

Universitatea Tehnică "Gh. Asachi" Iași  
Facultatea de Automatică și Calculatoare  
Specializarea Tehnologia Informației

Lucrare de licență

# **Identificarea obiectelor bazată pe puncte de interes, în secvențe video**

Absolvent

**Carata Lucian**

Coordonator științific  
Prof. Dr. Ing. Vasile Manta

Iași, 2009



Programs must be written for people to read, and only incidentally for machines to execute.

*Programele trebuie scrise mai întâi pentru a fi citite de oameni și doar apoi pentru a fi executate de mașini.*

— Abelson and Sussman



---

# Cuprins

---

<b>Cuprins</b>	<b>i</b>
<b>1 Introducere</b>	<b>1</b>
1.1 Recunoașterea automată a obiectelor . . . . .	1
1.2 Formularea și abordarea temei . . . . .	2
<b>2 Noțiuni teoretice</b>	<b>5</b>
2.1 Abordări ale temei în literatura de specialitate . . . . .	5
2.2 Identificarea trăsăturilor . . . . .	7
2.2.1 Detectorul Harris . . . . .	7
2.2.2 Detectorul SIFT . . . . .	9
2.3 Identificarea trăsăturilor în timp real . . . . .	14
<b>3 Chapter's title</b>	<b>15</b>
<b>Bibliografie</b>	<b>17</b>
<b>A Anexa 1</b>	<b>21</b>
<b>Lista de simboluri și prescurtări</b>	<b>23</b>
<b>Listă de figuri</b>	<b>24</b>
<b>Listă de tabele</b>	<b>25</b>
<b>Glosar</b>	<b>27</b>



# Capitolul 1

---

## Introducere

---

### 1.1 Recunoașterea automată a obiectelor

Identificarea și recunoașterea automată a unor obiecte, în imagini statice sau secvențe video, este îndelung studiată în grafica pe calculator, prezentând interes din perspectiva dificultăților întâmpinate în rezolvarea problemei de sistemele de calcul în comparație cu sistemele biologice, dar și datorită aplicabilității în domenii din cele mai diverse.

Astfel, primele aplicații s-au conturat în mediul industrial, pentru inspectarea automată a produselor de pe liniile de fabricație (de exemplu, identificarea defectelor unor plăci integrate - lipituri incorecte, plasări incorecte de componente etc). Totuși, mediul de recunoaștere în aceste cazuri este unul controlat, existând limite stricte între care recunoașterea se realizează cu succes. Mai mult, obiectele pentru care se realizează identificarea au caracteristici bine cunoscute. Pornind de aici, s-au dezvoltat metode care încearcă să elimine cât mai multe dintre restricții, și să permită recunoașterea în cazul general. Aceste abordări largesc gama de aplicații și în zona utilizatorilor obișnuiți, pentru îmbunătățirea următoarelor generații de motoare de căutare, programe de chat sau de supraveghere a locuințelor. Desigur, domenii precum medicina (recunoașterea sau numărarea celulelor de un anumit tip), robotica (dezvoltarea unor roboți care să interacționeze cu mediul înconjurător folosind "vederea artificială") utilizează și ele recunoașterea obiectelor ca subproblemă. Aplicații similare există în domeniul militar (recunoașterea unor dispozitive suspecte în aeroporturi, identificarea persoanelor pe baza înfățișării).

Problemele cele mai mari în identificarea și recunoașterea obiectelor apar datorită variațiilor din mediul în care obiectul este plasat (culoare, lumină, umbre, reflexii, ocluzionarea obiectului țintă de către alte obiecte). De asemenea, apar dificultăți și datorită diferențelor de poziționare și perspectivă între reprezenta-

rea inițială a obiectului care se dorește a fi identificat (de cele mai multe ori, o fotografie a respectivului obiect) și situația reală în care se încearcă identificarea acestuia (când el poate fi "privit" la o altă scală, rotit sau dintr-un punct de vedere diferit).

Au fost găsite mai multe abordări pentru rezolvarea acestor probleme, majoritatea detectând în fiecare imagine anumite zone caracteristice, invariante la modificări ale parametrilor de mediu/perspectivă. Se realizează apoi o potrivire între ele și o bază de date în care au fost anterior reținute caracteristicile obiectelor căutate. Diferențele între metode se referă la modalitatea de detecție a zonelor, la forma lor (puncte, arii din imagine) și la informațiile reținute pentru fiecare zonă în parte astfel încât ea să poată fi regăsită într-o nouă imagine și atribuită ca aparținând obiectului căutat.

La momentul actual, tehnicile de recunoaștere nededicate permit o detecție cu un procentaj de reușită și repetabilitate a rezultatelor suficient de mare (tipic peste 80%) pentru a fi considerate aplicabile cu succes în aplicații de orice tip. Totuși, se pune problema selectării unor metode cât mai eficiente, care să poată fi aplicate "în timp real".

## 1.2 Formularea și abordarea temei

Recunoașterea unor obiecte (furnizate ca imagini, drept date de intrare), într-o secvență video live sau filmată anterior (video salvat pe hdd), presupune detectarea existenței obiectelor și identificarea poziției acestora în fiecare frame, urmată de "adnotarea" frame-ului în zona obiectului recunoscut. Atât inițial, pentru imaginile ce definesc obiectele, cât și pentru secvența video, se aplică aceeași algoritmi de determinare a zonelor caracteristice. Apoi, se realizează o potrivire între rezultatele obținute pentru frame-ul curent și modelul determinat pentru fiecare dintre obiecte. În măsura în care există corespondențe (în frame-ul curent există zone similare cu cele ale obiectului), se stabilește prezența obiectului în frame, precum și poziția acestuia. Se dorește o variație cât mai mică a rezultatelor la schimbări de scală, rotații, modificări ale perspectivei 3D și a luminozității ambientale, urmărind în același timp o repetabilitate crescută a experimentelor. De asemenea, este de preferat să nu se impună restricții legate de modul în care sunt capturate imaginile sau de calibrarea anterioară a dispozitivelor (camere video, aparate foto digitale).

Dacă drept zone caracteristice se folosesc puncte de interes, există 2 pași generali în analiza fiecărei imagini prelucrate:

1. *Localizarea*: determinarea automată a poziției punctelor de interes (în două imagini ale aceluiași obiect, se dorește ca punctele de interes să fie localizate în aceeași zonă relativ la obiect - Figura 1.1)



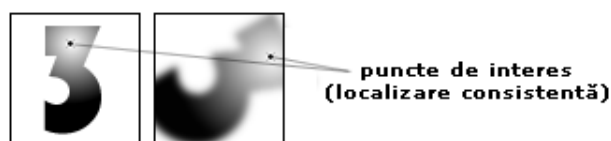


Figura 1.1: *Localizarea punctelor de interes*: rulând în mod independent algoritmul pe două imagini ale aceluiași obiect în situații diferite, se dorește ca punctele de interes să fie identificate în aceleași poziții relativ la obiect

2. *Descrierea*: fiecărui punct de interes determinat anterior îi sunt asociate o mulțime de date rezultate din analiza imaginii, astfel încât el să poată fi identificat cu un grad ridicat de individualitate în comparație cu restul punctelor de interes (Figura 1.2). În imagini diferite ale aceluiași obiect, vectorul obținut trebuie să fie invariant la modificări ale mediului (luminozitate) sau la transformări afine asupra obiectului (translații, rotații, scalări), pentru a asigura o recunoaștere adecvată (căutarea se realizează pe baza descriptorilor).

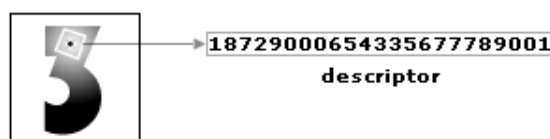


Figura 1.2: *Descrierea punctelor de interes*: Asocierea de informații pentru identificare, considerând vecinătatea punctului de interes.

Pentru fiecare dintre acești pași, există diferiți algoritmi, unii asigurând o ”acoperire” mai bună a obiectelor cu puncte de interes, alții concentrându-se pe stabilitatea trăsăturilor determinate sau pe eficiența computațională. Desigur, trebuie realizat un compromis astfel încât să se ajungă la o soluție acceptabilă pentru cât mai multe aplicații. De menționat că algoritmi de localizare și cei de descriere pot fi aleși în mod independent, dar o analiză a performanțelor nu poate fi realizată decât la nivelul efectului aplicării celor 2 pași, în mod secvențial.

Metodele ce au la bază trasături identificate prin puncte de interes pot fi folosite și în alte aplicații, nu doar cea a recunoașterii. Astfel de algoritmi pot fi aplicați, de exemplu, ca prim pas în reconstruirea unor modele 3D ale obiectelor din imagini sau video, pentru calibrarea aparatelor foto digitale sau pentru crearea de panorame din secvențe de imagini. Prin urmare, sunt de utilitate mare implementările cât mai generale, flexibile, care să poată fi utilizate într-o gamă largă de aplicații sau teste comparative. Lucrarea de față se referă la detaliile unei astfel de implementări.

*Capitolul 2* face o scurtă trecere în revistă a cercetărilor realizate în domeniu, punând accentul pe descrierea noțiunilor teoretice și a algoritmilor folosiți.

*Capitolul 3* prezintă detaliile legate de proiectarea aplicației propuse, descriind structura detaliată a modulelor unei platforme software pentru prelucrarea fluxurilor de imagini și a submodulelor ce implementează algoritmi pentru detecția de obiecte.

## Capitolul 2

---

# Noțiuni teoretice

---

### 2.1 Abordări ale temei în literatura de specialitate

În funcție de restricțiile aplicației practice în care este utilizată, recunoașterea obiectelor poate lua mai multe forme, de la simpla împărțire a imaginii în zone ce pot reprezenta obiecte (segmentare pe bază de culoare), la o recunoaștere completă, ce implică determinarea locației  $(x, y)$  a unui obiect, reconstituirea poziționării sale în spațiu (sau 2D în planul imaginii) și recunoașterea denumirii obiectului respectiv pe baza unor cunoștințe anterioare ale sistemului.

Oricare ar fi gradul de complexitate, la modul general se pune problema ca pornind de la o matrice de pixeli (imaginea), să fie determinată o submulțime a acestora care reprezintă un obiect. Fără a scădea din generalitate, considerăm că obiectul este dat de o regiune contiguă de pixeli din imaginea originală.

O abordare directă a problemei, presupunând că deținem o imagine a obiectului, ar fi căutarea tuturor pixelilor săi într-o altă imagine dată. Îmbunătățiri ale acestei metode, caracterizată de potrivirea unor "tipare" reprezentând obiectul în scene care îl conțin, au reprezentat începutul cercetării în domeniu (Figura 2.1a). Soluția (template matching), în forma ei inițială, este ineficientă computațional și sensibilă atât la modificări ale mediului în care dorim să realizăm recunoașterea (luminozitate, reflexii) cât și la ocluzionări parțiale ale obiectelor. Pentru obținerea unei oarecare invarianțe, a fost propusă corelarea nivelurilor de gri din diverse zone ale imaginii reprezentând obiectul, cu zone din imagini care se presupune că îl conțin. [Ballard and Brown, 1982, Goshtasby et al., 1984] Aceste studii sunt făcute în contextul sistemelor de stereo-vizualizare, unde scena este fotografiată simultan din perspective diferite și se dorește determinarea unor corespondențe între imagini, evitând o calibrare anterioară sau cunoașterea geometriei epipolare a sistemului. Mai recent, există variante care propun modificări ale metodei pentru a o putea rula în timp real [Cole et al., 2004].

template matching

Pentru a depăși o parte din problemele metodei anterioare, se pleacă de la observația că pentru recunoașterea unui obiect nu este nevoie de toți pixelii săi, ci doar de o parte din aceștia, ce definesc forma specifică a obiectului sau caracteristici importante ale acestuia. Se realizează o sintetizare a informației din imaginea originală, făcându-se primul pas înspre reprezentarea respectivului obiect într-un mod abstract. Abordarea recunoașterii obiectelor prin potrivirea unor astfel de trăsături abstracte (feature matching) este cea de-a doua direcție de cercetare în domeniu. (Figura 2.1b)

feature matching

Aplicarea algoritmilor de acest tip presupune macarea trăsăturilor din imagine ca puncte de interes, având ca informație minimală locația,  $(x_t, y_t)$ . Există desigur și posibilitatea stocării unor date suplimentare precum orientarea sau scala caracteristicii determinate.

În mod tradițional, trăsăturile alese pentru identificare și potrivire sunt muchii, colțuri sau contururi [Cheng and Huang, 1984, Ullman, 1979]. În momentul de față, sunt propuși algoritmi care realizează și identificarea unor alte structuri, precum petele luminoase sau întunecate (eng. blob) [Lowe, 2003, Bay et al., 2006].



Figura 2.1: (a): Potrivire bazată pe tipare (template matching), (b): Potrivire bazată pe trăsături (feature matching)

Avantajul acestei metode este că necesită mai puțină putere de calcul (operând asupra unui număr relativ restrâns de pixeli) și poate fi prin urmare aplicată cu ușurință în timp real. În plus, datorită faptului că se lucrează cu o reprezentare intermediară a obiectului (trăsături), metodele pot fi proiectate pentru a obține un grad ridicat de invarianță la anumiți parametri de mediu sau la transformări afine aplicate obiectului. Abordarea eșuează însă dacă nu se reușește o determinare repetabilă și consistentă a trăsăturilor unui obiect în imagini diferite. În acest caz, potrivirea nu are loc și obiectul nu este detectat.

Datorită flexibilității crescute și a rezultatelor foarte bune obținute în practică de către abordarea potrivirii bazate pe trăsături, lucrarea de față utilizează

această metodă pentru recunoașterea obiectelor.

## 2.2 Identificarea trăsăturilor

În identificarea trăsăturilor, se disting 2 metode mai importante, utilizate, cu unele adaptări, în cele mai multe dintre aplicațiile practice curente: Detectorul Harris, și SIFT (Scale Invariant Feature Transform). Aplicația propusă în lucrare utilizează o variantă îmbunătățită a algoritmului SIFT, adaptată pentru procesarea fluxurilor de imagini, în timp real. O parte a ideilor propuse inițial de Harris și Stephens pentru detectorul Harris sunt reluate în algoritmul SIFT, prin urmare considerăm utilă prezentarea ambelor metode.

### 2.2.1 Detectorul Harris

Prima abordare, propusă de Harris și Stephens, identifică în imagine colțurile și muchiile [Harris and Stephens, 1988]. Cei doi pornesc de la o observație anterioară a lui Moravec, care definește un colț ca fiind un pixel care nu se aseamănă cu pixelii din vecinătatea sa. Astfel, pe o suprafață uniformă, un pixel va avea valori apropiate de cele ale vecinilor săi; pe o muchie, în vecinătatea pixelului se vor identifica modificări mari relativ la valorile vecinilor perpendiculari pe direcția muchiei, dar modificări mici în direcția muchiei. Însă dacă pixelul aparține unei trăsături cu variații în toate direcțiile (un colț), atunci nici una dintre vecinătăți nu va fi similară pixelului. În [Harris and Stephens, 1988], formalizând matematic aceste observații, se definește noțiunea de autocorelație. Funcția de autocorelație măsoară modificările locale ale semnalului 2D (imaginea), folosind ferestre deplasate pe distanțe mici în vecinătatea punctului considerat. Fiind dată o deplasare  $(\Delta x, \Delta y)$  și un punct  $(x, y)$ , funcția de autocorelație este

$$E(\Delta x, \Delta y) = \sum_{x,y} w(x, y) [I(x + \Delta x, y + \Delta y) - I(x, y)]^2 \quad (2.1)$$

unde  $w(x, y)$  reprezintă funcția fereastră (și poate fi aleasă ca fiind o funcție nucleu rectangulară sau, pentru a reduce influența zgomotului, un nucleu Gaussian) iar  $I(\cdot, \cdot)$  este funcția imagine.

Imaginea din fereastra deplasată este aproximată prin dezvoltarea în serie Taylor, trunchiată la primii termeni,

$$I(x + \Delta x, y + \Delta y) = I(x, y) + \begin{bmatrix} I_x & I_y \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix} \quad (2.2)$$

$I_x$  și  $I_y$  fiind derivatele parțiale pe direcția  $x$ , respectiv  $y$ .

Înlocuind 2.2 în 2.1 și considerând  $\Delta x$  și  $\Delta y$  suficient de mici, obținem o ecuație de forma:

$$E(\Delta x, \Delta y) \cong \begin{bmatrix} \Delta x & \Delta y \end{bmatrix} M \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix}$$

M fiind o matrice  $2 \times 2$  calculată din derivatele locale parțiale ale imaginii,

$$M = \sum_{x,y} w(x,y) \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix}$$

Matricea M descrie structura locală a imaginii în vecinătatea pixelului considerat. Fie  $\lambda_1, \lambda_2$  valorile proprii ale acestei matrici. Există 3 cazuri care trebuie considerate:

1. Dacă atât  $\lambda_1$  cât și  $\lambda_2$  au valori mici, astfel încât funcția de autocorelație este plată (schimbări mici ale lui  $E(\Delta x, \Delta y)$  în orice direcție), zona din fereastra considerată este aproximativ uniformă.
2. Dacă o valoare proprie este mare iar cealaltă este mică, astfel încât funcția de autocorelație are forma unei trepte, atunci deplasările ferestrei într-o direcție (de-a lungul treptei) produc modificări mici ale lui E, iar deplasările pe o direcție ortogonală primeia produc modificări mari. Acest lucru indică o muchie.
3. Dacă valorile proprii sunt ambele mari, deplasările în orice direcție vor produce modificări mari ale lui E, indicând un colț.

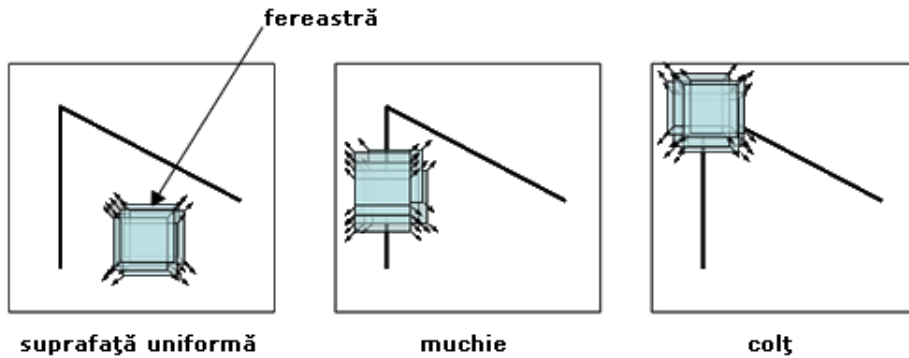


Figura 2.2: Detectorul Harris (muchii și colțuri)

Intuitiv, modul de operare al detectorului Harris este prezentat în Figura 2.2. Performanțele sale au fost analizate în detaliu [C.Schmid et al., 2000]. Concluzia studiului este că detectorul Harris este unul robust, putând fi aplicat cu succes

chiar și pe imagini zgomotoase și fiind invariant la rotații sau schimbări ale luminozității ambientale. Totuși, repetabilitatea rezultatelor sale scade drastic la schimbări ale perspectivei. O altă problemă a detectorului este că nu este invariant la modificările de scală ale obiectelor considerate. Acest poate fi observat cu ușurință în Figura 2.3.

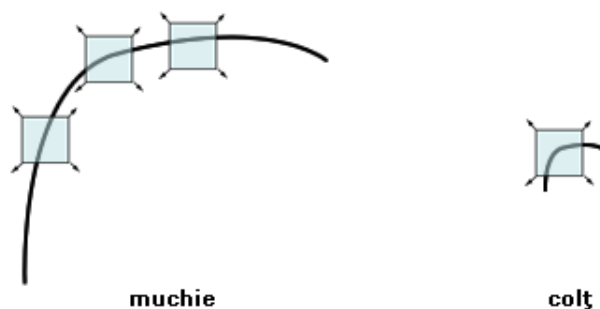


Figura 2.3: Modificarea scalei imaginii poate duce la rezultate diferite în cazul detectorului Harris

Au fost propuse și variante care să fie invariante la scalări (Harris-Laplacian), acestea fiind similare ca abordare cu cel de-al doilea detector important, SIFT.

### 2.2.2 Detectorul SIFT

SIFT (Scale Invariant Feature Transform) este un algoritm propus de Lowe [Lowe, 2003], care include ca pas intermediar detecția unor puncte de interes asimilate unor trăsături de tip "zonă luminoasă" sau "zonă întunecată". Prin construcția algoritmului, aceste zone sunt determinate pentru a fi invariante la scalări, rotații și parțial invariante la modificări ale luminozității și la transformări afine. Metoda aplică o filtrare în cascadă, pentru a asigura calitatea punctelor de interes determinate, dar și pentru a aplica operațiile intensive computațional doar acelor zone care trec unele teste inițiale. Pe lângă determinarea locației punctelor de interes, algoritmul SIFT propune și metode de descriere a acestora în mod individual, astfel încât să poată fi identificate cu probabilitate mare în imagini noi. Practic, fiecărui punct de interes îi este asociat un descriptor (vector caracteristic), calculat pe baza informațiilor imaginii în vecinătatea punctului de interes.

Aceste caracteristici fac SIFT ideal pentru aplicarea în zona recunoașterii obiectelor. Pentru aceasta, mai întâi se extrag trăsăturile SIFT pentru un set de imagini de referință ce reprezintă obiectele, stocând descriptorii rezultați într-o bază de date. Unei imagini noi, în care se dorește identificarea unuia dintre obiectele existente în baza de date, i se aplică același algoritm, iar descriptorii

punctelor de interes rezultate sunt comparați individual cu descriptorii din baza de date. Potrivirile între descriptori se fac pe baza distanței Euclidiene între vectori (nu se caută doar potriviri exacte). Totuși, într-o imagine aglomerată, multe trăsături din fundal nu vor avea corespondenți în baza de date, dând potriviri false, pe lângă cele corecte. Potrivirile corecte pot fi însă filtrate prin identificarea unor submulțimi de puncte de interes care sunt consistente cu aceeași localizare, scală și orientare a obiectului în noua imagine. Determinarea acestor clustere poate fi realizată eficient folosind transformata Hough.

### Localizarea punctelor de interes

Primul pas în determinarea punctelor de interes SIFT îl reprezintă detectarea locațiilor din imagine care sunt invariante la scalări, prin căutarea trăsăturilor stabile, folosind o funcție de scală cunoscută sub denumirea de spațiu al scalarilor (eng. scale space). Spațiul scalarilor pentru o imagine este definit de funcția  $L(x, y, \sigma)$ , obținută prin convoluția unui nucleu Gaussian  $G(x, y, \sigma)$  cu imaginea,  $I(x, y)$ . Pentru a obține scalări diferite, se variază  $\sigma$ :

$$L(x, y, \sigma) = G(x, y, \sigma) * I(x, y),$$

unde  $*$  reprezintă operația de convoluție, iar nucleul Gaussian  $G$  este dat de formula:

$$G(x, y, \sigma) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$$

Pentru a detecta punctele de interes stabile în spațiul scalarilor, Lowe propune determinarea extremelor locale ale funcției "diferență de nucleu Gauss cu scalări diferite", în convoluție cu imaginea,  $D(x, y, \sigma)$ . Aceasta poate fi calculată din diferența a două scalări separate de un factor  $k$ :

$$\begin{aligned} D(x, y, \sigma) &= (G(x, y, k\sigma) - G(x, y, \sigma)) * I(x, y) \\ &= L(x, y, k\sigma) - L(x, y, \sigma) \end{aligned}$$

Există mai multe motive pentru care a fost aleasă această funcție în mod particular. În primul rând, imaginile pentru care se aplică filtrul Gaussian (convoluție), trebuie oricum calculate în procesul de creare al spațiului scalarilor,  $D$  calculându-se în mod eficient prin scăderea imaginilor din două scale adiacente. În al doilea rând, diferența nucleelor Gauss (Difference of Gaussian, DOG) aproximează foarte bine Laplacianul Gaussian-ului,  $\sigma^2 \nabla^2 G$ . S-a demonstrat că extremele acestei funcții reprezintă trăsături foarte stabile ale imaginii, în comparație cu trăsăturile determinate cu alte funcții precum gradientul, Hessian-ul sau detectorul Harris.

Pentru a detecta extremele locale ale lui  $D$ , se realizează o eșantionare a funcției atât spațial  $(x, y)$ , cât și pentru parametrul de scală  $(\sigma)$ . Frecvența aleasă pentru eșantionare reprezintă un compromis între precizia localizării extremelor



și puterea de calcul necesară pentru determinarea lor. Astfel, o eșantionare cu frecvență mare duce la costuri mari din punct de vedere computațional, iar o frecvență mică duce la o precizie scăzută a algoritmului.

Fiecare punct eșantionat este comparat cu cei 8 vecini ai săi din imaginea curentă, și cei 9 vecini din scalările adiacente celei curente (Figura 2.4). Punctul este selectat doar dacă este mai mare sau mai mic comparativ cu toți vecinii săi. Această abordare se dovedește eficientă pentru că majoritatea punctelor sunt eliminate după doar câteva comparații.

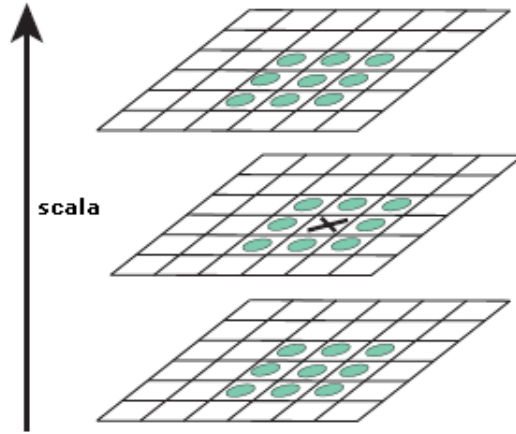


Figura 2.4: *SIFT: Detectarea minimelor și maximelor locale*; punctul central este comparat cu toți vecinii marcați ([Lowe, 2003])

O precizie crescută a localizării punctelor de interes se poate obține folosind o metodă de aproximare a poziționării maximului, prin interpolare. Astfel, se încearcă aproximarea punctelor eșantionate cu o funcție cuadrică, 3D. Practic, se realizează o dezvoltare în serie Taylor până la termenii de grad 2, a funcției  $D(x, y, \sigma)$ , translată astfel încât punctul eșantionat să fie în origine:

$$D(\mathbf{x}) = D + \frac{\partial D^T}{\partial \mathbf{x}} \mathbf{x} + \frac{1}{2} \mathbf{x}^T \frac{\partial^2 D}{\partial \mathbf{x}^2} \mathbf{x} \quad (2.3)$$

unde  $D$  și derivatele sale sunt evaluate în punctul de eșantionare iar  $\mathbf{x} = (x, y, \sigma)^T$  este deplasarea față de acest punct. Localizarea precisă a extremului,  $\hat{\mathbf{x}}$  este determinată prin derivarea ecuației 2.3 în raport cu  $\mathbf{x}$  și egalarea cu zero, rezultând

$$\hat{\mathbf{x}} = -\frac{\partial^2 D^{-1}}{\partial \mathbf{x}^2} \frac{\partial D}{\partial \mathbf{x}} \quad (2.4)$$

Pentru a elimina punctele care sunt maxime locale dar se află într-o regiune cu un contrast slab (fiind deci instabile), se vor reține doar acelea pentru care  $D(\hat{\mathbf{x}})$  este mai mare decât o valoare prag (Lowe alege valoarea de prag 0.03 pentru experimentele sale).

Totuși, pentru o stabilitate crescută, nu e suficientă îndepărtarea trăsăturilor cu un contrast slab. Funcția "diferență de nuclee Gauss" va avea un răspuns puternic de-a lungul muchiilor, chiar dacă locația respectivă este determinată imprecis, sensibilă la zgomotele din imagine. Pentru eliminarea acestor răspunsuri, se folosește o abordare bazată pe o matrice Hessiană  $2 \times 2$ , calculată în poziția și pentru scala punctului de interes:

$$\mathbf{H} = \begin{bmatrix} D_{xx} & D_{xy} \\ D_{xy} & D_{yy} \end{bmatrix} \quad (2.5)$$

Derivatele se estimează prin diferențele față de punctele eșantionate din vecinătate. Pentru eliminarea răspunsurilor de-a lungul muchiilor, se impune ca raportul valorilor proprii ale acestei matrici să fie sub o valoare prag (Lowe alege valoarea 10). Pentru că eliminarea se face în funcție de raportul valorilor proprii, nu este necesară calcularea individuală a acestora. În loc, se folosesc determinantul și urma matricii  $\mathbf{H}$ . Dacă notăm valorile proprii cu  $\lambda_1$  și  $\lambda_2$ , atunci:

$$\begin{aligned} Tr(\mathbf{H}) &= D_{xx} + D_{yy} = \lambda_1 + \lambda_2 \\ Det(\mathbf{H}) &= D_{xx}D_{yy} - (D_{xy})^2 = \lambda_1\lambda_2 \end{aligned}$$

Considerăm arbitrar  $\lambda_1 > \lambda_2$  și notăm raportul valorilor proprii cu  $r$ , astfel încât  $\lambda_1 = r\lambda_2$ . Atunci, avem:

$$\frac{Tr(\mathbf{H})^2}{Det(\mathbf{H})} = \frac{(\lambda_1 + \lambda_2)^2}{\lambda_1\lambda_2} = \frac{(r\lambda_2 + \lambda_2)^2}{r\lambda_2^2} = \frac{(r+1)^2}{r} \quad (2.6)$$

Prin urmare, pentru a impune  $r$  ca valoare prag, trebuie verificată doar condiția:

$$\frac{Tr(\mathbf{H})^2}{Det(\mathbf{H})} < \frac{(r+1)^2}{r} \quad (2.7)$$

### Descrierea punctelor de interes

După stabilirea precisă a locației unei trăsături, se dorește asocierea unui vector caracteristic (descriptor), astfel încât ea să poată fi identificată și în alte imagini.

Primul pas constă în atribuirea unei orientări fiecărui punct de interes, astfel încât descriptorul să poată fi reprezentat relativ la orientarea sa locală. Pentru operațiile care urmează, se alege imaginea filtrată cu nucleu Gaussian având scala cât mai apropiată de cea determinată prin interpolare pentru punctul de interes. Folosind această imagine, se calculează norma și orientarea gradientului într-un număr de puncte din vecinătatea punctului de interes. Valorile obținute sunt organizate într-o histogramă a orientărilor, cu 36 de intervale. Fiecare vector gradient este adăugat în intervalul corespunzător orientării sale și ponderat cu valoarea normei. Vârfurile din histogramă corespund orientărilor dominante ale gradientilor locali. Cel mai mare vârf este ales ca orientare a punctului de interes.

Dacă cel de-al doilea vârf al histogramei este comparabil ca mărime, atunci în aceeași poziție din imagine se va crea un al doilea punct de interes, care să aibă orientarea acestui al doilea vârf.

Parametrii de poziție, scală și orientare determinați până acum stabilesc un sistem de coordonate 2D, local punctului de interes, față de care se realizează descrierea acestuia.

În vecinătatea determinată de sistemul local de coordonate al punctului de interes se realizează o eșantionare, iar în punctele alese se calculează norma și orientarea gradientului, relativ la orientarea punctului de interes (Figura 2.5). Normele sunt ponderate de o funcție Gaussiană, (cercul din figură) cu  $\sigma$  de 1.5 ori mai mare decât dimensiunea vecinătății considerate (în experimente  $16 \times 16$  pixeli). Vecinătatea este împărțită apoi într-un număr de subregiuni care nu se suprapun (16 regiuni de  $4 \times 4$  pixeli). Pentru fiecare subregiune, valorile gradientelor sunt acumulate într-o histogramă, similară celei folosite anterior. Dacă o histogramă discretizează unghiurile de orientare în 8 valori posibile, descriptorul punctului de interes va conține  $4 \times 4 \times 8 = 128$  elemente, obținute prin concatenarea valorilor din toate histogramele. (Figura 2.5).

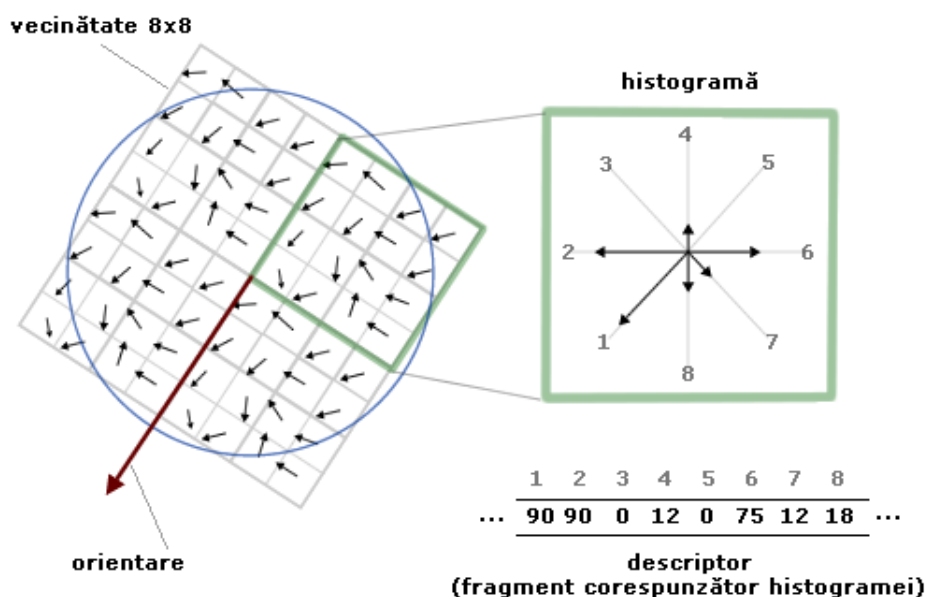


Figura 2.5: *SIFT: Procesul de determinare al descriptorului*; pentru claritatea reprezentării, a fost aleasă o vecinătate 8x8 a punctului de interes. Algoritmul folosește vecinătăți 16x16.

### 2.3 Identificarea trăsăturilor în timp real

Dintre cele 2 metode prezentate, SIFT se remarcă datorită invarianței la un număr mare de parametri precum și datorită stabilității punctelor de interes determinate. Totuși, este evident că aplicarea algoritmului SIFT implică un număr mult mai mare de operații în comparație cu alți detectori (Harris). Deoarece majoritatea aplicațiilor îl vor utiliza doar ca pas intermediar, se pune problema unei post-procesări a punctelor de interes (de exemplu, pentru a identifica obiecte) și se dorește ca ansamblul algoritmilor de procesare să ruleze în timp real. În forma prezentată, SIFT poate prelucra în jur de 5 frame-uri (de dimensiune  $650 \times 315$ ) pe secundă. Prin urmare, se justifică o căutare a unor îmbunătățiri care să determine o scădere a timpului de prelucrare, fără a afecta calitatea rezultatelor finale.

SURF (Speeded-Up Robust Features) este una dintre soluțiile propuse în acest sens, fiind și metoda utilizată de aplicația descrisă în această lucrare. Deoarece majoritatea pașilor urmați sunt identici cu cei ai algoritmului SIFT, vom prezenta în continuare doar elementele noi pe care le aduce în comparație cu acesta.

## Capitolul 3

---

# Chapter's title

---

... some text ...

Some reference



---

# Bibliografie

---

- [Ballard and Brown, 1982] Ballard, D. and Brown, C. (1982). *Computer Vision*. Prentice-Hall, Englewood Cliffs. [cited at p. 5]
- [Bay et al., 2006] Bay, H., Tuytelaars, T., and Gool, L. V. (2006). Surf: Speeded-up robust features. In *ECCV*, pages 404–417. [cited at p. 6]
- [Cheng and Huang, 1984] Cheng, J. and Huang, T. (1984). Image registration by matching relational structures. *Pattern Recognition*, 17(1):149–159. [cited at p. 6]
- [Cole et al., 2004] Cole, L., Austin, D., and Cole, L. (2004). Visual object recognition using template matching. In *Proceedings of Australian Conference on Robotics and Automation*. [cited at p. 5]
- [C.Schmid et al., 2000] C.Schmid, Mohr, R., and C.Bauckhage (2000). Evaluation of interest point detectors. *International Journal of Computer Vision*, 37(2):151–172. [cited at p. 8]
- [Goshtasby et al., 1984] Goshtasby, A., Gage, S., and Bartholic, J. (1984). A two-stage cross correlation approach to template matching. *IEEE Transactions, Pattern Analysis & Machine Intelligence*, 6(3):374–378. [cited at p. 5]
- [Harris and Stephens, 1988] Harris, C. and Stephens, M. (1988). A combined corner and edge detection. In *4th Alvey Vision Conference*, pages 147–151. [cited at p. 7]
- [Lowe, 2003] Lowe, D. G. (2003). Distinctive image features from scale-invariant keypoints. [cited at p. 6, 9, 11, 24]
- [Ullman, 1979] Ullman, S. (1979). *The Interpretation of Visual Motion*. MIT Press, Cambridge, MA. [cited at p. 6]





# Anexe



## **Anexa A**

---

## **Anexa 1**

---

... some text ...



---

## Lista de simboluri și prescurtări

---

Prescurtare	Descriere	Definiție
SIFT	Scale Invariant Feature Transform	page 9
DOG	Difference of Gaussian	page 10
SURF	Speeded-Up Robust Features	page 14

---

## Listă de figuri

---

1.1	<i>Localizarea punctelor de interes:</i> rulând în mod independent algoritmul pe două imagini ale aceluiași obiect în situații diferite, se dorește ca punctele de interes să fie identificate în aceleași poziții relativ la obiect . . . . .	3
1.2	<i>Descrierea punctelor de interes:</i> Asocierea de informații pentru identificare, considerând vecinătatea punctului de interes. . . . .	3
2.1	(a): Potrivire bazată pe tipare (template matching), (b): Potrivire bazată pe trăsături (feature matching) . . . . .	6
2.2	Detectorul Harris (muchii și colțuri) . . . . .	8
2.3	Modificarea scalei imaginii poate duce la rezultate diferite în cazul detectorului Harris . . . . .	9
2.4	<i>SIFT: Detectarea minimelor și maximelor locale;</i> punctul central este comparat cu toți vecinii marcați ([Lowe, 2003]) . . . . .	11
2.5	<i>SIFT: Procesul de determinare al descriptorului;</i> pentru claritatea reprezentării, a fost aleasă o vecinătate 8x8 a punctului de interes. Algoritmul folosește vecinătăți 16x16. . . . .	13

---

## Listă de tabele

---





---

# Glosar

---

blob, 6

detector Harris, 7, 10  
DOG, 10

feature matching, 6

nucleu gaussian, 10

scale space, 10  
SIFT, 9, 14  
spațiul scalărilor, 10  
SURF, 14

template matching, 5  
text, 15