

# RAPORT STIINTIFIC FINAL

Proiect PN-III-P4-PCE-2021-1959

Acronim: [FIGHS](#)

*Metode Detaliat de Reconstructie Tridimensională a Persoanelor*

Institutul de Matematica 'Simion Stoilow' al Academiei Romane

*Grupul de Cercetari in Vedere Artificiala si Invatare Computationala*

Director de proiect: Prof-univ. dr. Cristian Sminchisescu

## Impactul Sintetic al Rezultatelor Obținute:

Impactul estimat al rezultatelor obținute, cu sublinierea celui mai semnificativ rezultat obținut.

- Articolul asociat temei centrale a proiectului acceptat cu revizii minore la IEEE Transactions on Pattern Analysis and Machine Intelligence, 2025, IF=20.8.
- Website cu datele unice continand interactiuni umane si metodologia dezvoltate in proiect <https://ci3d.imar.ro/>. Peste 400 de grupuri de cercetare majore din Europa, Statele Unite si Asia s-au inregistrat si folosesc activ datele noastre pentru cercetare si publicatii.

## Obiective:

Obiectivul proiectului "Metode Detaliat de Reconstructie Tridimensională a Persoanelor" (FIGHS) este acela de a avansa metodologia fundamentala pentru analiza tridimensională detaliata a mai multor persoane in imagini si video. Desi s-au facut progrese semnificative in domeniul perceptiei 3d a persoanei, nici un sistem existent nu poate analiza la nivel de detaliu scene continand mai multe persoane, potential acoperindu-se unele pe altele in imagine si implicate in interactiuni. In acest scop, propunerea noastra tinteste (a) sa creeze baze de date state-of-the-art pe scara larga cu imagini ground-truth sincronizate cu reprezentarea tridimensională a formei si posturii corpului pentru antrenarea si evaluarea modelelor pentru multiple persoane implicate in interactiuni, ca si (b) tehnici de meta-invatare auto-supervizata, modelare si optimizare care sa sprijine reconstructia in detaliu a scenelor complexe care contin oameni interactionand. Un asemenea sistem de succes de perceptie 3d a persoanelor ar fi folositor unor audiente si domenii largi din sanatate, comunicatii, divertisment, analiza performantei sportivilor, robotics si masini care se conduc singure, sau productie. De asemenea, un astfel de sistem ar sprijini o analiza detaliata extensiva si precisa a masurilor de distantare sociala in timpul unei pandemii.

## **Activitati stiintifice si realizari:**

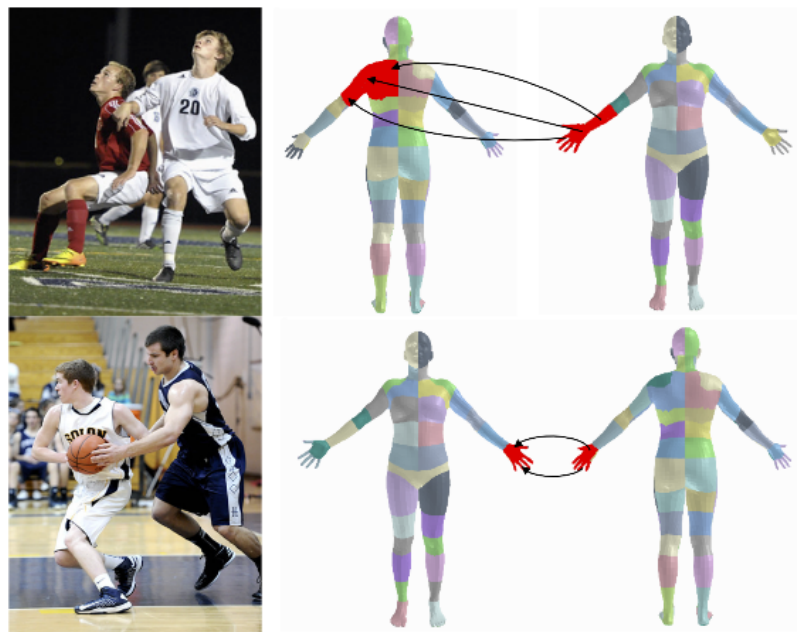
Înțelegerea tridimensională a interacțiunilor umane este fundamentală pentru analiza detaliată a scenei 3d și pentru modelarea comportamentului uman. Cele mai multe dintre modelele existente se concentrează pe analiza unei singure persoane izolate, iar acelea care procesează mai multe persoane pun accentul pe rezolvarea asocierii între multiple persoane, mai degrabă decât să înțeleagă interacțiunile dintre ele. Acest lucru duce la estimări care, chiar și atunci când sunt impresionante în ceea ce privește postura și forma plauzibile de la distanță, ratează esența evenimentului la un scrutin de detaliu, atunci când, de exemplu, două reconstrucții nu reușesc să surprindă contactul în timpul unei strângeri de mână, o atingere pe umăr, sau o îmbrățișare. Astfel de interacțiuni sunt deosebit de dificil de rezolvat, deoarece efectele lor se compun: pe de o parte, incertitudinea privind adâncimea și forma corpului ar putea duce la compensare prin împingerea membrelor în față sau mai departe de poziția lor de adevăr la sol, atunci când se face inferența 3d din imaginile monoculare; pe de altă parte, ocluzia parțială și detaliul (rezoluția) relativ limitate pentru zonele de contact din imagini, tipice multor interacțiuni umane, pot face dovezile vizuale neconcludente. Astfel, aceste modele conduc la estimări 3d incorecte, nerealiste care pierd din vedere aspectele subtile ale contactului uman și sunt prea puțin utile înțelegerii comportamentului uman din imagini.

Într-un articol aflat în evaluare la IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), aliniat cu obiectivele FIGHS, sunt investigate și propuse următoarele contribuții: (1) modele pentru estimarea semnăturii interacțiunilor care cuprind detectia contactului, segmentarea lui și predicția semnăturii 3d a contactului; (2) demonstrarea utilității unor asemenea modele pentru a produce funcții de loss extinse pentru a asigura consistența reprezentării contactului uman în procesul de reconstrucție 3d a persoanelor; (3) propunerea unei metodologii pentru recuperarea posturii și formei ground-truth a persoanelor într-un mediu controlat; (4) prezentarea unor modele și algoritmi baseline pentru a ilustra felul în care estimarea contactului uman sprijină reconstrucția 3d superioară a persoanelor atunci când sunt capturate interacțiuni esențiale între ele.

Articolul propune un prim set de elemente metodologice pentru a aborda reconstrucția persoanelor care interacționează, într-o manieră mai principială, bazându-se pe recunoaștere, segmentare, mapare și reconstrucție 3d. Mai precis, problema furnizării de reconstrucții 3d veridice ale persoanelor care interacționează este descompusă în (a) detectia contactului, (b) segmentarea binară a regiunilor de contact pe suprafețele corespunzătoare asociate oamenilor care interacționează; (c) predicția semnăturii contactului pentru a produce estimări ale hărții potențiale de corespondență many-to-many dintre regiunile aflate în contact; și (d) reconstrucție 3d folosind funcții de loss extinse construite folosind

constrângeri suplimentare de contact de suprafață, atunci când exista o semnătură de contact. Cu ajutorul unui set extins de experimente, s-au evaluat toate componentele sistemului și s-au oferit comparații cantitative și calitative care arată modul în care abordarea propusă poate surprinde interacțiunile umane reprezentate 3d în mod realist.

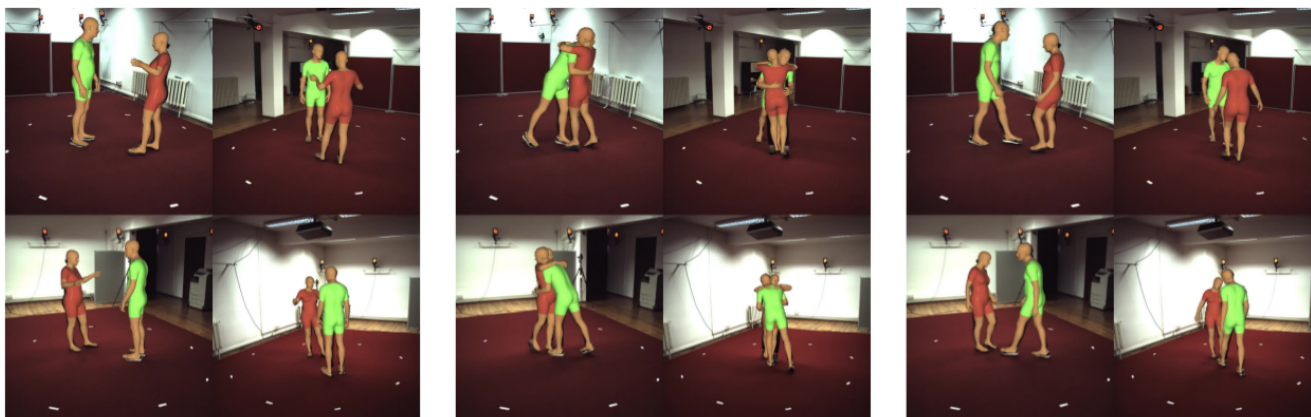
Pentru primele trei task-uri de mai sus (a-c) s-au folosit metode de învățare bazate pe rețele neuronale profunde care primesc ca date de intrare bounding box-urile a doua persoane care interacționează împreună cu posturile 2d detectate. Clasificarea contactului presupune estimarea existenței unui contact fizic între cele două persoane din imagine, care se face cu o rețea de clasificare binară care prezice existența contactului sau a absentei lui. Segmentarea contactului și predicția semnăturii de contact se fac la nivel de regiune și sunt antrenate împreună o rețea mulți-task care folosește adnotari ground-truth despre contact și segmentarea sa, reprezentabile ca în figura de mai jos în care se poate observa cum adnotatorii marchează regiuni ale corpurilor celor două persoane aflate în contact. Corespondențele între regiuni pe care ei le marchează pe mesh-urile 3d ale persoanelor din imagine constituie o segmentare automată a regiunilor aflate în contact.



Odată estimată semnatura contactului dintre două persoane aflate într-o imagine, această informație poate fi folosită pentru a reconstrui postura 3d și parametrii mesh-urilor celor două persoane prin extinderea unui framework de reconstrucție 3d a mai multor persoane dintr-o imagine cu funcții de loss care folosesc semnatura de contact. În primul rând, având în vedere că mâinile sunt cel mai adesea

implicate în interacțiuni de contact, s-a folosit modelarea lor pentru a augmenta funcțiile de loss care cuantifica eroarea de reproiectie 2d fata de evidenta vizuala extrasa din imagine sub forma posturii 2d și a etichetarii semantice a partilor corpului, respectiv funcția de regularizare a posturii și formei mesh-ului. Suplimentar a fost introdus un nou termen de loss specific informației de contact care măsoară alinierea geometrica a regiunilor puse în corespondenta. Mai exact, termenul încearcă sa minimizeze suma distantelor dintre toate perechile de regiuni aflate în contact și, respectiv, sa forteze alinierea orientarilor suprafetelor acestor regiuni.

Seturile de date folosite au fost extinse prin adnotarea întregii perioade temporale a contactului fizic în fiecare secvență video. De asemenea, metodologia propusa pentru a obține postura și forma ground-truth a persoanelor care interacționează depășește simpla extrapolare a pozițiilor markerilor 3d la parametrii unui model de forma a corpului, deoarece: (a) doar unui singur subiect din fiecare videoclip ii este urmărita mișcarea cu ajutorul sistemului de captura a miscarii și (b) subiecții inregistrați nu sunt doar apropiați unul de celalalt, dar sunt, de cele mai multe ori, chiar în contact fizic. În plus, sistemul de captura a miscarii nu permite urmarirea articulatiei mainii în cea mai mare parte a ei. Prin valorificarea informațiilor de la senzorii de detectie a mișcarii, a celor de la mai multe camere RGB și respectiv a celor de la un scanner 3d, dar și din adnotări de contact, informații a priori despre postura corpului și constrângeri fizice, s-au realizat reconstrucții 3d similare dpdv al corectitudinii cu cele rezultate din reprezentarile ground-truth (vezi fig. de mai jos).



Pentru atingerea acestui tel a fost nevoie de un model de optimizare care sa ia în considerare suma unor termeni de loss care sa cuantifice diversele constrangeri mentionate mai sus. În primul rând, scanuri 3d ale subiecților au fost folosite pentru a potrivi un model parametrizat (GHUM) la ele. Diferența dintre reconstrucția 3d folosind markeri ground-truth și cea rezultata din GHUM a fost folosită doar pentru subiectul urmărit de sistemul mocap. În 2d, funcția de loss aferenta a cuantificat pentru fiecare dintre cele 4 viewpoint-uri RGB ale scenei eroarea dintre proiectia 2d a articulatiilor 3d estimate ale modelului GHUM și reprezentarile 2d ale articulatiilor estimate direct din imagine (pentru subiectul cu markeri,

estimarea se face într-o zonă decupată strâns în jurul subiectului, zona obținută prin proiectarea markerilor 3d ground-truth în imagine folosind parametrii cunoscuți ai camerei; pentru celălalt subiect s-a folosit un detector de keypoints pe întreaga imagine). Pentru prevenirea auto-intersecțiilor și respectiv a interpenetrării a două mesh-uri s-au folosit teste de tip generalized winding numbers pentru a stabili relația vertex-urilor interioare unui mesh față de cele mai apropiate vertex-uri vecine (nearest neighbors) fie din propriul mesh (self-collision) fie din celălalt mesh (interpenetration). De asemenea, s-a exercitat modelarea corectului contact cu solul ca și constrangerea impusă de o funcție de loss care cuantifică semnatura de contact bazată pe interpretarea fatetelor în contact ale celor două mesh-uri.

Pentru a avea o reprezentare calitativă a importanței modelării interacțiunilor între persoane atașăm mai jos reconstrucții 3d ale unor persoane aflate în interacțiune conform unor varii scenarii. Prima coloană prezintă imaginile RGB, urmată de coloana a doua care prezintă reconstrucția 3d fără informație de contact (interacțiune). Coloanele 3 și 4 prezintă rezultatele reconstrucției 3d folosind informație de contact la nivel de regiuni (abordare coarse-grain) la diferite granularități ((35 și respectiv 75 de regiuni utilizate). Ultima coloană prezintă rezultatele reconstrucției 3d atunci când se folosește informație de contact la nivel fine-grain, mai exact la nivel de fatete de mesh 3d. După cum se observă, chiar dacă folosirea constrangerilor la nivel de fateta produce cele mai bune rezultate, rezultate rezonabile se obțin și la un nivel de utilizare coarse-grain al informației de contact.



De asemenea, s-au adnotat mișcările de interacțiune cu descrieri text deschizând astfel calea pentru antrenarea și evaluarea modelelor care generează mișcări de interacțiune 3d din text. Fiecarei perechi de mesh-uri în interacțiune i-au fost asociate adnotări text cu o lungime medie de 12 cuvinte în engleză care oferă și posibilitatea unor augmentări suplimentare (de ex, schimbarea ordinii mișcărilor celor două persoane fără a schimba adnotarea textuală). Un exemplu de adnotare poate fi văzut mai jos:

*A man holds his right arm around somebody's shoulder and raises his left hand for a picture.*



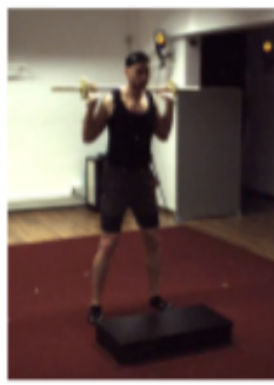
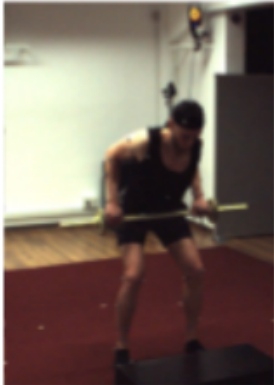
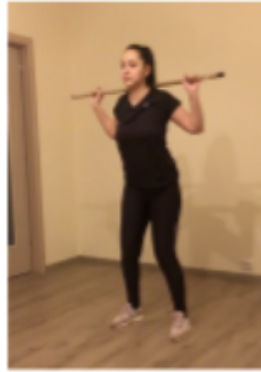
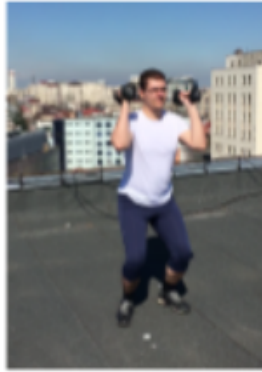
Pe lângă secvențele comune 3d, articolul face publice secvențele de mișcare ground-truth atât în format GHUM cât și SMPLX și pune la dispoziția comunității de cercetare un server de evaluare cu set de test ascuns, împreună cu un benchmark public, cu scopul de a avansa stadiul cunoașterii în reconstrucția 3d a persoanelor aflate în interacțiune.

Dezvoltarea recentă spectaculoasă a modelelor lingvistice pe scară largă (Large Language Models, LLM) de tipul GPT (Generative Pre-trained Transformer) și a modelelor derivate text-to-image de tip DALL-E care se bazează pe ele au motivat așa cum spuneam mai sus efortul de augmentare a înregistrărilor mișcărilor de interacțiune 3d ale persoanelor cu text descriptiv, cu intenția de a putea genera mișcări de interacțiune 3d din descrieri text. În acest context, am trimis spre publicare la IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) un articol care prezintă modelarea reconstrucției 3d a persoanelor care fac fitness. Articolul prezintă un sistem automat care poate fi folosit acasă, la sala sau în aer liber pentru a reconstrui postura, forma și mișcarea persoanei, segmentează repetarea exercițiilor cu acuratețe și identifică în timp real deviațiile mișcărilor persoanei față de mișcările standard recomandate de antrenori calificați (instructori de fitness). Ca rezultat, sistemul oferă feedback cantitativ pentru execuția corectă a exercițiilor ceea ce reduce riscul de accidentare și permite îmbunătățirea continuă a performanței. Pentru a sprijini cercetarea și evaluarea, articolul introduce un set de date pe scară largă cu peste 3 milioane de imagini și configurațiile aferente de captură a mișcării pentru 37 de exerciții de

fitness executate deopotrivă de amatori/elevi și instructori. Sistemul oferă asistența unui antrenor virtual statistic care evaluează critic performanța elevilor pe baza unui parametru care permite atât adaptarea recomandărilor critice nivelului de pregătire al elevului (incepator, avansat, expert) cât și ajustarea acuratetei estimate a metodei de reconstrucție 3d a posturii umane. Antrenorul virtual oferă feedback în limbaj natural sprijinit de evidența spatio-temporală în imaginile video ale exercițiilor elevului.

Metodologia propusă în articol include evaluarea monoculară, multi-view pe scară largă a reconstrucției 3d a posturii umane, modele pentru identificarea automată a repetițiilor unui exercițiu în video, ca și metode de comparație între performanțele elevilor și cele ale instructorilor potrivit unor politici de evaluare statistice. În funcție de parametrul care guvernează gradul de exigență critică al antrenorului virtual cu privire la performanța elevului, parametru care variază între 0 și 1 (0 însemnând atitudine foarte critică a instructorului în vreme ce 1 reprezintă o evaluare relaxată a performanței elevului), sistemul de feedback bazat pe estimarea posturilor 3d atinge o acuratețe înaltă, comparabilă cu aceea a unui asemenea sistem bazat pe o reconstrucție 3d bazată pe captura mișcării cu markeri ground-truth.

Imaginea de mai jos oferă o ilustrare a feedback-ului oferit de antrenorul virtual atunci când sistemul operează asupra unor video-uri capturate în medii reale cu camera video a unui telefon mobil. În rândul de sus este prezentată imaginea elevului care execută o mișcare de fitness în mod greșit, reconstrucția sa 3d pe rândul următor urmata de imaginile execuției corecte a exercițiului de către un instructor calificat și feedback-ul textual care punctează eroarea de execuție. Această descriere text este punctul de plecare pentru un sistem generativ care pe baza indicațiilor text ale antrenorului cu privire la îmbunătățirea performanței elevului să genereze interacțiuni 3d între elev și antrenor care să reflecte vizual sfaturile antrenorului. Ca un exemplu concret, inspirat din imaginea de mai jos, se poate genera mesh-ul antrenorului într-o interacțiune virtuală cu mesh-ul elevului încercând să-i îndrepte coloana vertebrală atunci când nu execută corect o mișcare (v. exercițiul din ultima coloană) pornind de la feedback-ul text („you should keep your back straight”). Un astfel de sistem generativ va fi antrenat pe baza setului de date de captură a mișcărilor de interacțiune 3d a persoanelor augmentate cu descrieri text ale interacțiunilor, așa cum s-a menționat anterior.



**Barbell Row  
Repetition 4**

**You bend your  
elbows too much.**

**Overhead  
Extension Thruster  
Repetition 1**

**You bend your right  
elbow too much by  
15 degrees.**

**Squat  
Repetition 1**

**You should keep  
your elbows  
lower.**

**One Arm Row  
Repetition 3**

**You should  
keep your  
back straight.**



Un rezultat teoretic dezvoltat în cadrul și aliniat cu obiectivele proiectului, urmează să fie publicat în conferința NeurIPS 2024, în domeniul optimizării metodelor de învățare automată. Articolul vizează metode fracționale care extind metoda de optimizare bazată pe gradient descent prin incorporarea derivatelor de ordin fracțional. Această abordare permite o mai mare flexibilitate în navigarea peisajelor complexe de optimizare și oferă avantaje în anumite tipuri de probleme, în special cele care implică neliniarități și comportament dinamic haotic. Cu toate acestea, provocarea reglării fine a parametrilor de ordin fracțional rămâne nerezolvată. Lucrarea demonstrează că este posibil să se antreneze o rețea neuronală pentru a prezice în mod eficient ordinea gradientului.

În optimizarea convențională de ordinul întâi, funcția țintă este de obicei aproximată ca local liniară folosind o expansiune Taylor. Există posibilitatea de a beneficia de aproximări neliniare care surprind comportamentul funcției pe o vecinătate mai mare, oferind o reprezentare mai precisă decât aproximații liniare locale. Metode de gradient descent fracțional au fost dezvoltate pentru a profita de un astfel de aproximări. În literatura de specialitate se arată că acestea pot îmbunătăți considerabil rata de convergență a algoritmului de gradient descent în cazul convex. Aceste metode se bazează pe conceptul de derivate fracționale. Derivata fracționată poate fi gândită ca o „interpolare” între două derivate convenționale. De exemplu, semiderivata (adică ordinul fracțional  $\alpha = 0,5$ ) reprezintă o interpolare între funcția  $f$  însăși și prima sa derivată. Cu toate acestea, lipsește înțelegerea fină a felului în care se poate determina ordinea fracțională optimă pentru o anumită problemă. Au fost dezvoltate metode adaptive dar acestea depind de hiper-parametri suplimentari (de exemplu limite, puncte terminale). Între timp, în domeniul optimizării metodelor de învățare automată s-au realizat îmbunătățiri ale reglării fine a optimizatoarelor expresive. Lucrarea ilustrează modul în care optimizarea învățată poate fi folosită pentru a regla fin ordinul fracțional.

Pentru a învăța optimizarea unor funcții clasice, considerate ca fiind parametrizate de un vector de stare  $X_t$ , s-a antrenat o rețea neuronală care primește ca date de intrare starea curentă, gradientii normalizați ai funcției, magnitudinea lor și feature-uri Fourier și generează ordinul fracțional al derivatei și magnitudinea pasului de actualizare al metodei de gradient descent care vor fi folosite pentru calculul următoarei stări  $X_{t+1}$  folosind o aproximare Taylor de ordinul întâi. Rețeaua folosește optimizatorul AdamW pentru antrenament și o funcție de loss de tipul:

$$L_\theta = \log(f(X_{t+1})) - \log(f(X_t))$$

Rețeaua este antrenată atât supervizat cât și nesupervizat.

Rezultatele obținute au fost încurajatoare pentru sisteme de dimensionalitate mică și deși domeniul de aplicare al metodei este momentan limitat la meta-training, continuarea cercetării în această direcție poate conduce la soluții viabile pentru probleme de dimensionalitate mare.