

# Detecția mâinii într-o scenă încărcată de obiecte, în scopul recunoașterii gesturilor

Delia Mitrea  
Universitatea Tehnică Cluj-Napoca,  
Facultatea de Automatică și  
Calculatoare,  
str. G. Baritiu nr.26-28, Cluj-Napoca  
delia.mitrea@cs.utcluj.ro

Sergiu Nedevschi  
Universitatea Tehnică Cluj-Napoca,  
Facultatea de Automatică și  
Calculatoare,  
str. G. Baritiu nr.26-28, Cluj-Napoca  
sergiu.nedevschi@cs.utcluj.ro

Dorian Gorgan  
Universitatea Tehnică Cluj-Napoca,  
Facultatea de Automatică și  
Calculatoare  
str. G. Baritiu nr.26-28, Cluj-Napoca  
dorian.gorgan@cs.utcluj.ro

## REZUMAT

Detecția vizuală a mâinii, într-o scenă încărcată de obiecte, este o etapă importantă în scopul recunoașterii gesturilor, în contextul unei interfețe multimodale ce implica mijloace de interacțiune naturală cu utilizatorul. Problema este aceea de a distinge, în timp real, obiectul cautat fața de fond și fața de celelalte obiecte din scenă, indiferent de poziția sa și de condițiile de iluminare, respectiv de acoperirile cu alte obiecte. Se va adopta o metodă robustă bazată pe caracteristici de textură, culoare și formă.

## Categorii și descriptori ai subiectelor

**H.1.2. [User/Machine Systems]** Human Factors, Human Information Processing

**I.4 [Image Processing and Computer Vision]:** Segmentation - *edge and feature detection, pixel classification*, Feature Measurement - *texture*, Scene Analysis - *shape, object recognition*

<http://www.acm.org/class/1998/>

## Termeni generali

*Documentation, Theory, Human Factors, Design, Algorithms, Performance*

## Cuvinte-cheie

decor multi-modal, scenă naturală, invarianță vizuală, detecția mâinii, gesturi

## 1. INTRODUCERE

Scopul acestei lucrări este acela de a dezbate problema detecției vizuale a mâinii într-o scenă încărcată de obiecte și de a propune o metodă robustă pentru soluționarea acesteia.

Detecția vizuală a mâinii prezintă o deosebită importanță în contextul recunoașterii gesturilor, constituind prima etapă a acestui proces complex care implică, în primul rând, localizarea mâinii, urmând ca abia apoi să se facă demersurile necesare pentru identificarea unor poziții și expresii specifice ale mâinii, respectiv

urmărirea mișcării acesteia și recunoașterea traiectoriei parcurse – ceea ce conferă o caracterizare completă a gestului. Identificarea și localizarea mâinii pe cale vizuală, în comparație cu alte modalități de localizare a mâinii (senzori, termo-viziune artificială) prezintă următoarele avantaje:

1. metoda este apropiată de percepția umană, procesul de recunoaștere umană a gesturilor fiind unul pur vizual
2. costul echipamentelor hardware necesare (cameră video, PC) este mai redus decât în celelalte cazuri.

Detecția mâinii într-o scenă încărcată de obiecte este o instanță a problemei generale de detecție a obiectelor în scenă. Se vor discuta aspectele legate de această problemă, dificultățile cu care se confruntă literatura internațională de specialitate în a o soluționa, apoi se va trece la cazul particular al detecției mâinii. Se vor introduce și defini termeni noi, specifici, întâlniți în literatura de specialitate, cărora li se va găsi un corespondent adecvat în limba română – decor multimodal, invarianță vizuală, scenă încărcată de obiecte. Se vor trece apoi în revistă metodele existente în literatura de specialitate pentru detecția obiectelor, respectiv a mâinii și se vor identifica principalele direcții ale acestor metode. Metoda propusă va urmări găsirea acelor caracteristici (trăsături) ale obiectelor care sunt invariante la schimbarea condițiilor precum poziția obiectului, direcția de iluminare, distanța față de observator; de asemenea, se va adopta un algoritm care să furnizeze rezultatul așteptat în timp real – aceasta din urmă fiind o necesitate imperativă în cazul interfețelor om-calculator bazate pe interacțiune naturală. Se dorește și obținerea unui anumit grad de generalitate a metodei, pentru a face posibilă utilizarea acesteia atât în cazul unei aplicații comandată prin gesturi, cât și în cazul unui translator al limbajului mimico-gestual.

## 2. DEFINIRE CONCEPTE

**Fundal:** planul, sau suprafața generată de mai multe plane, care mărginește, în extremitatea opusă observatorului, spațiul scenei.

**Decor:** obiectele din planul secundar, care diferă față de obiectul de interes; termenul originar este acela de *background* -care desemnează atât fundalul, cât și decorul;

**Decor multi-modal:** termenul originar fiind acela de *multi-modal background* [5], conceptul se referă la un decor în care se

petrec mișcări repetitive ale obiectelor, precum mișcările repetate ale unei ape sau mișcările crengilor arborilor datorate vântului;

**Scenă:** desemnează spațiul vizual ce constă din fundal, decor și obiectul de interes;

**Scenă naturală:** desemnează scena obișnuită a obiectelor într-o situație reală, fiind caracterizată prin decor multimodal, având loc atât mișcări periodice sau izolate ale elementelor prezente permanent în scenă, cât și apariția și dispariția unor obiecte din scenă, respectiv acoperiri între obiectele scenei; termenul de proveniență este acela de *cluttered scene (background)* [2],[8] – scenă încărcată de obiecte;

**Detectia obiectelor** implică formularea unui răspuns pozitiv sau negativ privind prezența într-o scenă naturală a unui obiect aparținând unei anumite categorii. Spre deosebire de recunoaștere, care presupune identificarea obiectului căutat – fața unei anumite persoane, mâna într-o poziție specifică, un anumit automobil-detectia obiectelor presupune găsirea unui obiect aparținând unei clase mai largi – fețe umane, mâini, automobile.

**Invariantă vizuală:** capacitatea unor elemente (textură, formă, culoare) sau părți componente ale obiectelor de a nu-și modifica proprietățile vizuale în momentul modificării unor condiții precum schimbarea orientării obiectului relativ la observator, schimbarea direcției de iluminare, modificarea distanței obiectului față de observator sau a rezoluției imaginii, acoperirile cu alte obiecte; termenul originar este adjectivul *view-invariant* [6]

**Textură:** Se referă la aranjamentul regulat al elementelor vizuale dintr-o anumită regiune a scenei, aranjament ce poate fi caracterizat prin mărimi statistice.

**Texton:** microstructura (elementul vizual) fundamental prin repetarea căruia rezultă textura

### 3. DETECȚIA OBIECTELOR PRIN TEHNICI VIZUALE. PARTICULARIZARE: DETECȚIA MĂINII

#### 3.1 Problema generală a detecției obiectelor prin tehnici vizuale

Întrucât, așa cum rezultă din definiția dată în paragraful anterior, detecția obiectelor în imagini înseamnă căutarea în imagine a acelor obiecte aparținând unor categorii precizate la un moment dat; astfel, problema detecției implică clasificare, recunoaștere, și, deci, un anumit grad de complexitate.

Demersul necesar detecției obiectelor poate fi împărțit în două etape:

- separarea regiunilor distincte din scenă (segmentarea scenei)
- localizarea obiectului căutat

Prima etapă, de segmentare a scenei, implică sesizarea granițelor dintre regiunile distincte ale imaginii care pot să corespundă obiectelor distincte din scenă. Aici intervin, de obicei, metodele bazate pe detecție de muchii, respectiv de contururi. Se vor pune în evidență, așadar, trecerile de la regiuni cu anumite caracteristici, la regiuni cu caracteristici diferite. Complexitatea scenelor din lumea reală și varietatea tipurilor de obiecte impun

metode sofisticate de detecție a granițelor între obiecte - simple treceri de la un nivel de intensitate la un alt nivel de intensitate, treceri de la o textură la o alta textură, contururi iluzorii.

Cea de-a doua etapă implică localizarea obiectului căutat, despre care se știe că are anumite caracteristici: culoare, formă, textură, sau grupare de regiuni cu caracteristici bine precizate. De obicei, această etapă este precedată de una de învățare, în care se rețin caracteristicile definitorii ale obiectului căutat, dintr-un set de imagini de antrenament.

Complexitatea procesului de separare rezidă în discontinuitatea conturilor, datorată acoperirilor dintre obiecte, schimbarea aparenței obiectului căutat la modificarea poziției, a condițiilor de iluminare, a distanței față de observator, dinamicii scenelor naturale, decorurilor multimodale. Percepția umană se bazează pe combinarea caracteristicilor de textură, formă, culoare, pe memorie, capacitate deductivă. Această combinație de factori trebuie integrată într-un sistem computațional, astfel încât acesta să funcționeze în timp real. Așadar, se vor căuta acele caracteristici ale obiectelor care sunt invariante la schimbarea condițiilor externe, iar pe de altă parte, se vor implementa metode și algoritmi performanți, cât mai rapizi, în scopul recunoașterii obiectului în condiții variate.

Metodele generale de detecție a obiectelor indiferent de clasa de apartenență, implică, până în momentul actual, o complexitate computațională deosebit de ridicată și o dimensiune foarte mare a setului de antrenament. De aceea, în cazul metodelor robuste, se preferă căutarea doar a unui anumit obiect specific, după un număr cât mai mic, dar suficient, de caracteristici, invariante la modificarea condițiilor externe.

În ceea ce privește detecția mâinii, se pot remarca următoarele caracteristici invariante: textura și culoarea pielii, forma 3D specifică, dependentă de câteva poziții de bază ale acesteia. Forma, mărimea, precum și uniformitatea culorii și a texturii, lipsa altor elemente componente precum ochi, ochelari, nas, gura, mustața, păr, contribuie la distincția mâinii față de alte obiecte având aceeași textură și culoare – cea a pielii, precum fața umană.

#### 3.2 Metode existente in literatura de specialitate

Pentru detecția muchiilor, conturilor și regiunilor distincte din imagine, metodele existente variază, de la detecția de muchii pe baza nivelurilor de intensitate - metoda Gradientului sau a Laplacianului [13], la metode mai complexe, precum prăguirea (*thresholding*) pe baza nivelelor de intensitate ale pixelilor [9], extragerea fondului [9], metode bazate pe textură [3], iar metodele pentru extragerea conturului, de la transformata Hough [13], la modele active [7].

În ceea ce privește localizarea obiectelor aparținând unor categorii specificate, există următoarele modalități de soluționare a problemei: detecție de contur, potrivirea formelor 2D ale acestor contururi [1], metode robuste, de timp real, bazate pe determinarea distanței dintre semnături ce surprind caracteristici de textură, culoare și formă a obiectului căutat [14], metode ce considera obiectul ca o constelație de părți și utilizează clasificatori bazați pe reguli structurale, precum și reprezentări probabilistice care țin cont de scală, formă și ocluziile dintre obiecte [4], detecția, independentă de orientare și scală, a unor obiecte aparținând unor clase multiple [14]. Multe metode

generice, ce detectează, de obicei, mai multe clase de obiecte, utilizează clasificatori multipli de tip ada-boost, ce determină trăsăturile cele mai relevante pentru clasificare [14]. O altă categorie de metode, ce urmăresc detecția obiectului indiferent de poziția sa, sau – în cazul fețelor umane, de expresia feței, utilizează așa-numitele modele de aparență, ce surprind variațiile posibile ale acestor obiecte și utilizează clasificatori de tipul rețelelor neuronale, lanțurilor Markov ascunse, clasificatorul Bayesian. Aceste metode implică o complexitate computațională deosebit de ridicată.

Un interes deosebit îl prezintă metodele robuste de localizare a obiectelor în mișcare, precum metoda [5], ce utilizează împărțirea imaginii în blocuri și compararea blocurilor corespondente din cadre consecutive ale unei secvențe video, pe baza unei trăsături statistice, ce caracterizează textura aceluia bloc, anume parametrul LBP (*Local Binary Pattern*):

$$LBP(x_c, y_c) = \sum_{p=0}^{P-1} s(g_p - g_c)^2 p \quad (1)$$

De un interes aparte este și metoda dezvoltată de renumitul sistem „Pffinder” – de urmărire în timp real a persoanelor, care definește regiuni cu caracteristici eterogene numite „blobs” și calculează probabilitatea de apariție a unei astfel de regiuni prin formula gaussiană de probabilitate. [9]

În ceea ce privește metodele specifice de detecție a mâinii, se preferă metode robuste, de tipul celor ce utilizează trăsături spectrale, extrase după calculul transformelor Fourier, Gabor sau Wavelet ale imaginii [8]; metode bazate pe histograme de culori sau statistici ce caracterizează textura pielii [2], metode ce consideră mâna ca o constelație de părți, metode asemănătoare cu cele potrivite pentru detecția fețelor: bazate pe trăsături, pe potrivirea șabloanelor, pe modele de aparență, pe culori, pe detecția mișcării [6].

### 3.3 Metoda propusă pentru detecția mâinii

Metoda propusă urmărește localizarea mâinii în imagini reprezentând scene naturale de tipuri variate, cu decor multimodal, adecvate situațiilor specifice pentru recunoașterea gesturilor, care conțin mâna umană ca element esențial, dar și alte obiecte din decorul care completează imaginea – cel mai frecvent întâlnite fiind elementele din partea superioară a corpului uman.

Elementul vizual de bază care se ia în considerare în segmentare este textura, care, din punct de vedere tehnic, reprezintă aranjamentul regulat al nivelurilor de intensitate ale pixelilor dintr-o anumită regiune. Dincolo de această definiție clasică, textura 3D poate fi privită și ca o configurație bine definită de microstructuri numite textoni, acestea reprezentând elementele fundamentale ale texturii. Aceste microstructuri sunt de următoarele tipuri: muchie (*edge*), creastă (*ridge*), pată (*spot*), undă (*wave*), undișoară (*ripple*). Conform acestei viziuni, caracterizarea texturii se va face prin histograma textonilor; procesul de construire a histogramei textonilor implică următorii pași [11]:

- calculul unui vector de trăsături pentru fiecare pixel, ce va conține rezultatele aplicării nucleelor de convoluție Laws [3] respectiv nucleul de convoluție Laplacian al Gaussianului [13]

- formarea textonilor, prin gruparea pixelilor cu caracteristici similare, folosind metoda numita *k-means clustering* [13]

- marcarea fiecărui pixel cu eticheta textonului căruia îi aparține
- construirea histogramei textonilor, care memorează numărul de etichete aparținând fiecărei clase din regiunea pentru care se face caracterizarea texturii

Pentru ca analiza să fie independentă de orientarea texturii, respectiv de condițiile de iluminare, se va considera aceeași regiune de material în imagini diferite, reprezentând condiții diferite de orientare, respectiv de iluminare. Vectorul de trăsături pentru un pixel va rezulta din concatenarea vectorilor de trăsături ai pixelului respectiv calculat în fiecare astfel de imagine. Se va obține, în acest mod, o caracterizare a texturii general valabilă, independentă de orientarea, respectiv direcția de iluminare a materialului.

În scopul segmentării scenei în obiecte, se va urmări determinarea granițelor dintre texturile care compun scena în felul următor:

- se va împărți imaginea în blocuri de dimensiuni mai mici decât 10 pixeli;

- se va calcula histograma textonilor pentru fiecare bloc

- se vor calcula distanțele dintre histogramele blocurilor

$$\chi^2(h_1, h_2) = \frac{1}{2} \sum_{n=1}^{\#bins} \frac{(h_1(n) - h_2(n))^2}{h_1(n) + h_2(n)} \quad (2)$$

învecinate, folosind distanța  $\chi^2$ :

- se va considera că există o trecere de la o textură la o altă textură, respectiv de la fundal la obiectele scenei și invers doar dacă distanța  $\chi^2$  este mai mare decât un prag stabilit în felul următor:

$$\text{Prag} = (\chi_{\min}^2 + \chi_{\max}^2) / 2 + \sigma_{\chi}^2 \quad (3)$$

unde  $\chi_{\min}^2$  și  $\chi_{\max}^2$  reprezintă valorile minimă, respectiv maximă a distanțelor calculate între toate blocurile adiacente ale imaginii, de la stânga spre dreapta,  $\sigma_{\chi}^2$  reprezintă abaterea medie pătratică a acestor distanțe.

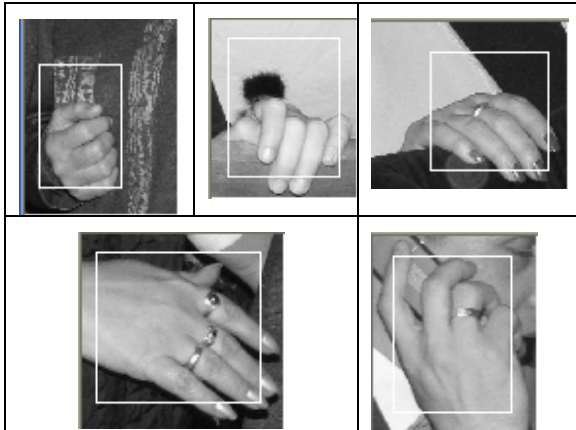
Se vor marca cu 0 regiunile aparținând fundalului și cu 1 regiunile din interiorul obiectelor.

Se vor căuta doar acele obiecte pentru care textura din interiorul blocurilor marcate corespunde cu textura din setul de antrenament, în cazul de față textura mâinii. Recunoașterea texturii se va face tot prin intermediul metodei bazate pe textoni, în varianta independentă de orientare și condiții de iluminare, utilizând ca metrică distanța dintre histograme.

Pentru a perfecționa rata de recunoaștere, se vor lua în considerare și caracteristici de formă și culoare – histograma de culori sau semnături RGB ale imaginii, pentru caracterizarea culorii, respectiv spectrul de formă ISS (Image Shape Spectrum) [12], a cărui valoare într-un punct caracterizează curbura locală a suprafeței de intensitate [12]. Astfel, setul de antrenament va conține imagini reprezentând mâna în câteva poziții de bază ale acesteia.

## 4. REZULTATE EXPERIMENTALE

Tabelul 1: Cazuri de detecție a mâinii



Tabelul de mai sus ilustrează câteva rezultate experimentale obținute în urma aplicării metodei descrise în secțiunea 3.3 – detecția mâinii pe baza informațiilor de textură. Mâna este surprinsă în diverse poziții ale sale, prin intermediul unui dreptunghi. Decorul este compus din elemente de îmbrăcăminte umană sau din accesorii precum inelele, alte părți ale corpului precum fața, obiecte ținute în mână (telefon). S-au utilizat imagini cu niveluri de gri reprezentate în formatul 8 biți/pixel, preluate la o rezoluție de 5 MP, fiind împărțite în blocuri de dimensiuni 10, 7, 5, respectiv 3 pixeli.

## 5. CONCLUZII

Metoda este potrivită pentru detecția mâinii în scene naturale, fiind capabilă de a localiza obiectul de interes într-un decor variat. Pentru perfecționarea ratei de succes, ne propunem utilizarea și a unor alte informații precum cele de formă și culoare, precum și detecția exactă a conturului mâinii prin intermediul unor modele active. În ceea ce privește viteza de execuție, metoda va fi comparată cu cea descrisă în [5], special concepută pentru detecția, pe baza texturii, a obiectelor aflate în mișcare, în timp real, dar care are drept inconvenient lipsa capacității de a detecta obiectele ce nu-și modifică poziția între două cadre succesive ale unei secvențe video.

## 6. REFERINȚE

- [1] Belongie, S., Malik, J. and Puzicha, J., *Shape matching and object recognition using shape contexts*, IEEE PAMI, vol. 24, no. 24, April 2002
- [2] de Campos, Teofilo Emidio, University of Oxford, *3D Hand and Object Tracking for Intention Recognition*, Transfer Report, 2003
- [3] Davis, Larry, Department of Computer Sciences, University of Texas at Austin, Austin, *Image Texture Analysis Techniques – A Survey*
- [4] Fergus, R; Perona, P., Zisserman, A., *Object class recognition by unsupervised scale-invariant learning*, CVPR, 2003
- [5] Heikkila, M., Pietikainen M., J. Heikkila, Department of Electrical and Information Engineering, University of Oulu, Finland, *A Texture-based Method for Detecting Moving Objects*, 2004
- [6] Hsuan, Ming Yang, Honda Research Institute, Mountain View, California, USA, *Advances in Face Processing: Detection*, ICPR 2004
- [7] Kass, M., Witkin, A., Terzopoulos, D., *Snakes: Active contour models*, International Journal of Computer Vision. v. 1, n. 4, pp. 321-331, 1987
- [8] Kolsch, M., Turk, M., Department of Computer Science, University of California, Santa Barbara *Robust Hand Detection*, 2003
- [9] Lombardi, Paolo, Dipartimento di Informatica e Sistemistica, Università di Pavia, Italy, *A survey on pedestrian detection for autonomous driving systems*, 2001
- [10] Mitrea, Delia, Gorgan, Dorian, *Comportament Adaptiv pentru Interfete Utilizator Avansate prin Intermediul Gesturilor*, Volumul de lucrari al Primei Conferințe Nationale de Interacțiune Om-Calculator –RoCHI 2004, Universitatea Politehnica Bucuresti 23-24 Septembrie 2004, Editura Printech, Bucuresti, 2004, ISBN 973-718-053-4, pg. 33-42
- [11] Mitrea Delia, Nedeveschi Sergiu, Technical University of Cluj-Napoca, *Road Quality Evaluation and Road Material Recognition using 3D textons*, 8-th IEEE International Conference of Intelligent Engineering Systems, INES 2004, Cluj-Napoca, Romania, septembrie 2004 Proceedings, Editura U.T.Press Cuj-Napoca, pg.236-241
- [12] Nastar, Chahab, Mitschke, Matthias, INRIA, France, *Real-Time Face Recognition using Feature Combination*, Proceedings of the Third IEEE International Conference on Automatic Face and Gesture Recognition, 14-16 April 1998, Nara, Japan
- [13] Nedeveschi, Sergiu, Editura Albastra, Cluj-Napoca, 1998: *Prelucrarea imaginilor si Recunoasterea Formelor*
- [14] Torralba, A., Murphy, K.P., and Freeman, W.T., *Sharing visual features for multiclass and multiview object detection*, CVPR'2004