorial*. [http://www.4i2i.com/h263\\_video\\_codec.htm](http://www.4i2i.com/h263_video_codec.htm), 2006.
- [Accame 98] M. Accame, F.G.B. De Natale & D.D.Giusto. *High Performance Hierarchical Block-based Motion Estimation for Real-Time Video Coding*. Real-Time Imaging, vol. 4, pag. 67–79, 1998.
- [Acosta 02] E. Acosta, L. Torres, A. Albiol & E. Delp. *An Automatic Face Detection and Recognition System for Video Indexing Applications*. IEEE International Conference on Acoustics, Speech and Signal Processing, vol. 4, pag. 3644–3647, 2002.
- [Adames 02] B. Adames, C. Dorai & S. Venkatesh. *Towards Automatic Extraction of Expressive Elements of Motion Pictures: Tempo*. IEEE Transactions on Multimedia, vol. 4, nr. 4, pag. 472–481, decembrie 2002.
- [Adjero 01] D.A. Adjero & M.C. Lee. *On Ratio-Based Color Indexing*. IEEE Transactions on Image Processing, vol. 10, nr. 1, pag. 36–48, ianuarie 2001.
- [Agoston 87] G.A. Agoston. *Color Theory and Its Application in Art and Design*. Optical Sciences, Springer-Verlag, 1987.
- [Aigrain 95] P. Aigrain, P. Jolly & V. Longueville. *Medium Knowledge-Based Macro-Segmentation into Sequences*.

- Working notes of IJCAI Workshop on Intelligent Multimedia Information Retrieval, pag. 5–14, Montreal, Canada 1995.
- [Akutsu 92] A. Akutsu, Y. Tonomura, H. Hashimoto & Y. Ohba. *Video Indexing Using Motion Vectors*. SPIE Visual Communications Image Processing, vol. 1818, 1992.
- [Alatan 01] A.A. Alatan, A.N. Akasu & W. Wolf. *Multimodal Dialog Scene Detection Using Hidden Markov Models for Content-Based Multimedia Indexing*. Multimedia Tools and Applications, vol. 14, nr. 2, pag. 137–151, 2001.
- [Alattar 93] A.M. Alattar. *Detecting and Compressing Dissolve Regions in Video Sequences with a DVI Multimedia Image Compression Algorithm*. IEEE International Symposium on Circuits and Systems, vol. 1, pag. 13–16, mai 1993.
- [Alattar 97] A.M. Alattar. *Detecting Fade Regions in Uncompressed Video Sequences*. IEEE International Conference on Acoustics, Speech and Signal Processing, vol. 4, pag. 3025–3028, 1997.
- [Allen 83] J.F. Allen. *Maintaining Knowledge About Temporal Intervals*. Communications of the ACM, vol. 26, nr. 11, noiembrie 1983.
- [Aner 01] A. Aner & J.R. Kender. *Mosaic-Based Clustering of Scene Locations in Videos*. IEEE Workshop on Content-based Access of Image and Video Libraries, decembrie, Hawaii, USA 2001.
- [ARGOS 06] ARGOS. *Campagne d'Evaluation d'Outils de Surveillance de Contenus Vidéo*. <http://www.irit.fr/argos>, 2006.
- [Ariki 03] Y. Ariki, M. Kumano & K. Tsukada. *Highlight Scene Extraction in Real Time From Baseball Live Video*. ACM 5th International Workshop on Multimedia Information Retrieval, pag. 209–214, 2003.
- [Arman 93a] F. Arman, A. Hsu & M.Y. Chiu. *Feature Management for Large Video Databases*. SPIE Storage and Retrieval

- for Image and Video Databases, vol. 1908, pag. 2–12, februarie 1993.
- [Arman 93b] F. Arman, A. Hsu & M.Y. Chiu. *Image Processing on Compressed Data for Large Video Database*. ACM International Conference on Multimedia, pag. 267–272, august, Anaheim, USA 1993.
- [A.R.Weeks 95] A.R. Weeks, G.E.Hague & H.R. Myler. *Histogram Equalization of 24-bits Color Images in the Color Difference (c-y) Color Space*. Journal of Electronic Imaging, vol. 4, nr. 1, pag. 15–22, ianuarie 1995.
- [Babaguchi 02] N. Babaguchi, Y. Kawai & T. Kitahashi. *Event Based Indexing of Broadcasted Sports Video by Intermodal Collaboration*. IEEE Transactions on Multimedia, vol. 4, nr. 1, 2002.
- [Barnard 03] K. Barnard, P. Duygulu, N. de Freitas, D. Forsyth, D.M. Blei & M.I. Jordan. *Matching Words and Pictures*. Journal of Machine Learning Research, vol. 3, pag. 1107–1135, 2003.
- [Barron 94] J.L. Barron, D.J. Fleet & S.S. Beauchemin. *Performance of Optical Flow Techniques*. International Journal of Computer Vision, vol. 12, nr. 1, pag. 43–77, februarie 1994.
- [Beaver 94] F. Beaver. *Dictionary of Film Terms*. New York: Twayne, 1994.
- [Ben-Yacoub 99] S. Ben-Yacoub, B. Fasel & J. Luettin. *Fast Face Detection Using MLP and FFT*. Second International Conference on Audio and Video-Based Biometric Person Authentication, pag. 31–36, 1999.
- [Benavente 04] R. Benavente & M. Vanrell. *Fuzzy Colour Naming Based on Sigmoid Membership Functions*. The Second European Conference on Colour Graphics, Imaging and Vision, pag. 135–139, aprilie 2004.
- [Benitez 01] A. Benitez, S.-F. Chang & J.R. Smith. *IMKA: A Multimedia Organization System Combining Perceptual and*

- [Benoit 07] H. Benoit & M. Bernard. *Interactive Video*. Signals and Communication Technology, Springer Berlin Heidelberg, vol. 1, pag. 27–42, 2007.
- [Berlin 91] B. Berlin & P. Kay. *Basic Color Terms: Their Universality and Evolution*. University of California Press, Berkeley 1991.
- [Bimbo 99] A. Del Bimbo. *Visual Information Retrieval*. Morgan Kaufmann Publishers, San Francisco, USA 1999.
- [Birren 69] F. Birren. *Principles of Color - A Review of Past Traditions and Modern Theories of Color Harmony*. New York: Reinhold, 1969.
- [Boccignone 00] G. Boccignone, M. De Santo & G. Percannella. *Automated Threshold Selection for the Detection of Dissolves in Mpeg Video*. IEEE International Conference on Multimedia and Expo, vol. 3, pag. 1535–1538, 2000.
- [Boreczky 98] J.S. Boreczky & L.D. Wilcox. *A Hidden Markov Model Framework for Video Segmentation Using Audion and Image Features*. IEEE International Conference on Acoustics, Speech, and Signal Processing, Seattle, USA 1998.
- [Bouthemy 98] P. Bouthemy & R. Fablet. *Motion Characterization from Temporal Cooccurrences of Local Motion-Based Measures for Video Indexing*. Pattern Recognition, vol. 1, pag. 905–908, august 1998.
- [Bretl 99] W. Bretl & M. Fimoff. *MPEG Tutorial*. [http://www.zenith.com/sub\\_hdtv/mpeg\\_tutorial/](http://www.zenith.com/sub_hdtv/mpeg_tutorial/), 1999.
- [Calic 02a] J. Calic & E. Izquierdo. *Efficient Key-Frame Extraction and Video Analysis*. International Conference on Information Technology: Coding and Computing, pag. 28–33, 2002.
- [Calic 02b] J. Calic & E. Izquierdo. *A Multiresolution Technique for Video Indexing and Retrieval*. IEEE International

- Conference on Image Processing, vol. 1, pag. 952–955, 2002.
- [Calic 04] J. Calic & E. B.T. Thomas. *Spatial Analysis in Key-Frame Extraction Using Video Segmentation*. Workshop on Image Analysis for Multimedia Interactive Services, Lisboa, Portugal 2004.
- [Campisi 99] P. Campisi, A. Longari & A. Neri. *Automatic Key Frame Selection Using a Wavelet Based Approche*. SPIE, vol. 3813, pag. 861–872, 1999.
- [Chang 98] S.-F. Chang, W. Chen, H. Meng, H. Sundara & D. Zhong. *A Fully Automatic Content-Based Video Search Engine Supporting Multi-Object Spatio-Temporal Queries*. IEEE Transactions on Circuits and Systems for Video Technology, vol. 8, pag. 602–615, 1998.
- [Chang 99] S.-F. Chang, W. Chen, H.J. Meng, H. Sundaram & D. Zhong. *Evaluation of Texture Segmentation Algorithms*. International Conference on Computer Vision and Pattern Recognition, vol. 1, pag. 294–299, 1999.
- [Chanussot 98] J. Chanussot. *Approches Vectorielles ou Marginales pour le Traitement d'Images Multi-composantes*. Teză de doctorat, Université de Savoie, Annecy-France, 1998.
- [Chen 99] Y. Chen, E.K. Wong, M.M. Yeunh, Y. Boon-Lock & A.C. Charles. *Augmented Image Histogram for Image and Video Similarity Search*. SPIE Conference on Storage and Retrieval for Image and Video Database VII, vol. 3656, pag. 523–532, 1999.
- [Chen 05] S.-C. Chen, M.-L. Shyu & N. Zhao. *An Enhanced Query Model for Soccer Video Retrieval Using Temporal Relationships*. IEEE International Conference on Data Engineering, aprilie 2005.
- [Chen 06] J.-Y. Chen, C. Taskiran, E.J. Delp & C.A. Bouman. *ViBE: A New Paradigme for Video Database Browsing and Search*. <http://stargate.ecn.purdue.edu/~ips/ViBE/>, 2006.

- [Choi 98] H. Choi & R. Baraniuk. *Multiscale Texture Segmentation using Wavelet-Domain Hidden Markov Models*. 32st Asilomar Conference on Signals, Systems and Computers, vol. 2, pag. 1692–1697, 1998.
- [CICA 06] CICA. *Centre International du Cinéma d'Animation*. <http://www.annecy.org>, 2006.
- [Coldefy 04] F. Coldefy & P. Bouthemy. *Unsupervised Soccer Video Abstraction Based on Pitch, Dominant Color and Camera Motion Analysis*. ACM Multimedia, pag. 268–271, Ney York, USA 2004.
- [Colombo 99] C. Colombo, A. Del Bimbo & P. Pala. *Semantics in Visual Information Retrieval*. IEEE Multimedia, vol. 6, nr. 3, pag. 38–53, 1999.
- [Cooper 02] M. Cooper & J. Foote. *Summarizing Video Using Non-Negative Similarity Matrix Factorization*. IEEE Workshop on Multimedia Signal Processing, pag. 25–28, St. Thomas, US Virgin Islands 2002.
- [Corridoni 95] J.M. Corridoni & A. Del Bimbo. *Film Semantic Analysis*. Proceedings of Computer Architectures for Machine Perception, pag. 202–209, septembrie, Como-Italy 1995.
- [Corridoni 99] J.M. Corridoni, A. Del Bimbo & P. Pala. *Retrieval in Paintings Using Effects Induced by Color Features*. IEEE Multimedia, vol. 6, nr. 3, pag. 38–53, 1999.
- [Cox 00] I.J. Cox, M. Miller, T.P. Minka, T.V. Papathomas & P.N. Yianilos. *The Bayesian Image Retrieval System, PicHunter: Theory, Implementation and Psychophysical Experiments*. IEEE Transactions on Image Processing, vol. 9, pag. 20–37, 2000.
- [CSAIL 06] CSAIL. *Color Name Dictionaries*. <http://swiss.csail.mit.edu/~jaffer/Color/Dictionaries.html>, 2006.
- [Dagtas 00] S. Dagtas, W. Al-Khatib, A. Ghafoor & R.L. Kashyap. *Models for Motion-Based Video Indexing and Retrieval*. IEEE Transactions on Image Processing, vol. 9, nr. 1, pag. 88–101, ianuarie 2000.

- [Dempster 77] A.P. Dempster, N.M. Laird & D.B. Rubin. *Maximum Likelihood from Incomplete Data via the EM Algorithm.* Journal of the Royal Statistical Society, vol. 39, nr. 1, pag. 1–38, 1977.
- [Detyniecki 03] M. Detyniecki & C. Marsala. *Discovering Knowledge for Better Video Indexing Based on Colors.* IEEE International Conference on Fuzzy Systems, vol. 2, pag. 1177–1181, Paris, France 2003.
- [Délibéré 89] M. Délibéré. *La Couleur.* Presses Universitaires de France, 6ème édition, collection: Que sais-je?, vol. 220, 1989.
- [Donderler 04] M.E. Donderler, O. Ulusoy & U. Gudukbay. *Rule-Based Spatio-Temporal Query Processing for Video Databases.* International Journal on Very Large Data Bases, vol. 13, nr. 1, ianuarie 2004.
- [Dorado 04] A. Dorado, J. Calic & E. Izquierdo. *A Rule-Based Video Annotation System.* IEEE Transactions on Circuits and Systems for Video Technology, vol. 14, nr. 5, pag. 622–633, mai 2004.
- [Doulamis 98] N.D. Doulamis, A.D. Doulamis, Y.S. Avrithis & S.D. Kollias. *Video Content Representation Using Optimal Extraction of Frames and Scenes.* IEEE International Conference on Image Processing, vol. 1, pag. 875–879, Chicago, USA 1998.
- [Doulamis 00a] A.D. Doulamis, N. Doulamis & S. Kollias. *Non-Sequential Video Content Representation Using Temporal Variation of Feature Vectors.* IEEE Transactions on Consumer Electronics, vol. 46, nr. 3, pag. 758–768, 2000.
- [Doulamis 00b] A.D. Doulamis, N.D. Doulamis & S.D. Kollias. *A Fuzzy Video Content Representation for Video Summarization and Content-Based Retrieval.* Signal Processing, vol. 80, nr. 6, pag. 1049–1067, iunie 2000.
- [Doulamis 00c] N.D. Doulamis, A.D. Doulamis, Y.S. Avrithis, K.S. Ntalianis & S.D. Kollias. *Efficient Summarization of Stereoscopic Video Sequences.* IEEE Transactions on

- Circuits and Systems for Video Technology, vol. 10, nr. 4, pag. 501–517, iunie 2000.
- [Drew 00] M.S. Drew, Z.N. Li & X. Zhong. *Video Dissolve and Wipe Detection Via Spatio-Temporal Images of Chromatic Histograms Differences*. IEEE International Conference on Image Processing, vol. 3, pag. 929–932, 2000.
- [Duan 06] L.-Y. Duan, J.S. Jin & C.-S. Xu Q. Tian. *Nonparametric Motion Characterization for Robust Classification of Camera Motion Patterns*. IEEE Transactions on Multimedia, vol. 8, nr. 2, pag. 323–340, aprilie 2006.
- [Dufaux 00] F. Dufaux. *Key Frame Selection to Represent a Video*. IEEE International Conference on Multimedia and Expo, vol. 2, pag. 275–278, 2000.
- [Dunn 73] J.C. Dunn. *A Fuzzy Relative of the ISODATA Process and Its Use in Detecting Compact Well-Separated Clusters*. Journal of Cybernetics, vol. 3, pag. 32–57, 1973.
- [Eick 04] C.F. Eick, N. Zeidat & Z. Zhao. *Supervised Clustering - Algorithms and Benefits*. 16th IEEE International Conference on Tools with Artificial Intelligence, pag. 774 – 776, 2004.
- [Eidenberger 04] H. Eidenberger. *A Video Browsing Application Based on Visual MPEG-7 Descriptors and Self-Organising Maps*. TFSA International Journal of Fuzzy Systems, vol. 6, nr. 3, pag. 122–135, septembrie 2004.
- [Elad 99] M. Elad, P. Teo & Y. Hel-Or. *Optimal Filters for Gradient-Based Motion Estimation*. IEEE International Conference on Computer Vision, pag. 559–565, Corfu, Greece 1999.
- [Electric 05] Mitsubishi Electric. *MERL - Timetunnel Interface for Video Browsing*. <http://www.merl.com/projects/timetunnel2/>, 2005.
- [Erol 00] B. Erol & F. Kossentini. *Automatic Key Video Object Plane Selection Using the Shape Information in the MPEG-4 Compressed Domain*. IEEE Transactions on Multimedia, vol. 2, pag. 129–138, 2000.

- [Erol 03] B. Erol & D.-S.H.J. Lee. *Multimodal Summarization of Meeting Recordings*. IEEE International Conference on Multimedia and Expo, vol. 3, pag. 25–28, 2003.
- [Estrela 04] V. Estrela, L.A. Rivera & M.H.S. Bassani. *Pel-Recursive Motion Estimation Using the Expectation-Maximization Technique and Spatial Adaptation*. International Conference in Central Europe on Computer Graphics, Visualization and Computer Vision, pag. 47–54, februarie, Plzen, Czech Republic. 2004.
- [Fablet 02] R. Fablet, P. Bouthemy & P. Pérez. *Nonparametric Motion Characterization Using Causal Probabilistic Models for Video Indexing and Retrieval*. IEEE Transactions on Image Processing, vol. 11, nr. 4, pag. 393–407, 2002.
- [Fan 01] J. Fan, W.G. Aref, A.K. Elmagamid, M.-S. Hadid, M.S. Marzouk & X. Zhu. *Multi-Level Video Content Representation and Retrieval*. Journal of Electronic Imaging, Special Issue on Multimedia Database, vol. 10, nr. 4, pag. 895–908, 2001.
- [Fan 04] J. Fan, H. Luo & A.K. Elmagarmid. *Concept-Oriented Indexing of Video Databases: Toward Semantic Sensitive Retrieval and Browsing*. IEEE Transactions on Image Processing, vol. 13, nr. 7, pag. 974–991, 2004.
- [Farnebäck 00] G. Farnebäck. *Fast and Accurate Motion Estimation using Orientation Tensors and Parametric Motion Models*. 15th International Conference on Pattern Recognition, vol. 1, pag. 135–139, septembrie, Barcelona, Spain 2000.
- [Faugeras 79] O.D. Faugeras. *Digital Color Image Processing within the Framework of a Human Visual Model*. IEEE Transactions on Acoustics, Speech, and Signal Processing, vol. 27, nr. 4, pag. 380–393, 1979.
- [Fauvet 04] B. Fauvet, P. Bouthemy, P. Gros & F. Spindler. *A Geometrical Key-Frame Selection Method Exploiting Dominant Motion Estimation in Video*. International Conference on Image and Video Retrieval, vol. 3115, pag. 419–427, Dublin, Ireland 2004.

- [Ferecatu 01] M. Ferecatu, N. Boujemaa & S. Boughorbel. *Local Color Structure Signatures for Image Retrieval*. International Workshop on Multimedia Content-Based Indexing and Retrieval, Rocquencourt, France 2001.
- [Ferman 99] A.M. Ferman & A.M. Tekalp. *Probabilistic Analysis and Extraction of Video Content*. IEEE International Conference on Image Processing, vol. 2, pag. 91–95, octombrie, Kobe, Japan 1999.
- [Fernando 99] W.A.C. Fernando, C.N. Canagarajah & D.R.Bull. *Fade and Dissolve Detection in Uncompressed and Compressed Video Sequence*. IEEE International Conference on Image Processing, vol. 3, pag. 299–303, octombrie, Kobe, Japan 1999.
- [Fernando 01] W.A.C. Fernando, C.N. Canagarajah & D.R.Bull. *Scene Change Detection Algorithms for Content-Based Video Indexing and Retrieval*. IEE Electronics and Communication Engineering Journal, pag. 117–126, iunie 2001.
- [Fix 51] E. Fix & J. Hodges. *Discriminatory Analysis. Nonparametric Discrimination: Consistency Properties*. Technical Report 11, USAF School of Aviation Medicine, Randolph Field, Texas, USA 1951.
- [Flickner 95] M. Flickner, H. Sawhney, W. Niblack, J. Ashley, Q. Huang, B. Dom, M. Gorkani, J. Hafner, D. Lee, D. Patkovic, D. Steele & P. Yanker. *Query by Image and Video Content: The QBIC System*. IEEE Computer, vol. 28, nr. 9, pag. 23–32, septembrie 1995.
- [Folimage 06] Studio Folimage. *Présentation du studio et de ses productions*. <http://www.folimage.com>, 2006.
- [Fraley 96] C. Fraley. *Algorithms for Model-Based Gaussian Hierarchical Clustering*. Technical Report No.311, University of Washington, Department of Statistics, Seattle, USA 1996.
- [Furht 95] B. Furht, S.W. Smoliar & H. Zhang. *Video and Image Processing in Multimedia Systems*. Norwell, Kluwer 1995.

- [Furnkraz 99] J. Furnkraz. *Separate-and-Conquer Rule Learning*. Artificial Intelligence Review, vol. 13, pag. 3–54, 1999.
- [Gargi 00] U. Gargi, R. Kasturi & S.H. Strayer. *Performance Characterization of Video-Shot-Change Detections Methods*. IEEE Transactions on Circuits, Systems for Video Technology, vol. 10, nr. 1, pag. 1–13, februarie 2000.
- [Ge 02] J. Ge & G. Mirchandani. *A New Hybrid Block-Matching Motion Estimation Algorithm*. IEEE International Conference on Acoustics, Speech and Signal Processing, vol. 4, pag. 4190, 2002.
- [Gibson 02] N.C.D. Gibson & B. Thomas. *Visual Abstraction of Wildlife Footage Using Gaussian Mixture Models*. International Conference on Vision Interface, pag. 814–817, 2002.
- [Gilvarry 99] J. Gilvarry. *Extraction of Motion Vectors from an MPEG Stream*. Technical report Dublin City University, <http://www.cdvp.dcu.ie/Papers/MVector.pdf>, 1999.
- [Gimel'farb 96] G.L. Gimel'farb & A.K. Jain. *On Retrieving Textured Images from an Image Database*. Pattern Recognition, vol. 29, nr. 9, pag. 1461–1483, 1996.
- [Gong 00] Y. Gong & X. Liu. *Generating Optimal Video Summaries*. IEEE International Conference on Multimedia and Expo, vol. 3, pag. 1559–1562, New York, USA 2000.
- [Gong 03] Y. Gong & X. Liu. *Video Summarization and Retrieval Using Singular Value Decomposition*. ACM Multimedia Systems Journal, vol. 9, pag. 157–168, 2003.
- [Graffigne 95] C. Graffigne. *Approche Région: Méthodes Markoviennes*. Analyse d'Images: Filtrage et Segmentation, Masson, Paris, pag. 281–304, octombrie 1995.
- [Gu 97] L. Gu, K. Tsui & D. Keightley. *Dissolve Detection in MPEG Compressed Video*. IEEE International Conference on Intelligent Processing Systems, vol. 2, pag. 1692–1696, 1997.

- [Guillaume 01] S. Guillaume. *Induction de Règles Floues Interprétables*. Teză de doctorat, INSA Toulouse-France, noiembrie 2001.
- [Guimaraes 03] S.J.F. Guimaraes, M. Couprie, A. de A. Araujo & N.J. Leite. *Video Segmentation Based on 2D Image Analysis*. Pattern Recognition Letters, nr. 24, pag. 947–957, 2003.
- [H. Sundaram 00] S. Chang H. Sundaram. *Determining Computable Scenes in Films and Their Structures using Audio-Visual Memory Models*. ACM Multimedia, pag. 95 – 104, Marina del Rey, California, US 2000.
- [Haering 00] N. Haering, R. Qian & I. Sezan. *A Semantic Event-Detection Approach and Its Application to Detecting Hunts in Wildlife Video*. IEEE Transactions on Circuits and Systems for Video Technology, vol. 10, nr. 6, pag. 857–868, 2000.
- [Hampapur 95] A. Hampapur, R. Jain & T.E. Weymouth. *Production Model Based Digital Video Segmentation*. Multimedia Tools and Applications, vol. 1, pag. 9–45, 1995.
- [Hampson 00] F.J. Hampson & J.-C. Pesquet. *Motion Estimation in the Presence of Illumination Variations*. Signal Processing: Image Communication, vol. 16, pag. 373–381, 2000.
- [Han 02a] J. Han & K.-K. Ma. *Fuzzy Color Histogram and Its Use in Color Image Retrieval*. IEEE Transactions on Image Processing, vol. 11, nr. 8, pag. 944–952, 2002.
- [Han 02b] M. Han, W. Hua, W. Xu & Y. Gong. *An Integrated Baseball Digest System Using Maximum Entropy Method*. ACM Multimedia, Juan-les-Pins, France 2002.
- [Hanjalic 97] A. Hanjalic, M. Ceccarelli, R.L. Lagendijk & J. Biemond. *Automation of Systems Enabling Search on Stored Video Data*. SPIE Storage and Retrieval for Image and Video Databases V, vol. 3022, pag. 427–438, februarie 1997.
- [Hanjalic 02] A. Hanjalic. *Shot-Boundary Detection: Unraveled and Resolved?* IEEE Transactions on Circuits and Systems

- for Video Technology, vol. 12, nr. 2, pag. 90–105, februarie 2002.
- [Hauptmann 98] A.G. Hauptmann & M.J. Witbrock. *Story Segmentation and Detection of Commercials in Broadcast News Video*. Advances in Digital Libraries, pag. 168–179, aprilie, Santa Barbara, USA 1998.
- [He 99] L. He, E. Sanocki, A. Gupta & J. Grudin. *Auto-Summarization of Audio-Video Presentations*. ACM Multimedia, pag. 489–498, Orlando, USA 1999.
- [Heckerman 96] D. Heckerman. *A Tutorial on Learning With Bayesian Networks*. Microsoft Research Advanced Technology Division, Technical Report MSR-TR-95-06, 1996.
- [Heng 99] W.J. Heng & K.N. Ngan. *Post Shot Boundary Detection Technique: Flashlight Scene Determination*. Signal Processing and Its Applications, vol. 1, pag. 447–450, august 1999.
- [Heng 01] W.J. Heng & K.N. Ngan. *Enhanced Shot Boundary Refinement for Post-Shot Boundary Detection*. IEEE International Conference on Electrical and Electronic Technology, vol. 1, pag. 259–263, 2001.
- [Hinneburg 00] A. Hinneburg & D.A. Keim. *Clustering Techniques for Large Data Sets From the Past to the Future*. IEEE International Conference on Bioinformatics and Biomedical Engineering, pag. 43–49, 2000.
- [Horn 81] B. Horn & B. Schunck. *Determining optical flow*. Artificial Intelligence, vol. 17, pag. 185–203, 1981.
- [Houten 03] Y.Van Houten, M.Van Setten & J.-G. Schuurman. *Patch-Based Video Browsing*. Human-Computer Interaction INTERACT, 2003.
- [Hsu 02] C.-T. Hsu & S.-J. Teng. *Motion Trajectory Based Video Indexing and Retrieval*. IEEE International Conference on Image Processing, vol. 1, pag. 605–608, septembrie, New York, USA 2002.

- [Huang 98] J. Huang & Y. Wang Z. Liu. *Integration of Audio and Visual Information for Content-Based Video Segmentation*. IEEE International Conference on Image Processing, vol. 3, pag. 526–530, octombrie, Chicago, USA 1998.
- [Ionescu 05a] B. Ionescu, D. Coquin, P. Lambert & V. Buzuloiu. *Analysis and Characterization of Animation Movies*. ORA-SIS Journées Francophones des Jeunes Chercheurs en Vision par Ordinateur, mai, Fournois, Puy-de-Dôme-France 2005.
- [Ionescu 05b] B. Ionescu, D. Coquin, P. Lambert & V. Buzuloiu. *An Approach to Scene Detection in Animation Movies and Its Applications*. Buletin Științific Universitatea Politehnica din București, vol. C (67), nr. 2, pag. 45–57, 2005.
- [Ionescu 05c] B. Ionescu, P. Lambert, D. Coquin & L. Drlea. *Color-Based Semantic Characterization of Cartoons*. IEEE International Symposium on Signals, Circuits and Systems, Special Issue on Statistical Models in Image Processing, vol. 1, pag. 223–226, iulie, Iași, România 2005.
- [Ionescu 06a] B. Ionescu, V. Buzuloiu, P. Lambert & D. Coquin. *Improved Cut Detection for the Segmentation of Animation Movies*. IEEE International Conference on Acoustic, Speech and Signal Processing, mai, Toulouse-France 2006.
- [Ionescu 06b] B. Ionescu, P. Lambert, D. Coquin, L. Ott & V. Buzuloiu. *Animation Movies Trailer Computation*. ACM Multimedia, octombrie, Santa Barbara, CA, USA 2006.
- [Ionescu 07a] B. Ionescu. *Caractérisation Symbolique de Séquences d'Images: Application aux Films d'Animation*. Teză de doctorat, Université de Savoie, LISTIC, <http://www.listic.univ-savoie.fr/>, Annecy-France 2007.
- [Ionescu 07b] B. Ionescu, P. Lambert, D. Coquin & V. Buzuloiu. *Caractérisation du Mouvement dans les Films d'Animation*. Buletin Științific Universitatea Politehnica din București, vol. C (69), nr. 2, 2007.

- [Ionescu 07c] B. Ionescu, P. Lambert, D. Coquin & V. Buzuloiu. *The Cut Detection Issue in the Animation Movie Domain*. Academy Publisher Journal of Multimedia, vol. 2, nr. 4, pag. 10–19, 2007.
- [Ionescu 08] B. Ionescu, D. Coquin, P. Lambert & V. Buzuloiu. *A Fuzzy Color-Based Approach for Understanding Animated Movies Content in the Indexing Task*. Eurasip Journal on Image and Video Processing, Special Issue on Color in Image and Video Processing, 2008.
- [Irani 95] M. Irani, P. Anandan & S. Hsu. *Mosaic Based Representations of Video Sequences and Their Applications*. Computer Vision, pag. 605–611, iunie 1995.
- [IRISA 05] IRISA. *Motion2D*. <http://www.irisa.fr/vista/Motion2D/>, Rennes, France 2005.
- [Itten 61] J. Itten. *The Art of Color: The Subjective Experience and Objective Rationale of Color*. New York: Reinhold, 1961.
- [Jain 91] J.R. Jain & A.K. Jain. *Displacement measurement and its application in interframe image coding*. IEEE Transactions on Communications, vol. 29, nr. 12, pag. 1799–1806, decembrie 1991.
- [Jain 99] A.K. Jain, M.N. Murty & P.J. Flynn. *Data Clustering: A Review*. ACM Computing Surveys, vol. 31, nr. 3, pag. 264–323, septembrie 1999.
- [Jeannin 01] S. Jeannin & A. Divakaran. *MPEG-7 Visual Motion Descriptors*. IEEE Transactions on Circuits and Systems for Video Technology, vol. 11, nr. 6, pag. 720–724, iunie 2001.
- [Jin 02] R. Jin, Y. Qi & A. Hauptmann. *A Probabilistic Model for Camera Zoom Detection*. The Sixteenth Conference of the International Association for Pattern Recognition, nr. 3, pag. 859–862, august, Quebec-Canada 2002.
- [Johnson 67] S.C. Johnson. *Hierarchical Clustering Schemes*. Psychometrika, vol. 2, pag. 241–254, 1967.

- [Kang 01] H.-B. Kang. *A Hierarchical Approach to Scene Segmentation*. IEEE Workshop on Content-Based Access of Image and Video Libraries, 2001.
- [Kay 03] P. Kay & T. Regier. *Resolving the Question of Color Naming Universals*. National Academy of Sciences, vol. 100, nr. 15, 2003.
- [Kelly 76] K.L. Kelly & D.B. Judd. *Color: Universal Language and Dictionary of Names*. National Bureau of Standards, 1976.
- [Kim 00a] C. Kim & J.N. Hwang. *An Integrated Scheme for Object-Based Video Abstraction*. ACM Multimedia, pag. 303 – 311, 2000.
- [Kim 00b] E.Y. Kim, K.I. Kim, K. Jung & H.J. Kim. *A Video Indexing System Using Character Recognition*. IEEE International Conference on Consumer Electronics, pag. 358–359, Los Angles, CA, USA 2000.
- [Kim 02] S.H. Kim & R.H. Park. *Robust Video Indexing for Video Sequences with Complex Brightness Variation*. IASTED International Conference on Signal and Image Processing, pag. 410–414, Kauai, Hawaii 2002.
- [Kim 04] J.-G. Kim, H.S. Chang, J. Kim & H.-M. Kim. *Threshold-Based Camera Motion Characterization of MPEG Video*. Electronics and Telecommunications Research Institute Journal, vol. 26, nr. 3, pag. 269–272, iunie 2004.
- [Klir 95] G. J. Klir & B. Yuan. *Fuzzy Sets and Fuzzy Logic: Theory and Applications*. Prentice Hall, New Jersey, 1995.
- [Kobla 99] V. Kobla, D. DeMenthon & D. Doermann. *Special Effect Edit Detection Using Video Trails: A Comparison with Existing Techniques*. SPIE Storage and Retrieval for Image and Video Databases VII, vol. 3656, pag. 302–313, 1999.
- [Kobla 00] V. Kobla, D. DeMenthon & D. Doermann. *Identifying Sports Video Using Replay, Text and Camera Motion Features*. SPIE Storage and Retrieval for Media Database, vol. 3972, pag. 332–343, 2000.

- [Kosch 01] H. Kosch, L. Boszorményi, A. Bachlechner, B. Dorflingera, C. Hanin, C. Hofbauer, M. Lang, C. Riedler & R. Tusch. *SMOOTH - A Distributed Multimedia Database System*. International Conference on Very Large Data Bases, Rome, Italy 2001.
- [Kotsiantis 07] S.B. Kotsiantis. *Supervised Machine Learning: A Review of Classification Techniques*. Informatica, vol. 31, pag. 249–268, 2007.
- [Kramer 05] P. Kramer & J.B. Pineau. *Camera Motion Detection in the Rough Index Paradigm*. TREC Video Retrieval Evaluation Online Proceedings, TRECVID, noiembrie 2005.
- [Kwok 04] S.H. Kwok, A.G. Constantinides & W.-C. Siu. *An Efficient Recursive Shortest Spanning Tree Algorithm Using Linking Properties*. IEEE Transactions on Circuits and Systems for Video Technology, vol. 14, nr. 6, pag. 852–863, 2004.
- [Kyungpook 06] National University Kyungpook. *Artificial Intelligence Laboratory*. <http://ailab.kyungpook.ac.kr/vindex/video-view.html>, 2006.
- [Laganière 08] R. Laganière, R. Bacco, A. Hocevar, P. Lambert, G. Pays & B. Ionescu. *Video Summarization from Spatio-Temporal Features*. ACM International Conference on Multimedia, Trecvid BBC Rushes Summarization Workshop, octombrie, Vancouver-Canada 2008.
- [Latecki 01] L.J. Latecki, D.D. Widldt & J. Hu. *Extraction of Key Frames from Videos by Optimal Color Composition Matching and Polygon Simplification*. Multimedia Signal Processing Conference, pag. 245–250, Cannes, France 2001.
- [Lay 04] J.A. Lay & L. Guan. *Retrieval for Color Artistry Concepts*. IEEE Transactions on Image Processing, vol. 13, nr. 3, pag. 125–129, martie 2004.
- [Lecce 99] V. Di Lecce, G. Dimauro, A. Guerriero, S. Impedovo, G. Pirlo & A. Salzo. *Image Basic Features Indexing*

- Techniques for Video Skimming.* International Conference on Image Analysis and Processing, pag. 715–720, 1999.
- [Lee 01] M.-S. Lee, Y.-M. Yang & S.-W. Lee. *Automatic Video Parsing Using Shot Boundary Detection and Camera Operation Analysis.* Pattern Recognition, vol. 34, pag. 711–719, 2001.
- [Lee 02a] H. Lee & S.-D. Kim. *Rate-Driven Key Frame Selection Using Temporal Variation of Visual Content.* Electronics Letters, vol. 38, nr. 5, pag. 217–218, februarie 2002.
- [Lee 02b] S. Lee & M.H. Hayes. *Real-Time Camera Motion Classification for Content-Based Indexing and Retrieval using Templates.* IEEE International Conference on Acoustics, Speech, and Signal Processing, vol. 4, pag. 3664–3667, 2002.
- [Lescieux 06] M. Lescieux. *Introduction à la Logique Floue: Plan du Cours.* [http://auto.polytech.univ-tours.fr/automatique/AUA/ressources/Introduction\\_logique\\_floue.ppt](http://auto.polytech.univ-tours.fr/automatique/AUA/ressources/Introduction_logique_floue.ppt), 2006.
- [Li 01] Y. Li, T. Zhang & D. Tretter. *An Overview of Video Abstraction Techniques.* HP Laboratories, HPL-2001-191, 2001.
- [Li 03] Y. Li, S. Narayanan & C.-C.J. Kuo. *Movie Content Analysis, Indexing and Skimming via Multimodal Information.* Video Mining, Chapter 5, Eds. Kluwer Academic Publishers, 2003.
- [Li 04] Z. Li, G. Schuster, A.K. Katsaggelos & B. Gandhi. *Optimal Video Summarization With a Bit Budget Constraint.* IEEE International Conference on Image Processing, vol. 1, pag. 617–620, Singapore 2004.
- [Lienhart 97] R. Lienhart, C. Kuhmunch & W. Effelsberg. *On the Detection and Recognition of Television Commercials.* IEEE Conference on Multimedia Computing and Systems, pag. 509–516, Ottawa, Canada 1997.

- [Lienhart 99a] R. Lienhart. *Comparison of Automatic Shot Boundary Detection Algorithms*. SPIE Storage and Retrieval for Still Image and Video Databases VII, vol. 3656, pag. 290–301, 1999.
- [Lienhart 99b] R. Lienhart, S. Pfeiffer & W. Effelsberg. *Scene Determination Based on Video and Audio Features*. IEEE International Conference on Multimedia, Computing and Systems, vol. 1, pag. 685–690, iunie, Florence-Italy 1999.
- [Lienhart 00] R. Lienhart. *Dynamic Video Summarization of Home Video*. SPIE Storage and Retrieval for Media Databases, vol. 3972, pag. 378–389, ianuarie 2000.
- [Lienhart 01a] R. Lienhart. *Reliable Dissolve Detection*. SPIE Storage and Retrieval for Media Databases, vol. 4315, pag. 219–230, ianuarie 2001.
- [Lienhart 01b] R. Lienhart. *Reliable Transition Detection in Videos: A Survey and Practitioner's Guide*. MRL, Intel Corporation, [http://www.lienhart.de/Publications/IJIG\\_AUG2001.pdf](http://www.lienhart.de/Publications/IJIG_AUG2001.pdf), august, Santa Clara, USA 2001.
- [Lim 01] S.H. Lim & A. El Gamal. *Optical Flow Estimation Using High Frame Rate Sequences*. IEEE International Conference on Image Processing, vol. 2, pag. 925–928, octombrie 2001.
- [Lin 98] C.-W. Lin, Y.-J. Chang & Y.-C. Chen. *Hierarchical Motion Estimation Algorithm Based on Pyramidal Successive Elimination*. International Computer Symposium, octombrie 1998.
- [Lin 02] W.-H. Lin & A.G. Hauptmann. *News Video Classification Using SVM-Based Multimodal Classifiers and Combination Strategies*. ACM Multimedia, pag. 323–326, Juan-les-Pins, France 2002.
- [Liu 02a] C.-C. Liu & A.L.P. Chen. *3D-List: A Data Structure for Efficient Video Query Processing*. IEEE Transactions on Knowledge and Data Engineering, vol. 14, nr. 1, pag. 106–122, ianuarie-februarie 2002.

- [Liu 02b] T. Liu & J.R. Kender. *An Efficient Error-Minimizing Algorithm for Variable-Rate Temporal Video Sampling*. IEEE International Conference on Multimedia and Expo, vol. 1, pag. 413–416, 2002.
- [Liu 02c] T. Liu & J.R. Kender. *Optimization Algorithms for the Selection of Key Frames Sequences of Variable Length*. European Conference on Computer Vision, vol. 2353, pag. 403–417, London, UK 2002.
- [Liu 03] T. Liu, H.-J. Zhang & F. Qi. *A Novel Video Key-Frame Extraction Algorithm Based on Perceived Motion Energy Model*. IEEE Transactions on Circuits and Systems for Video Technology, vol. 13, nr. 10, pag. 1006–1013, octombrie 2003.
- [Liu 04] T. Liu, X. Zhang, J. Freg & K. Lo. *Shot Reconstruction Degree: A Novel Criterion for Keyframe Selection*. Pattern Recognition Letter, vol. 25, nr. 12, pag. 1451–1457, septembrie 2004.
- [Lu 03] S. Lu, I. King & M. Lyu. *Video Summarization Using Greedy Method in a Constraint Satisfaction Framework*. 9th International Conference on Distributed Multimedia Systems, pag. 456–461, Miami, Florida, USA 2003.
- [Lundmark 01] A. Lundmark. *Non-Redundant Search Patterns in Log-Search Motion Estimation*. Swedish Society for Automated Image Analysis Symposium - SSAB, 2001.
- [Lupatini 98] G. Lupatini, C. Saraceno & R. Leonardi. *Scene Break Detection: A Comparison*. Research Issues in Data Engineering, Workshop on Continuous Media Databases and Applications, pag. 34–41, Orlando, FL, USA 1998.
- [Ma 01] Y.F. Ma, J. Sheng, Y. Chen & H.J. Zhang. *MSR-Asia at TREC-10 Video Track: Shot Boundary Detection Task*. 10th Text Retrieval Conference, pag. 371, 2001.
- [MacQueen 67] J.B. MacQueen. *Some Methods for Classification and Analysis of Multivariate Observations*. 5th Berkeley Symposium on Mathematical Statistics and Probability, University of California Press, vol. 1, pag. 281–297, 1967.

- [Maillet 03] S.M. Maillet. *Content-Based Video Retrieval: An Overview.* <http://viper.unige.ch/~marchand/CBVR/>, 2003.
- [Marichal 98] X. Marichal. *Motion Estimation and Compensation for Very Low Bitrate Video Coding.* Teză de doctorat, UCL - Université Catholique de Louvain, Laboratoire de Télécommunications et Télédétection, Louvain-la-Neuve, Belgique 1998.
- [Mazière 00] M. Mazière, F. Chassaing, L. Garrido & P. Salem-bier. *Segmentation and Tracking of Video Objects for a Content-Based Video Indexing Context.* IEEE International Conference on Multimedia Computing and Systems, vol. 2, pag. 1191–1194, New York, USA 2000.
- [Mehtre 97] B.M. Mehtre, M.S. Kankanhalli & W.F. Lee. *Shape Measures for Content Based Image Retrieval: A Comparison.* Information Processing and Management, vol. 33, nr. 3, pag. 319–337, mai, 1997.
- [Meng 95] J. Meng, Y. Juan & S.F. Chang. *Scene Change Detection in a MPEG Compressed Video Sequence.* SPIE Symposium, vol. 2419, pag. 14–25, februarie 1995.
- [Miene 01] A. Miene, A. Dammeyer, T. Hermes & O. Herzog. *Advanced and Adaptive Shot Boundary Detection.* ECDL WS Generalized Documents, pag. 39–43, 2001.
- [Miura 03] K. Miura, R. Hamada, I. Ide, S. Sakai & H. Tanaka. *Motion Based Automatic Abstraction of Cooking Videos.* IPSJ Transactions on Computer Vision and Image Media, vol. 44, 2003.
- [Mojsilovic 00] A. Mojsilovic, J. Kovacevic, R.J. Safranek J. Hu & S.K. Ganapathy. *Matching and Retrieval Based on the Vocabulary and Grammar of Color Patterns.* IEEE Transactions on Image Processing, vol. 9, nr. 1, pag. 38–54, 2000.
- [Morphing 08] Morphing. *Introduction to Media Computation.* <http://coweb.cc.gatech.edu/mediaComp-plan/65>, 2008.

- [Nagasaka 92] A. Nagasaka & Y. Tanaka. *Automatic Video Indexing and Full-Video Search for Object Appearances*. Visual Database Systems II, pag. 113–127, Amsterdam, Netherlands 1992.
- [Nam 98] J. Nam, M. Alghoniemy & A.H. Tewfik. *Audio-Visual Content-Based Violent Scene Characterization*. IEEE International Conference on Image Processing, vol. 1, pag. 353–357, Chicago, USA 1998.
- [Nam 00] J. Nam & A.H. Tewfik. *Dissolve Transition Detection Using B-Spline Interpolation*. IEEE International Conference on Multimedia and Expo, vol. 3, pag. 1349–1352, iulie 2000.
- [Naphade 01a] M.R. Naphade. *A Probabilistic Framework for Mapping Audio-Visual Features to High-Level Semantics in Terms of Concepts and Context*. Teză de doctorat, Department of Electrical and Computing Engineering, University of Illinois, 2001.
- [Naphade 01b] M.R. Naphade & T.S. Huang. *A Probabilistic Framework for Semantic Video Indexing, Filtering and Retrieval*. IEEE Transactions on Multimedia, vol. 3, nr. 1, pag. 141–151, 2001.
- [Naphade 02] M.R. Naphade & T.S. Huang. *Extracting Semantics from Audiovisual Content: The Final Frontier in Multimedia Retrieval*. IEEE Transactions on Neural Networks, vol. 13, nr. 2, pag. 793–810, 2002.
- [Netravali 79] A.N. Netravali & J.D. Robbins. *Motion-compensated television coding: Part I*. BELL System Technical Journal, vol. 58, nr. 3, pag. 631–670, martie 1979.
- [Ngo 00] C.-W. Ngo, T.-C. Pong, H.-J. Zhang & R.T. Chin. *Motion Characterization by Temporal Slices Analysis*. IEEE International Conference on Computer Vision and Pattern Recognition, vol. 2, pag. 768–773, 2000.
- [Ngo 03] C.-W. Ngo, Y.-F. Ma & H.-J. Zhang. *Automatic Video Summarization by Graph Modeling*. IEEE International Conference on Computer Vision, vol. 1, pag. 104, Nice, France 2003.

- [Niculescu 06] R.S. Niculescu, T.M. Mitchell & R.B. Rao. *Bayesian Network Learning with Parameter Constraints*. Journal of Machine Learning Research, vol. 7, pag. 1357–1383, 2006.
- [Ohta 80] Y.I. Ohta, T. Kanade & T. Sakai. *Color Information for Region Segmentation*. Computer Graphics and Image Processing, vol. 13, pag. 222–241, 1980.
- [Otsuji 91] K. Otsuji, Y. Tonomura & Y. Ohba. *Video Browsing Using Brightness Data*. SPIE Visual Communications and Image Processing, vol. 1606, pag. 980–989, 1991.
- [Ott 07] L. Ott, P. Lambert, B. Ionescu & D. Coquin. *Animation Movie Abstraction: Key Frame Adaptative Selection based on Color Histogram Filtering*. International Conference on Image Analysis and Processing, Computational Color Imaging Workshop, septembrie, Modena, Italy 2007.
- [Pan 00] C. Pan & S. Ma. *3D Motion Estimation of Human by Genetic Algorithm*. 15th International Conference on Pattern Recognition, vol. 1, pag. 1159, 2000.
- [Pan 01] H. Pan, P. Beek & M. Sezan. *Detection of Slow-Motion Replay Segments in Sports Video for Highlights Generation*. IEEE International Conference on Acoustics, Speech and Signal Processing, vol. 3, pag. 1649–1652, Salt Lake City, Utah, SUA 2001.
- [Papoulis 91] Papoulis. *Probability, Random Variables, and Stochastic Processes*. Mc Graw Hill, Inc., New-York, 3rd edition, 1991.
- [Pfeiffer 96] S. Pfeiffer, R. Lienhart, S. Fisher & W. Effelsberg. *Abstracting Digital Movies Automatically*. Journal of Visual Communication and Image Representation, vol. 7, nr. 4, decembrie 1996.
- [Pilu 97] M. Pilu. *On Using Raw MPEG Motion Vectors To Determine Global Camera Motion*. HP - Hewlett Packard, <http://www.hpl.hp.com/techreports/97/HPL-97-102.pdf>, august 1997.

- [Pineau 05] J.B. Pineau. *Extraction des Objets Couleur en Mouvement des Séquences Vidéo*. LABRI UMR CNRS 5800, <http://www.labri.fr/ImageetSon/AIV>, 2005.
- [Porter 00] S.V. Porter, M. Mirmehdi & B.T. Thomas. *Video Cut Detection Using Frequency Domain Correlation*. 15th International Conference on Pattern Recognition, vol. 3, pag. 413–416, Barcelona, Spain 2000.
- [Porter 01] S.V. Porter, M. Mirmehdi & B.T. Thomas. *Detection and Classification of Shot Transitions*. 12th British Machine Vision Conference, pag. 73–82, septembrie 2001.
- [Porter 03] S.V. Porter, M. Mirmehdi & B.T. Thomas. *A Shortest Path Representation for Video Summarization*. 12th International Conference on Image Analysis and Processing, pag. 460–465, septembrie 2003.
- [QBIC 03] IBM QBIC. *Hermitage Museum - Query by Image Content*. <http://www.hermitage-museum.org>, 2003.
- [Qian 99] R. Qian, N. Haering & I. Sezan. *A Computational Approach to Semantic Event Detection*. Computer Vision and Pattern Recognition, vol. 1, pag. 206, iunie 1999.
- [Radhakrishnan 04] R. Radhakrishnan, A. Divakaran & Z. Xiong. *A Time Series Clustering Based Framework for Multimedia Mining and Summarization Using Audio Features*. ACM International Workshop on Multimedia Information Retrieval, pag. 157–164, New York, USA 2004.
- [Rasheed 03] Z. Rasheed & M. Shah. *Scene Detection in Hollywood Movies and TV Shows*. IEEE Computer Vision and Pattern Recognition, vol. 2, pag. 343–351, Wisconsin, USA 2003.
- [Ren 03] W. Ren & S. Singh. *Video Transition: Modeling and Prediction*. Pattern Analysis and Neural Networks, PANN, [http://www.dcs.ex.ac.uk/research/pann/pdf/pann\\_SS\\_089.PDF](http://www.dcs.ex.ac.uk/research/pann/pdf/pann_SS_089.PDF), 2003.
- [Reoxiang 94] L. Reoxiang, Z. Bing & M.L. Liou. *A New Three-Step Search Algorithm for Block Motion Estimation*. IEEE

- Transactions on Circuits and Systems for Video Technology, vol. 4, nr. 4, pag. 438–442, august 1994.
- [Research 05] Compaq Corporate Research. *Audio Search Using Speech Recognition*. <http://speechbot.research.compaq.com>, 2005.
- [Rivlin 95] E. Rivlin & I. Weiss. *Local Invariants for Recognition*. IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 17, nr. 3, pag. 226–238, 1995.
- [Rong 04] J. Rong, W. Jin & L. Wu. *Key Frame Extraction Using Inter-Shot Information*. IEEE International Conference on Multimedia and Expo, vol. 1, Taipei, Taiwan 2004.
- [Rosenblatt 62] F. Rosenblatt. *Principles of Neurodynamics*. Spartan, New York, 1962.
- [Rowe 01] L.A. Rowe, D. Harley, P. Pletcher & S. Lawrence. *BIBS: A Lecture Webcasting System*. Berkley Multimedia Research Center, TR 2001-160, iunie 2001.
- [Rumelhart 86] D.E. Rumelhart, G.E. Hinton & R.J. Williams. *Learning Internal Representations by Error Propagation*. Parallel Distributed Processing: Explorations in the Microstructure of Cognition, MIT Press, Cambridge, vol. 1, pag. 318–362, 1986.
- [Sanson 81] H. Sanson. *Motion Affine Models Identification and Application to Television Image Coding*. Visual Communication and Image Processing, vol. 1605, nr. 2, pag. 570–581, noiembrie 1981.
- [Saraceno 98] C. Saraceno & R. Leonardi. *Identification of Story Units in AV Sequencies by Joint Audio and Video Processing*. IEEE International Conference on Image Processing, vol. 1, pag. 363–367, octombrie, Chicago, USA 1998.
- [Saur 97] D.D. Saur, Y.P. Tan, S.R. Kulkarni & P.J. Ramadge. *Automated Analysis and Annotation of Basketball Video*. SPIE Symposium on Storage and Retrieval for Image and Video Databases V, vol. 3022, pag. 176–187, 1997.

- [Scaringella 06] N. Scaringella, G. Zoia & D. Mlynek. *Automatic Genre Classification of Music Content*. IEEE Signal Processing Magazine, vol. 23, nr. 2, pag. 133–141, martie 2006.
- [Schmid 97] C. Schmid & R. Mohr. *Local Grayvalue Invariants for Image Retrieval*. IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 19, nr. 5, pag. 530–535, 1997.
- [Schneiderman 00] H. Schneiderman & T. Kanade. *A Statistical Method for 3D Object Detection Applied to Faces and Cars*. IEEE Computer Vision and Pattern Recognition, vol. 1, pag. 746–751, Hilton Head Island, SC, USA 2000.
- [Shahrray 95] B. Shahrray. *Scene Change Detection and Content-Based Sampling of Video Sequences*. SPIE Conference on Digital Video Compression: Algorithms and Technologies, vol. 2419, pag. 2–13, februarie 1995.
- [Shen 97] B. Shen. *HDH Based Compressed Video Cut Detection*. Visual 97, pag. 149–156, decembrie, San Diego, USA 1997.
- [Smeulders 00] A.W.M. Smeulders, M. Worring, S. Santini, A. Gupta & R. Jain. *Content-Based Image Retrieval at the End of the Early Years*. IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 22, nr. 12, pag. 1349–1380, decembrie 2000.
- [Smith 98] M.A. Smith & T. Kanade. *Video Skimming and Characterization Through the Combination of Image and Language Understanding*. International Workshop on Content-Based Access of Image and Video Databases, Bombay, India 1998.
- [Smith 04] P. Smith, T. Drummond & R. Cipolla. *Layered Motion Segmentation and Depth Ordering by Tracking Edges*. IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 26, nr. 4, pag. 479–492, aprilie 2004.
- [Smith 99] J.R. Smith. *VideoZoom Spatial-Temporal Video Browsing*. IEEE Transactions on Multimedia, vol. 1, iunie 99.

- [Snoek 05a] C.G.M. Snoek & M. Worring. *Multimedia Event Based Video Indexing Using Time Intervals*. IEEE Transactions on Multimedia, vol. 7, nr. 4, august 2005.
- [Snoek 05b] C.G.M. Snoek & M. Worring. *Multimodal Video Indexing: A Review of the State-of-the-art*. Multimedia Tool and Applications, vol. 25, nr. 1, pag. 5–35, 2005.
- [Song 02] H.S. Song, I.K. Kim & N.I. Cho. *Scene Change Detection by Feature Extraction from Strong Edge Blocks*. SPIE Visual Communications and Image Processing, vol. 4671, pag. 784–792, 2002.
- [Stricker 95] M. Stricker & M. Orengo. *Similarity of color images*. SPIE Conference on Storage and Retrieval for Image and Video Databases, vol. 2420, pag. 381–392, 1995.
- [Su 05a] C.-W Su, H.-Y.M. Liao, H.-R. Tyan, K.-C. Fan & L.-H. Chen. *A Motion-Tolerant Dissolve Detection Algorithm*. IEEE Transactions on Multimedia, vol. 7, nr. 6, pag. 1106–1113, 2005.
- [Su 05b] C.W. Su, H.R. Tyan, H.Y. Mark Liao & L.H. Chen. *A Motion-Tolerant Dissolve Detection Algorithm*. IEEE Transactions on Multimedia, vol. 7, nr. 6, pag. 1106–1113, 2005.
- [Sun 00] X. Sun & M.S. Kankanhalli. *Video Summarization Using R-Sequences*. Journal of Real Time Imaging, vol. 6, pag. 449–459, 2000.
- [Sundaram 02] H. Sundaram & S.-F. Chang. *Video Skims: Taxonomies and an Optimal Generation Framework*. IEEE International Conference on Image Processing, vol. 2, pag. 21–24, Rochester, USA 2002.
- [Taniguchi 95] Y. Taniguchi, A. Akutsu, Y. Tonomura & H. Hamada. *An Intuitive and Efficient Access Interface to Real-Time Incoming Video Based on Automatic Indexing*. ACM Multimedia, pag. 25–33, San Francisco, California, United States 1995.
- [Tardini 05] G. Tardini, C. Grana, R. Marchi & R. Cucchiara. *Shot Detection and Motion Analysis for Automatic MPEG-7*

- Annotation of Sports Videos.* 13th International Conference on Image Analysis and Processing, pag. 653–660, septembrie 2005.
- [Timoner 01] S.J. Timoner & D.M. Freeman. *Multi-image gradient-based algorithms for motion estimation.* SPIE Optical Engineering, vol. 40, nr. 9, pag. 2003–2016, 2001.
- [Tong 01] S. Tong & E. Chang. *Support Vector Machine Active Learning for Image Retrieval.* ACM Multimedia Conf., vol. 9, Ottawa, Canada 2001.
- [Trecvid 08] Trecvid. *Video Retrieval Evaluation.* <http://www-nlpir.nist.gov/projects/trecvid/>, 2008.
- [Trémeau 04] A. Trémeau, C. Fernandez-Maloigne & P. Bonton. *Image Numérique Couleur: De l'Acquisition au Traitement.* DUNOD ISBN 2-10-006843-1, 2004.
- [Truong 00a] B.T. Truong & C. Dorai. *Automatic Genre Identification for Content-Based Video Categorization.* IEEE International Conference on Pattern Recognition, vol. 4, pag. 230–233, Barcelona, Spain 2000.
- [Truong 00b] B.T. Truong, C. Dorai & S. Venkatesh. *New Enhancements to Cut, Fade, and Dissolve Detection Processes in Video Segmentation.* ACM Multimedia, pag. 219–227, noiembrie 2000.
- [Truong 07] B.T. Truong & S. Venkatesh. *Video Abstraction: A Systematic Review and Classification.* ACM Transactions on Multimedia Computing, Communications and Applications, vol. 3, nr. 1, 2007.
- [Turaga 98] D. Turaga & M. Alkanhal. *Search Algorithms for Block-Matching in Motion Estimation.* [http://www.ece.cmu.edu/~ee899/project/deepak\\_mid.htm](http://www.ece.cmu.edu/~ee899/project/deepak_mid.htm), 1998.
- [Vasconcelos 00] N. Vasconcelos & A. Lippman. *Statistical Models of Video Structure for Content Analysis and Characterization.* IEEE Transactions on Image Processing, vol. 9, pag. 3–19, ianuarie 2000.

- [Vendriga 01] J. Vendriga & M. Worring. *Evaluation of Logical Story Unit Segmentation in Video Sequences*. IEEE International Conference on Multimedia and Expo, pag. 1092–1095, Tokyo, Japan 2001.
- [Venkataraman 99] A. Venkataraman. *Decision Trees*. <http://www.speech.sri.com/people/anand/771/html/node28.html>, 1999.
- [Vermaak 02] J. Vermaak, P. Prez, M. Gangnet & A. Blake. *Rapid Summarization and Browsing of Video Sequences*. British Machine Vision Conference, vol. 1, pag. 424–433, Cardiff, UK 2002.
- [Vertan 04] C. Vertan, B. Ionescu, I. Stefan & M. Ciuc. *Osteoporosis detection in calcaneum X-ray images by digital image processing: fractal and statistical approaches*. Interdisciplinary Applications of Fractal and Chaos Theory, Eds. R. Dobrescu, C. Vasilescu, Editura Academiei Române (ISBN 973-27-1070-5), 2004.
- [Vertan 08] C. Vertan. *Prelucrarea și Analiza Imaginilor Color*. Suport de curs, laboratorul LAPI - Laboratorul de Analiza și Prelucrarea Imaginilor, Universitatea Politehnica din București, <http://alpha.imag.pub.ro/ro/cursuri/paic/index.html>, 2008.
- [Visibone 06] Visibone. *Webmaster Palette*. <http://www.visibone.com/colorlab>, 2006.
- [Vogl 99] S. Vogl, K. Manske & M. Muhlhäuser. *A VRML Approach to Web Video Browsing*. Multimedia Computing and Networking, pag. 276–285, San Jose, USA 1999.
- [Wallin 01] H. Wallin, C. Christopoulos & F. Furesjo. *Robust parametric motion estimation for image mosaicing in the MPEG-7 standard*. IEEE International Conference on Image Processing, vol. 2, pag. 961–964, octombrie 2001.
- [Wang 92] L.-X. Wang. *Fuzzy Systems are Universal Approximators*. IEEE Conference on Fuzzy Systems, pag. 1163–1170, San Diego, USA 1992.

- [Wang 00] Y. Wang, Z. Liu & J.-C. Huang. *Multimedia Content Analysis Using Both Audio and Visual Clues*. IEEE Signal Processing Magazine, vol. 17, nr. 6, pag. 12–36, noiembrie 2000.
- [Wang 01] J.Z. Wang, J. Li & G. Wiederhold. *SIMPLICITY: Semantic-Sensitive Integrated Matching for Picture Libraries*. IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 23, nr. 9, pag. 947–963, 2001.
- [Welling 05] M. Welling. *Support Vector Machines*. Note de curs, University of Toronto, Department of Computer Science, Canada, [http://www.ics.uci.edu/~welling/classnotes/papers\\_class/SVM.pdf](http://www.ics.uci.edu/~welling/classnotes/papers_class/SVM.pdf), 2005.
- [Wikipedia 08] Wikipedia. *The Free Online Encyclopedia*. <http://www.wikipedia.org/>, 2008.
- [Witten 05] I.H. Witten & E. Frank. *Data Mining - Practical Machine Learning Tools and Techniques*. Elsevier, Morgan Kaufman Publishers, second edition, pag. 265–270, 2005.
- [Wolf 96] W. Wolf. *Key Frame Selection by Motion Analysis*. IEEE International Conference on Acoustics, Speech and Signal Processing, vol. 2, pag. 1228–1231, 1996.
- [Wyszecki 82] G. Wyszecki & W.S. Stiles. *Color Science: Concepts and Methods, Quantitative Data and Formulae*. John Wiley and sons, second edition, 1982.
- [Xiong 97] W. Xiong, J. C.-M. Lee & R.-H. Ma. *Automatic Video Data Structuring Through Shot Partitioning and Key Frame Computing*. Machine Vision and Applications, vol. 10, nr. 2, pag. 51–65, 1997.
- [Xiong 03] Z. Xiong, R. Radhakrishnan & A. Divakaran. *Generation of Sports Highlights Using Motion Activity in Combination With a Common Audio Feature Extraction Framework*. IEEE International Conference on Image Processing, vol. 1, Barcelona, Spain 2003.
- [Yahiaoui 01] L. Yahiaoui, B. Merialdo & B. Huet. *Automatic Video Summarization*. CBMIR, 2001.

- [Yao 01] A. Yao & J. Jin. *The Developing of a Video Metadata Authoring and Browsing System in XML*. ACM International Conference - Pan-Sydney Workshop on Visual Information Processing, vol. 2, Sydney, Australia 2001.
- [Yeo 95] B.-L. Yeo & B. Liu. *Rapid Scene Analysis on Compressed Video*. IEEE Transactions on Circuits, Systems and Video Technology, vol. 5, pag. 533–544, decembrie 1995.
- [Young 02] T. Young. *On the Theory of Light and Colors*. Philosophical Transactions of the Royal Society, nr. 91, 1802.
- [Yu 97] H. Yu, G. Bozdagı & S. Harrington. *Feature-Based Hierarchical Video Segmentation*. IEEE International Conference on Image Processing, vol. 2, pag. 498–501, 1997.
- [Yu 03] B. Yu, W.-Y. Ma, K. Nahrstedt & H.-J. Zhang. *Video Summarization Based on User Log Enhanced Link Analysis*. ACM Multimedia, pag. 382–391, Berkeley, USA 2003.
- [Yu 04] X.-D. Yu, L. Wang, Q. Tian & P. Xue. *Multi-level Video Representation With Application to Keyframe Extraction*. International Conference on Multimedia Modeling, pag. 117–121, Brisban, Australia 2004.
- [Yuan 95] Y. Yuan & M.J. Shaw. *Induction of Fuzzy Decision Trees*. Fuzzy Sets and Systems, vol. 69, pag. 125–139, 1995.
- [Zabih 95] R. Zabih, J. Miller & K. Mai. *A Feature-Based Algorithm for Detecting and Classifying Scene Breaks*. ACM Multimedia, pag. 189–200, noiembrie, San Francesco, USA 1995.
- [Zabih 99] R. Zabih, J. Miller & K. Mai. *A Feature-Based Algorithm for Detecting and Classification Production Effects*. Multimedia Systems, vol. 7, pag. 119–128, 1999.
- [Zadeh 65] L.A. Zadeh. *Fuzzy Sets*. Information and Control, vol. 8, nr. 3, pag. 338–353, 1965.

- [Zahariadis 96] T. Zahariadis & D. Kalivas. *A Spiral Search Algorithm For Fast Estimation Of Block Motion Vectors*. Signal Processing VIII, Theories and Applications. Eighth European Signal Processing Conference, vol. 2, pag. 1079–1082, 1996.
- [Zaim 01] M. Zaim, A. El Ouaazizi & R. Benslimane. *Genetic Algorithms Based Motion Estimation*. Vision Interface Annual Conference, iunie, Ottawa, Canada 2001.
- [Zeng 02] W. Zeng, W. Gao & D. Zhao. *Video Indexing by Motion Activity Maps*. IEEE International Conference on Image Processing, vol. 1, pag. 912–915, 2002.
- [Zhang 93] H. Zhang, A. Kankanhalli & S.W. Smoliar. *Automatic Partitioning of Full-Motion Video*. Multimedia Systems, vol. 1, nr. 1, pag. 10–28, 1993.
- [Zhang 94] H. Zhang, C.Y. Low, Y. Gong & S.W. Smoliar. *Video Parsing Using Compressed Data*. SPIE Image and Video Processing II, vol. 2182, pag. 142–149, februarie 1994.
- [Zhang 97] H.J. Zhang, J. Wu, D. Zhong & S.W. Smoliar. *An Integrated System for Content-Based Video Retrieval and Browsing*. Pattern Recognition, vol. 30, nr. 4, pag. 643–658, 1997.
- [Zhao 03] M. Zhao, J. Bu & C. Chen. *Audio and Video Combined for Home Video Abstraction*. IEEE International Conference on Acoustic, Speech and Signal Processing, vol. 5, pag. 620–623, Hong Kong, China 2003.
- [Zhong 96] D. Zhong, H. Zhang & C. Chang. *Clustering Methods for Video Browsing and Annotation*. SPIE Storage and Retrieval for Still Image and Video Databases IV, vol. 2670, pag. 239–246, 1996.
- [Zhong 97] D. Zhong & S.-F. Chang. *Spatio-temporal Video Search using the Object-Based Video Representation*. IEEE International Conference on Image Processing, vol. 1, pag. 1–12, 1997.

- [Zhou 02] W. Zhou, S. Dao & C.-C.J. Kuo. *On-Line Knowledge and Rule-Based Video Classification System for Video Indexing and Dissemination*. Information Systems, vol. 27, nr. 8, 2002.
- [Zhu 05] C.-Z. Zhu, T. Mei & X.-S. Hua. *Video Booklet - Natural Video Browsing*. ACM Multimedia, pag. 265 – 266, noiembrrie, Singapore 2005.
- [Zhuang 98] Y. Zhuang, Y. Rui, T.S. Huang & S. Mehrota. *Adaptive Key Frame Extraction Using Unsupervised Clustering*. IEEE International Conference on Image Processing, vol. 1, pag. 866–870, 1998.

