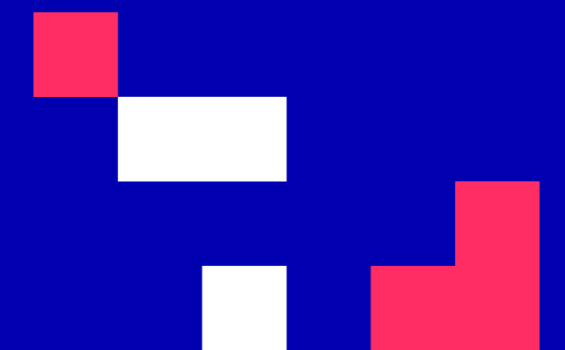University of Cyprus

# MAI645 - Machine Learning for Graphics and Computer Vision

**Andreas Aristidou, PhD**

Spring Semester 2025

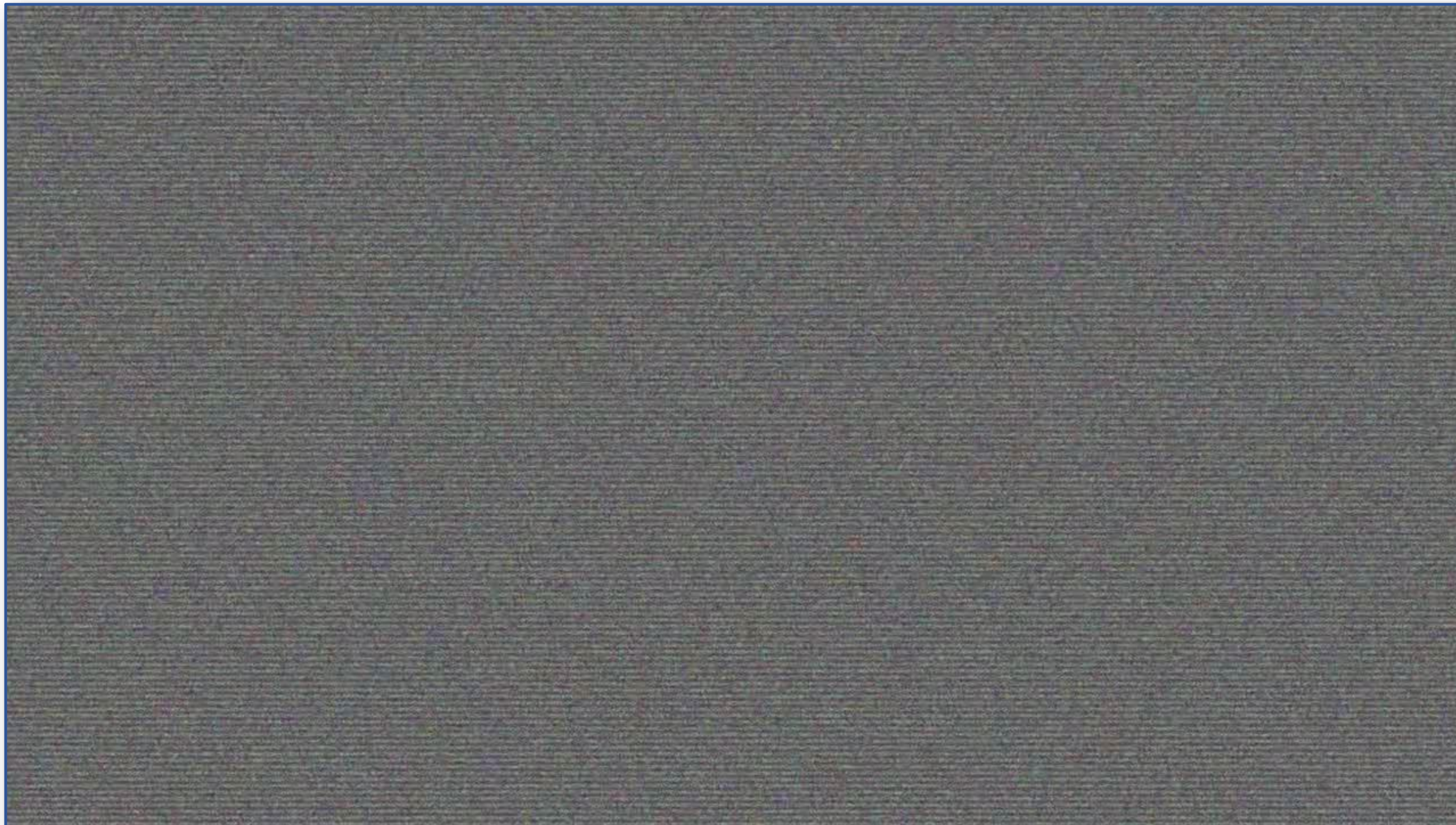## Character Animation

## How does the magic happen?

Uncharted 4: A Thief's End, © Sony Computer Entertainment

# How do we make this movie?



Uncharted 4: https://youtu.be/zL46dpNEPPA

**Modeling**

- Geometry
- Materials
- Lighting

**Animation**

- How do they move?

**Rendering**

- Shadows
- Camera
- Special effects
- Post-processing

MAI4CAREU

Master programmes in Artificial Intelligence 4 Careers in Europe

# MAI4CAREU

## Motivation

- Bring Animated characters to life
  - Animator analogous to film actors

- Many applications use **character** or **object** animation
  - Entertainment technology (e.g., films, games)
  - Virtual, or augmented reality
  - Simulations, demonstrations, or training systems

- Other forms of animation?
  - Trees, liquids, animals, clouds, etc.

- Other Important factors in character animation
  - Lighting, Rendering, etc.



© Pixar Animation





CD-MPM: Continuum Damage Material Point
Methods for Dynamic Fracture Animation
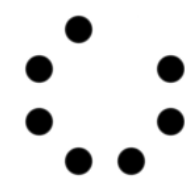Wolper, Fang, Li, Lu, Gao, Jiang

# Introduction to Animation

## Moving Picture & Animation

- The perception of motion is based on two optical illusions, the phi phenomenon and beta movement.
  - **phi** is an optical illusion whereas we perceive motion from fast luminous impulses in sequence. Our visual system "fills in" the missing information.
  - **beta** movement is the illusion of motion created when stimuli changes position in a sequence of images. Instead of being perceived as a series of images we perceive movement.
  - Quick succession of images (frames) causes this sensation of movement (1/25sec)
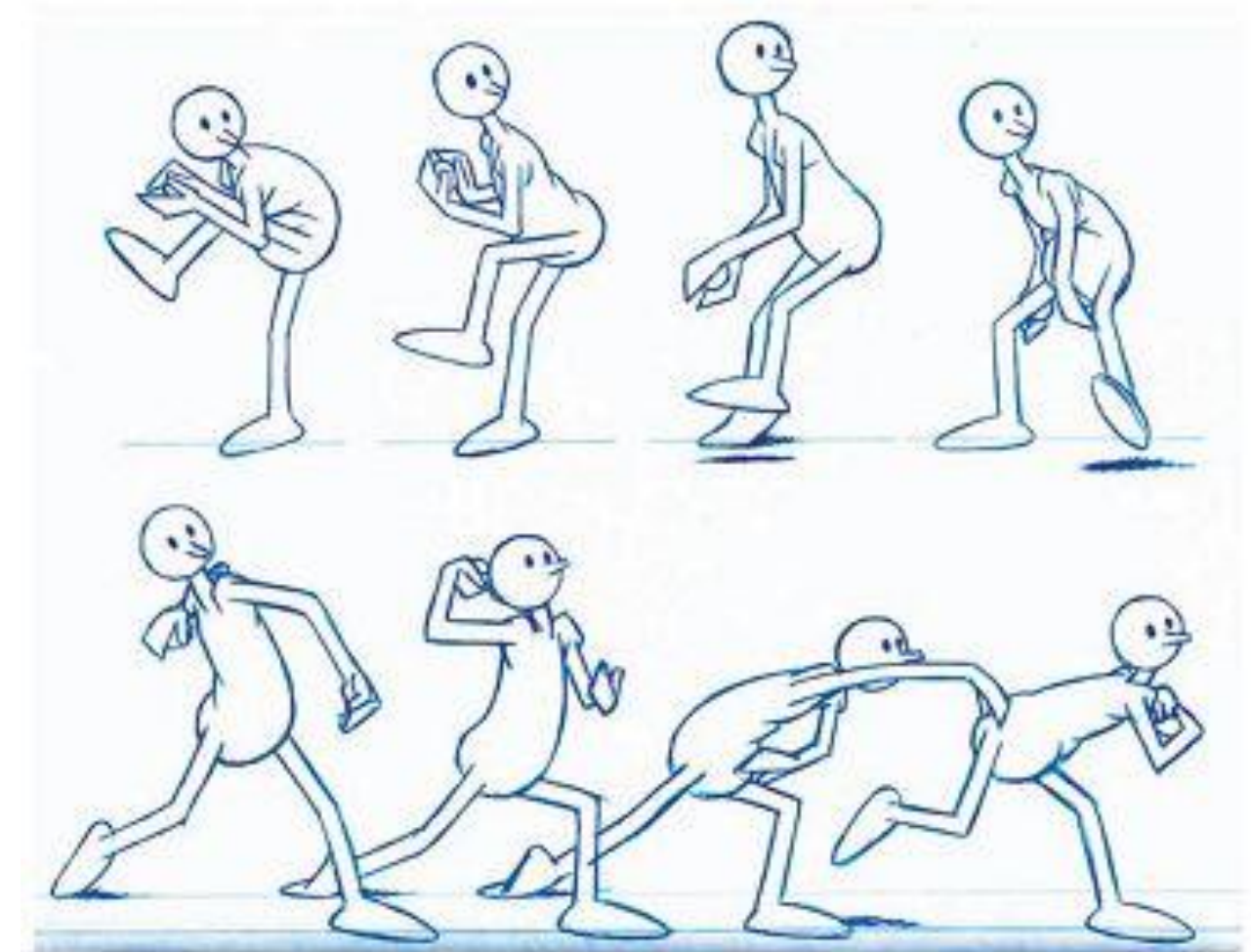
**Phi
Phenomenon**

**Beta
Movement**

## Moving Picture & Animation

- Below 16 images/sec flicker is observed.

- Movies play at 24 images/sec.

- ~10 images/sec still provide sensation of movement.

- Traditional animation was created "on twos"
  - A new image every second frame.

- Faster motions are executed "on ones"
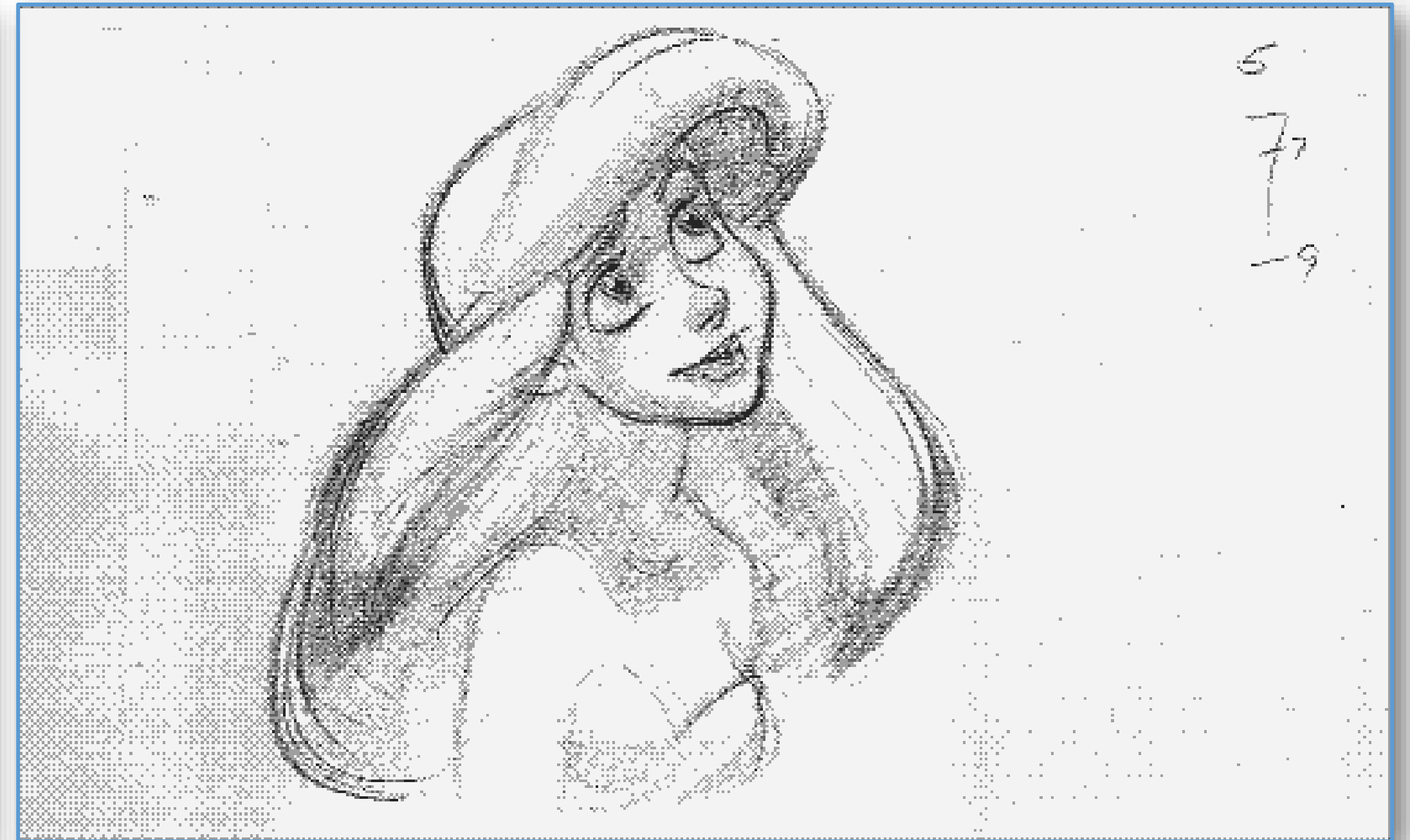  - A new image every frame.

## Moving Picture & Animation



AncientMagicToys.com

## Moving Picture & Animation

Co-financed by the European Union
Connecting Europe Facility

11

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423

## Moving Picture & Animation

- Senior animators draw keyframes (important/extreme shots)
- Junior animators (inbetweeners) fill in the in-between frames
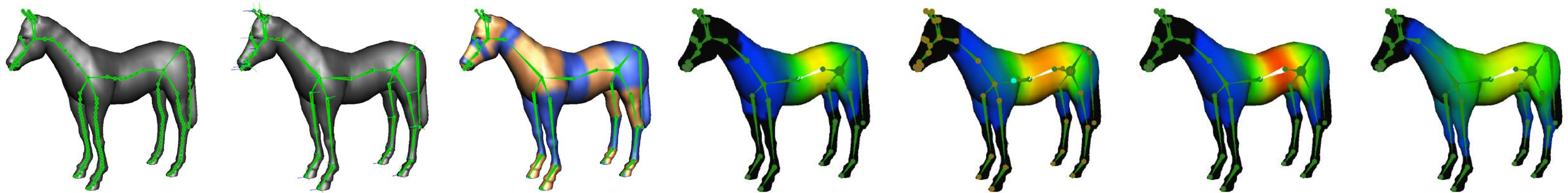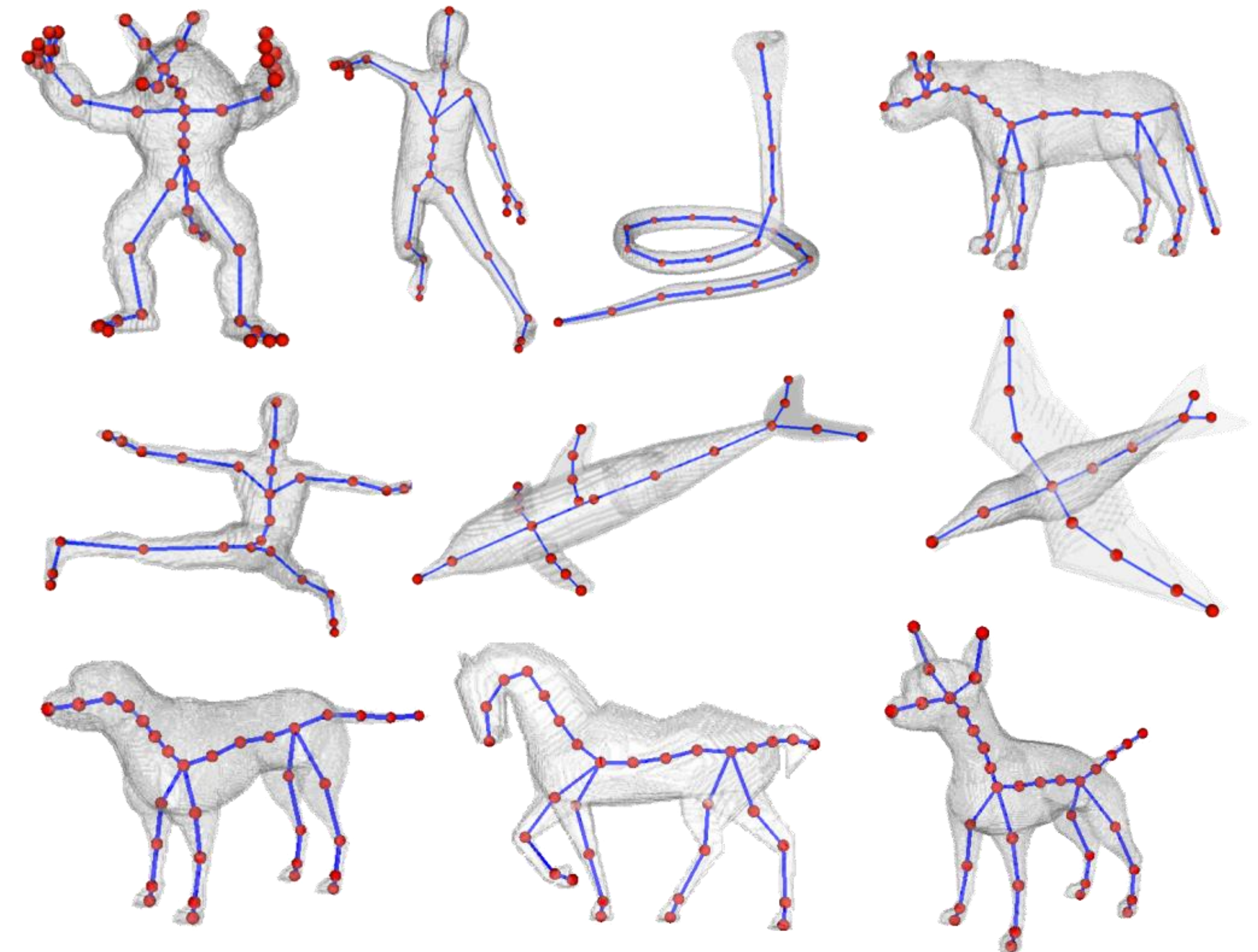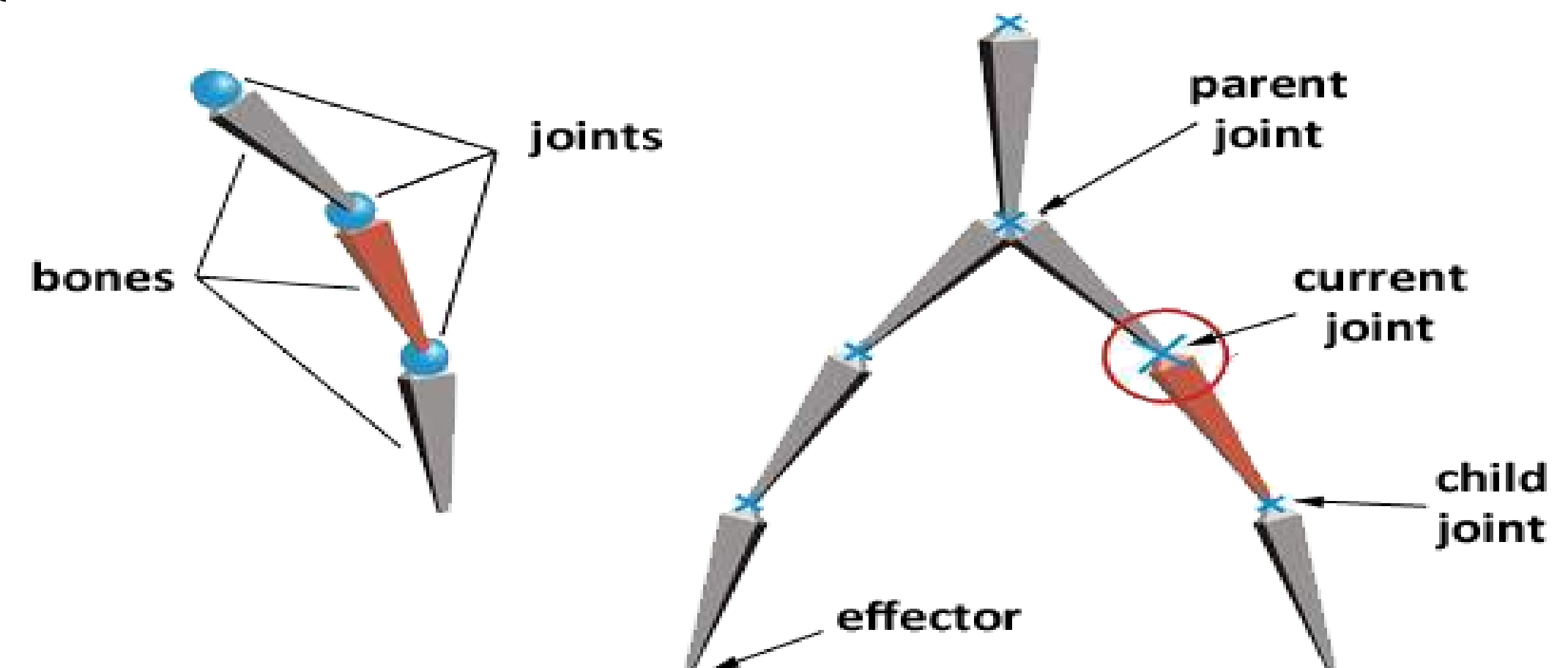


A cartoon animation may require <u>thousands</u> of hand-drawn images

## What is Computer Animation?

- **Computer Animation** is the branch of computer graphics interested in developing techniques for creating moving images.

- Computer Animation is the modernized brother of traditional animation.

# Character rigging & skinning

# Rigging

- 3D **rigging** is the process of creating a skeleton for a 3D model so it can move.

- A 'rig' has numerous degrees of freedom (DOFs) that can be used to control various properties.

- One character could have several rigs. One rig could control several characters…

## Rigging: *The Rig*

- A skeletal system (**rig**) is comprised of kinematic chains:

    - A hierarchical set of interconnected bones

    - A chain:

        - starts from a **root**,

        - it has multiple **bones**,

        - connected by **joints**, and

        - ends at the **end-effector**.

Co-financed by the European Union
Connecting Europe Facility

16

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423

## Rigging: *The Rig*

- A skeleton allows higher-level control of the character's animation.

- The skeleton is only a control mechanism – it is not rendered into the final image.

- Typically, there are many constraints.

# Rigging: *The Rig*

## Skinning

- Skinning.
  - Attach a mesh ("**skin**") to the skeletal system of the character.

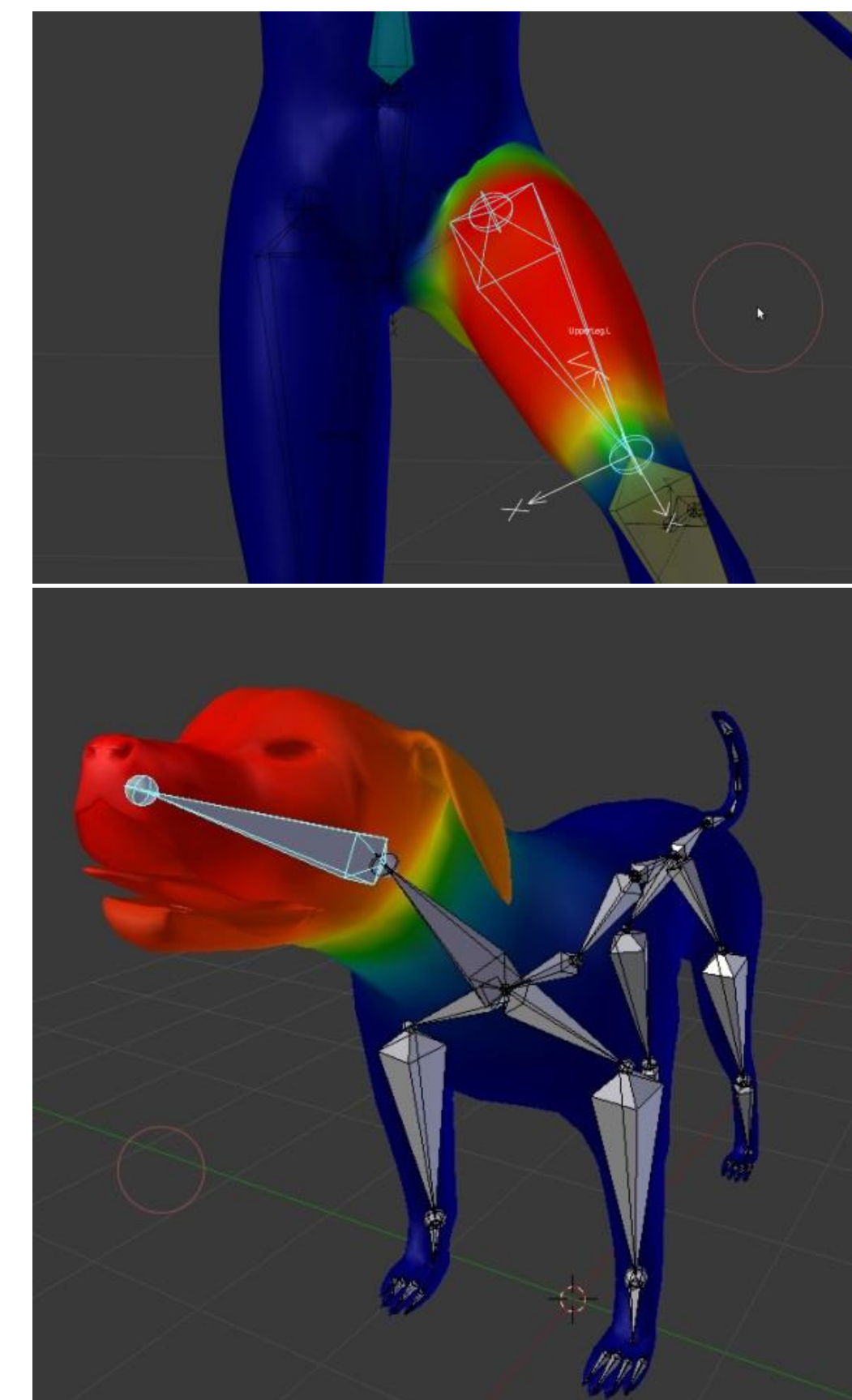- The **skin** is represented as a polygon mesh, e.g., a set of vertices, or a parametric surface

Co-financed by the European Union
Connecting Europe Facility

19

This Master is run under the context of Action No 2020-EU-IA-0087, co-financed by the EU CEF Telecom under GA nr. INEA/CEF/ICT/A2020/2267423

## Skinning: *The Skin*

- We **bind** the skeleton to the mesh when we first associate them.
  - The **T-pose** (or "**bind pose**") refer to the initial transformation matrices of the rig and skin when they are first associated.
  - The T-pose defines a coordinate system used later when animating the skin via the skeleton.
  - The T-pose is a convention used because:
    - modeling the mesh and the skeleton is easier, using symmetry.
    - rigging is much easier when the limbs are spread apart.

## Skinning: *The Skin*

- Each vertex is associated with a bone in the skeleton, and moves relative to that bone.

- Each vertex is multiplied by several "weighted" transformation matrices that provide the influence factor each bone has to the vertex, and the results are added together.
  - The skin's vertices can then be assigned weights.
    - Rigid skinning: 1 bone per vertex (weight = 1.0)
    - Smooth skinning: Multiple bones per vertex (weights != 1.0)

## Texture



Co-financed by the European Union
Connecting Europe Facility

22

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423

# Skeletal Animation

# What 3D character animation involves?

Animating characters can be broken down to:

- **Skeletal animation** – animating their main body parts.

## What 3D character animation involves?

Animating characters can be broken down to:

- **Skeletal animation** – animating their main body parts.
- **Facial animation** – animating their facial features.

## What 3D character animation involves?

Animating characters can be broken down to:

- **Skeletal animation** – animating their main body parts.
- **Facial animation** – animating their facial features.
- **Hair** (and **fur)** animation

## Skeletal Animation: *Keyframing*

We present a control system based on 3D muscle actuation

## Computer generated animation: *Motion Capture*

**Optical Motion Capture**

- The system combines the information of the tracked markers to describe the 3D position of the object

  - Repeat this operation several times per second the system can provide us the volumetric trajectory of the marker according to time and space (usually from 30Hz to 960Hz)

- Great naturalness and realism in the captured movements.

  - High quality recording

  - Capturing of both main and secondary movements

  - Ease of use (Skeletal geometry is given)

- Capture volume is the physical space where the cameras can combine their fields of view

## Computer generated animation: *Motion Capture*

### Optical Motion Capture

- Each person wears a suit with markers attached.

- Enters a space that is surrounded with cameras.

- Divided into two main categories: passive and active

# Motion Capture pipeline

**Motion Capture pipeline**

| Marker-based motion acquisition | Label Markers | Marker Data Clean-up | Convert to Joint Angles |

# Other popular motion capture systems

## Inertial Markers

- Micro-inertial sensors, biomechanical models and sensor fusion algorithms.
- Use a number of gyroscopes and accelerometers to measure rotational rates.
- These rotations are translated to a skeleton model.

## Depth-Based

- Use a combination of color cameras and depth sensors.
  - the subject's silhouette is captured from multiple angles.
- Reconstruct the object's volume (mesh) from the point clouds.
- Fit a skeleton into the 3D model to estimate motion.

## Vision-Based

- Use a singe or multiple RGB cameras
- Mainly based on deep-learning methods
- Use large amount of training motion data



**Fusion4D**

**Real-time Performance Capture of Challenging Scenes**

Mingsong Dou, Sameh Khamis, Yury Degtyarev, Philip Davidson*, Sean Ryan Fanello*,
Adarsh Kowdle*, Sergio Orts Escolano*, Christoph Rhemann*, David Kim,
Jonathan Taylor, Pushmeet Kohli, Vladimir Tankovich, Shahram Izadi

*equal contribution

**MICROSOFT RESEARCH**
contact: shahrami@microsoft.com

**MonoPerfCap:**
**Human Performance Capture from Monocular Video**

(with voiceover)

Weipeng Xu[1]  Avishek Chatterjee[1]  Michael Zollhöfer[1]  Helge Rhodin[2]
Dushyant Mehta[1]  Hans-Peter Seidel[1]  Christian Theobalt[1]

[1]Max Planck Institute for Informatics, Saarland Informatics Campus  [2]EPFL

## Motion Capture: *Current technological trends*

© Virtual Reality Lab, Department of Computer Science, University of Cyprus

## 3D scanning and animation

# More Examples

## Motion Capture Data

- Depending on the sensors used

- Popular file formats:
    - ASF/AMC (Acclaim's skeleton and motion capture files)
    - BVH (BioVision Hierarchy)
    - C3D (Coordinate 3D – biomechanics – C3D.org)

# Motion Capture Data: *BVH format*

```
HIERARCHY
ROOT Hips
{
    OFFSET 0.00000 0.00000 0.00000
    CHANNELS 6 Xposition Yposition Zposition Zrotation Yrotation Xrotat
    JOINT LHipJoint
    {
        OFFSET 0 0 0
        CHANNELS 3 Zrotation Yrotation Xrotation
        JOINT LeftUpLeg
        {
            OFFSET 3.13874 -1.57224 1.49786
            CHANNELS 3 Zrotation Yrotation Xrotation
            JOINT LeftLeg
            {
                OFFSET 2.10955 -5.79594 0.00000
                CHANNELS 3 Zrotation Yrotation Xrotation
                JOINT LeftFoot
                {
                    OFFSET 2.41843 -6.64458 0.00000
                    CHANNELS 3 Zrotation Yrotation Xrotation
                    JOINT LeftToeBase
                    {
                        OFFSET 0.04713 -0.12948 1.66229
                        CHANNELS 3 Zrotation Yrotation Xrotation
                        End Site
                        {
                            OFFSET 0.00000 -0.00000 0.85167
                        }
                    }
                }
            }
        }
    }
    JOINT RHipJoint
    {
        OFFSET 0 0 0
        CHANNELS 3 Zrotation Yrotation Xrotation
```

# Motion Capture Data: *BVH format*

## Motion Capture Data: *SMPL format*



SMPL Model

## Motion Capture: *Advantages*

- Great naturalness and realism in the movements that have been recorded.
  - High-quality recording
  - Recording of both primary and secondary movements
- Recording at a very high frequency
  - Up to 980 samples per second (e.g., birds)
- Ease of use
  - Geometry is a given.
  - Freedom of movement for users

## Motion Capture: *Limitations*



Animators could use more
than 750 controls to create
Shrek's performance. Some
controlled one joint or
muscle, others controlled
groups of several.

- Only realistic motion captured (movement that does not follow the laws of physics cannot be captured).
  - Cartoony or superhero animations are not possible to be captured.

- WYSIWYG (what you see is what you get).
  - Can't add more expression.
  - Continually need to recapture motion.

- What about muscles?

## Some home-built motion capture systems

Co-financed by the European Union
Connecting Europe Facility

45

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423

## Some home-built motion capture systems

## Other Challenges: *Motion Retargeting*

## Other Challenges: *Motion Retargeting*

# What is motion retargeting?

- A method to retarget animations onto models with different morphologies.

- A way to remap animations onto characters with very different animation-specific structures.

## Other Challenges: *Motion Retargeting*

# Why Motion Retargeting?

- Improves content reuse.
- Easy integration of procedurally generated animations.
- Sometimes is not possible to motion capture the subject (e.g. animal with human behavior, character does not exist – fiction movies).

## Other Challenges: *Motion Retargeting*

- Preserve angles or end-effector positions (flying).

- Foot-skating.

- Characters with different proportions may have body penetration.

# Our Dance Motion Capture Database



1st Antikristos

Zumba

# Our Dance Motion Capture Database

## Our Dance Motion Capture Database
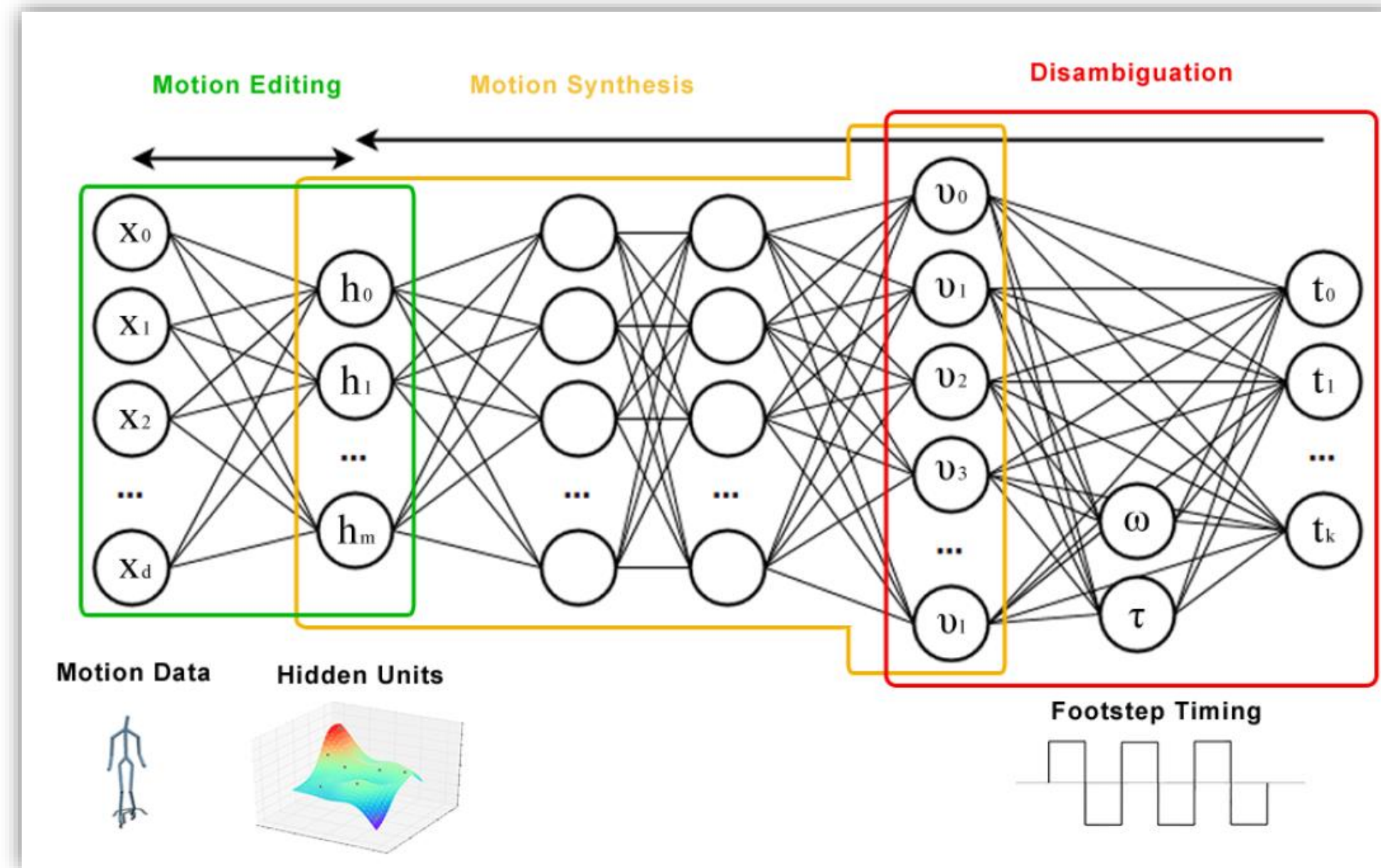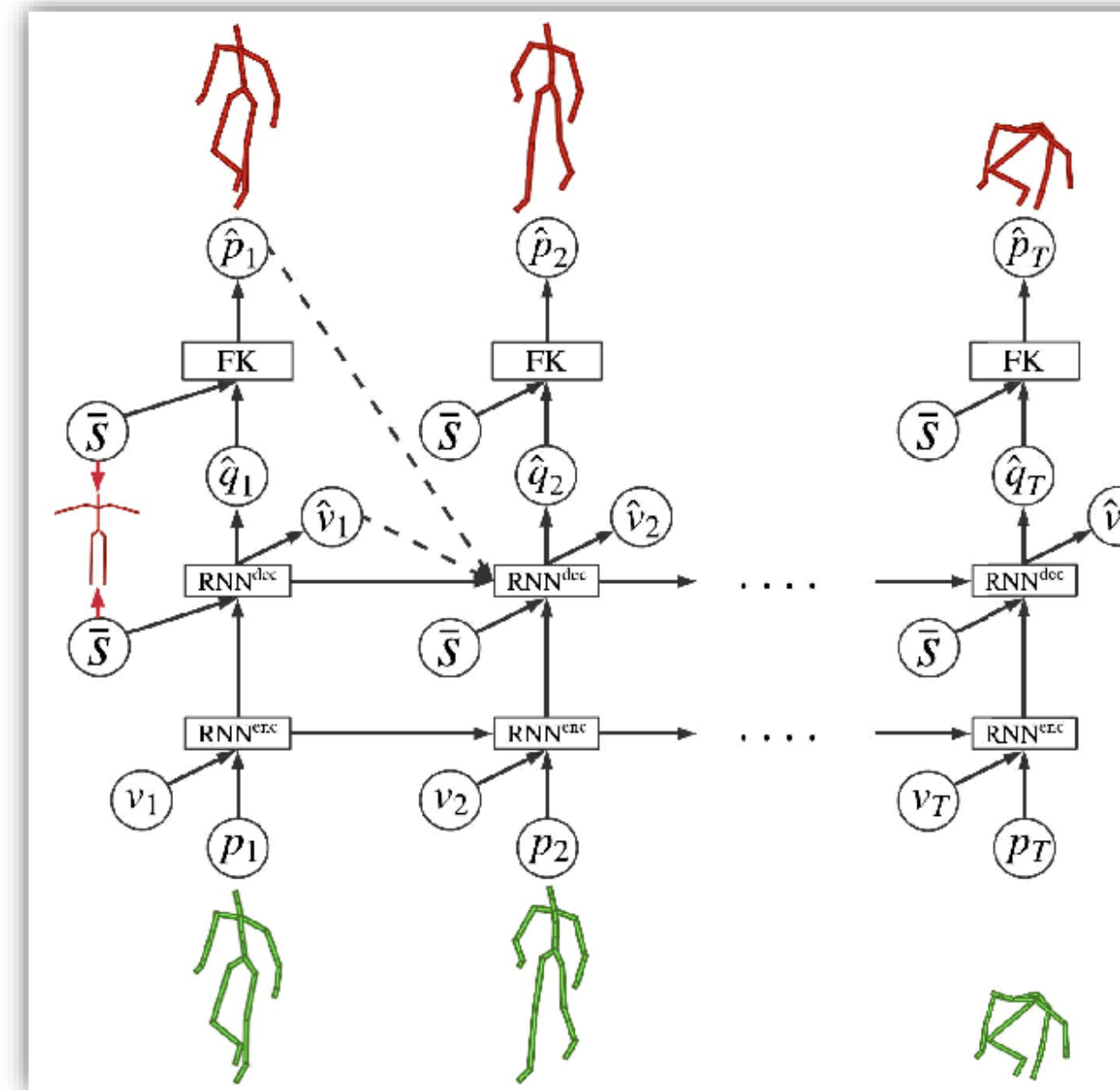
## Our Dance Motion Capture Database
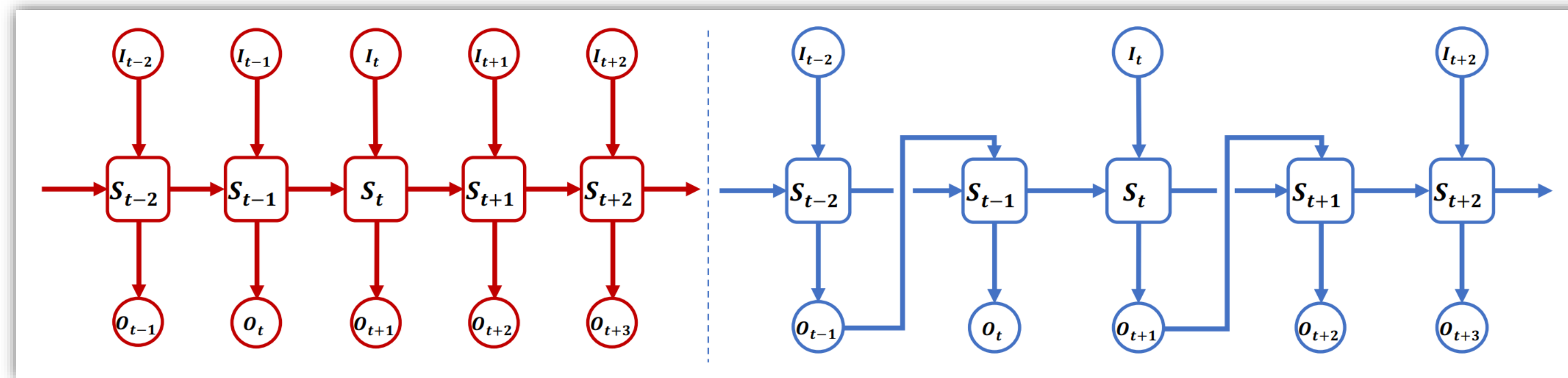
# Deep Character Animation
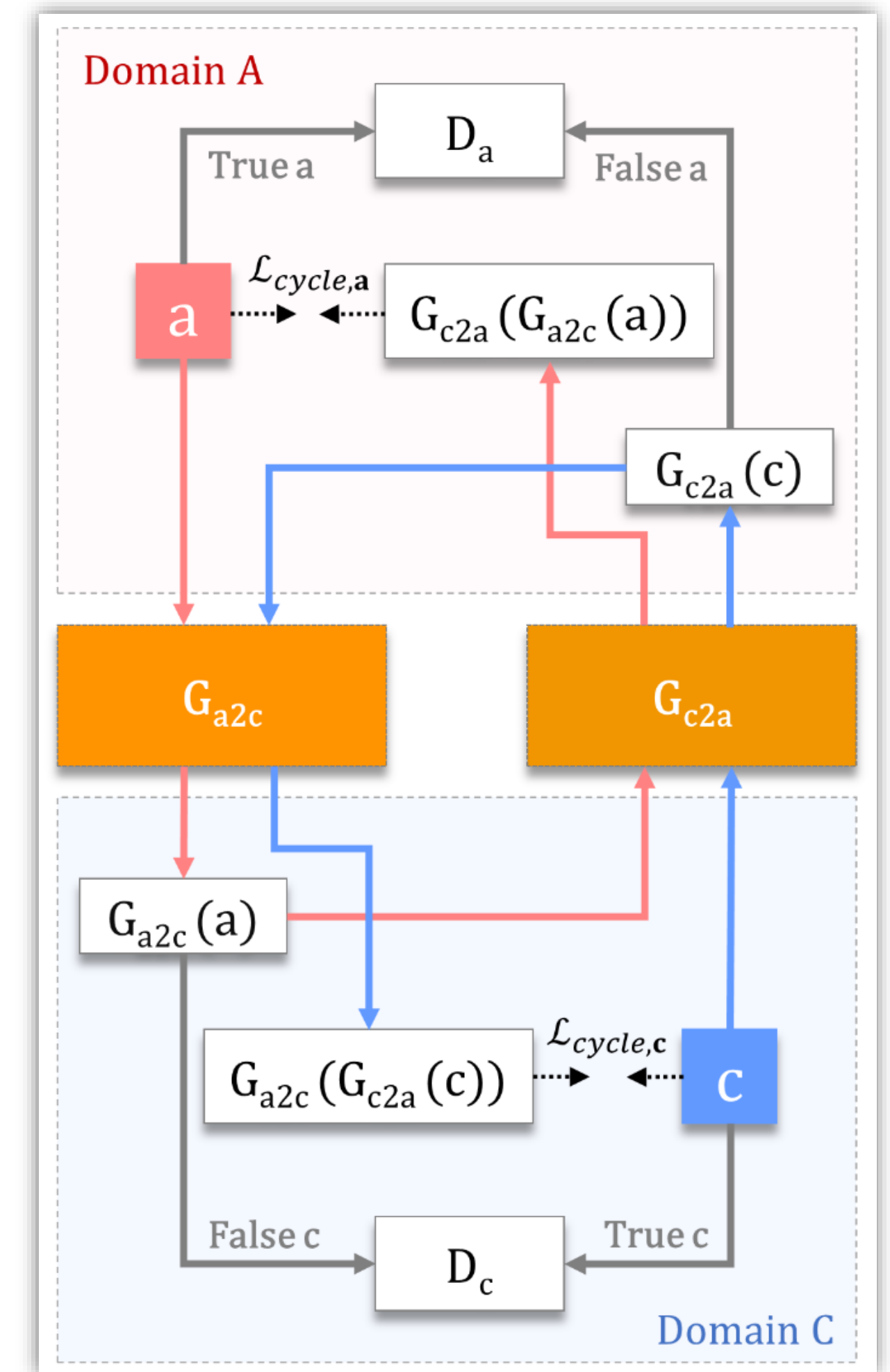
# Deep Neural Networks



Convolutional Neural Networks (CNNs)



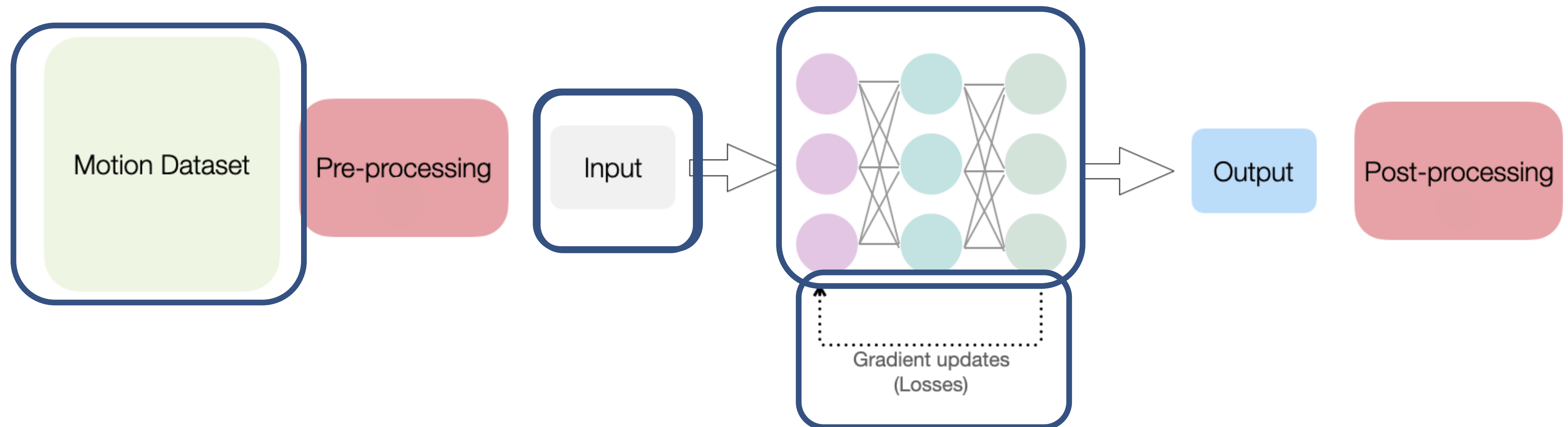Recurrent Neural Networks (RNNs)



auto-conditional LSTM



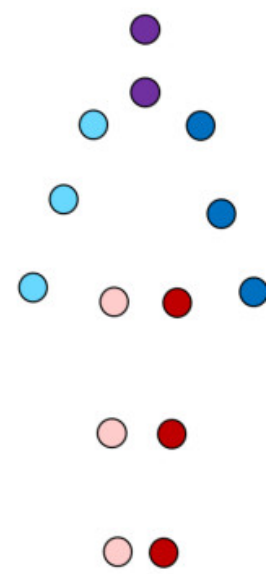Generative Adversarial Networks (GANs)

# Motivation

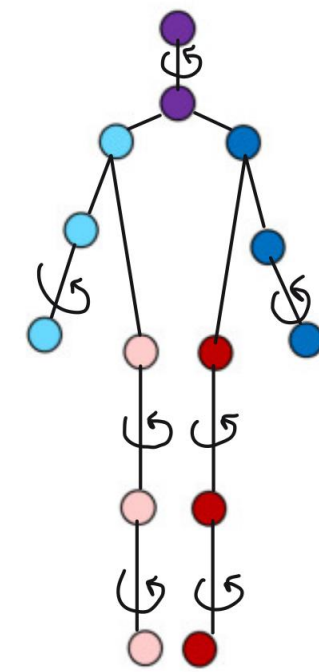*Character Animation with Deep Learning*
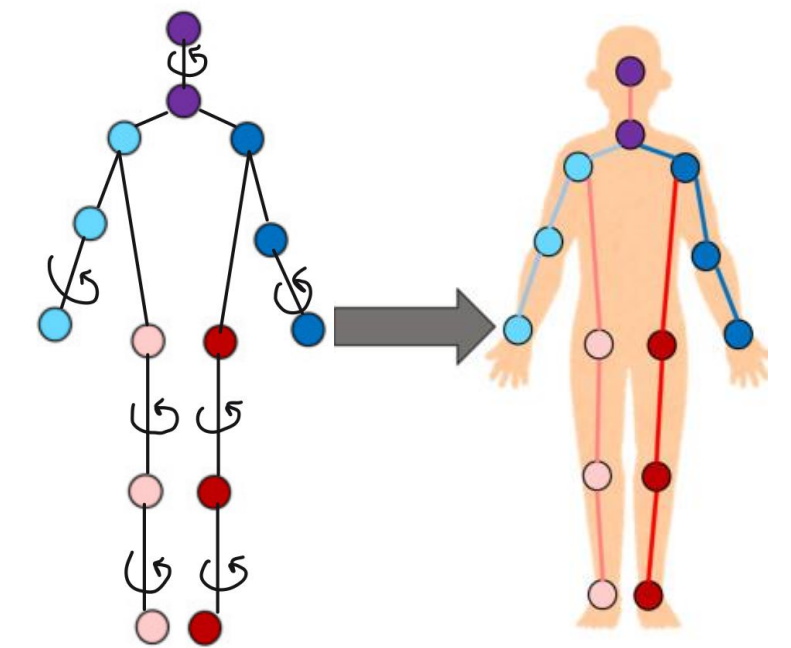
# Motivation

## *Common Pose Representations*



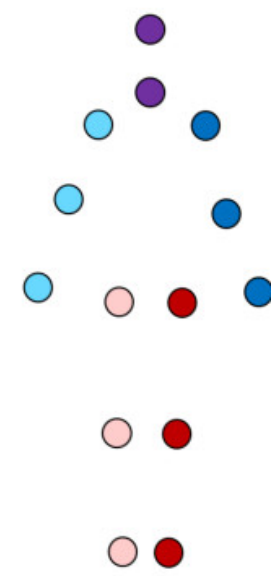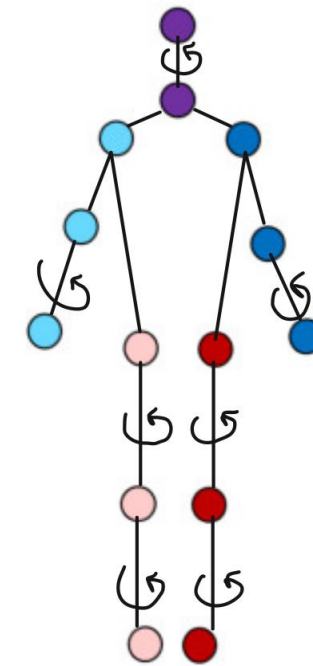Positional



Angular



Hybrid

# Motivation

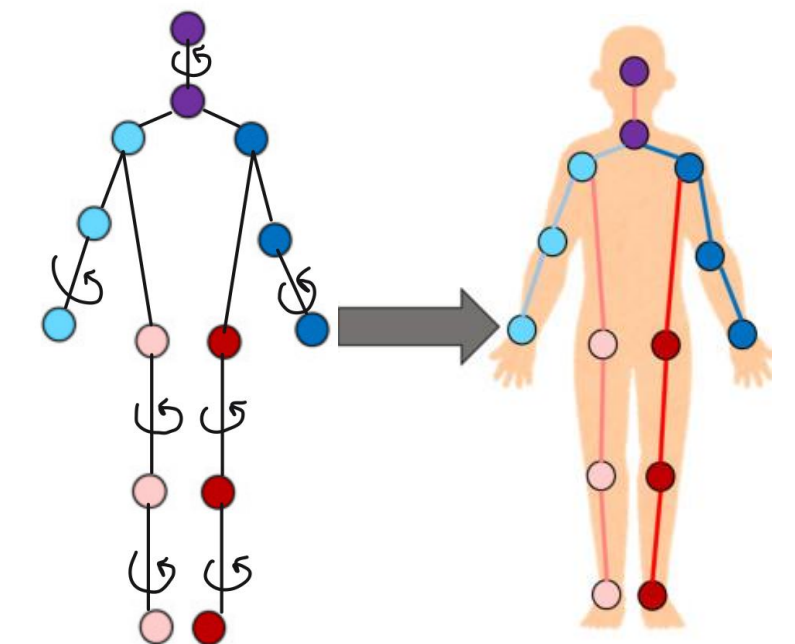## *Common Pose Representations*

**Positional**

- Euclidean joint locations [Zhou et al., 2018]
- Motion Capture markers [Zhang et al., 2020]

**Angular**

- Exponential Maps
- Quaternions
- Euler angles
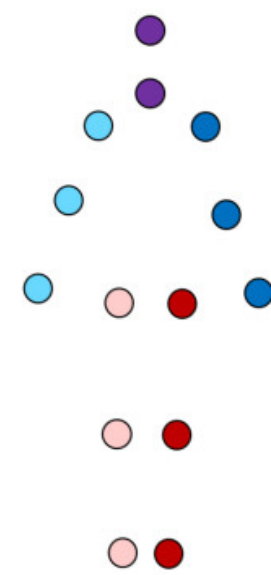- Rotation Matrices
- Ortho6D [Zhou et al., 2019]

**Hybrid**

- Positional and angular
- Joint velocities/accelerations [Holden et al., 2017]
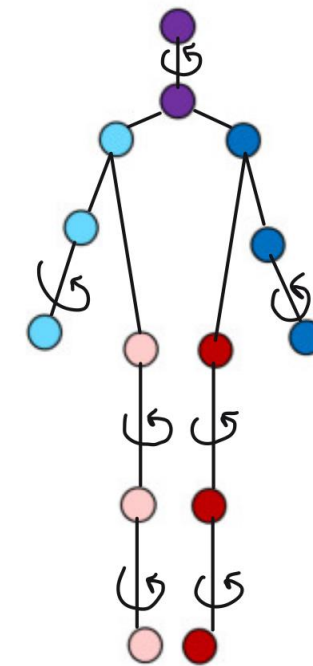* Angular representations with positional losses [Aberman et al., 2020]

# Motivation
## *Common Pose Representations*



Positional



Angular

✓ Intuitive
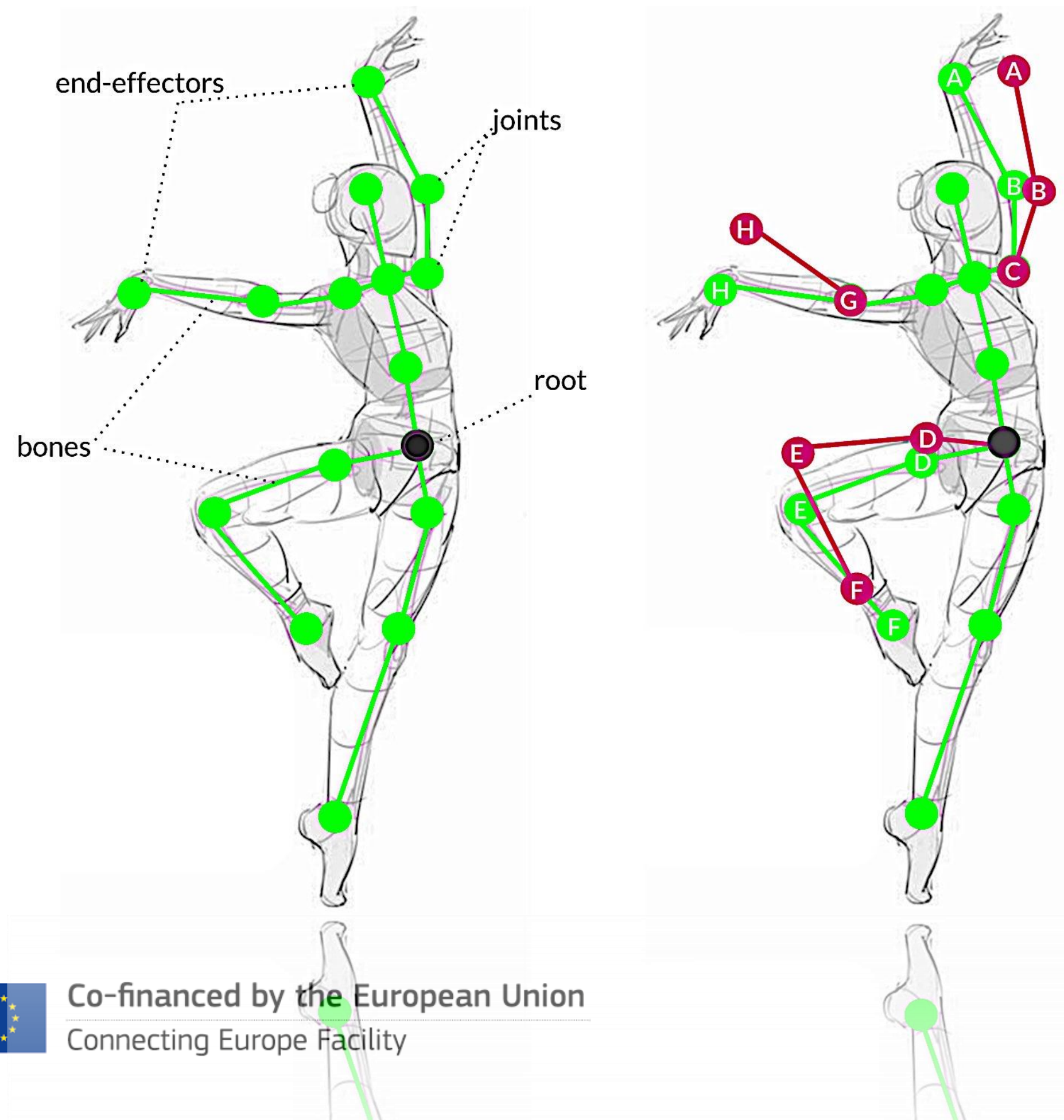✓ Visual result
✗ Not straightforward to apply to different characters

✓ Disentangle shape/skeletal proportions
✓ Convenient to work with
✗ Common rotation representations are discontinuous [Zhou et al., 2019]
✗ Error accumulation [Pavllo et al., 2018]
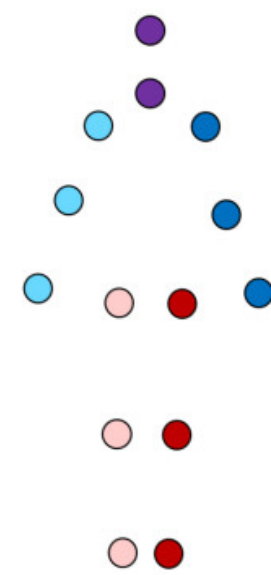
# Method

## *Error accumulation along kinematic tree*

**Problem:** Error accumulation along chain

- o Angular representation causes problems in optimization-based methods

- Angular representations are often paired with loss that averages errors over joints

- Skeleton is a connected graph

- Ignores the fact that prediction errors of different joints have varying impact on qualitative results
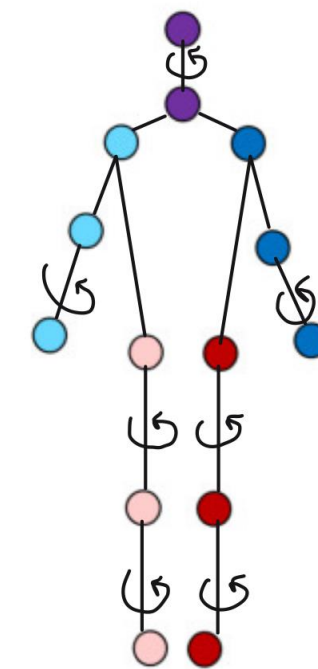
# Motivation

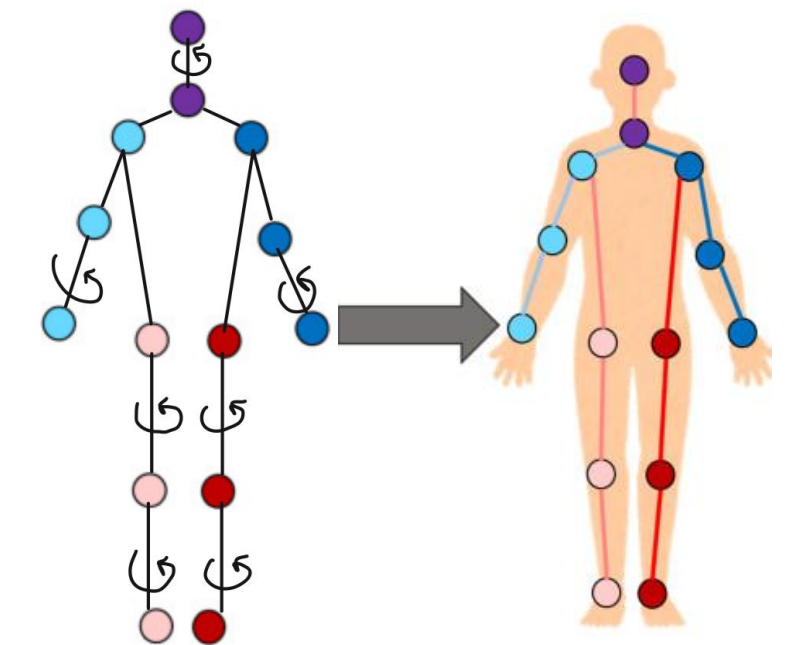## *Common Pose Representations*



**Positional**

✓ Intuitive
✓ Visual result
✗ Not straightforward to apply to different characters

**Angular**

✓ Disentangle shape/skeletal proportions
✓ Convenient to work with
✗ Common rotation representations are discontinuous [Zhou et al., 2019]
✗ Error accumulation [Pavllo et al., 2018]

**Hybrid**

✓ Combinations of positional + angular work better
✓ Angular representations can be paired with positional losses (requires FK)
✗ Excessive information
✗ Correspondence often ignored
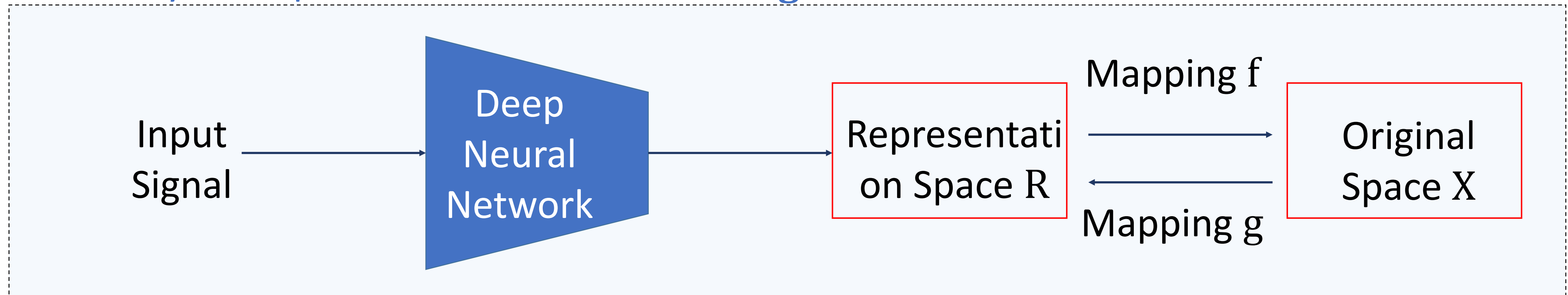✗ Positional losses hinder the rotational information

# Deep Neural Networks

**How to get a continuous representation in neural networks?**

Let's say that the

- mapping to the original space $\mathbf{f} : \mathbf{R} \to \mathbf{X}$, and

- mapping to the representation space $\mathbf{g} : \mathbf{X} \to \mathbf{R}$.

We can say $(\mathbf{f}, \mathbf{g})$ is a good *representation* if for every $\mathbf{x} \in \mathbf{X}$; $\mathbf{f}(\mathbf{g}(\mathbf{x})) = \mathbf{x}$, that is, $\mathbf{f}$ is a left inverse of $\mathbf{g}$.

We can say the representation is continuous if $\mathbf{g}$ is continuous.

Input Signal → Deep Neural Network → Representation Space R

Mapping f

Mapping g

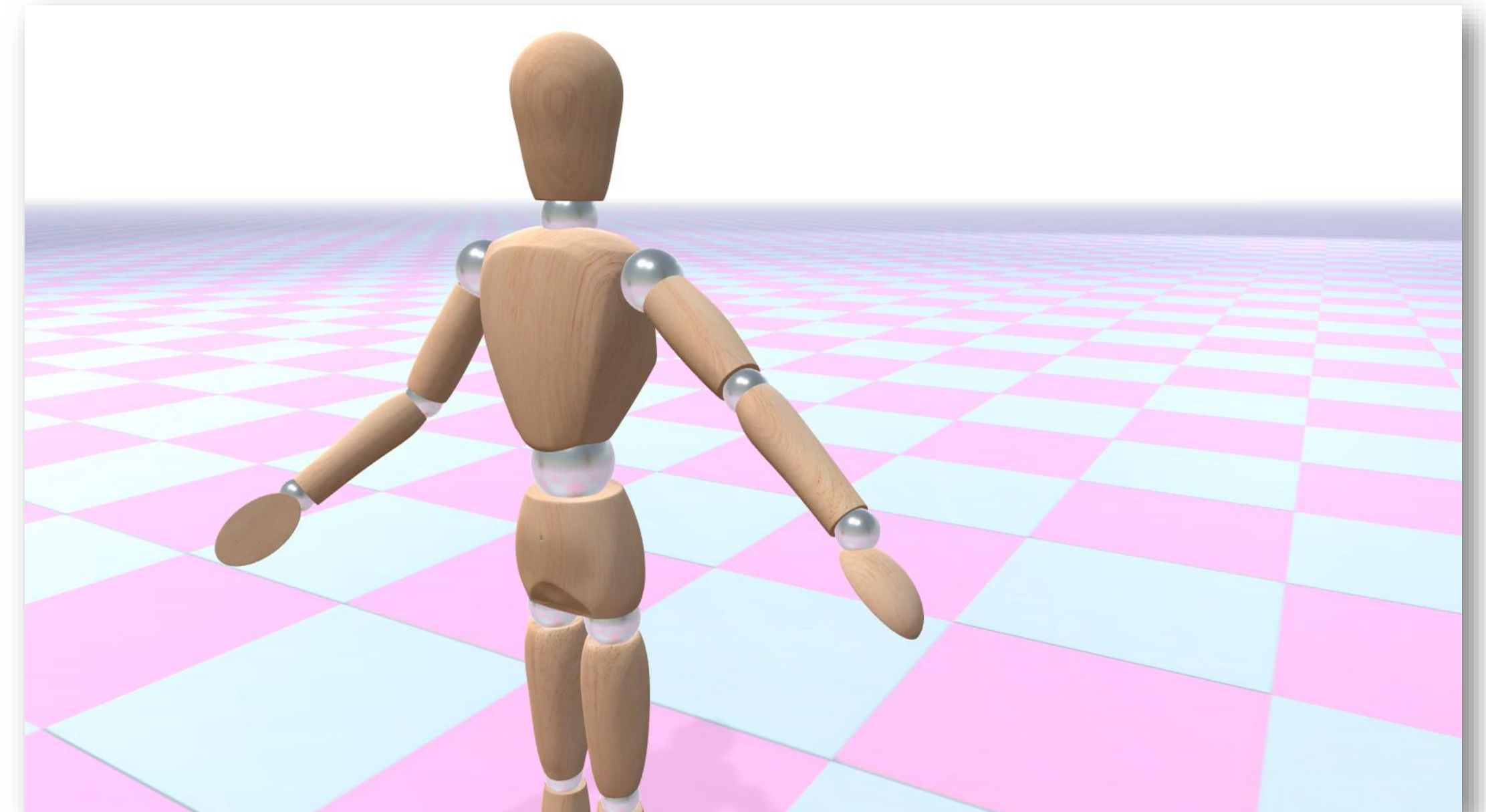Original Space X

*more details in [Zhou et al. 2018]

# Positional Data

Have been used on early machine learning approaches

- **Advantages:** Good in continuity

- **Disadvantages : (a)** Ambiguity problems → cannot describe the full human motion articulation, **(b)** Skeletal model violations
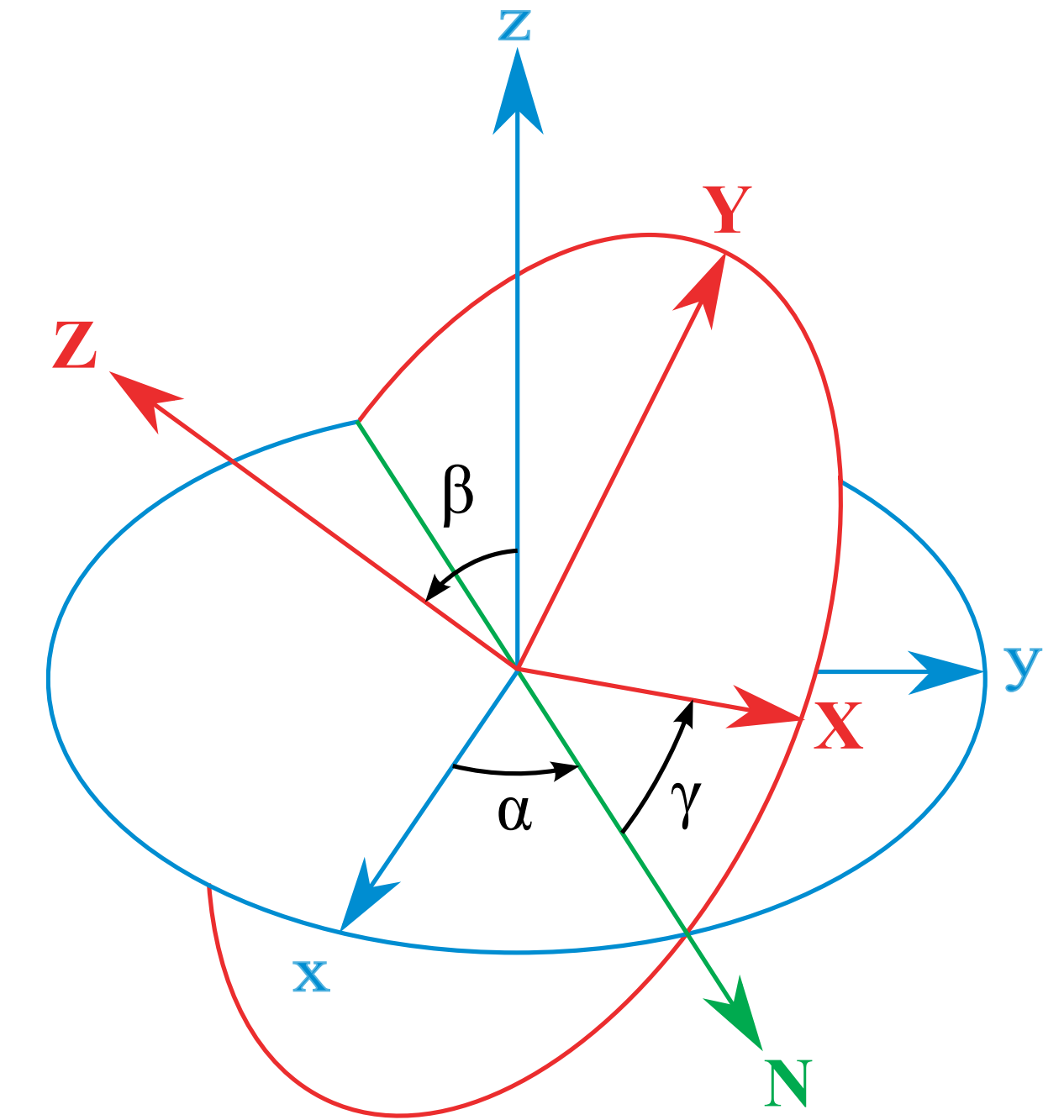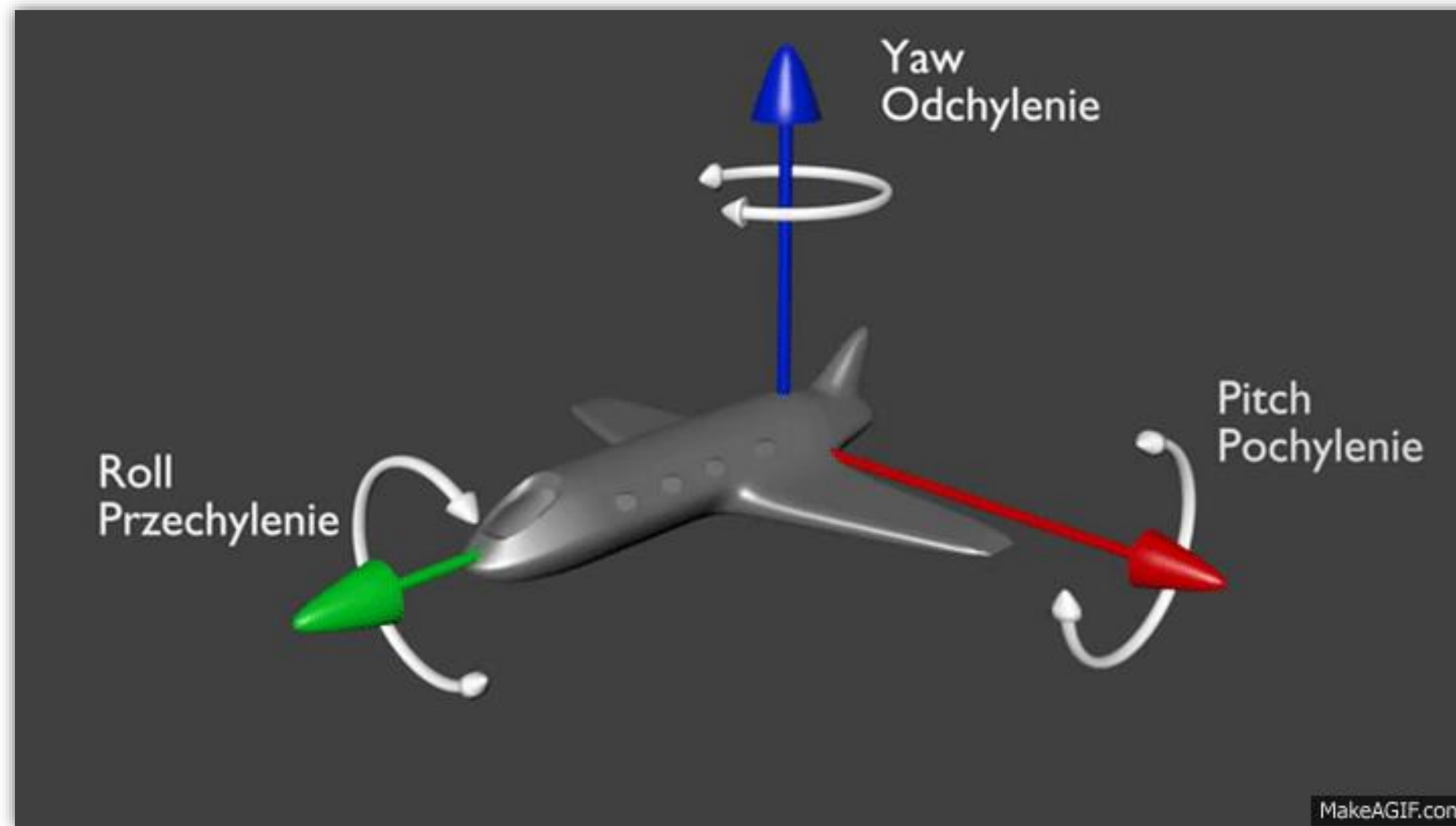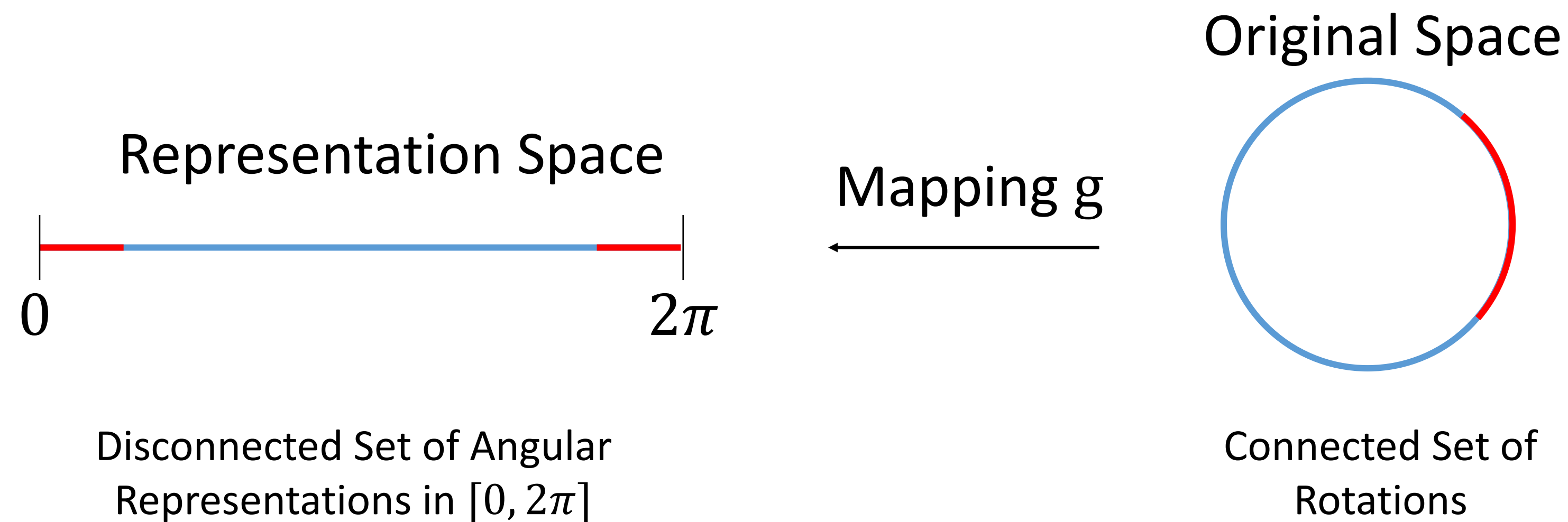


[Cao et al. 2018]

# Euler Angles

Rotate the angles of $\gamma, \beta$ and $\alpha$ along the $X, Y$ and $Z$ axes from the reference frame.

# Euler Angles: *Limitations*

- Gimbal Lock
- Discontinuity
- Singularities that cause learning problems

Original Space

Representation Space

Mapping g

$0$             $2\pi$

Disconnected Set of Angular
Representations in $[0, 2\pi]$
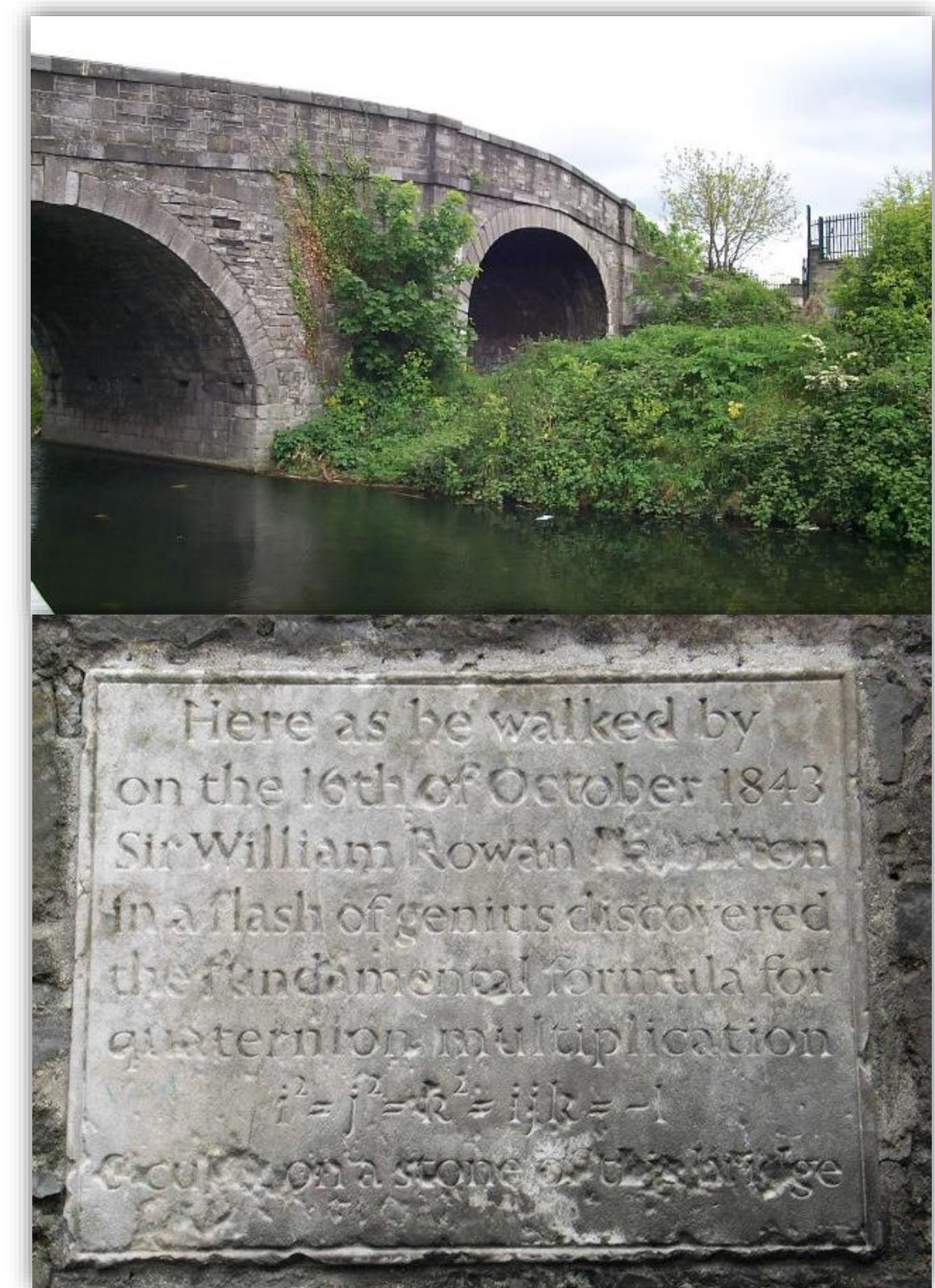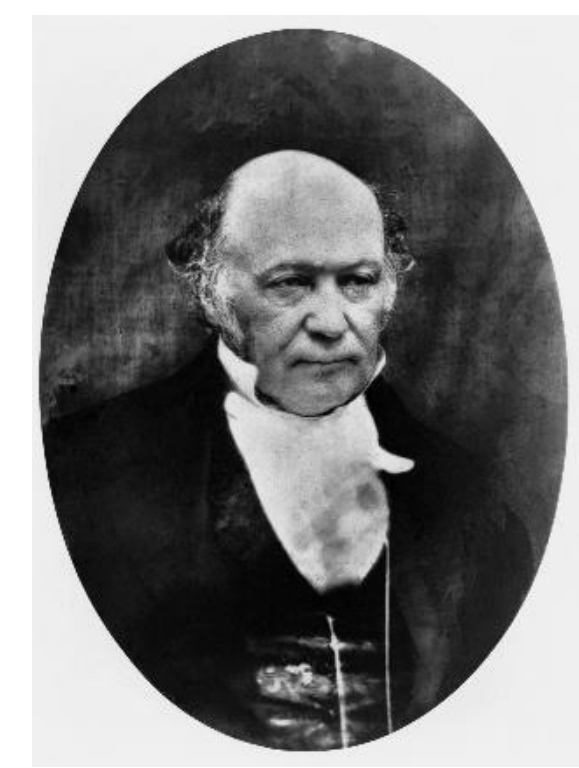
Connected Set of
Rotations

# Quaternions

- Mathematical abstractions alternative to Euler Angles

- Revised and Formulated by Sir William R. Hamilton in 1843

- 4-D complex numbers
  - With one real axis
  - And three imaginary axes, the basis vectors

$$\mathbf{i}, \mathbf{j}, \mathbf{k}$$

## How are quaternions represented?

$$\mathbf{q} = (w, \mathbf{V}) = w + x\mathbf{i} + y\mathbf{j} + z\mathbf{k}$$
$$\mathbf{q} = (q_0, \mathbf{V}) = q_0 + q_1\mathbf{i} + q_2\mathbf{j} + q_3\mathbf{k} \quad \text{or}$$
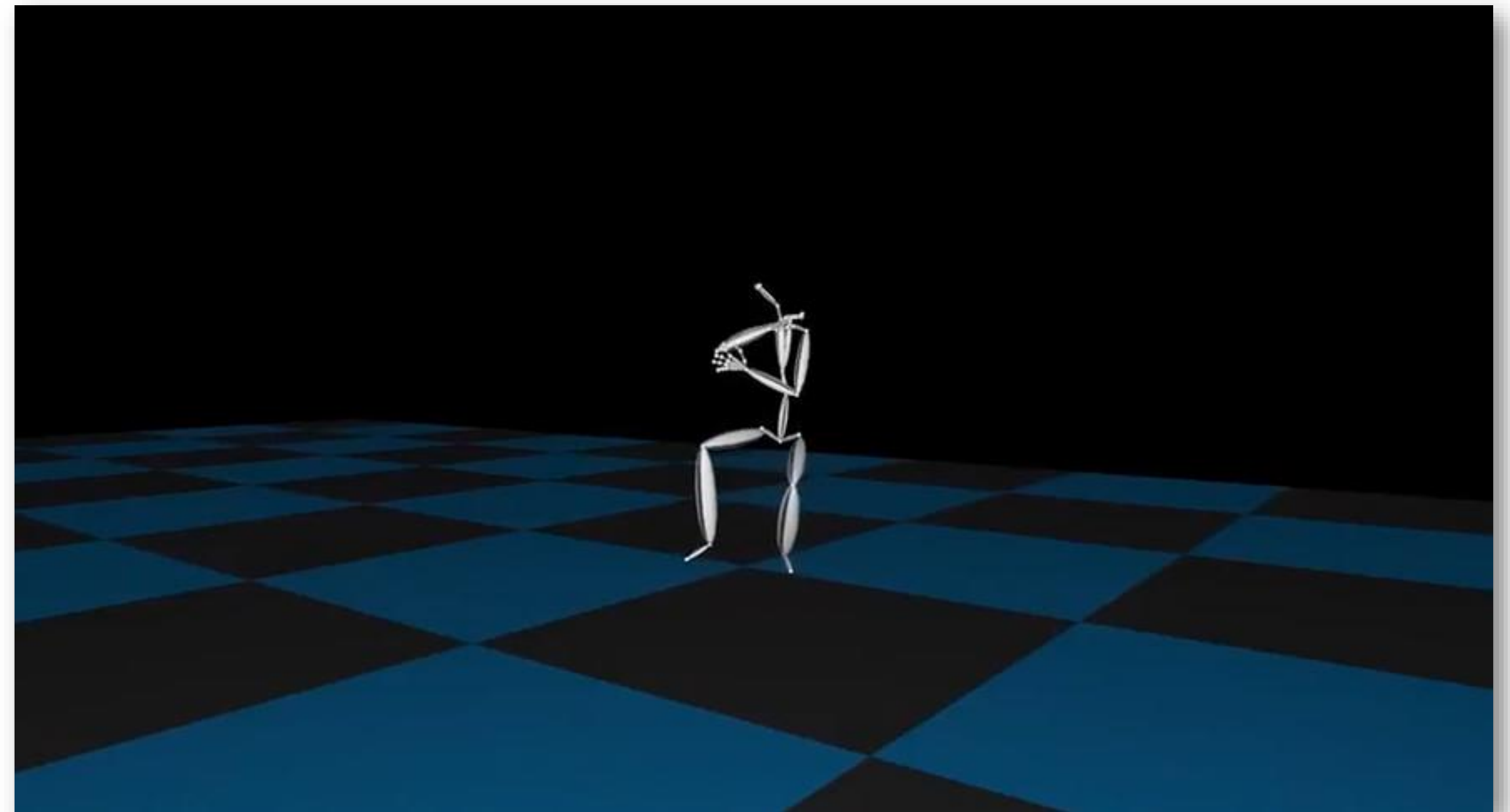
Hamilton Math Inst.,
Trinity College

# Motion representation in popular works

In an attempt to overcome these limitations, the character animation community proposed some alternatives/improvements:
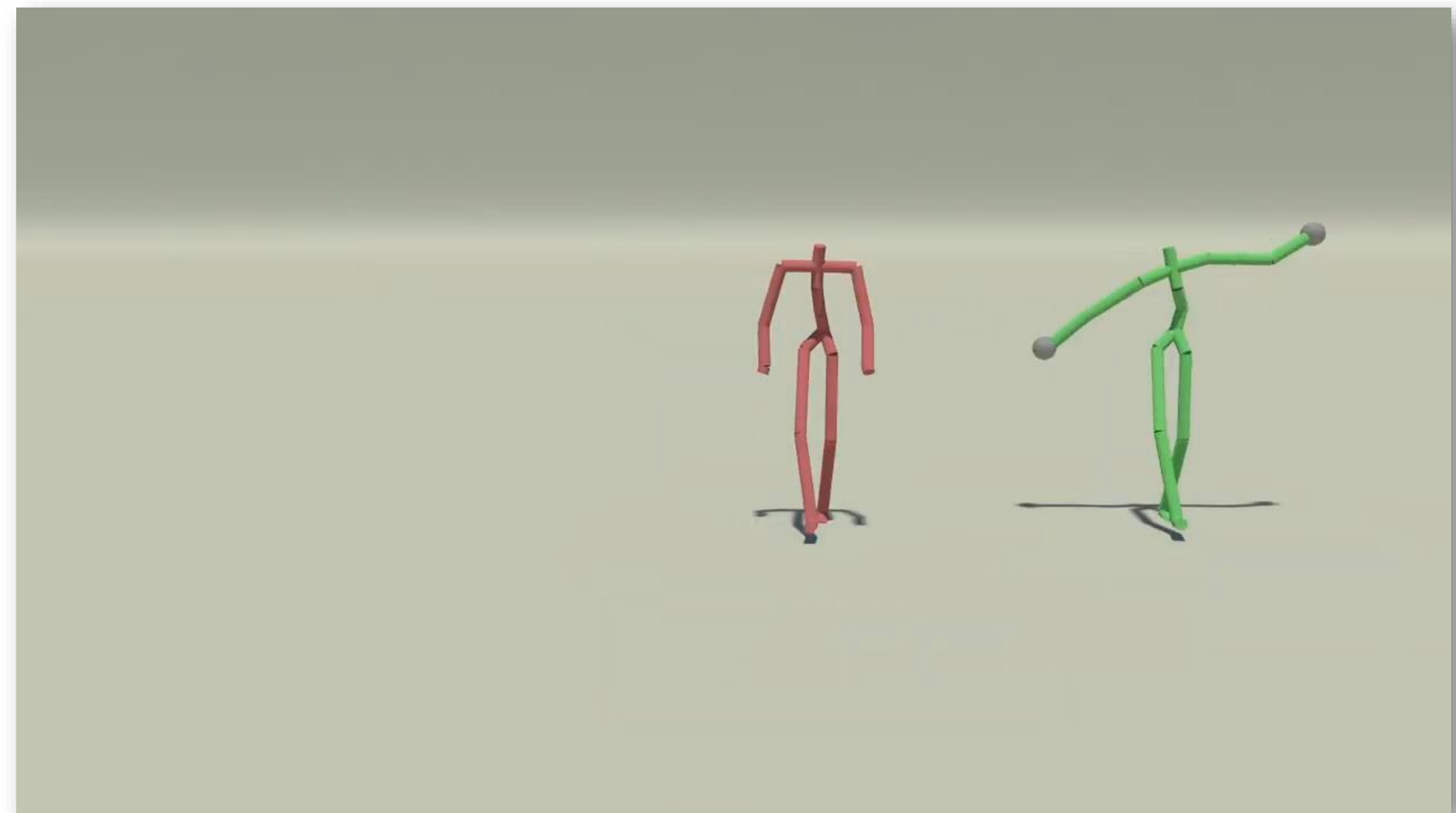
Training using only positional data:

- Zhou et al. 2018. Auto-Conditioned Recurrent Networks for Extended Complex Human Motion Synthesis. International Conference on Learning Representations

# Positional data, with bone length constraints:

- Holden et al. 2016. A deep learning framework for character motion synthesis and editing. ACM Trans. Graphics.

- Wang et al. 2021. Spatio-temporal manifold learning for human motions via long-horizon modelling. IEEE Trans. Visualization and Computer Graphics.

# Quaternions, with a Forward Kinematic layer so as to add a positional loss:

- Harvey et al. 2020. Robust Motion In-betweening. ACM Trans. Graphics.

- Aberman et al. 2020. Skeleton-Aware Networks for Deep Motion Retargeting. ACM Trans. Graphics.

Pavllo et al. 2018. QuaterNet: A Quaternion-based Recurrent Model for Human Motion. British Machine Vision Conference

# Quaternions, amended with positional data:

- Park et al. 2021. Diverse Motion Stylization for Multiple Style Domains via Spatial-Temporal Graph-Based Generative Model. ACM Comput. Graph. Interact. Tech.



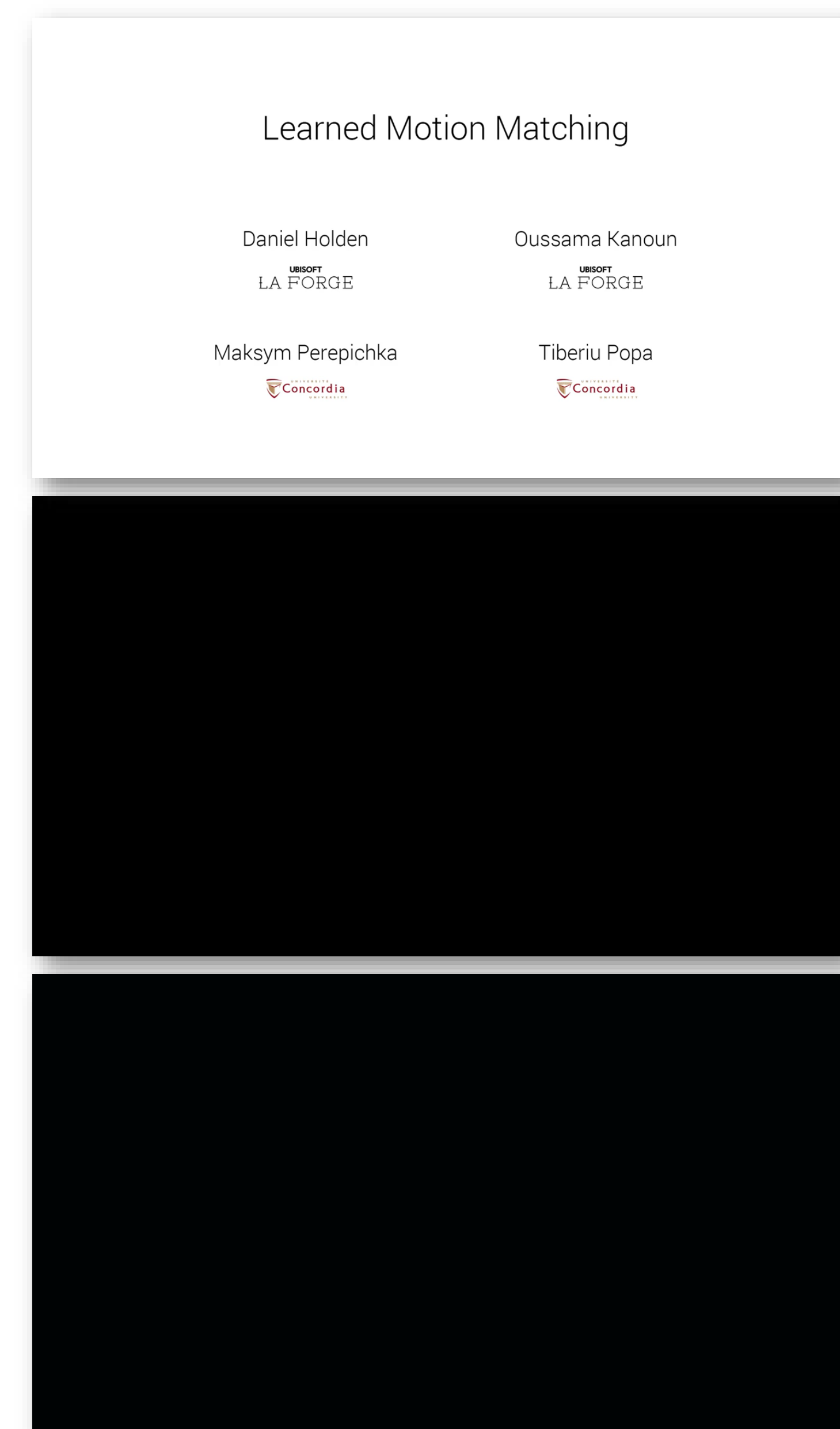Diverse Motion Stylization for Multiple Style Domains via Spatial-Temporal Graph-Based Generative Model

(Supplementary material)
Full demo

## Quaternions, amended with positional data, and motion dynamics:

- Holden et al. 2021. Learned Motion Matching. ACM Trans. Graphics

- Starke et al. 2021. Neural Animation Layering for synthesizing martial arts movements. ACM Trans. Graphics

- Starke et al. 2021. Neural state machine for character-scene interactions. ACM Trans. Graphics



Learned Motion Matching

Daniel Holden — UBISOFT LA FORGE
Oussama Kanoun — UBISOFT LA FORGE
Maksym Perepichka — Concordia
Tiberiu Popa — Concordia

# Method

## *Dual Quaternion Representation*

- Hybrid representation based on Dual Quaternions
- Unified entity

$$q = q_r + \epsilon q_d \text{ where } \epsilon^2 = 0$$

rotation        translation

- More compact than homogeneous transformation matrix (8 values per joint) and efficient [Kenwright et al., 2012]
- Well-established mathematical properties

# Method
*Dual Quaternion Representation*

- Can be defined in a root-centered coordinate system mitigating common problems such as error accumulation along the kinematic chain [Pavllo et al., 2018]

# Method
## *Losses*



Dual Quaternions

Motion Dataset

Pre-processing

Input

Gradient updates
(Losses)

Output

Post-processing

- Euclidean distance of joint positions/locations
- MSE on joint rotations
- **Offset loss** → maintain skeletal structure

# Deep character animation networks



Holden et al., 2016



Holden et al., 2017

## Deep character animation networks



Alexanderson et al., 2020

# Deep character animation networks



Ling et al., 2021

# Deep character animation networks



Frangiadaki et al. 2015



Zhou et al. 2018

Co-financed by the European Union
Connecting Europe Facility

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423

# Many challenges in Character Animation have been re-defined

Rigging/Skinning

Motion Synthesis

Motion in-betweening

Motion Control

Motion Retargeting

Style Transfer

Audio/music-driven synthesis

Text-to-animation

etc.

# Research in our lab

# MotionNet

**MotioNet:** 3D Human Motion Reconstruction from Monocular Video with Skeleton Consistency

by M. Shi, K. Aberman, A. Aristidou, T. Komura, D. Lischinski, D. Cohen-Or, B. Chen

ACM Transactions on Graphics

[Pavllo et al., CVPR 2019 ]

**Use IK to convert the 3d position to rotation**

[VNect, Mehta et al., SIGGRAPH 2017 ]

**Apply rigging to make the rotation to a consistent skeleton**

[HMR, Kanazawa et al., CVPR 2018 ]

Co-financed by the European Union
Connecting Europe Facility

86

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423

**The Key Idea:**

What is the common representation of motion in MoCap datasets?

BVH - the most used output format of MoCap system



Initial pose with a hierarchical structure

```
Frames: 2
Frame Time: 0.04166667
-9.533684    4.447926    -0.566564    -7.757381    -1.735414    89.207932    9.763572
             6.289016    -1.825344    -6.106647     3.973667    -3.706973   -6.474916
            -14.391472   -3.461282   -16.504230     3.973544    -3.805107   22.204674
             2.533497   -28.283911   -6.862538      6.191492     4.448771  -16.292816
             2.951538    -3.418231    7.634442      11.325822     5.149696  -23.069189
           -18.352753    15.051558   -7.514462      8.397663     2.953842   -7.213992
             2.494318    -1.543435    2.970936     -25.086460    -4.195537   -1.752307
             7.093068    -1.507532   -2.633332      3.858087      0.256802   7.892136
            12.803010   -28.692566    2.151862     -9.164188      8.006427   -5.641034
           -12.596124    4.366460
```

Time-framed joint information(rotation)

skeleton

$\mathbf{q}_0^t$   $\mathbf{q}_1^t$   $\mathbf{q}_3^t$

. . .

$\oplus$ Rotation = Human pose

Forward
Kinematics

motion

Joint rotations

$\oplus$: Forward kinematics

**MAI4CARE**

# MotioNet: 3D Human Motion Reconstruction from Monocular Video with Skeleton Consistency

MINGYI SHI, Shandong University, China, and AICFVE, Beijing Film Academy, China
KFIR ABERMAN, AICFVE, Beijing Film Academy, China, and Tel-Aviv University, Israel
ANDREAS ARISTIDOU, University of Cyprus and RISE Research Centre, Cyprus
TAKU KOMURA, Edinburgh University, Japan
DANI LISCHINSKI, Shandong University, China and The Hebrew University of Jerusalem, Israel and AICFVE, Beijing Film Academy, Israel
DANIEL COHEN-OR, Tel-Aviv University, Israel, and AICFVE, Beijing Film Academy, Israel
BAOQUAN CHEN, CFCS, Peking University, China, and AICFVE, Beijing Film Academy, China

Fig. 1. Given a monocular video of a performer, our approach, MotioNet, reconstructs a complete representation of the motion, consisting of a single symmetric skeleton, and a sequence of global root positions and 3D joint rotations. Thus, inverse kinematics is effectively integrated within the network and is data-driven rather than based on a universal prior. The images on the right were rendered from the output of our system after a simple rigging process.

We introduce *MotioNet*, a deep neural network that directly reconstructs the motion of a 3D human skeleton from a monocular video. While previous methods rely on either rigging or inverse kinematics (IK) to associate a consistent skeleton with temporally coherent joint rotations, our method is the first data-driven approach that directly outputs a kinematic skeleton, which is a complete, commonly used motion representation. At the crux of our approach lies a deep neural network with embedded kinematic priors, which decomposes sequences of 2D joint positions into two separate attributes: a single, symmetric skeleton encoded by bone lengths, and a sequence of 3D joint rotations associated with global root positions and foot contact labels. These attributes are fed into an integrated forward kinematics (FK) layer that outputs 3D positions, which are compared to a ground truth. In addition, an adversarial loss is applied to the velocities of the recovered rotations to ensure that they lie on the manifold of natural joint rotations. The key advantage of our approach is that it learns to infer natural joint rotations directly from the training data rather than assuming an underlying model, or inferring them from joint positions using a data-agnostic IK solver. We show that enforcing a single consistent skeleton along with temporally coherent joint rotations constrains the solution space, leading to a more robust handling of self-occlusions and depth ambiguities.

CCS Concepts: • **Computing methodologies → Motion processing**; **Neural networks**;

Additional Key Words and Phrases: Pose estimation, motion capturing, motion analysis

## 1 INTRODUCTION

Capturing the motion of humans has long been a fundamental task with a wide spectrum of applications in data-driven computer animation, special effects, gaming, activity recognition, and behavioral analysis. Motion is most accurately captured in a controlled setting using specialized hardware, such as magnetic
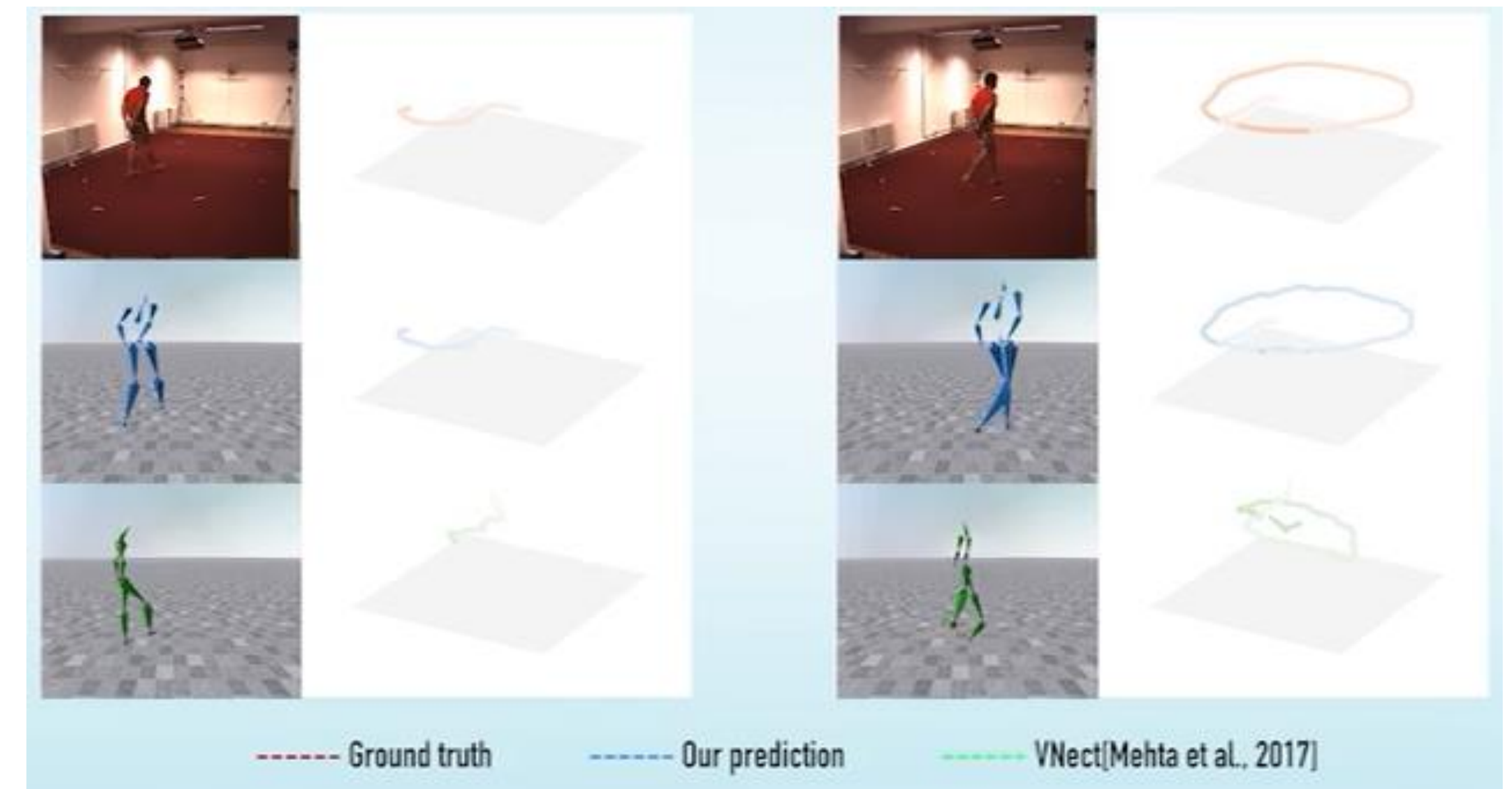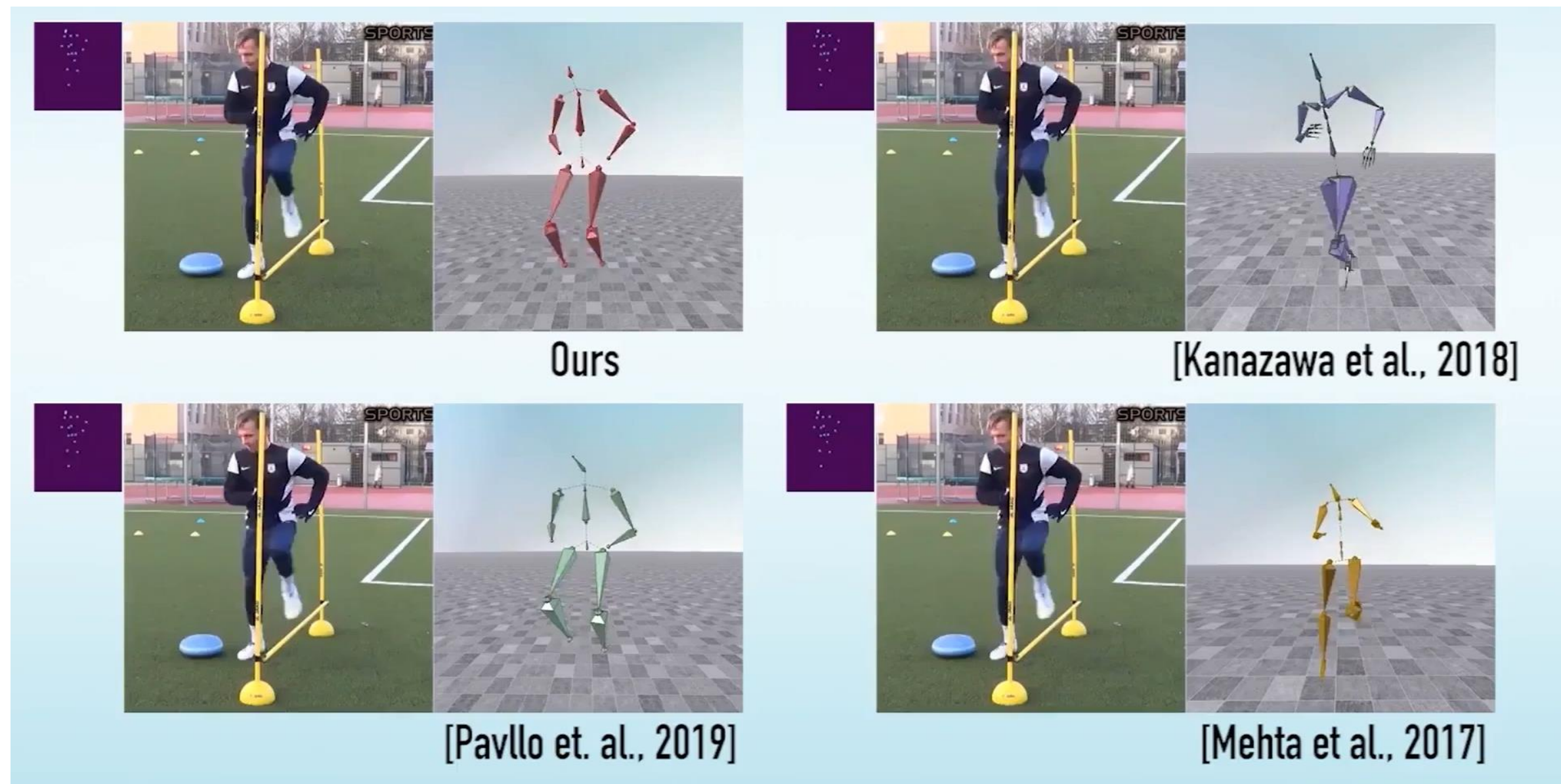
**confidence value simulation**

$T$

2D poses

+

confidence values

natural rotations

$\frac{\partial}{\partial t}$

$D$ → real / fake

3D poses

Co-financed by the European
Connecting Europe Facility

run under the context of Action
nanced by the EU CEF Telecom
INEA/CEF/ICT/A2020/2267423

# Results



Ours

[Kanazawa et al., 2018]

[Pavllo et. al., 2019]

[Mehta et al., 2017]



------- Ground truth       ------- Our prediction       ------- VNect[Mehta et al., 2017]

# Motion Analysis
# Emotion and Style



**Maori wedding (Haka)**
https://youtu.be/QUbx-AcDgXo

**Emotions through dance**
https://youtu.be/m0R-ftFBm38

# Human Motion Style



**Happy**

**Depressed**

**Style** is an **abstract** attribute

# Related Work

Model **style**, which is **not well-defined**, as some **hand-crafted representations**, such as,

- Difference in spectral domain.
- Physical parameters of human body.
- Low-level features based on the LMA theories on human analysis

# Related Work

**Data Driven**

Learn the mapping based on
**labeled & paired motion data**.

- limited to **structurally similar** motions in the dataset

- limited to a **pre-defined set** of styles in mocap data

- limited to **style recorded by MoCap** systems

# Inspiration from Image Style Transfer

Adaptive Instance Normalization **(AdaIN) layer -** spatially invariant, maintains geometry, manipulates style.



$$\text{AdaIN}(x, y) = \sigma(y) \left( \frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$

[Huang and Belongie, 2017]

style

conten t



**Geometry and shapes** are preserved

**Adapting to** motion?

# Adaptive Instance Normalization (AdaIN)



Deep Joint Rotations vs Time

**Shape of signal** is preserved

# Architecture

Temporally-invariant

Style Motion $\mathbf{y}^t$

(joint positions)

$E_S$

Content Motion $\mathbf{x}^s$

(joint rotations)

$E_C$ IN

# Loss Terms



Is $\tilde{\mathbf{x}}^t$ a real motion of style t?

Style Motion $\mathbf{y}^t$ (joint positions)

$E_S$

Style code $\mathbf{z}_S$

MLP

Content Motion $\mathbf{X}^s$ (joint rotations)

$E_C$  IN

$\mathbf{z}_c$ Temporal content code

$F$  AdaIN

Output Motion $\tilde{\mathbf{X}}^t$

# Loss Terms



Style Motion $\mathbf{y}^S$ (joint positions)

$E_S$

Style code $\mathbf{z}_S$

MLP

Multi-Style Adversarial Loss

$D$

Content Motion $\mathbf{x}^S$ (joint rotations)

$E_C$ IN

Temporal content code $\mathbf{z}_c$

$F$ AdaIN

Output Motion $\tilde{\mathbf{x}}^S$

Content Consistency Loss

$\tilde{\mathbf{x}}^S$ reconstructs $\mathbf{x}^S$ ?

# Loss Terms



Push embeddings of the same style closer to each other.

Style Motion $\mathbf{y}^t$
(joint positions)

$E_S$

Triplet Loss

Style code $\mathbf{z}_S$

MLP

Multi-Style Adversarial Loss

$D$

Content Motion $\mathbf{x}^s$
(joint rotations)

$E_C$ IN

$\mathbf{z}_c$

Temporal content code

$F$ AdaIN

Output Motion $\tilde{\mathbf{x}}^t$

Content Consistency Loss

# Results

Style Input (proud)



Content Input



Output



## Unpaired Motion Style Transfer from Video to Animation

KFIR ABERMAN*, AICFVE, Bejing Film Academy & Tel-Aviv University
YIJIA WENG*, CFCS, Peking University & AICFVE, Beijing Film Academy
DANI LISCHINSKI, The Hebrew University of Jerusalem & AICFVE, Beijing Film Academy
DANIEL COHEN-OR, Tel-Aviv University & AICFVE, Beijing Film Academy
BAOQUAN CHEN†, CFCS, Peking University & AICFVE, Beijing Film Academy

Transferring the motion style from one animation clip to another, while preserving the motion content of the latter, has been a long-standing problem in character animation. Most existing data-driven approaches are supervised and rely on paired data, where motions with the same content are performed in different styles. In addition, these approaches are limited to transfer of styles that were seen during training.

In this paper, we present a novel data-driven framework for motion style transfer, which learns from an unpaired collection of motions with style labels, and enables transferring motion styles not observed during training. Furthermore, our framework is able to extract motion styles directly from videos, bypassing 3D reconstruction, and apply them to the 3D input motion.

Our style transfer network encodes motions into two latent codes, for content and for style, each of which plays a different role in the decoding (synthesis) process. While the content code is decoded into the output motion by several temporal convolutional layers, the style code modifies deep features via temporally invariant adaptive instance normalization (AdaIN).

Moreover, while the content code is encoded from 3D joint rotations, we learn a common embedding for style from either 3D or 2D joint positions, enabling style extraction from videos.

Our results are comparable to the state-of-the-art, despite not requiring paired training data, and outperform other methods when transferring previously unseen styles. To our knowledge, we are the first to demonstrate style transfer directly from videos to 3D animations - an ability which enables one to extend the set of style examples far beyond motions captured by MoCap systems.

Fig. 1. Style transfer from video to animation. Our network, which is trained with unpaired motion sequences, learns to disentangle content and style. Our trained generator is able to produce a motion sequence that combines the content of a 3D sequence with the style extracted directly from a video.

CCS Concepts: • **Computing methodologies → Motion processing; Neural networks.**

Additional Key Words and Phrases: motion analysis, style transfer

## 1 INTRODUCTION

The style of human motion may be thought of as the collection of motion attributes that convey the mood and the personality of

*equal contribution
†corresponding author

Authors' addresses: Kfir Aberman, kfiraberman@gmail.com; Yijia Weng, halfsummer11@gmail.com; Dani Lischinski, danix3d@gmail.com; Daniel Cohen-Or, cohenor@gmail.com; Baoquan Chen, baoquan@pku.edu.cn.

# Challenges

- Data are not always available…



Russell from movie "Up"    Bonnie from "Toy Story 4"

102

# Challenges

- Given the difficulties of motion capturing children, can we just use adult mocap data on child characters?

- Can we convince the viewers that the motions are from children?

- Jain et al[2016] found that viewers can differentiate child motion from adult motion by viewing point light display videos.

# Key Ideas

- Adapt adult motions to child motions that captures both the postures and the timing of child motions.

- Achieve this goal without temporally aligned data given that adult motions and child motions can be drastically different.

# Overall Architecture

**Adversarial loss**

$$\mathcal{L}_{G_{c2a}} = 0.5 * \mathbb{E}_{c \sim p(c)} \left[ D_a(G_{c2a}(c)) - 1 \right]$$

$$\mathcal{L}_{G_{a2c}} = 0.5 * \mathbb{E}_{a \sim p(a)} \left[ D_c(G_{a2c}(a)) - 1 \right]$$

**Cycle loss**

$$\mathcal{L}_{cycle,c} = G_{a2c}(G_{c2a}(c)) - c$$

$$\mathcal{L}_{cycle,a} = G_{c2a}(G_{a2c}(a)) - a$$

**Coherence loss**

$$\mathcal{L}_{coherence,a} = \sum_t \sum_{DOF} ||G_{a2c}(a)(t) - G_{a2c}(a)(t-1)||$$

$$\mathcal{L}_{coherence,c} = \sum_t \sum_{DOF} ||G_{c2a}(c)(t) - G_{c2a}(c)(t-1)||$$

**Transition loss**

$$y = G_{c2a}(c)$$

$$\mathcal{L}_{transition,c} = \sum_t \sum_{DOF} ||y_i(t_{overlap:end}) - y_{i+1}(0 : t_{overlap})||$$

# Results: *Punch*

Input adult

Ours

Reference child

# Results: *Run as fast as you can*

Input adult        Ours        Reference child

Input adul...

Reference child

# Contextual Analysis

# Deep motifs and motion signatures

# Motion Words and Motifs



Deep Motion Motifs

motion motif

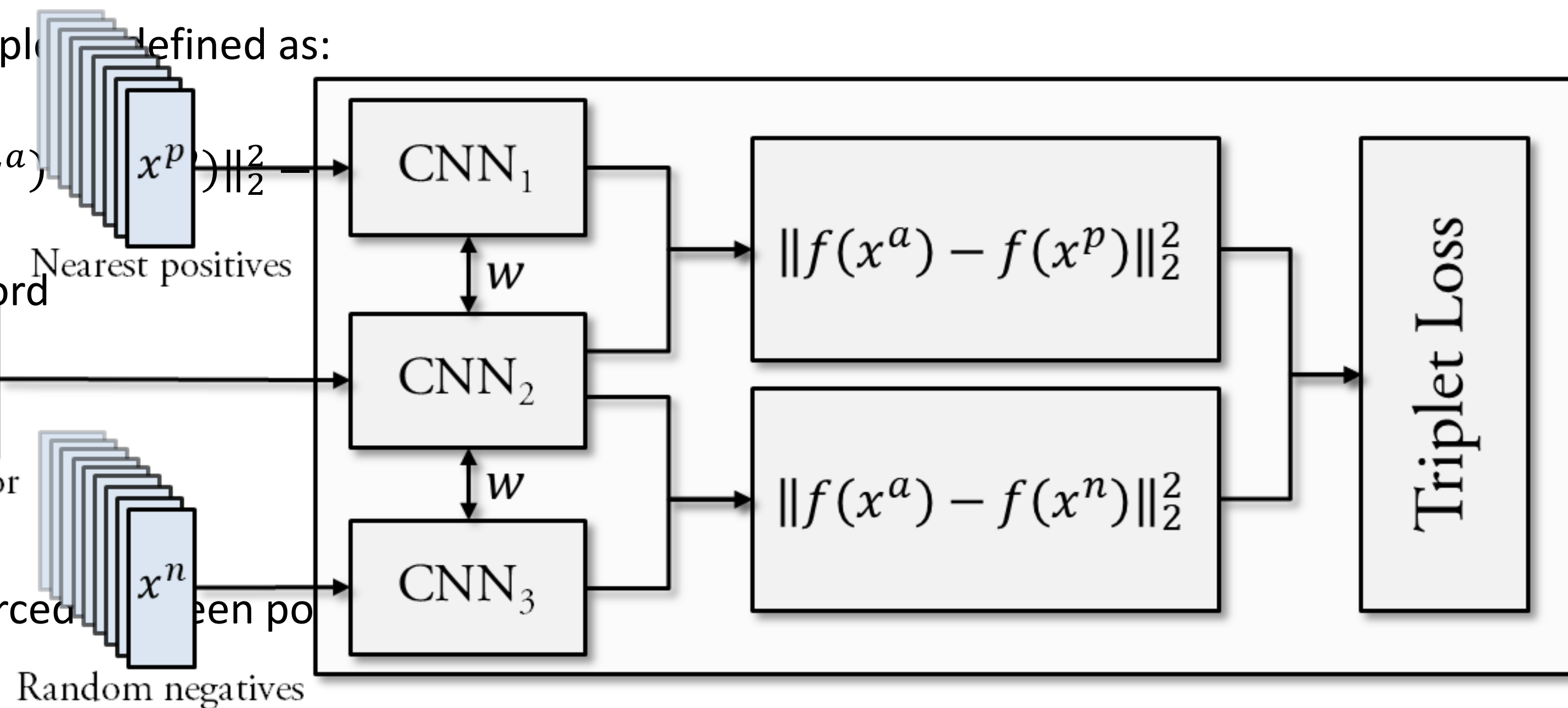motion motif

motion motif

# Triplet Loss Network

The loss for a single triplet is defined as:

$$L(x^a, x^p, x^n) = [\|f(x^a) - f(x^p)\|_2^2$$

$x^a$ - anchor motion word

$x^p$ - positive sample

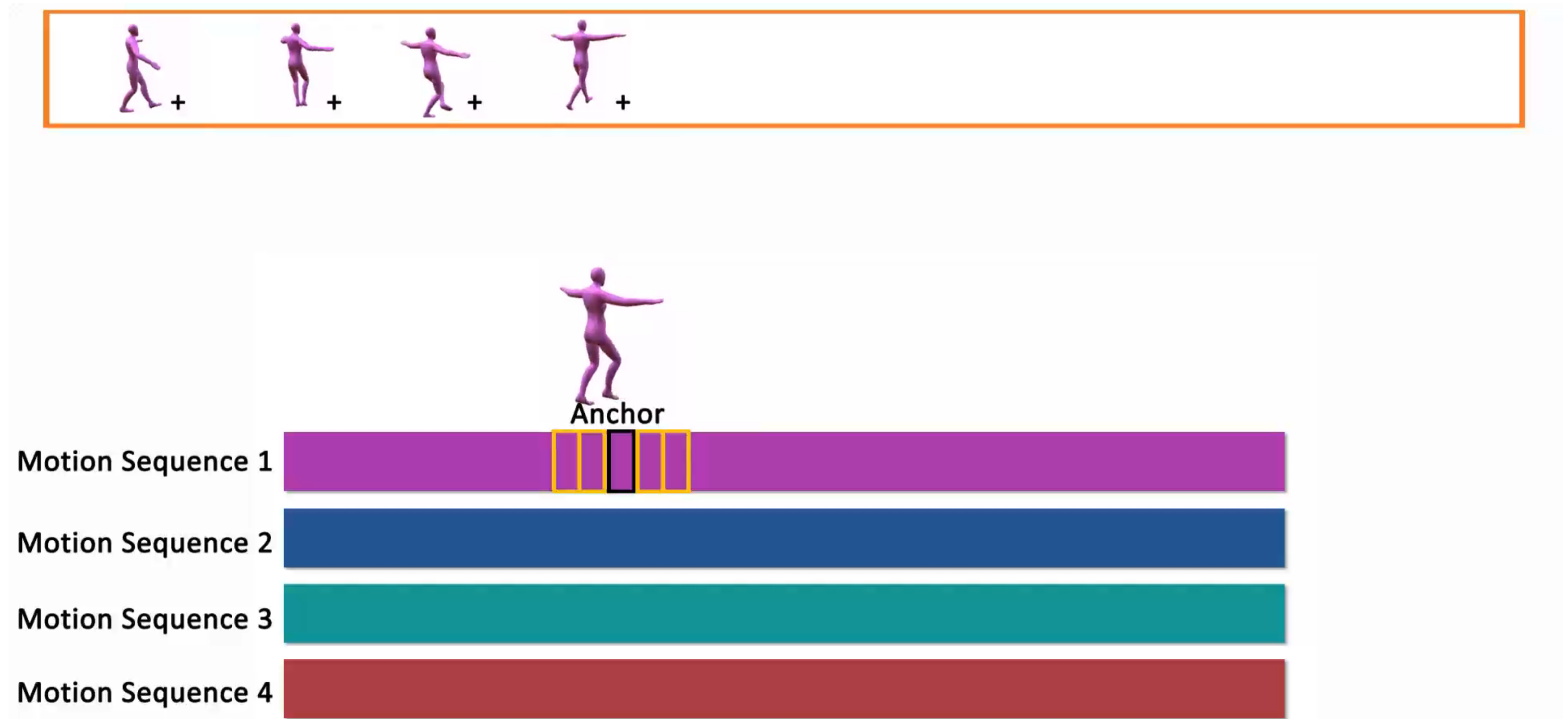$x^n$ - negative sample
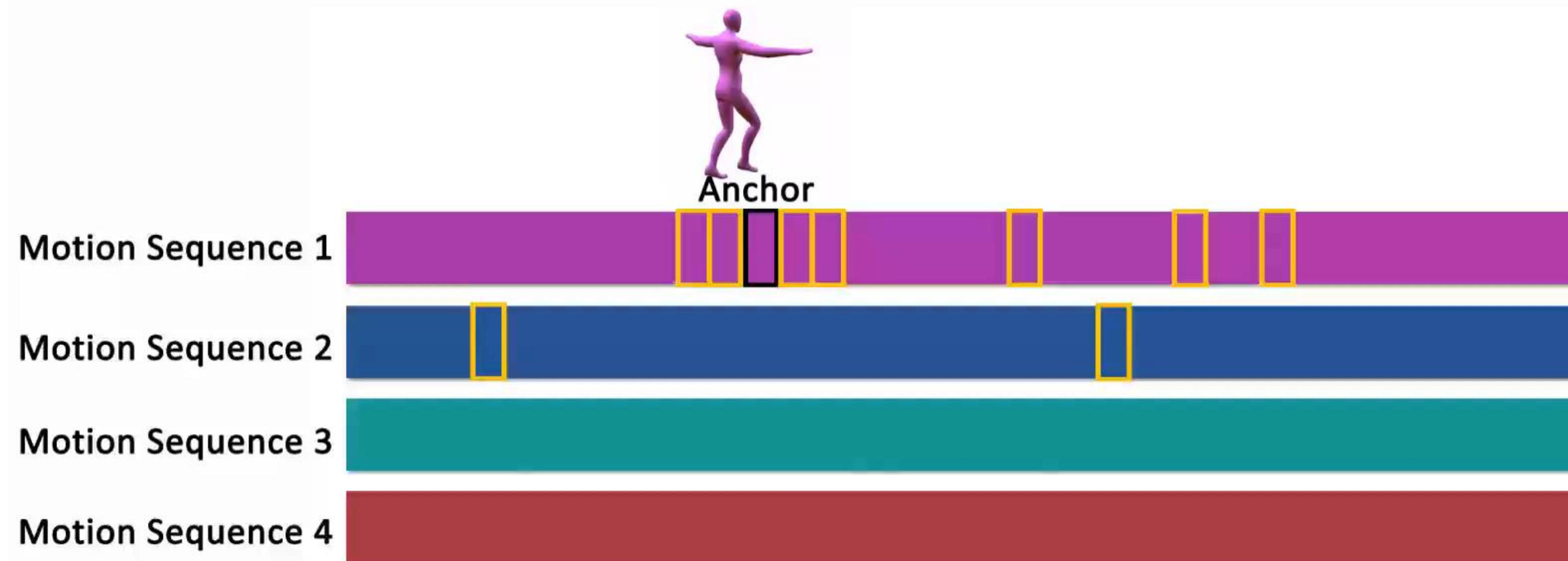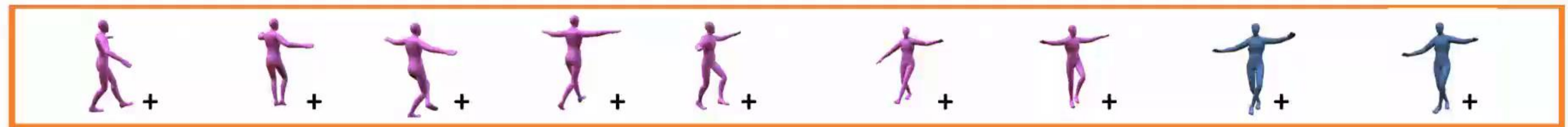
α - margin that is enforced between po

# Triplet Loss Network

# Triplet Loss Network

# Motion Words

# Motion Signatures

# Motion Signatures

# Motion Signatures



Salsa

Modern

Greek Folk

Indian Bollywood

# Deep motifs and motion signatures
# **Fine-grained details**



Data from the SNU Motion Database

Data from the CMU Motion Database

# Fine-grained details



Fighter A

Fighter B

Data from the SNU Motion Database

Leader

Follower

Data from the CMU Motion Database

# Motion Segmentation

# Organizing large collections: *Dance ethnography*

# Deep Motifs and Motion Signatures

ANDREAS ARISTIDOU, The Interdisciplinary Center
DANIEL COHEN-OR, Tel-Aviv University
JESSICA K. HODGINS, Carnegie Mellon University
YIORGOS CHRYSANTHOU, University of Cyprus & RISE Research Center
ARIEL SHAMIR, The Interdisciplinary Center

Fig. 1. Our motion signatures are defined using a deep analysis of motion words and selection of motion-motifs. Each signature is represented by a horizontal bar that shows the frequency of motion-motifs using color coding from red (high) through blue (low) to gray (zero). Note that the signatures represent distributions and not time evolution – the horizontal axis is not temporal. Three signatures of sequences are shown for each motion type – as can be seen, motions of similar type produce similar signatures where many motifs align. The rectangles in the sequence of motion to the left of the signatures illustrate motion words associated with the motifs shown by the corresponding arrow above the signature.

Many analysis tasks for human motion rely on high-level similarity between sequences of motions, that are not an exact matches in joint angles, timing, or ordering of actions. Even the same movements performed by the same person can vary in duration and speed. Similar motions are char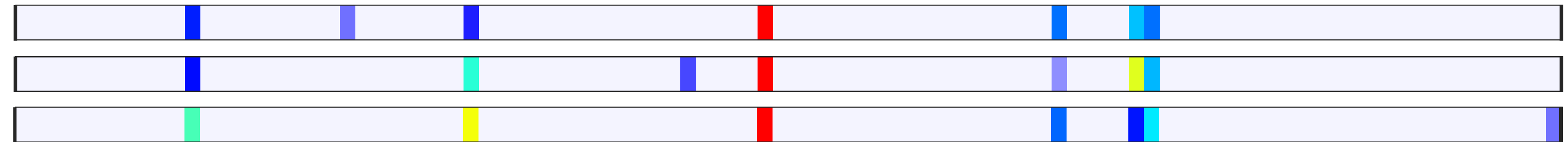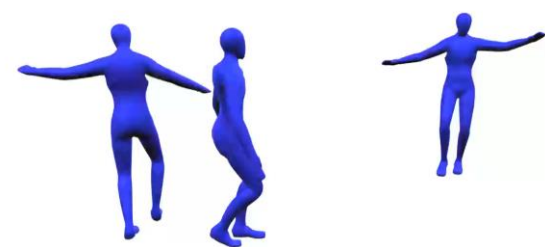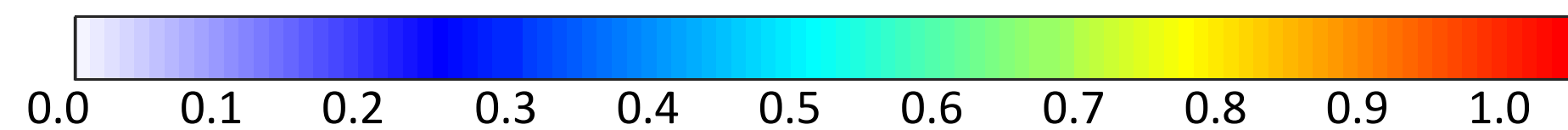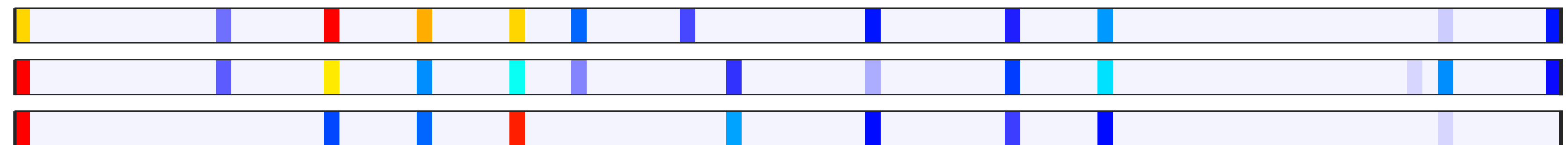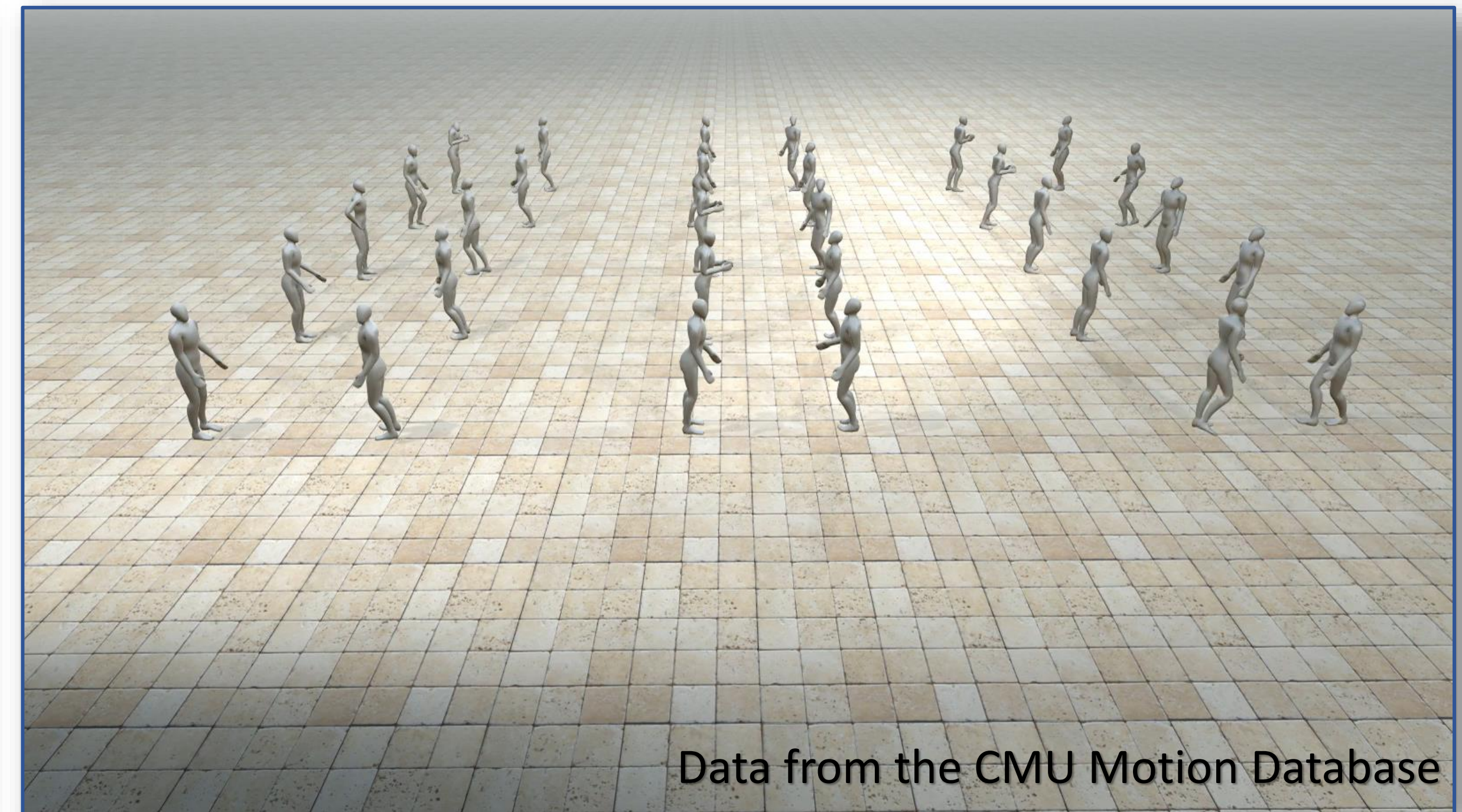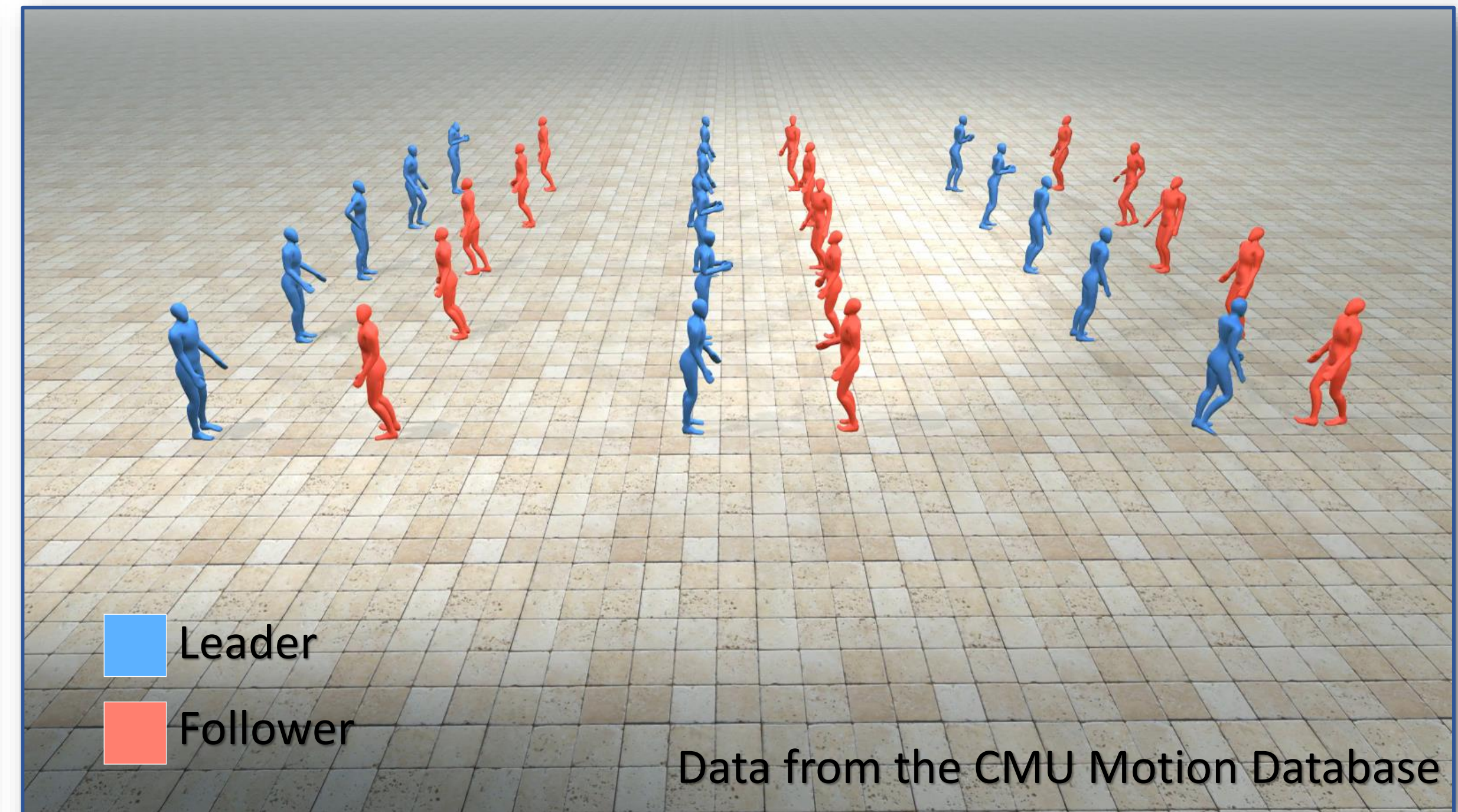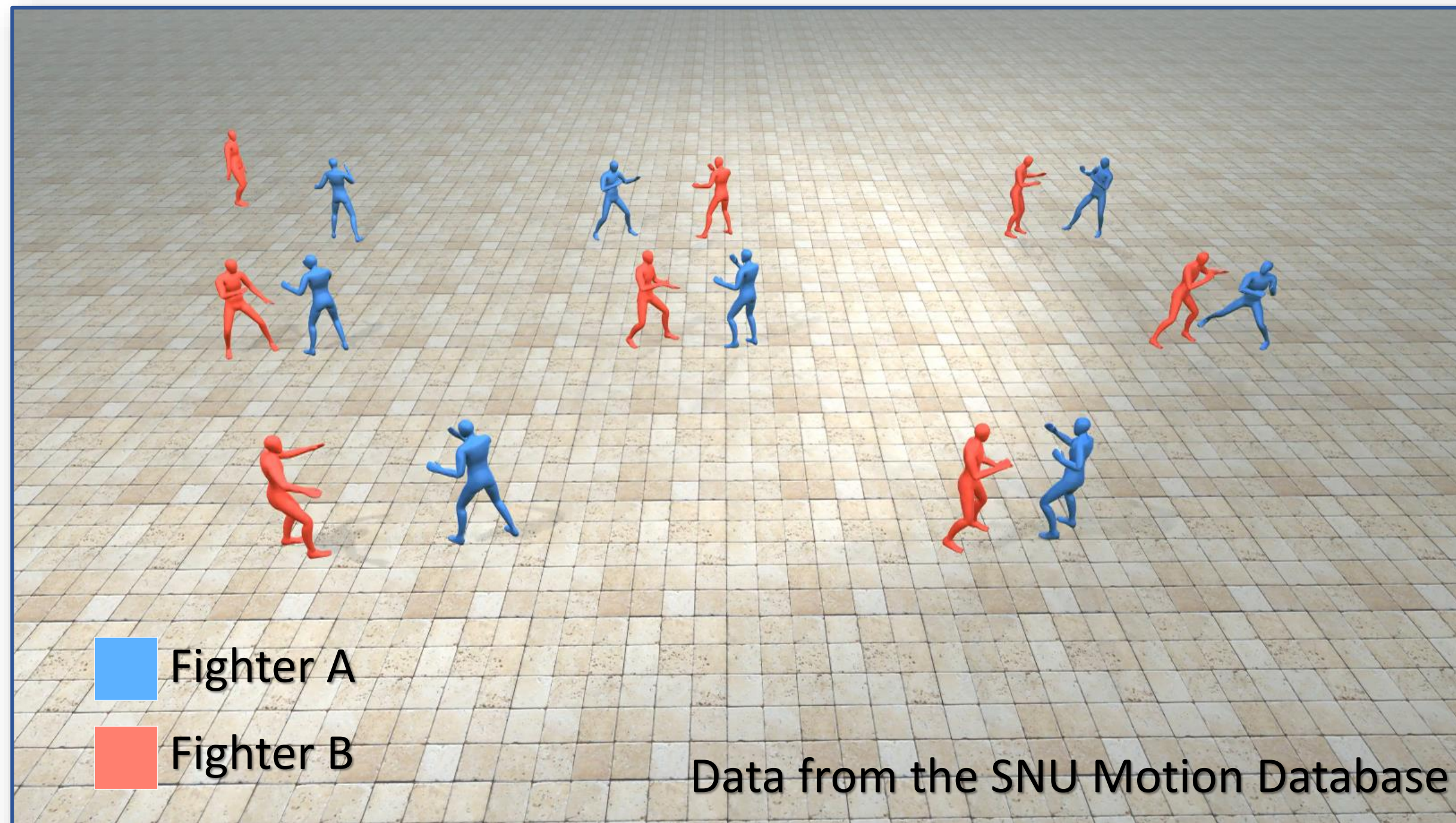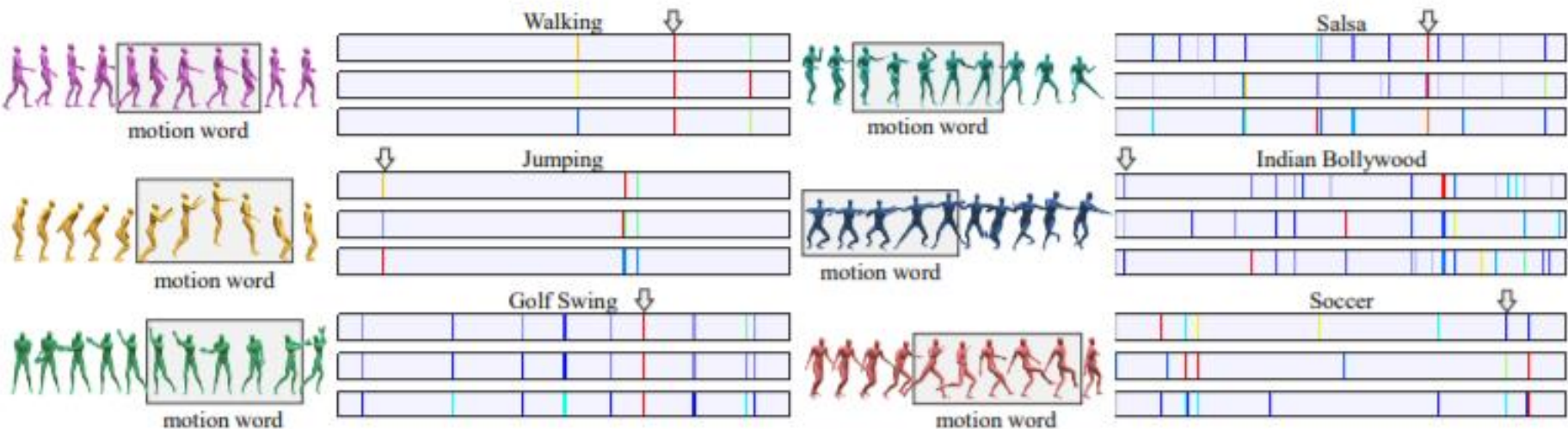acterized by similar sets of actions that appear frequently. In this paper we introduce *motion motifs* and *motion signatures* that are a succinct but descriptive representation of motion sequences. We first break the motion sequences to short-term movements called motion words, and then cluster the words in a high-dimensional feature space to find motifs. Hence, motifs are words that are both common and descriptive, and their distribution represents the motion sequence. To cluster words and find motifs, the challenge is to define an effective feature space, where the distances among motion words are semantically meaningful, and where variations in speed and duration are handled. To this end, we use a deep neural network to embed the motion words into feature space using a triplet loss function. To define a signature, we choose a finite set of motion-motifs, creating a bag-of-motifs representation for the sequence. Motion signatures are agnostic to movement order, speed or duration variations, and can distinguish fine-grained differences between motions of the same class. We illustrate examples of characterizing motion sequences by motifs, and for the use of motion signatures in a number of applications.

CCS Concepts: • **Computing methodologies** → **Motion capture**; Motion processing;

Additional Key Words and Phrases: Animation, Motion Word, Motif, Motion Signature, Convolutional Network, Triplet Loss.

**ACM Reference Format:**
Andreas Aristidou, Daniel Cohen-Or, Jessica K. Hodgins, Yiorgos Chrysanthou, and Ariel Shamir. 2018. Deep Motifs and Motion Signatures. *ACM Trans. Graph.* 37, 06, Article 187 (November 2018), 13 pages. https://doi.org/10.1145/3272127.3275038

Authors' addresses: Andreas Aristidou, The Interdisciplinary Center, Kanfei Nesharim, Herzliya, Israel, 4610101, a.aristidou@ieee.org; Daniel Cohen-Or, Tel-Aviv University, Ramat Aviv, Tel-Aviv, Israel, 6997801, cohenor@gmail.com; Jessica K. Hodgins, Carnegie Mellon University, 5000 Forbes Ave, Pittsburgh, USA, PA 15213, jkh@cs.cmu.edu; Yiorgos Chrysanthou, University of Cyprus & RISE Research Center, 75, Kallipoleos, Nicosia, Cyprus, 1678, yiorgos@cs.ucy.ac.cy; Ariel Shamir, The Interdisciplinary Center, Kanfei Nesharim, Herzliya, Israel, 4610101, arik@idc.ac.il.

## 1 INTRODUCTION

The availabilthy of human motion data in big repositories is growing with the emergence of simpler motion capture devices [Mehta et al. 2017; Pavlakos et al. 2017]. Content-based techniques and searching methods become essential to facilitate the use of such data. However, motion data is not always annotated or parameterized, hindering the semantic analysis of motions, the search in motion datasets, and the comparison between motion data. Working directly with the motion sequences is challenging due to the high-dimensional, temporal, nature of the motion, their large variations

**Belly Dance**
**Chinese Xin-Jian...**

# Introduction
# Contextual motion analysis



https://youtu.be/weSvQCGuTvU

Dance is "a performing-art form consisting of purposefully selected and controlled rhythmic sequences of human movements". These movements have aesthetic and often symbolic value.

S. H. Fraleigh, Dance and the Lived Body: A Descriptive Aesthetics. University of Pittsburgh Press, 1987.

Dance is "a performing-art form consisting of <span style="color:red">purposefully selected</span> and <span style="color:red">controlled rhythmic</span> sequences of human movements". These movements have aesthetic and often symbolic value.

S. H. Fraleigh, Dance and the Lived Body: A Descriptive Aesthetics. University of Pittsburgh Press, 1987.

The premise of our work
# Music-driven motion synthesis

* Raw Data

# Music-driven motion synthesis



**Our Architecture**

# Music-driven motion synthesis



**Our Architecture**

# Music-driven motion synthesis

# Music-driven motion synthesis



**Our Architecture**

# Music-driven motion synthesis

# Music-driven motion synthesis

# Music-driven motion synthesis

- Audio representation (Librosa Library [Ellis 2007]):
  - Rhythmic features $\mathbf{a}_r^t \in \mathbb{R}^4$
  - Spectral features $\mathbf{a}_s^t \in \mathbb{R}^{87}$

- Pose representation: $\mathbf{f}^t = \left[ \mathbf{f}_t, \mathbf{f}_q \right] \in \mathbb{R}^{3+4J}$
  - the root displacement $\mathbf{f}_t \in \mathbb{R}^3$
  - joint rotations in unit quaternions, $\mathbf{f}_q \in \mathbb{R}^{4J}$ , for $J = 31$ joints

- Motif representation:
  - $\mathbf{m}^t \in \mathbb{R}^d$, where $d = 184$ universal features
  - motion words are segmented on the beat; time-scaled to 13 frames

# Music-driven motion synthesis

- The input to the network at time $t$ is:

$$\mathbf{n}^t = [\mathbf{a}_r^t, \mathbf{m}^t, \mathbf{f}^t, \mathbf{c}^t] \in \mathbb{R}^{4+d+4J+2}$$

- where $\mathbf{c}^t \in \{0,1\}^2$ is a binary vector representing the left and the right foot contact labels

# Music-driven motion synthesis

- **Foot Sliding Cleaning (pose level)**
  - predict foot contact labels

- **Motion Diversity (motif level)**
  - AdaIN layer to inject style variation using $\mathbf{a}_s^t$

- **Motion Perceptual-Loss (motif level)**
  - controls the content of motion words

- **Motif Transition matrix (choreography level)**
  - describes probability of the temporal connectivity between consecutive motion motifs

- **Signature difference (choreography level)**
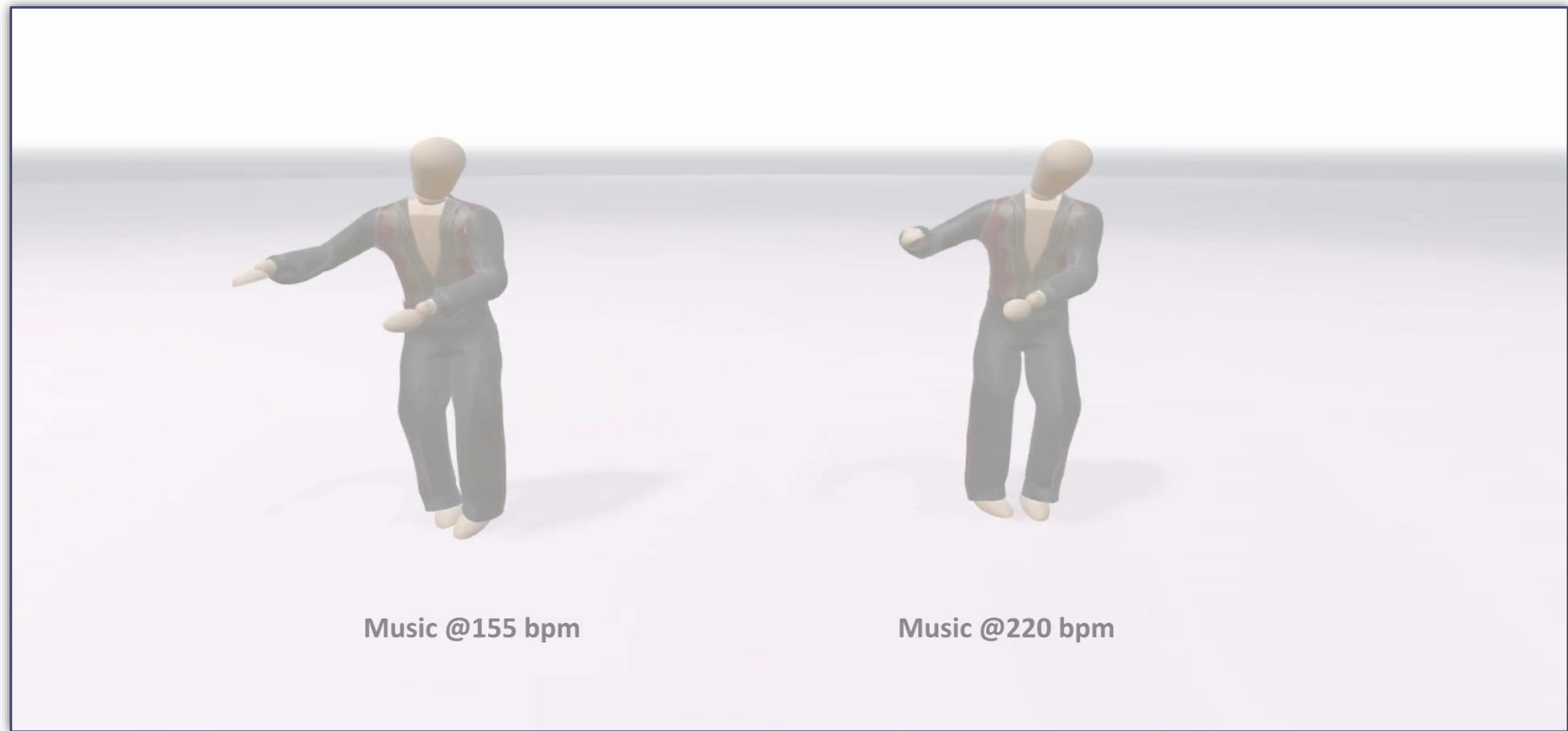  - compares the current signature to the target signature

# Music-driven motion synthesis



Target Signature

Current Signature

# Dance synthesis at different bpm



Music @155 bpm

Music @220 bpm

# Dance synthesis with variation

# Spectral audio for subtle variations



Motion generated without spectral features
Motion stylized with spectral audio (166 bpm) - Case A
Motion stylized with spectral audio (166 bpm) - Case B

# Music-driven motion synthesis



Modern Dance

# Recreate an existing dance



Target Motion

Generated Motion

## Other Applications

Co-financed by the European Union
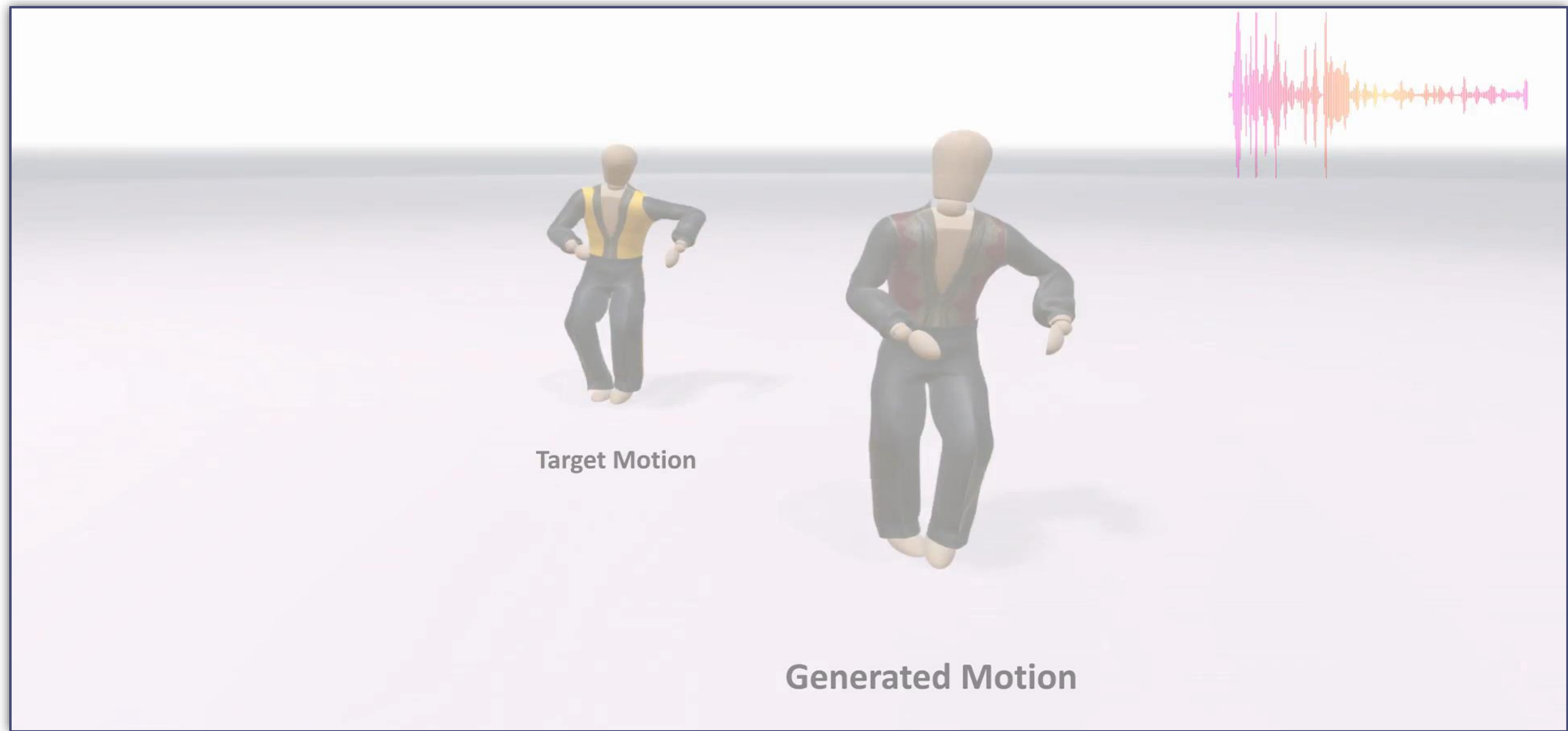Connecting Europe Facility

144

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423

# Other Applications



Salsa Style: Mambo

## Other Applications

SparsePoser: Real-time full-body motion from sparse data

Submission ID: 678

# Other Applications



Dance Central 4 - Shape of You



Just Dance® 2019 - Me Me Me

## Other Applications

Co-financed by the European Union
Connecting Europe Facility

148

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
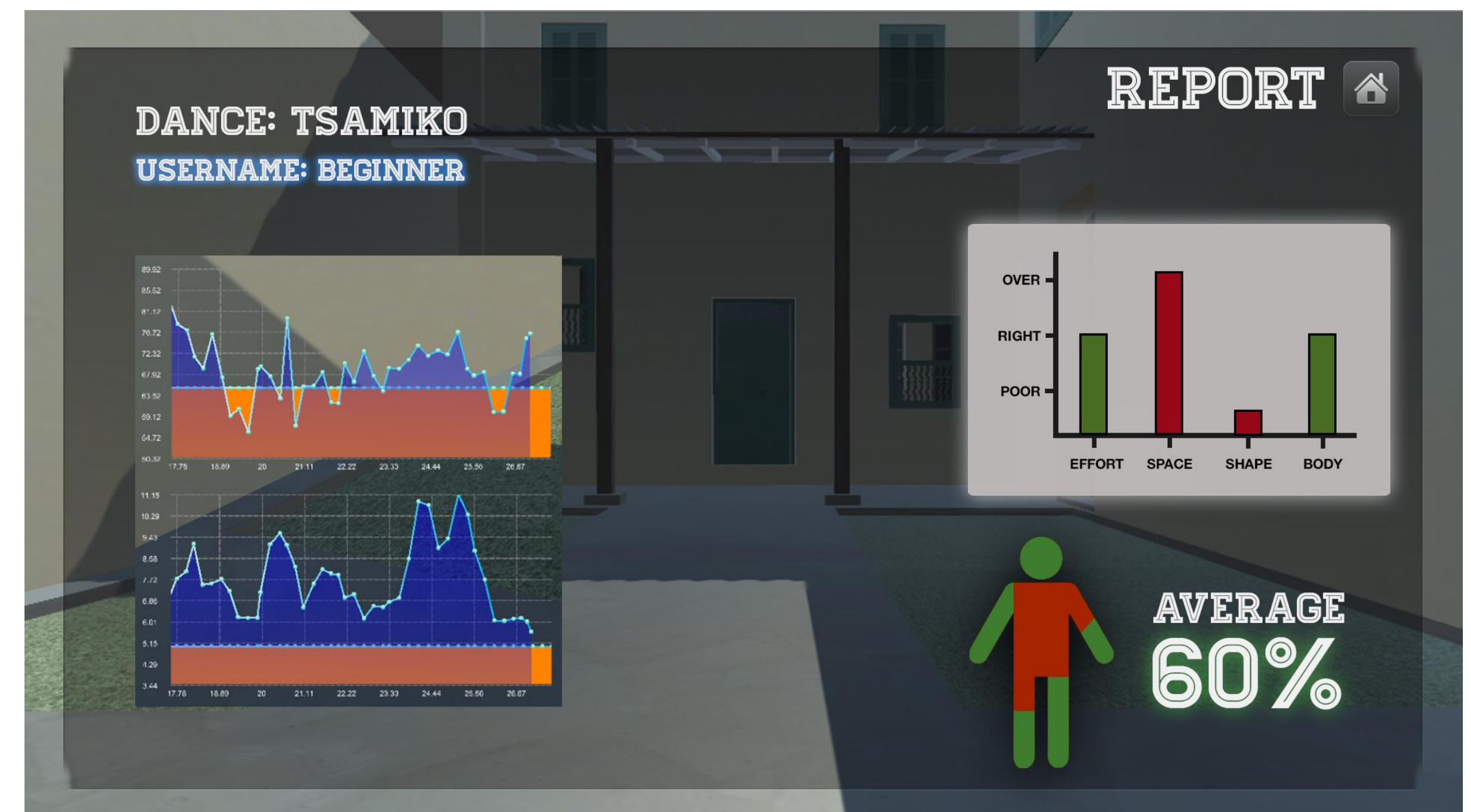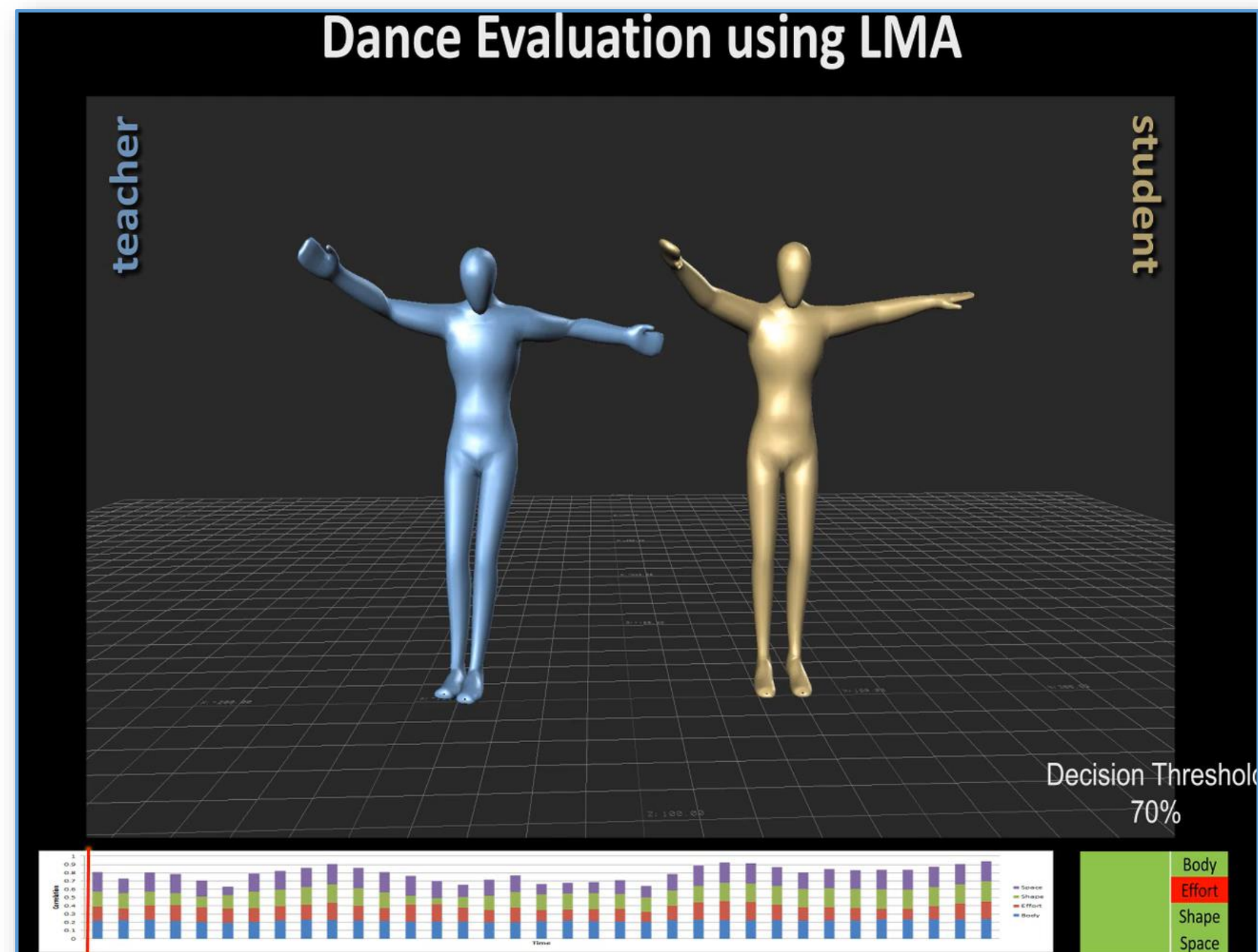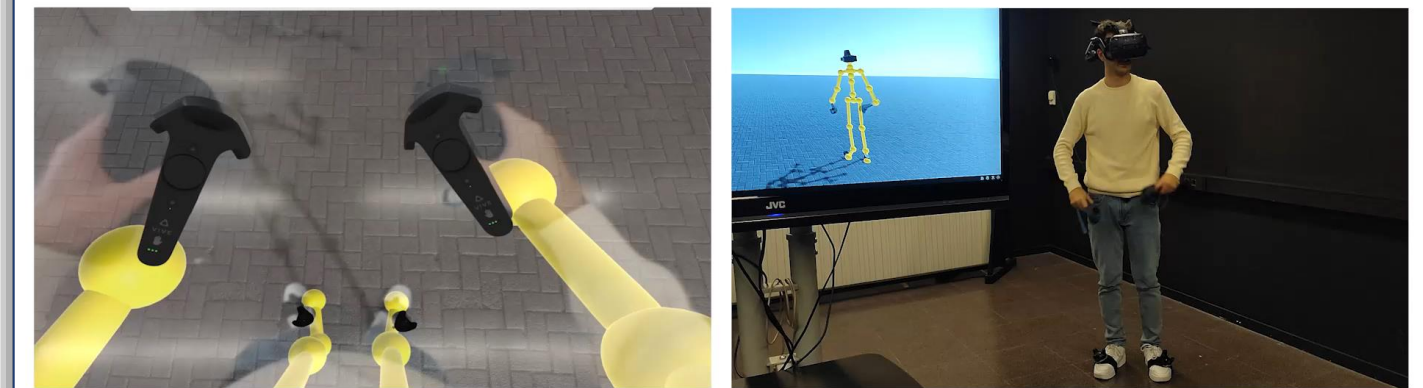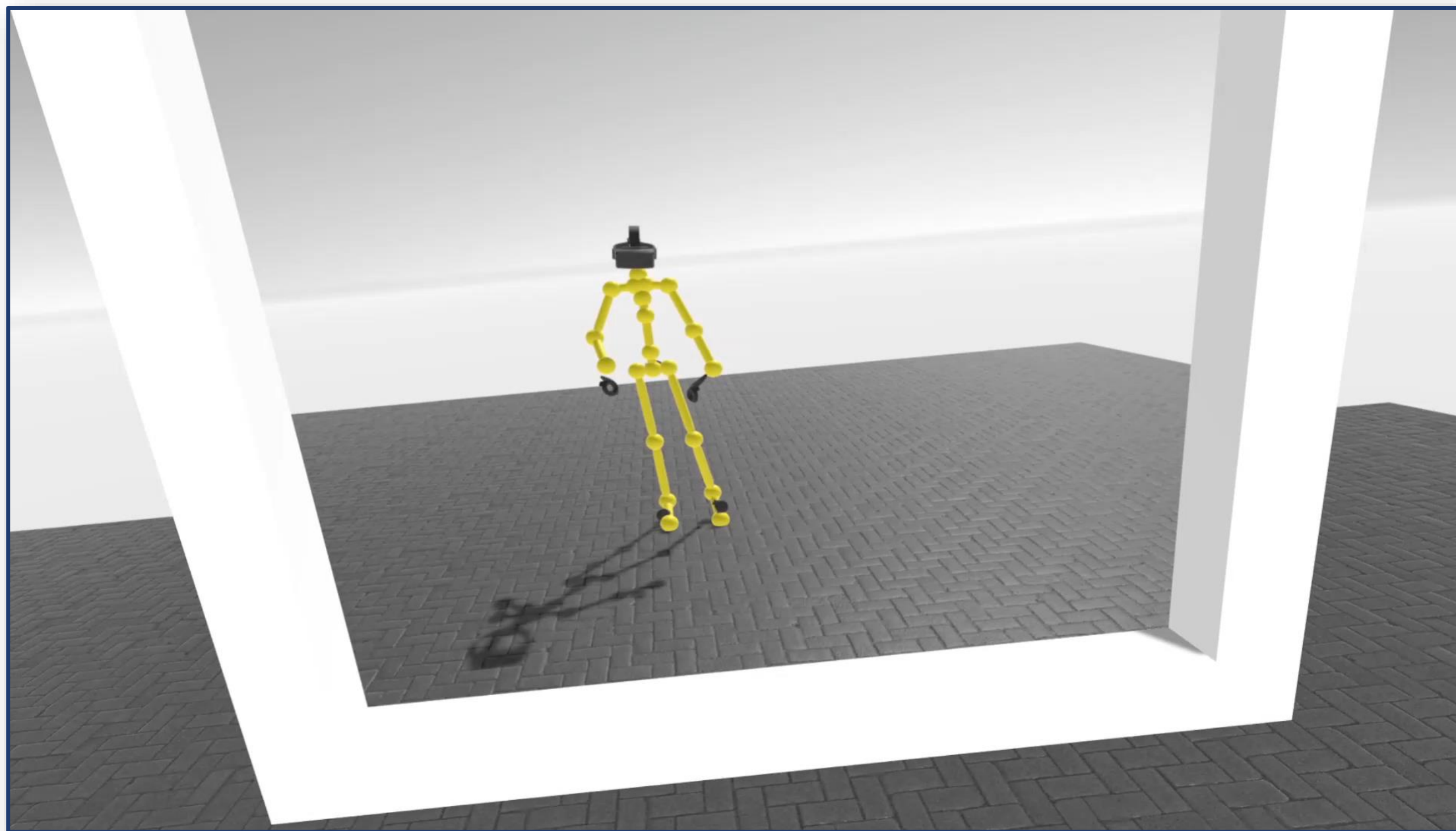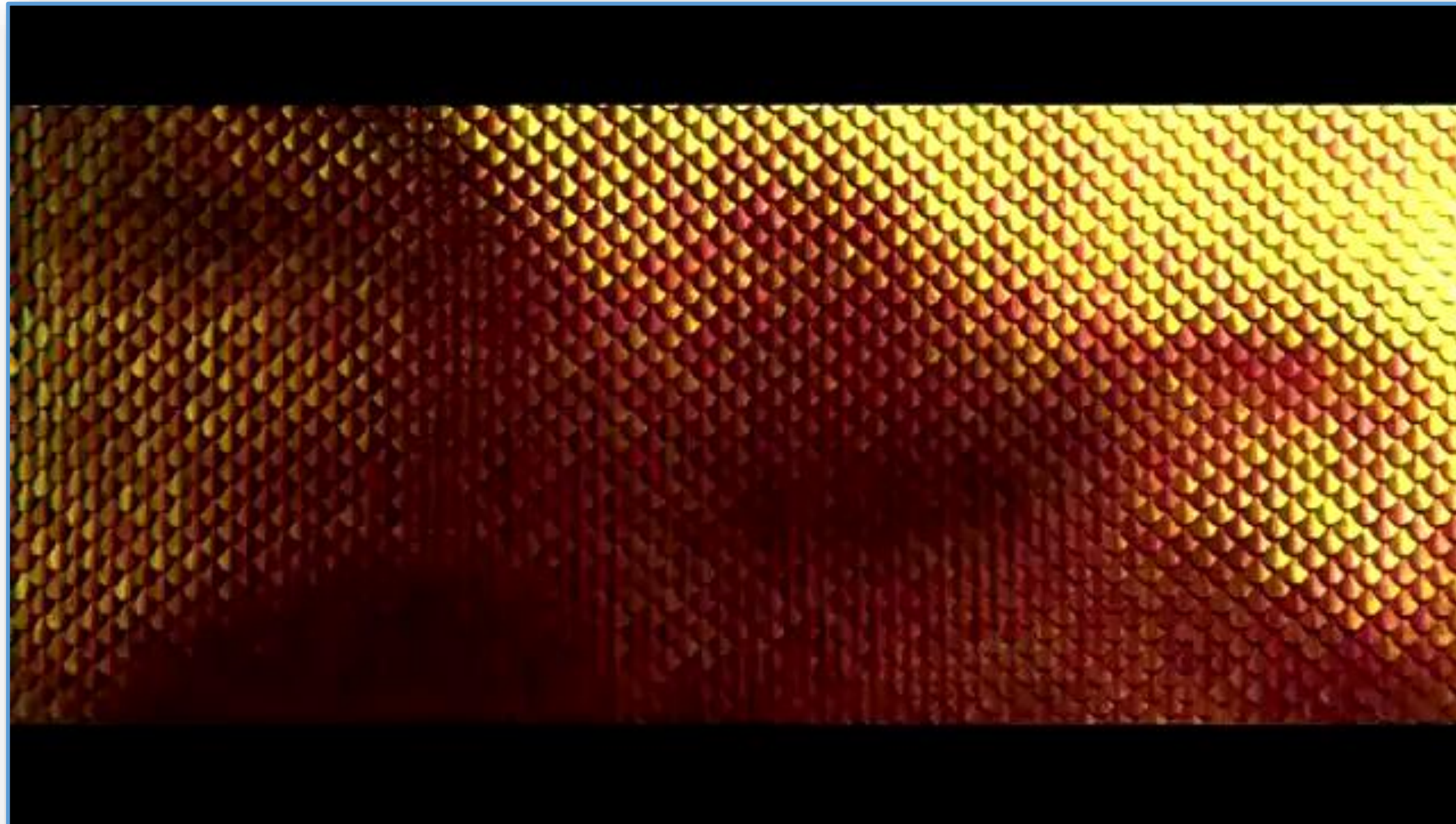under GA nr. INEA/CEF/ICT/A2020/2267423

## Other Applications



AICP sponsor reel by Method Studios
https://youtu.be/fd_9qwpzVBQ



Dancing Phantoms by Kiyan Forootan
https://youtu.be/Ig7A6fZrWyM

Co-financed by the European Union
Connecting Europe Facility

149

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423

# Other Applications



PARKER • Become the Fool
https://youtu.be/5oJUfpB4f90

## Join our team at the *Graphics & Extended Reality Lab*

# Andreas Aristidou
**Assistant Professor**

**Office**: FST01, Room B113

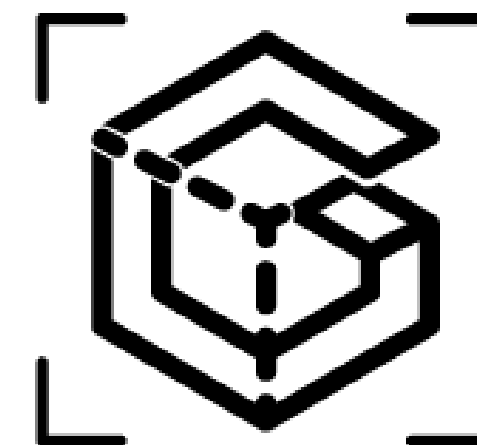**Office hours**: Only after appointment

**email**: andarist@ucy.ac.cy

**Research Interests:**

Machine Learning, Deep Learning and its applications in Computer Graphics and Character Animation, Virtual/Augmented Reality, Digital Heritage

https://www.cs.ucy.ac.cy/~andarist
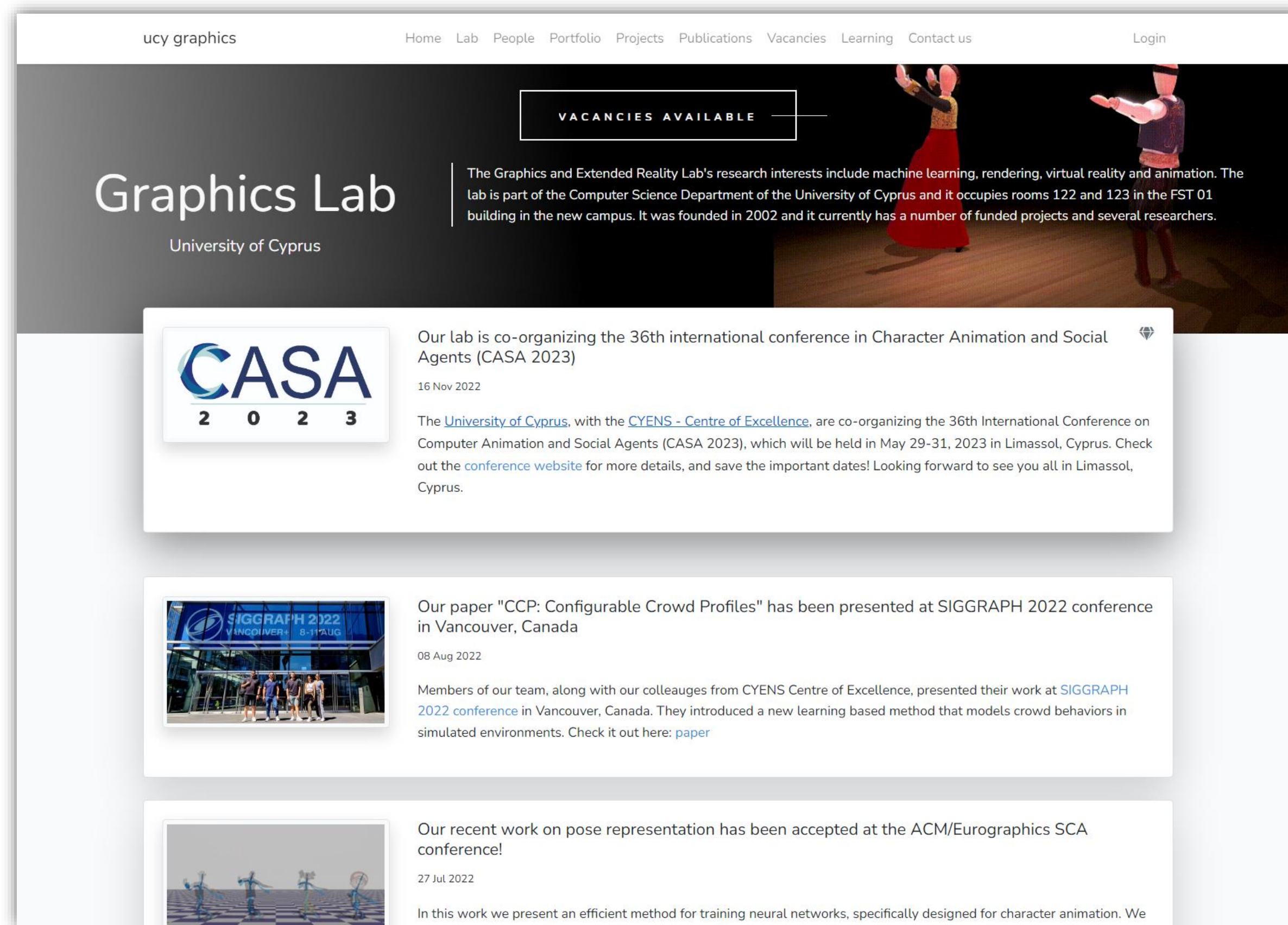
**GRAPHICS & EXTENDED REALITY**
L A B

Co-financed by the European Union
Connecting Europe Facility

155

This Master is run under the context of Action No 2020-EU-IA-0087, co-financed by the EU CEF Telecom under GA nr. INEA/CEF/ICT/A2020/2267423

# Join our team at the *Graphics & Extended Reality Lab*



The **Graphics and Extended Reality Lab** at the University of Cyprus, part of the Computer Science Department, conducts research in areas such as machine learning, rendering, virtual reality and animation.

Founded in 2002, the lab is located in rooms 122 and 123 of the FST 01 building on the university's new campus, and is staffed by two faculty members and twelve research associates. It also has several active, funded projects.

Website: https://graphics.cs.ucy.ac.cy/

# Thank you!

That's all folks!!!