

RAPORT STIINTIFIC

privind implementarea proiectului in perioada iulie 2022 – decembrie 2023

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

Acronim: FIGHS

Metode Detaliate de Reconstructie Tridimensionala 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

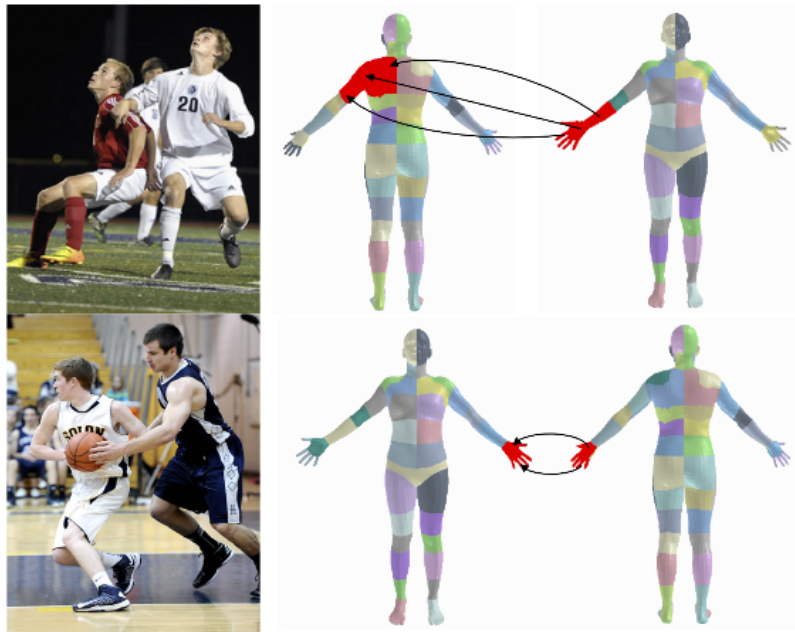
Activitati stiintifice si realizari:

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 există 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 două 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 absenței lui. Segmentarea contactului și predicția semnăturii de contact se fac la nivel de regiune și sunt antrenate împreună o rețea mulți-task care folosește adnotări 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 cantifică eroarea de reproiecție 2d față de evidența vizuală extrasă din imagine sub forma posturii 2d și a etichetării semantice a părților 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 geometrică a regiunilor puse în corespondență. Mai exact, termenul încearcă să minimizeze suma distanțelor dintre toate perechile de regiuni aflate în contact și, respectiv, să forțeze alinierea orientărilor suprafețelor acestor regiuni.

Seturile de date folosite au fost extinse prin adnotarea întregii perioade temporale a contactului fizic în fiecare secvență video. De asemenea, metodologia propusă pentru a obține postura și forma ground-truth a persoanelor care interacționează depășește simpla extrapolare a pozițiilor markerilor 3d la parametrii unui model de formă a corpului, deoarece: (a) doar unui singur subiect din fiecare videoclip îi este urmărită mișcarea cu ajutorul sistemului de captură a mișcării și (b) subiecții înregistrați nu sunt doar apropiați unul de celălalt, dar sunt, de cele mai multe ori, chiar în contact fizic. În plus, sistemul de captură a mișcării nu permite urmărirea articulației mâinii în cea mai mare parte a ei. Prin valorificarea informațiilor de la senzorii de detecție a mișcării, 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 reprezentările 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 subiectilor au fost folosite pentru a potrivi un model parametrizat (GHUM) la ele. Diferența dintre reconstructia 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 zona decupata strâns în jurul subiectului, zona obtinuta prin proiectarea markerilor 3d ground-truth în imagine folosind parametrii cunoscuți ai camerei; pentru celalalt subiect s-a folosit un detector de keypoints pe întreaga imagine). Pentru prevenirea auto-intersecțiilor și respectiv a interpenetrării a doua mesh-uri s-au folosit teste de tip generalized winding numbers pentru a stabili relația vertex-urilor interioare unui mesh fata 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 impusa de o funcție de loss care cuantifica semnatura de contact bazata pe interpretarea fatetelor în contact ale celor doua mesh-uri.

Pentru a avea o reprezentare calitativa a importantei modelarii interactiunilor între persoane atasam mai jos reconstructii 3d ale unor persoane aflate în interacțiune conform unor varii scenarii. Prima coloana prezinta imaginile RGB, urmata de coloana a doua care prezinta reconstructia 3d fără informație de contact (interacțiune). Coloanele 3 și 4 prezinta rezultatele reconstructiei 3d folosind informatie de contact la nivel de regiuni (abordare coarse-grain) la diferite granularitati ((35 și respectiv 75 de regiuni utilizate). Ultima coloana prezinta rezultatele reconstructiei 3d atunci când se folosește informație de contact la nivel fine-grain, mai exact la nivel de fatete de mesh 3d. După cum se observa, 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 informatiei de contact.



De asemenea, s-au adnotat mișcările de interacțiune cu descrieri text deschizand astfel calea pentru antrenarea și evaluarea modelor care genereaza miscari de interacțiune 3d din text. Fiecarei perechi de mesh-uri în interacțiune i-au fost asociate adnotari text cu o lungime medie de 12 cuvinte în engleza care ofera și posibilitatea unor augmentari suplimentare (de ex, schimbarea ordinii miscarilor celor doua persoane fără a schimba adnotarea textuala). 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.



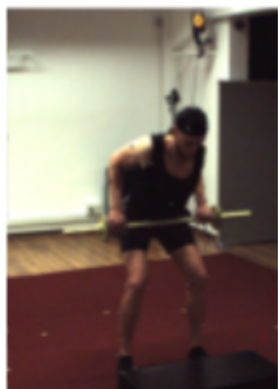
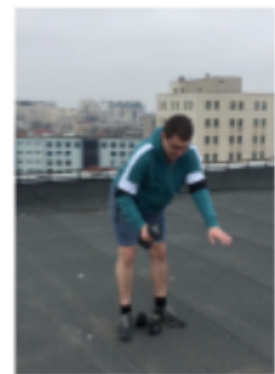
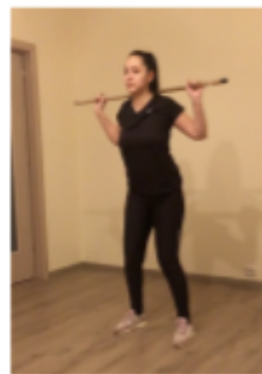
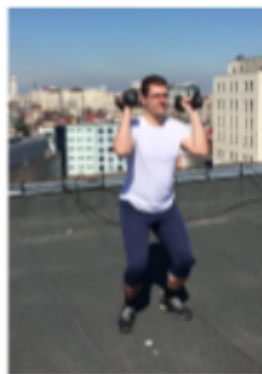
Pe lângă secvențele comune 3d, articolul face publice secvențele de mișcare ground-truth atât în format GHUM cât și SMPLX și pune la dispoziția comunității de cercetare un server de evaluare cu set de test ascuns, împreună cu un benchmark public, cu scopul de a avansa stadiul cunoașterii în reconstrucția 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 sală sau în aer liber pentru a reconstrui postura, forma și mișcarea persoanei, segmentează repetițiile 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 (începător, avansat, expert) cât și ajustarea acurateții 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 captură 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 prezentata imaginea elevului care executa o mișcare de fitness în mod greșit, reconstrucția sa 3d pe rândul următor urmata de imaginile executiei corecte a exercitiului de către un instructor calificat și feed-back-ul textual care puncteaza eroarea de execuție. Aceasta descriere text este punctul de plecare pentru un sistem generativ care pe baza indicatiilor text ale antrenorului cu privire la imbunatatirea performantei elevului sa genereze interactiuni 3d între elev și antrenor care sa 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 feed-back-ul text („you should keep your back straight”). Un astfel de sistem generativ va fi antrenat pe baza setului de date de captura a miscarilor de interacțiune 3d a persoanelor augmentat cu descrieri text ale interacțiunilor, așa cum s-a menționat anterior.



**Barbell Row
Repetition 4**

**Overhead
Extension Thruster
Repetition 1**

**Squat
Repetition 1**

**One Arm Row
Repetition 3**

**You bend your
elbows too much.**

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

**You should keep
your elbows
lower.**

**You should
keep your
back straight.**

Gradul de realizare al obiectivelor

Obiectivele etapelor de până acum au fost atinse în întregime după cum urmează. S-a realizat setup-ul colecției de date atât în contextul augmentării unor seturi de date existente cât și pentru folosirea noului sistem de captură a mișcării achiziționat. Metodologia folosită pentru organizarea și setup-ul datelor în ambele cazuri este principial aceeași, diferențele survenind exclusiv în privința metodei tehnice de achiziție a datelor care este, în mod evident, particulară echipamentului de captură a mișcării umane folosit.

S-au folosit aceste baze de date augmentate pentru evaluarea metodologiei state-of-the-art de estimare 3d a interacțiunii dintre persoane și, așa cum s-a menționat anterior, s-a extins această metodologie prin propunerea unui prim set de elemente metodologice care presupun o abordare principială a reconstrucției 3d a persoanelor care interacționează, bazată pe recunoaștere, segmentare, mapare și reconstrucție 3d. În plus, este în curs de dezvoltare metodologia de estimare a interacțiunilor umane în scene 3d din natură pe baza datelor înregistrate cu noul sistem de captură a mișcării achiziționat. Toate aceste dezvoltări sunt aliniate cu obiectivele proiectului și contribuie la elaborarea metodologiei de estimare a interacțiunilor umane în scene 3d din natură și realizarea bazei de date, activități care se extind până la finalul proiectului.

Primele rezultate ale proiectului sunt disponibile într-o variantă sumară (dezvoltările metodologice descrise mai sus apar în articole aflate în evaluare sau trimise spre publicare și nu pot fi publicate încă în detaliu pe site) la adresa de web <https://vision.imar.ro/fighs>

În aceste condiții, consideram obiectivele curente ale proiectului ca fiind îndeplinite.

Sumarul progresului înregistrat în proiect

1. setup-ul colecției de date disponibil
2. evaluarea și extinderea metodologiei de estimare 3d a interacțiunilor umane (incluzând și extinderea unor seturi de date existente)
3. demararea activităților de culegere a datelor cu noul sistem de captură a mișcării umane și dezvoltare a metodologiei necesare pentru estimarea interacțiunilor umane în scene 3d din natură
4. transmiterea unui articol care descrie progresele înregistrate spre publicare la IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), articol aflat momentan în plin proces de evaluare

5. transmiterea spre publicare a unui alt articol care va permite extinderea modelarii interacțiunilor umane către sisteme generative la IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)

Concluzii:

Activitățile proiectului s-au desfășurat conform calendarului stabilit și s-au adaptat condițiilor contractuale întâi prin reorientarea către folosirea unor baze de date existente pentru continuarea cercetării în domeniul analizei scenelor 3d cu mai multe persoane care interacționează între ele, urmata de demararea achiziționării de date cu noul sistem de captură a mișcării umane disponibil în proiect și începerea elaborării metodologiei specifice de estimare a interacțiunii persoanelor în scene 3d din natura. Ca rezultate concrete, s-au trimis spre publicare două articole la IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), revista cu unul dintre cei mai mari factori de impact din domeniul științei calculatoarelor (IF actual: 23.6).