



A New Era of Performance Capture with Machine Learning

Daniel Holden
Machine Learning Researcher, Ubisoft La Forge

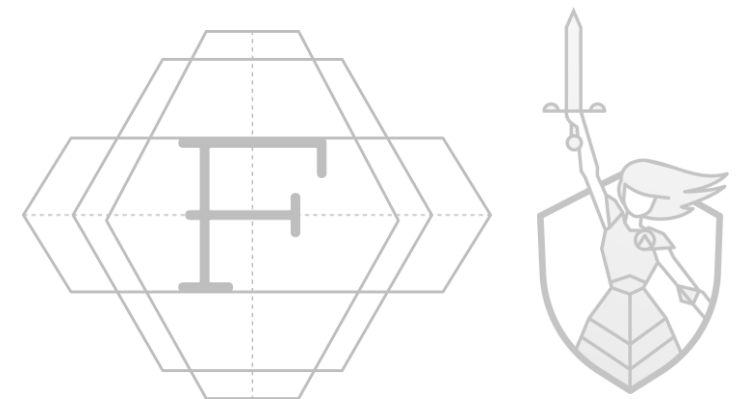
GAME DEVELOPERS CONFERENCE

MARCH 18–22, 2019 | #GDC19

History

Mocap Cleaning
Facial Tracking
Audio to Facial

The Future

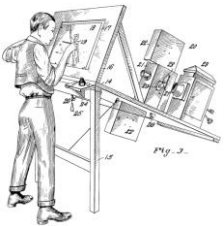


1939



1915

1978



1970

1980

1990

2000

2010

2020

Now

- http://graphics.stanford.edu/courses/cs448-09-spring/motion_capture.pdf
- <https://ca.ign.com/articles/2014/07/11/a-brief-history-of-motion-capture-in-the-movies>
- https://www.siggraph.org/education/materials/HyperGraph/animation/character_animation/motion_capture/history1.htm

GDC

GAME DEVELOPERS CONFERENCE

MARCH 18-22, 2019 | #GDC19

Era of the Rotoscope (1920s – 1980s)

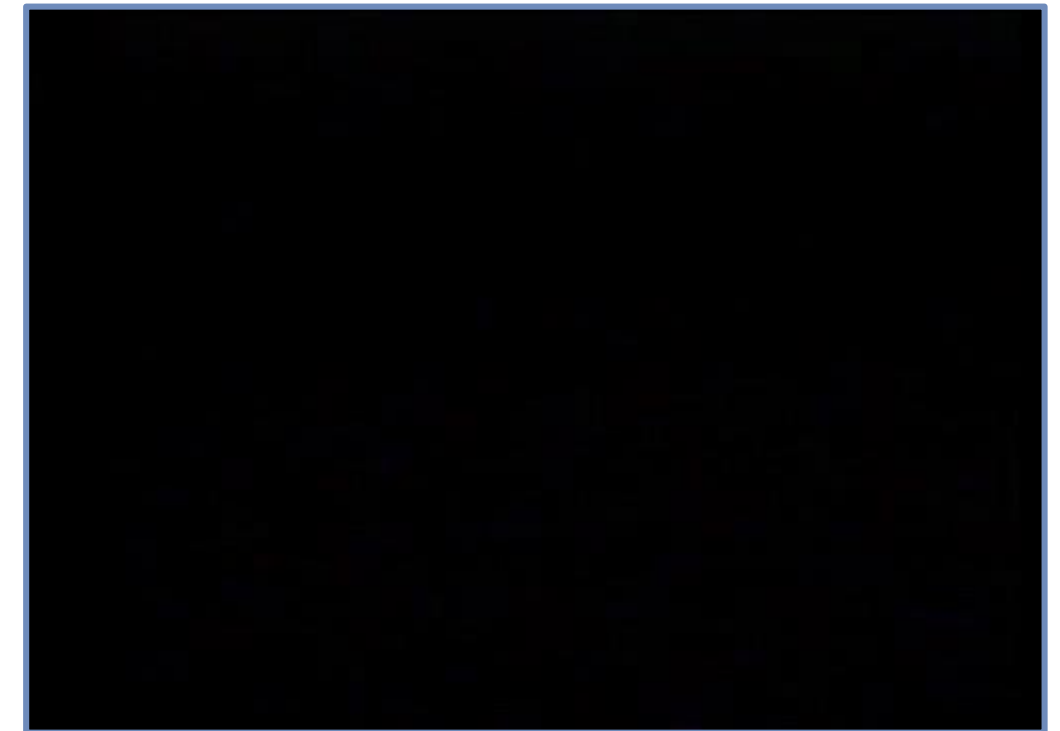
1992



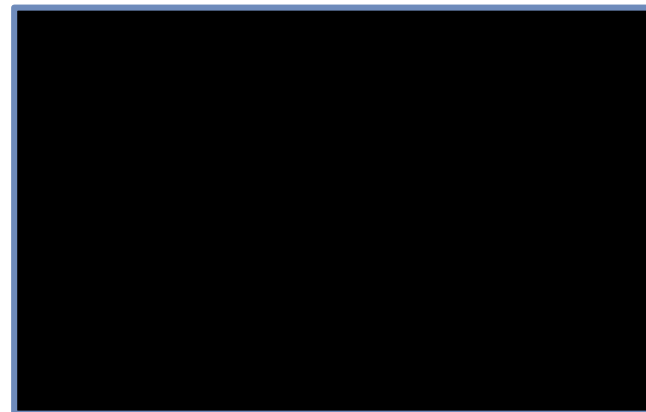
1939



1978



1989



- <https://imgur.com/gallery/IZkSR>
- <https://www.youtube.com/watch?v=kJAVgY8DMk>

1939

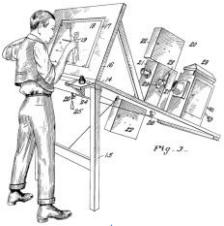


1988



1915

1978



1989



1970

1980

1990

2000

2010

2020

Now

- http://graphics.stanford.edu/courses/cs448-09-spring/motion_capture.pdf
- <https://ca.ign.com/articles/2014/07/11/a-brief-history-of-motion-capture-in-the-movies>
- https://www.siggraph.org/education/materials/HyperGraph/animation/character_animation/motion_capture/history1.htm

GDC

GAME DEVELOPERS CONFERENCE

MARCH 18-22, 2019 | #GDC19

Era of Experimentation (1980s)

1983



1988



1987



1939



1988



1999



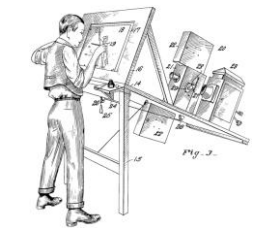
1989



1978



1915



1970

1980

1990

2000

2010

2020

Now

- http://graphics.stanford.edu/courses/cs448-09-spring/motion_capture.pdf
- <https://ca.ign.com/articles/2014/07/11/a-brief-history-of-motion-capture-in-the-movies>
- https://www.siggraph.org/education/materials/HyperGraph/animation/character_animation/motion_capture/history1.htm

GDC

GAME DEVELOPERS CONFERENCE

MARCH 18-22, 2019 | #GDC19

Era of Optical Capture (1990s)

1999



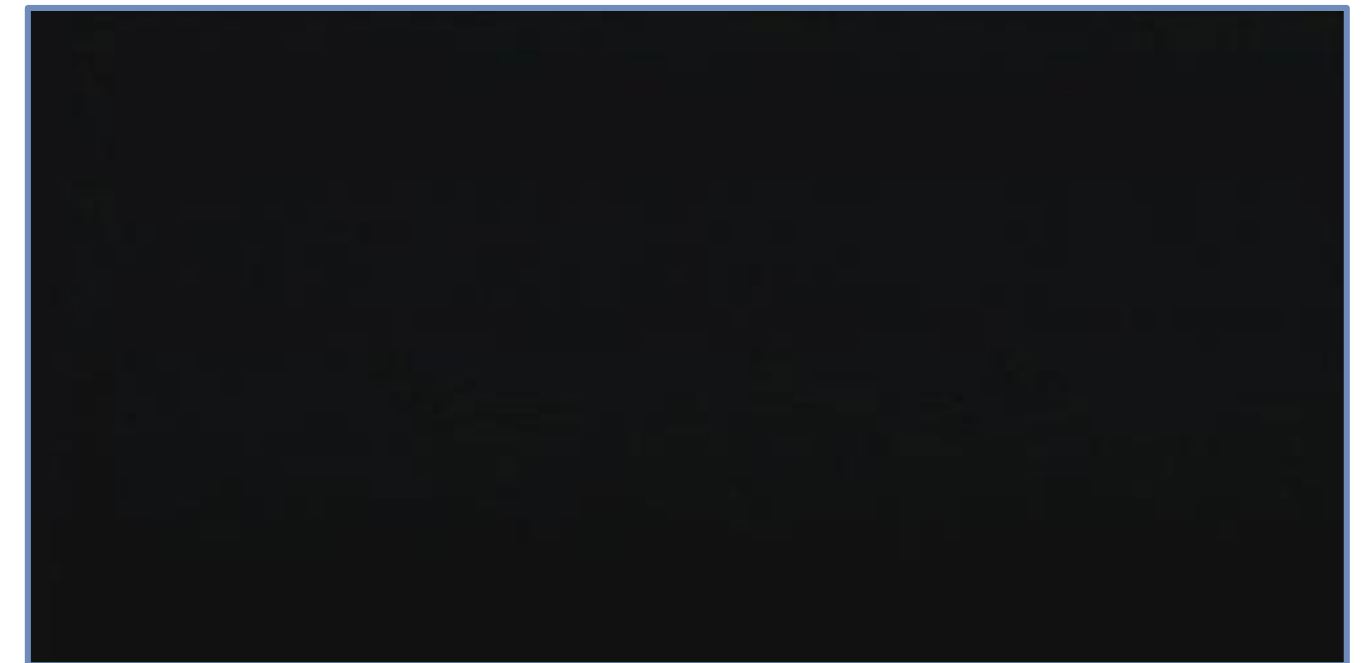
1993

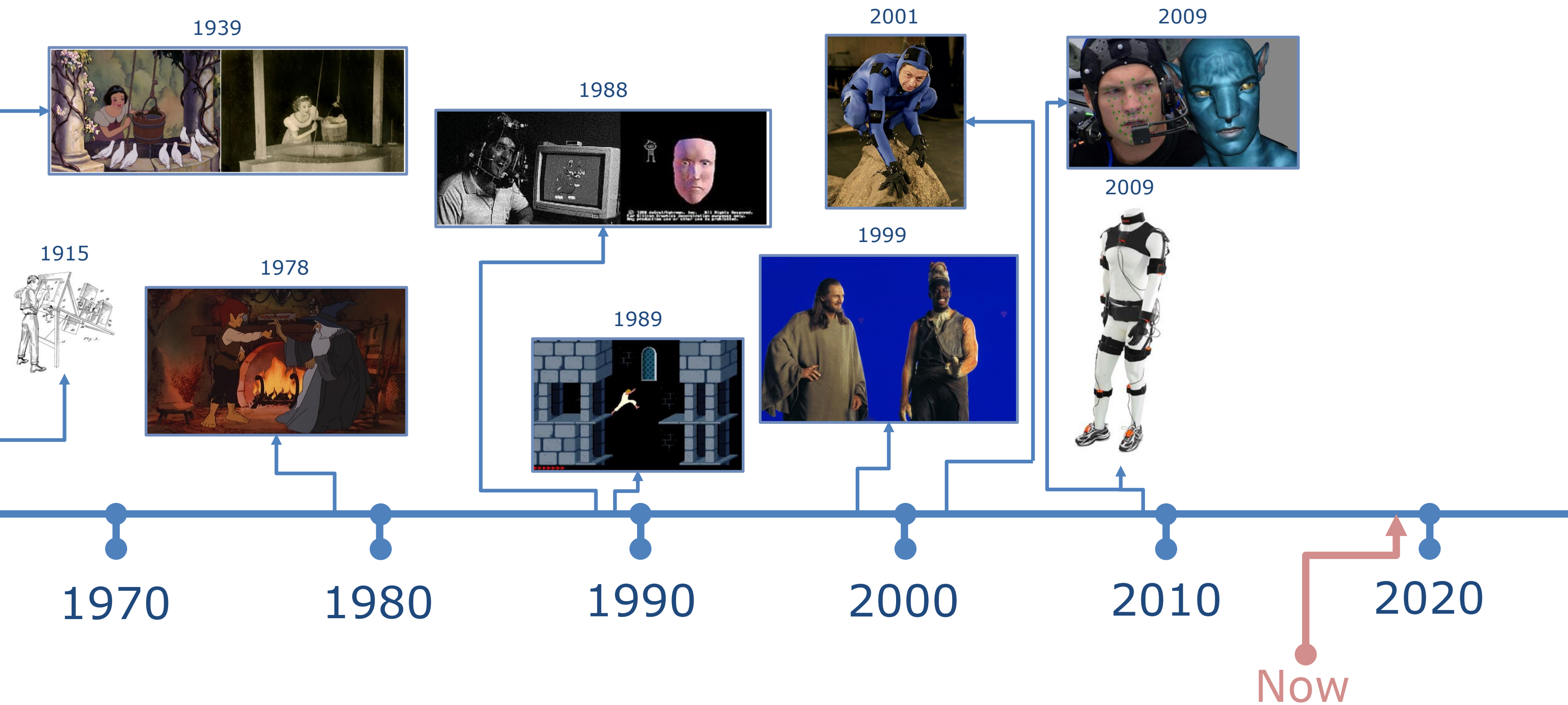


1999



2000





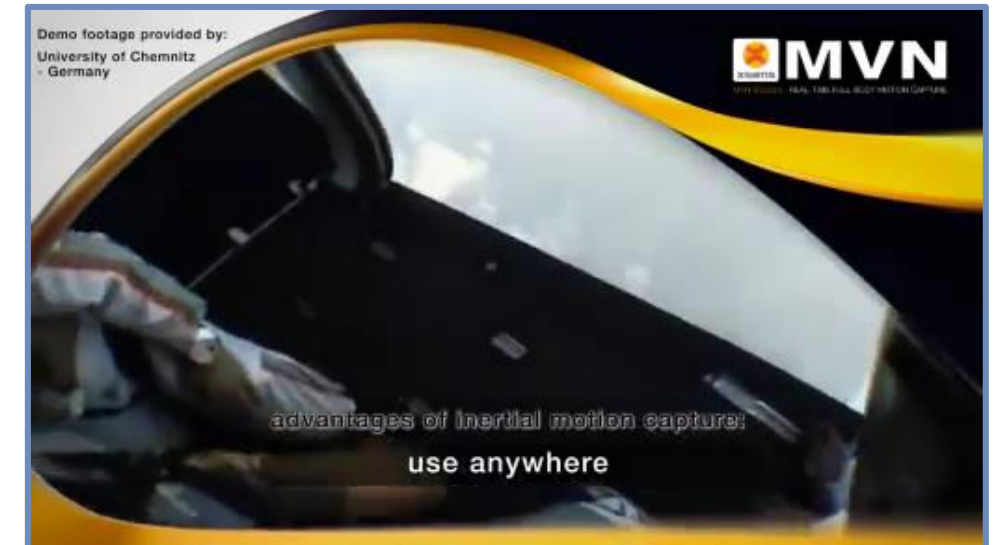
- http://graphics.stanford.edu/courses/cs448-09-spring/motion_capture.pdf
- <https://ca.ign.com/articles/2014/07/11/a-brief-history-of-motion-capture-in-the-movies>
- https://www.siggraph.org/education/materials/HyperGraph/animation/character_animation/motion_capture/history1.htm

Era of Competition (2000s)

2004

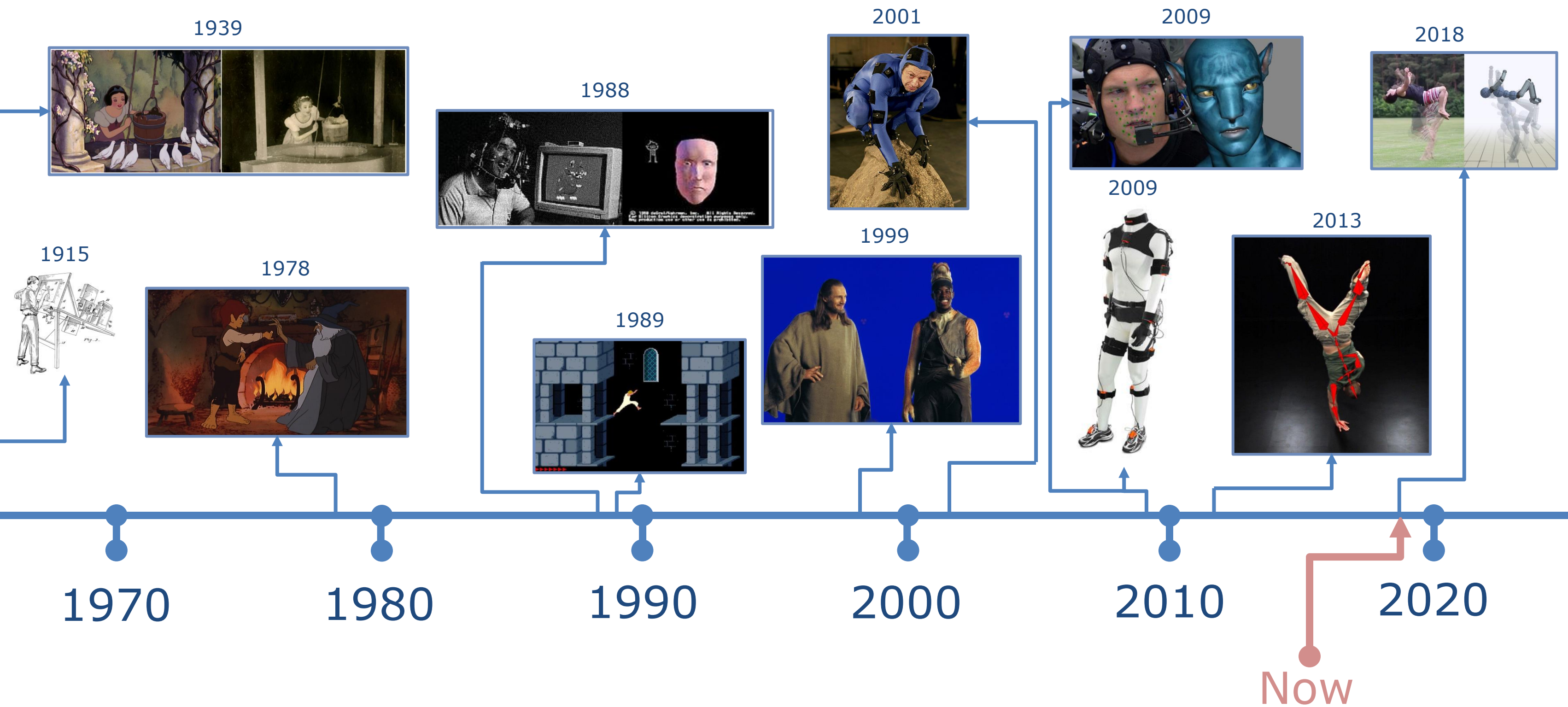


2009



2009

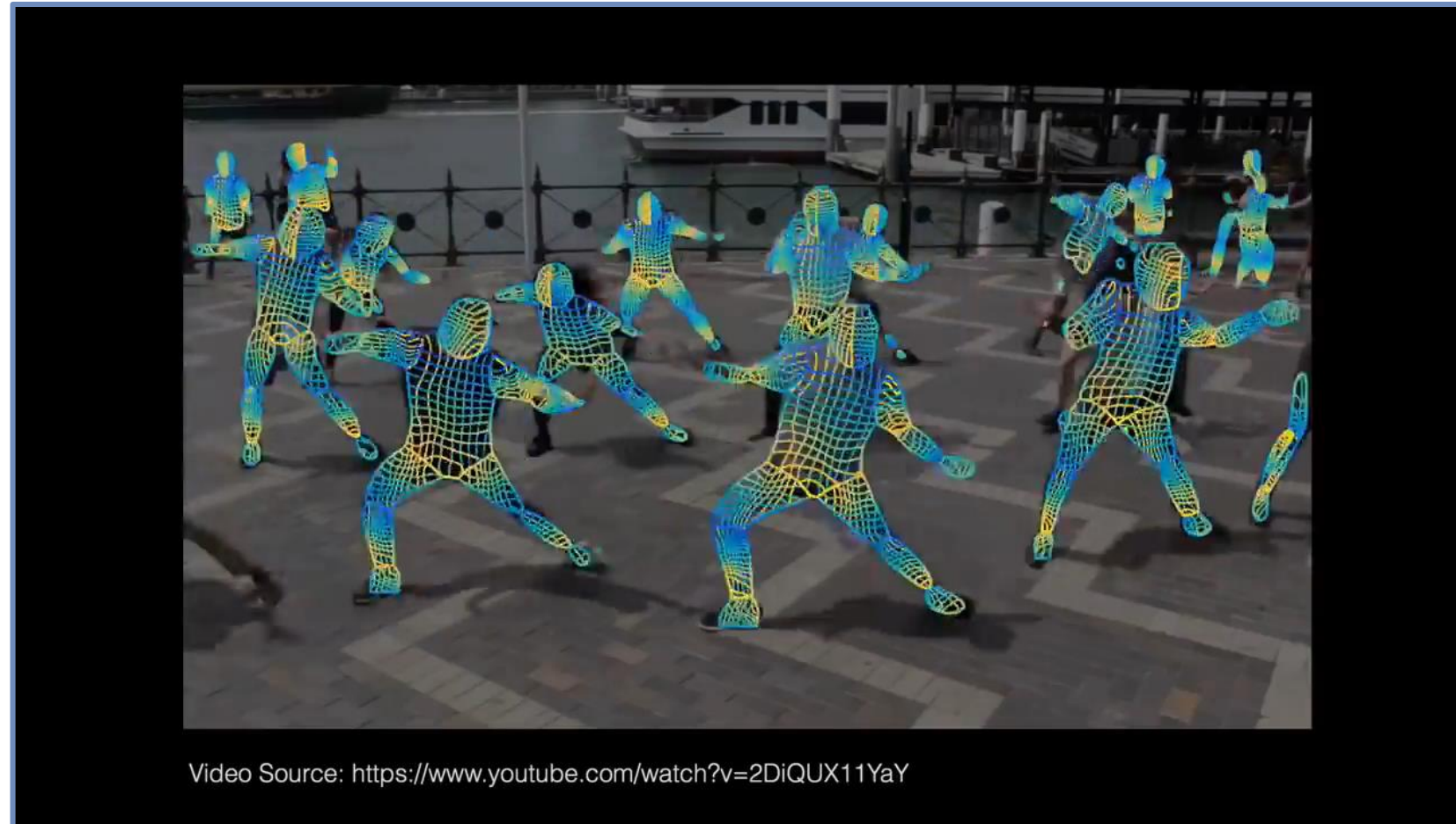




- http://graphics.stanford.edu/courses/cs448-09-spring/motion_capture.pdf
- <https://ca.ign.com/articles/2014/07/11/a-brief-history-of-motion-capture-in-the-movies>
- https://www.siggraph.org/education/materials/HyperGraph/animation/character_animation/motion_capture/history1.htm

Era of Machine Learning (2010s)

2018



2018

Results

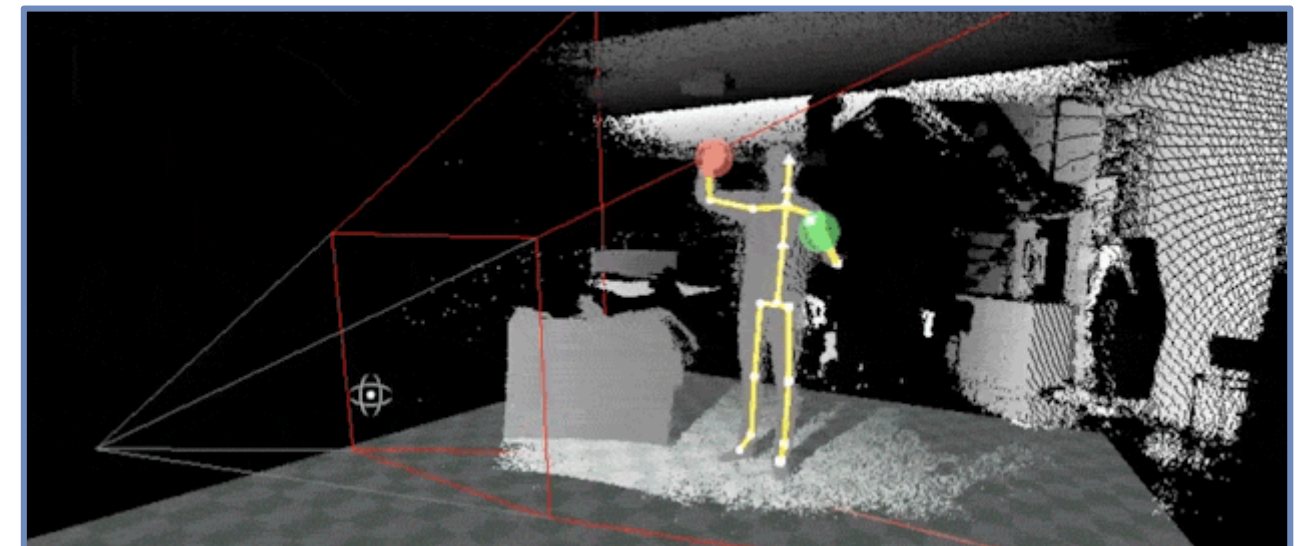


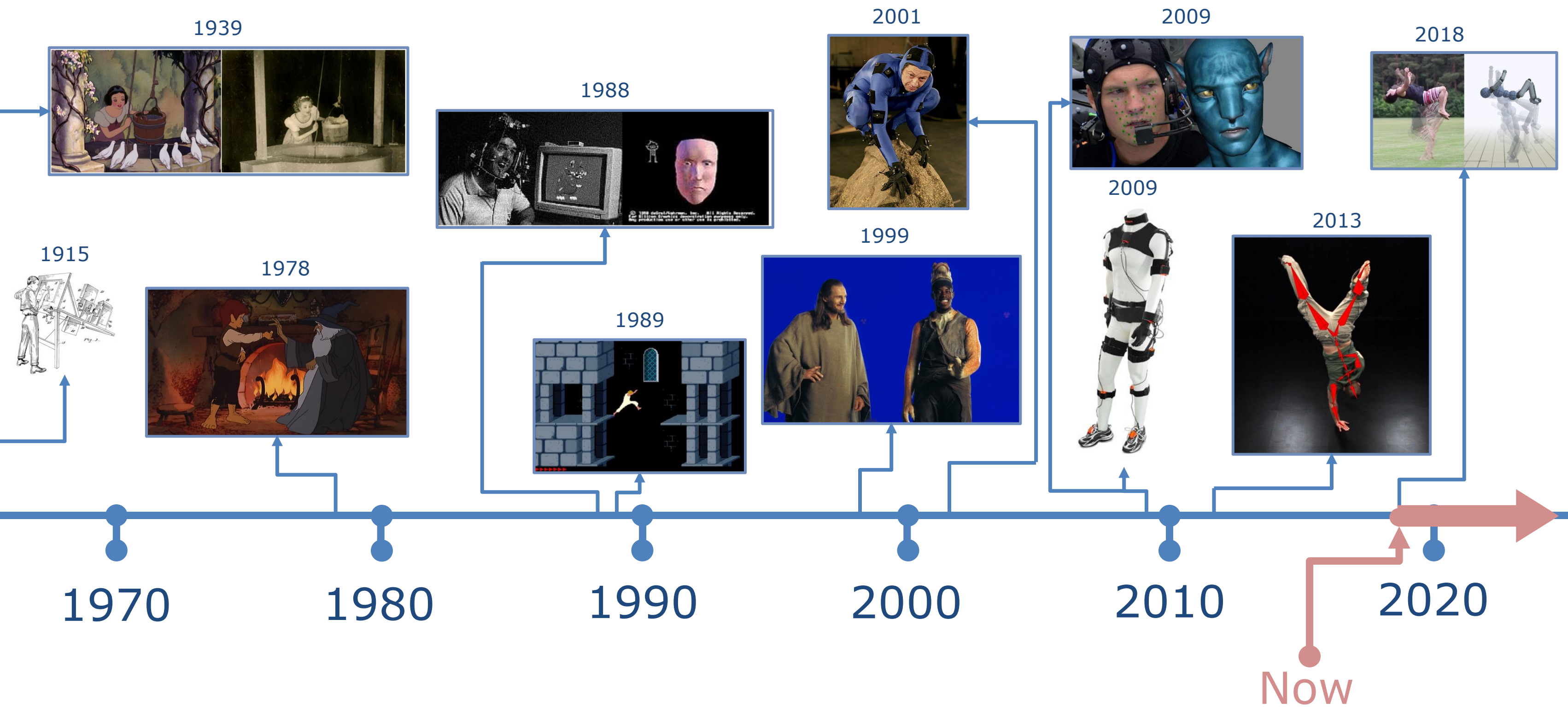
Video: Jumping Jack



Policy

2010



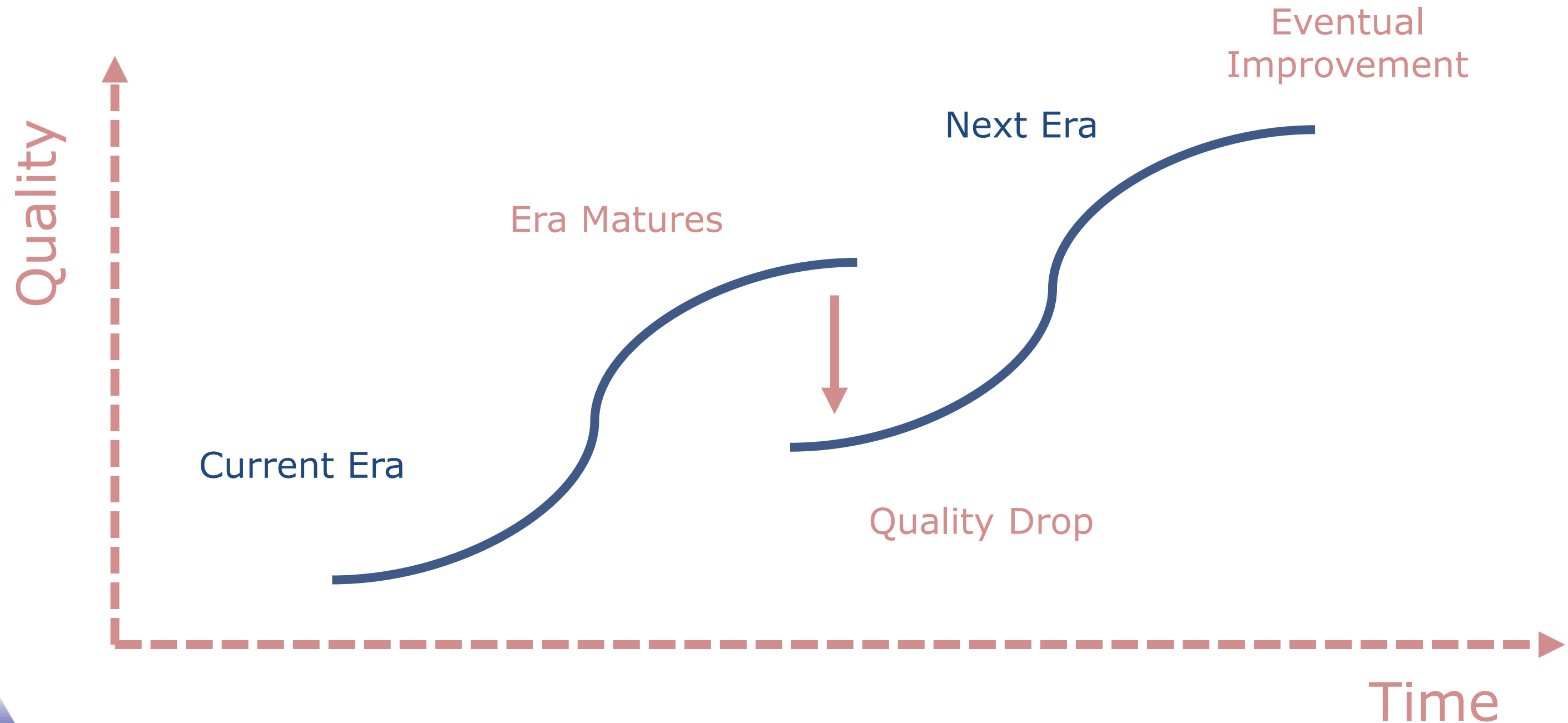


- http://graphics.stanford.edu/courses/cs448-09-spring/motion_capture.pdf
- <https://ca.ign.com/articles/2014/07/11/a-brief-history-of-motion-capture-in-the-movies>
- https://www.siggraph.org/education/materials/HyperGraph/animation/character_animation/motion_capture/history1.htm

A New Era

- Each era the previous era's technology matures.
- Machine Learning is starting to be used in production.

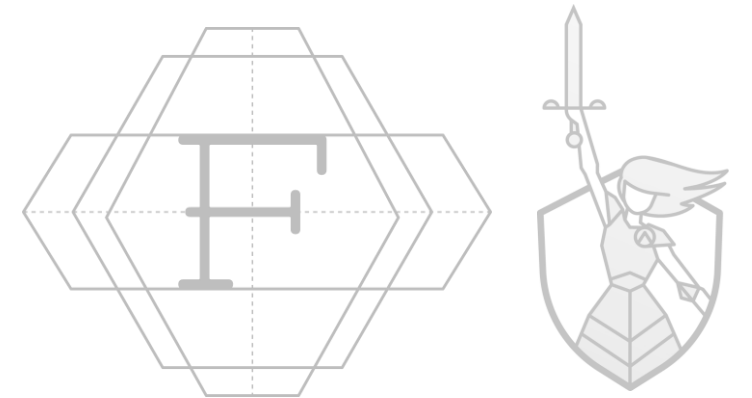
Progress



History

Mocap Cleaning
Facial Tracking
Audio to Facial

The Future

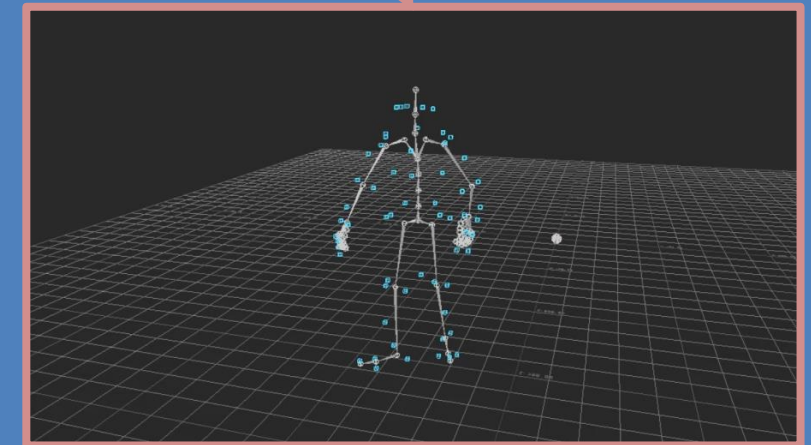
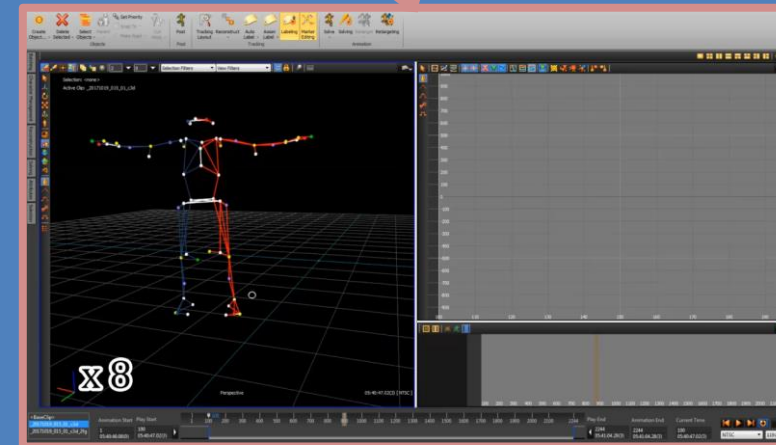
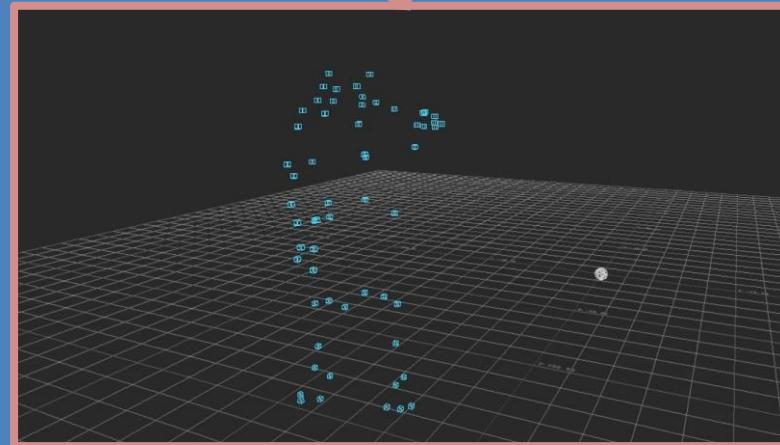


Motion Capture Pipeline

Tracking

Cleaning

Solving

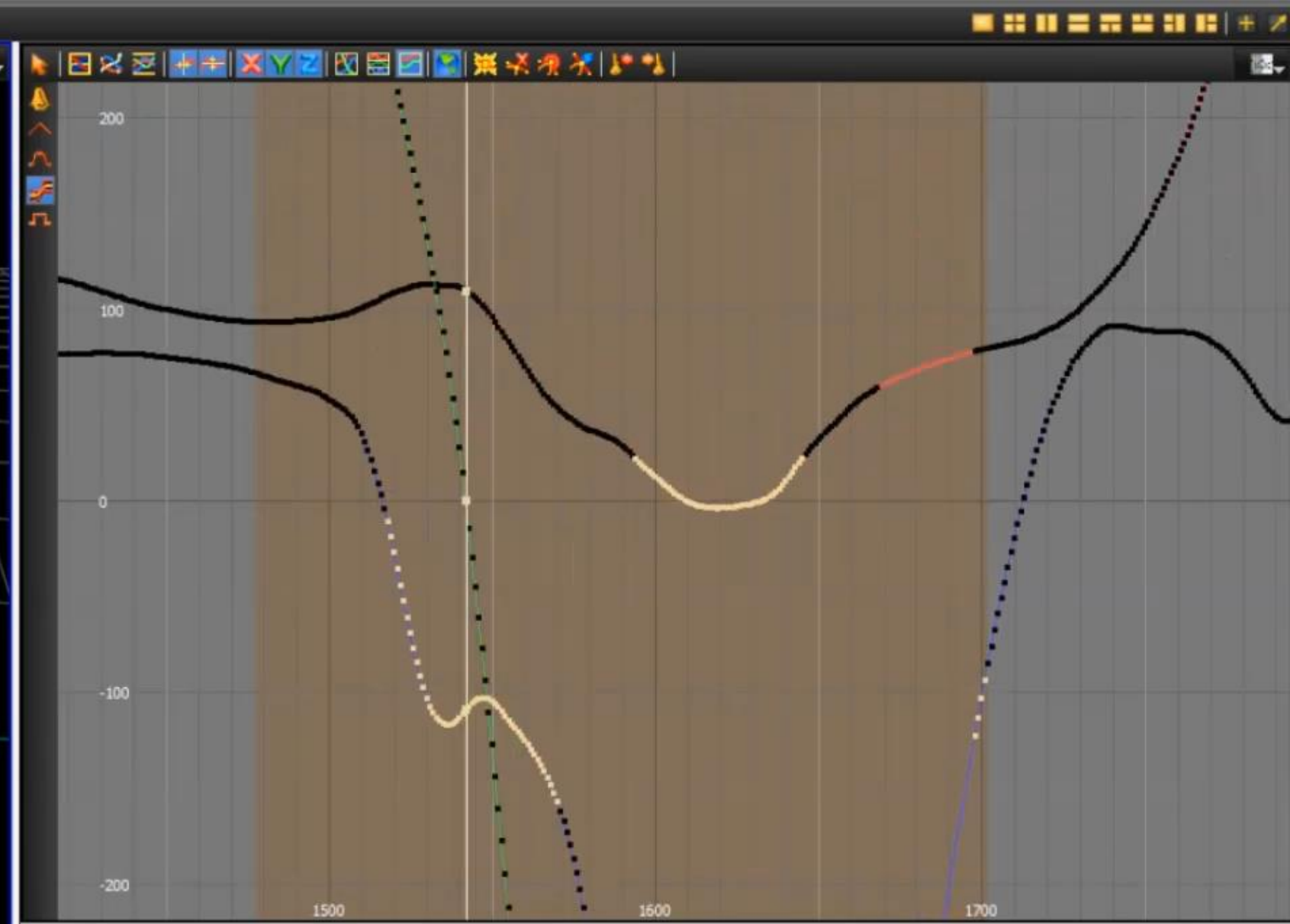




Occluded Markers

Top toolbar with icons for: Create Object..., Delete Selected, Select Objects, Parent Objects, Set Priority, Snap To, Make Rigid, Cut Keys, Post, Tracking Layout, Reconstruct, Auto Label, Axiom Label, Labeling, Marker Editing, Solve, Solving, Retarget, Retargeting, Animation.

3D Viewport showing a skeletal rig in perspective. The rig consists of white spheres (joints) connected by colored lines (bones). The background is a dark grid. Text in the top left of the viewport: Selection: RFWT, Active Clip: _20171019_015_01_c3d. A large white 'X8' logo is in the bottom left. The bottom right of the viewport shows 'Perspective' and '05:40:59:03(1) [NTSC]'.



Timeline window for 'JonathanGrondin_ROM_2'. It shows a horizontal timeline with markers and a playhead. The timeline is divided into segments, with the current time marked at 1542.

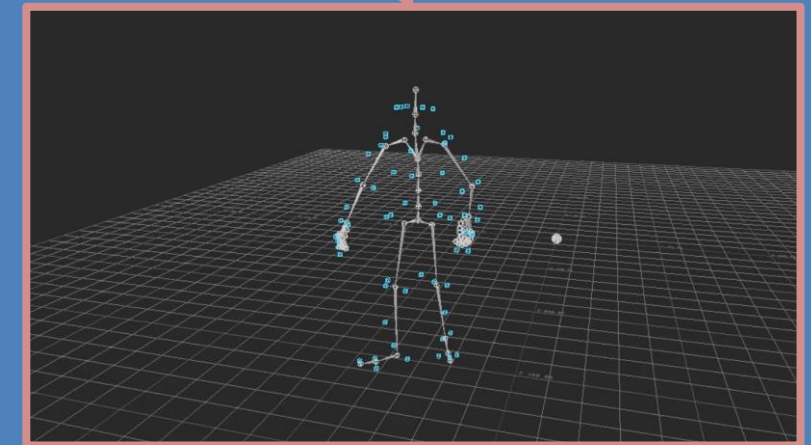
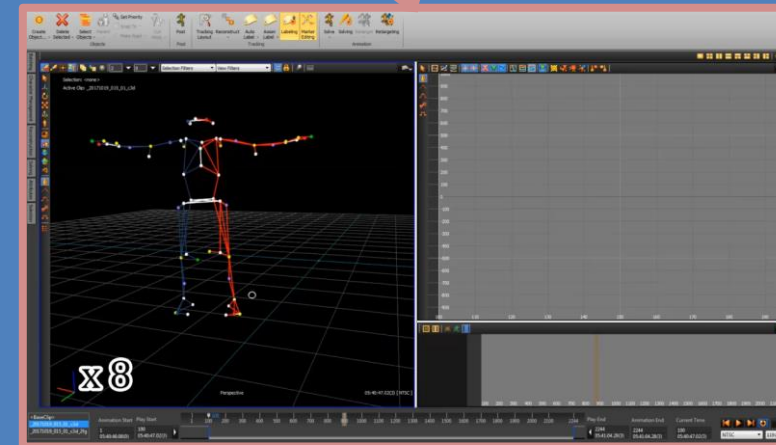
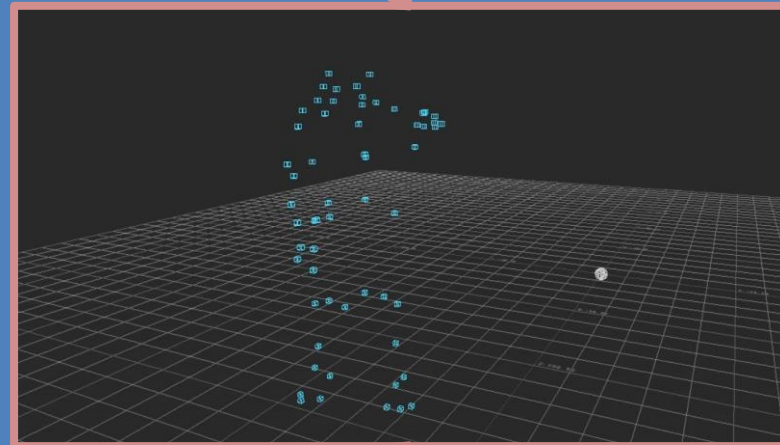
Bottom status bar and controls. On the left, a list of clips: <BaseClip>, 20171019_015_01_c3d, 20171019_015_01_c3d_2fg. In the center, a timeline with a playhead at 1542. On the right, playback controls (stop, play, next, previous, etc.) and a dropdown menu set to 'NTSC'.

Motion Capture Pipeline

Tracking

Cleaning

Solving

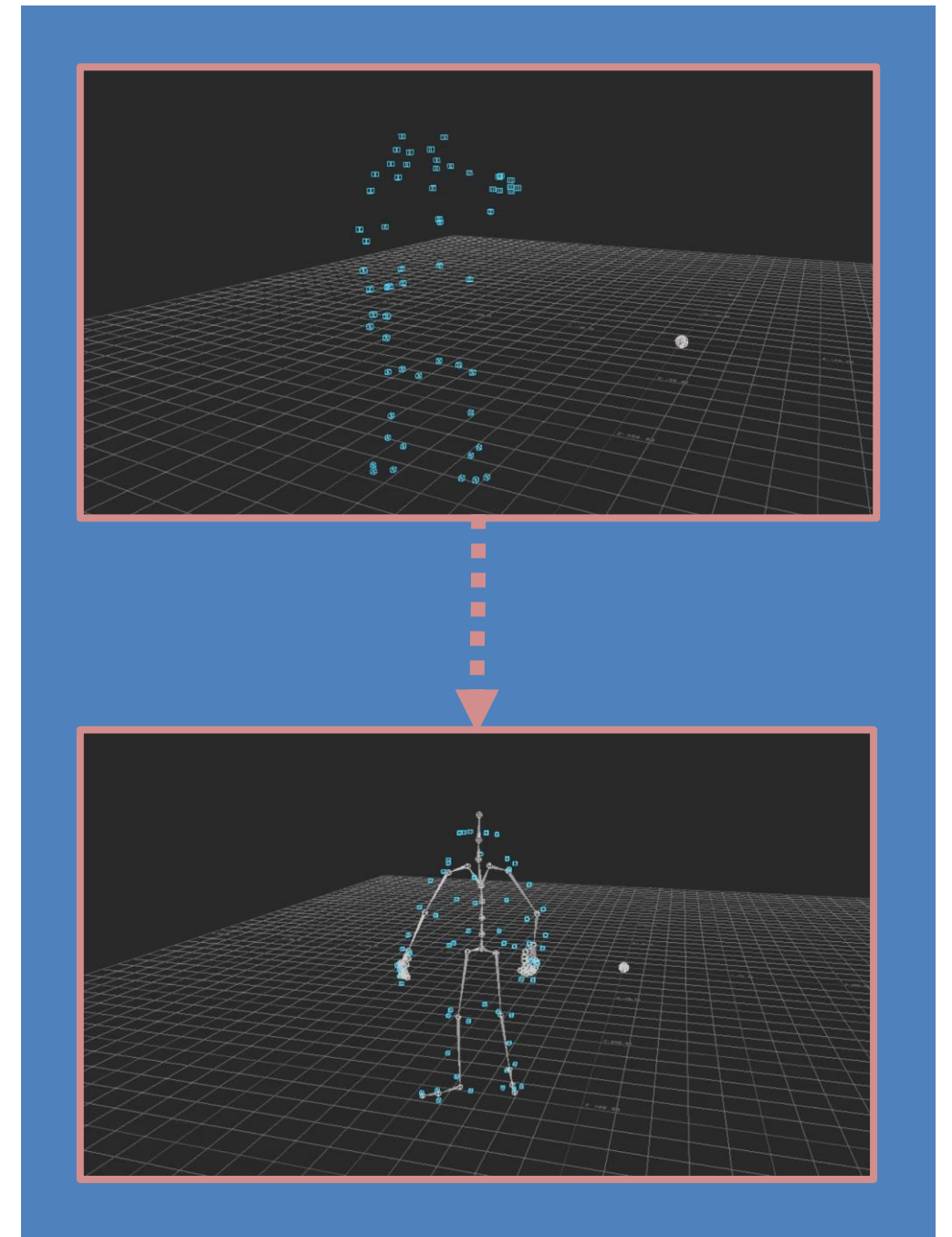


*What if we could go directly from
unclean marker data to joints?*

Robust Solving

Robust Solving

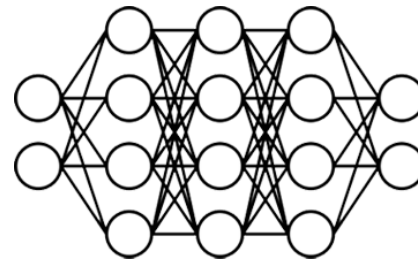
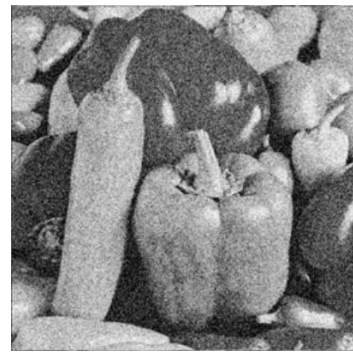
1. We want to train a **Neural Network** to learn a mapping from **Markers** to **Joints**.
2. We want to ensure the **Neural Network** is **Robust** to errors / noise in the markers.



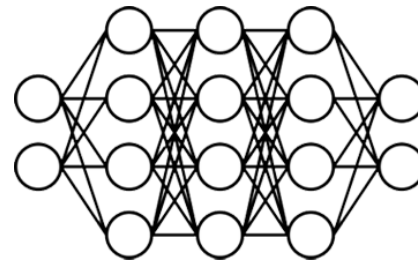
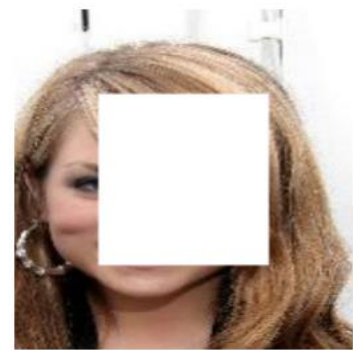
Robust Solving

How can we ensure the **Neural Network** is **Robust** to errors / noise in the markers?

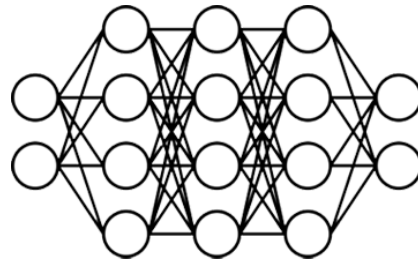
Denoising



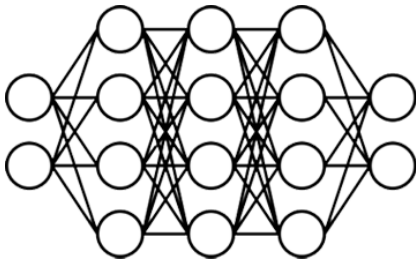
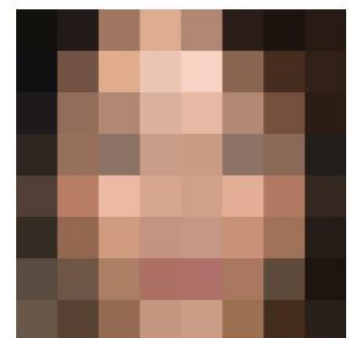
[Xie et al. 2012]



[Yu et al. 2018]

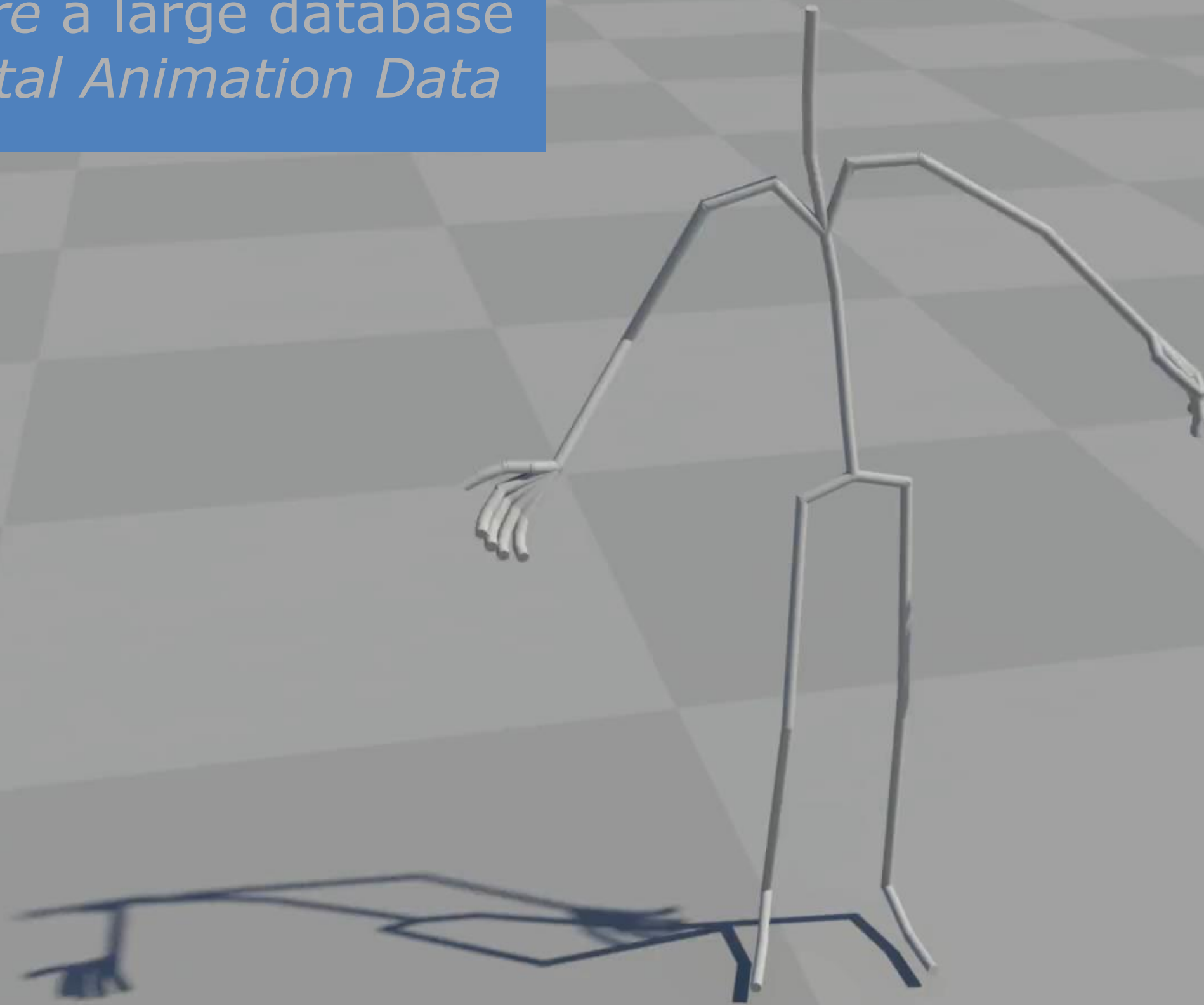


[Zhang et al. 2016]

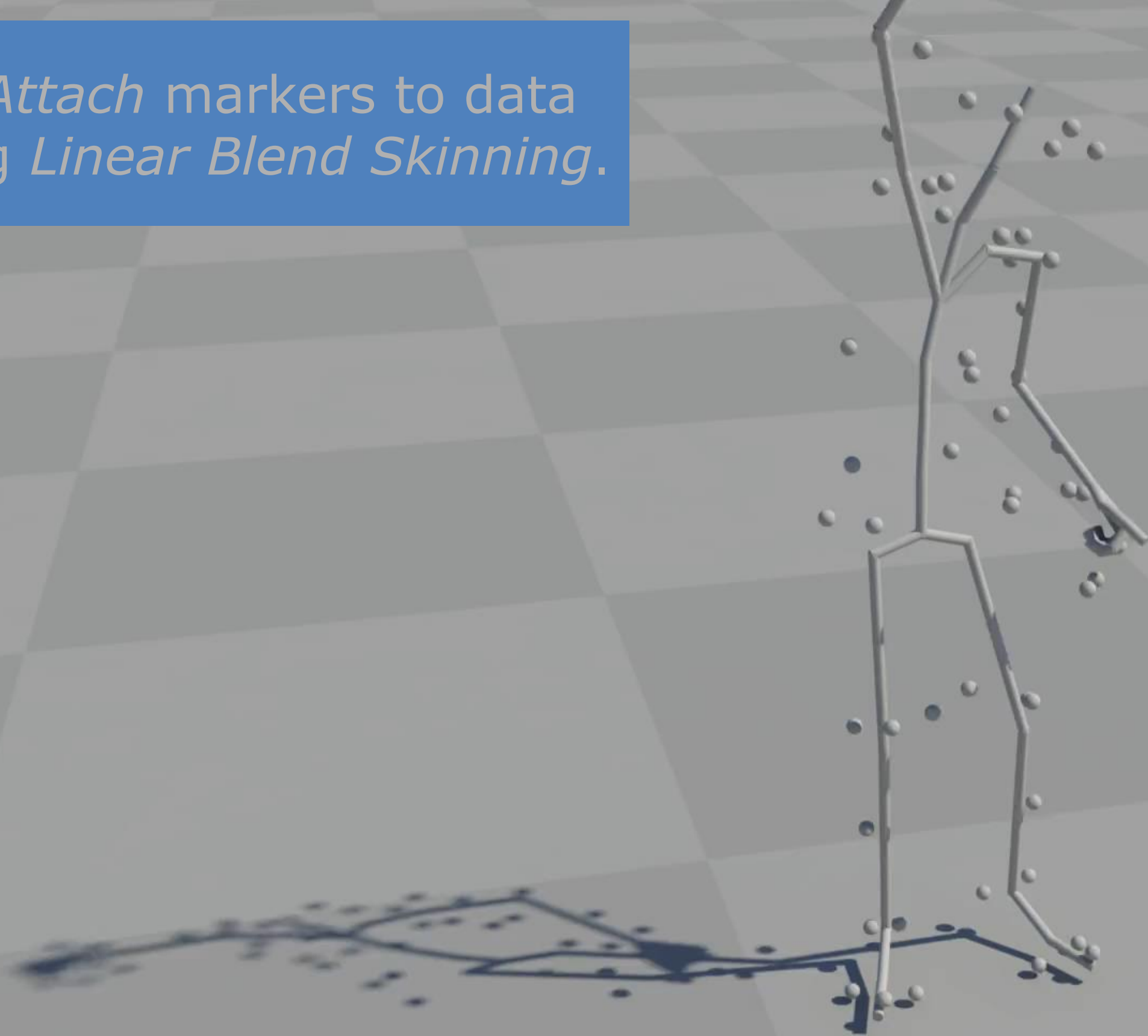


[Dahl et al. 2017]

1. *Acquire a large database of Skeletal Animation Data*



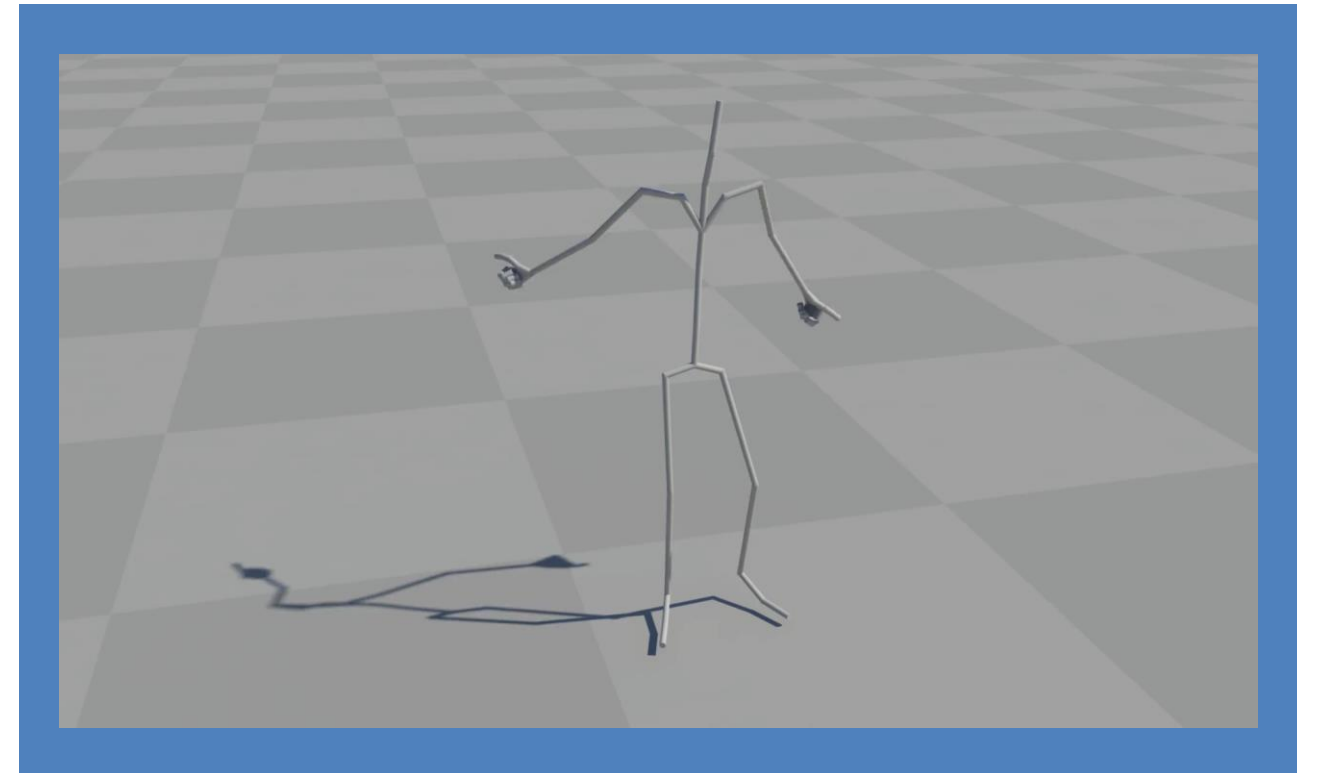
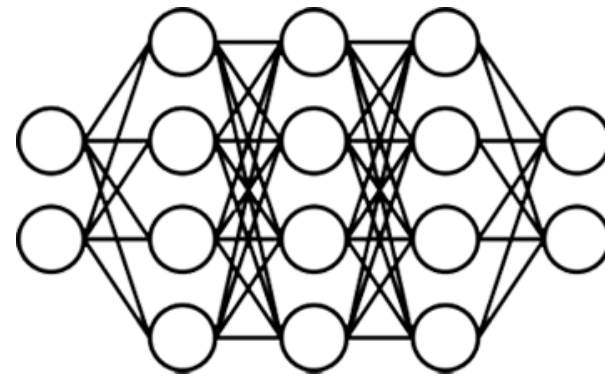
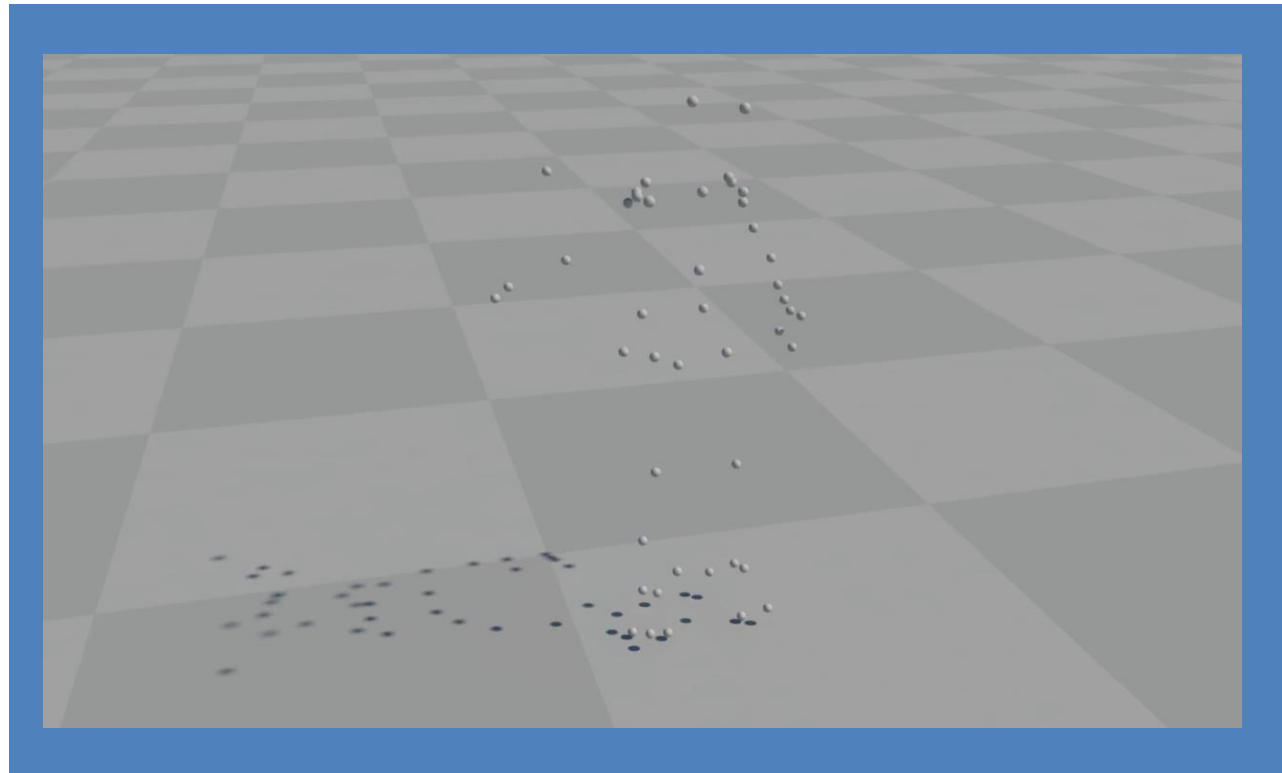
2. *Attach* markers to data using *Linear Blend Skinning*.



3. *Corrupt markers using a Custom Noise Function.*



Robust Solving by Denoising



4. Train a Neural Network to map from corrupted markers to original Motion.

Representation

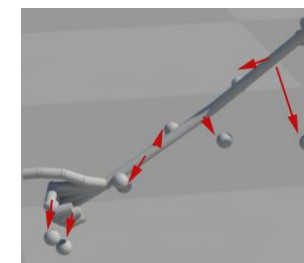
- x **Marker Positions** for a single pose flattened into a vector.
- y **Joint Transforms** for a single pose flattened into a vector.
- z **Local Offsets** from each marker to each associated joint.

Training Algorithm

Sample a set of marker configurations.

$$Z \sim \mathcal{N}(z^\mu, z^\Sigma)$$

z



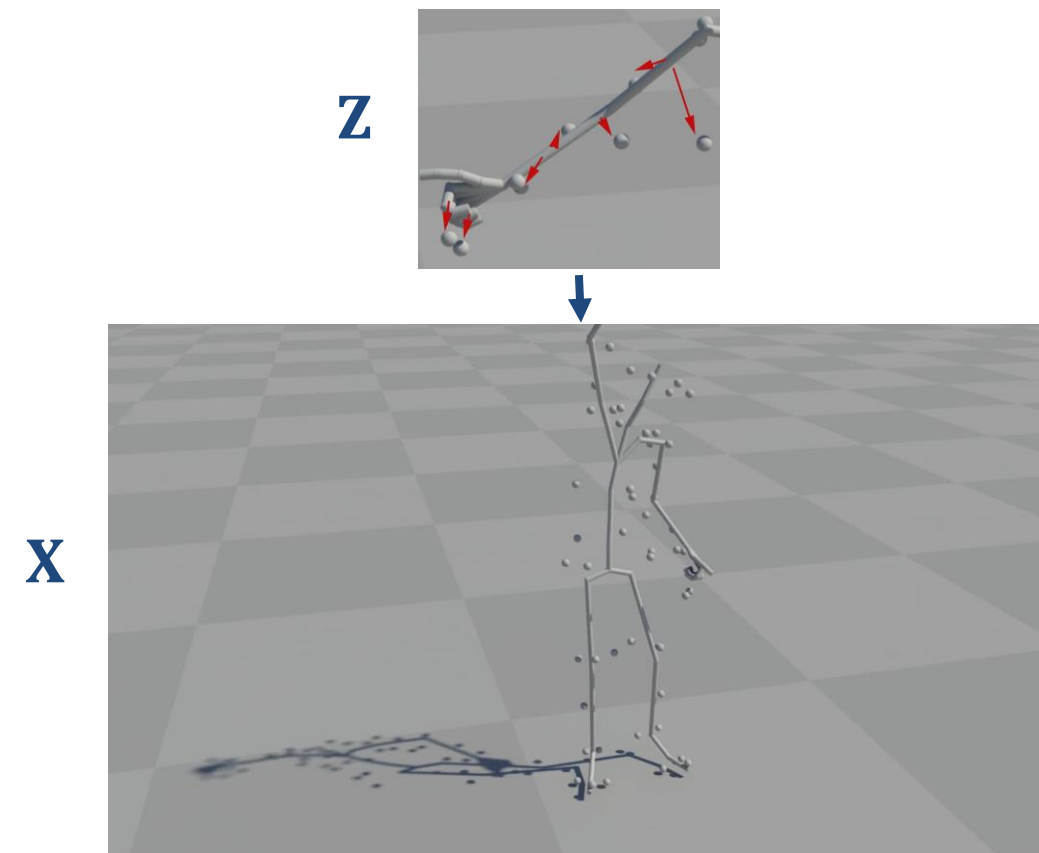
Training Algorithm

Sample a set of marker configurations.

$$Z \sim \mathcal{N}(z^\mu, z^\Sigma)$$

Compute marker positions via linear blend skinning.

$$X \leftarrow \text{LBS}(Y, Z)$$



Training Algorithm

Sample a set of marker configurations.

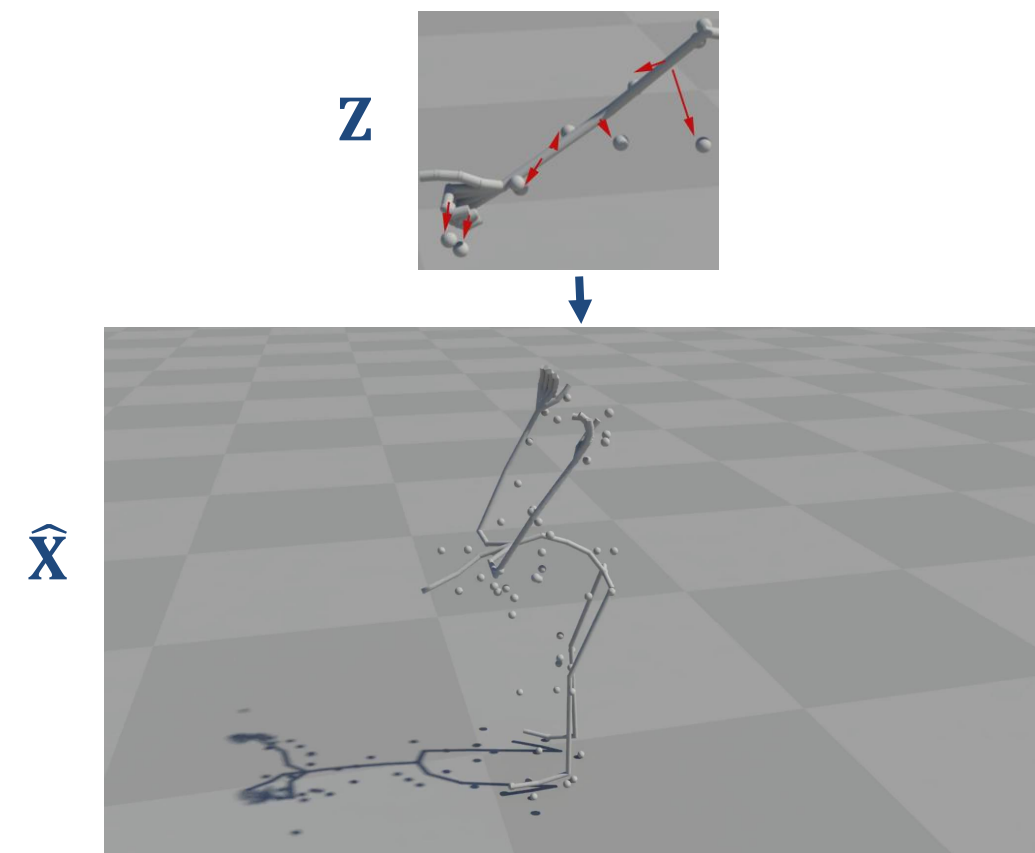
$$\mathbf{Z} \sim \mathcal{N}(\mathbf{z}^\mu, \mathbf{z}^\Sigma)$$

Compute marker positions via linear blend skinning.

$$\mathbf{X} \leftarrow \text{LBS}(\mathbf{Y}, \mathbf{Z})$$

Corrupt markers.

$$\hat{\mathbf{X}} \leftarrow \text{Corrupt}(\mathbf{X})$$



Training Algorithm

Sample a set of marker configurations.

$$\mathbf{Z} \sim \mathcal{N}(\mathbf{z}^{\mu}, \mathbf{z}^{\Sigma})$$

Compute marker positions via linear blend skinning.

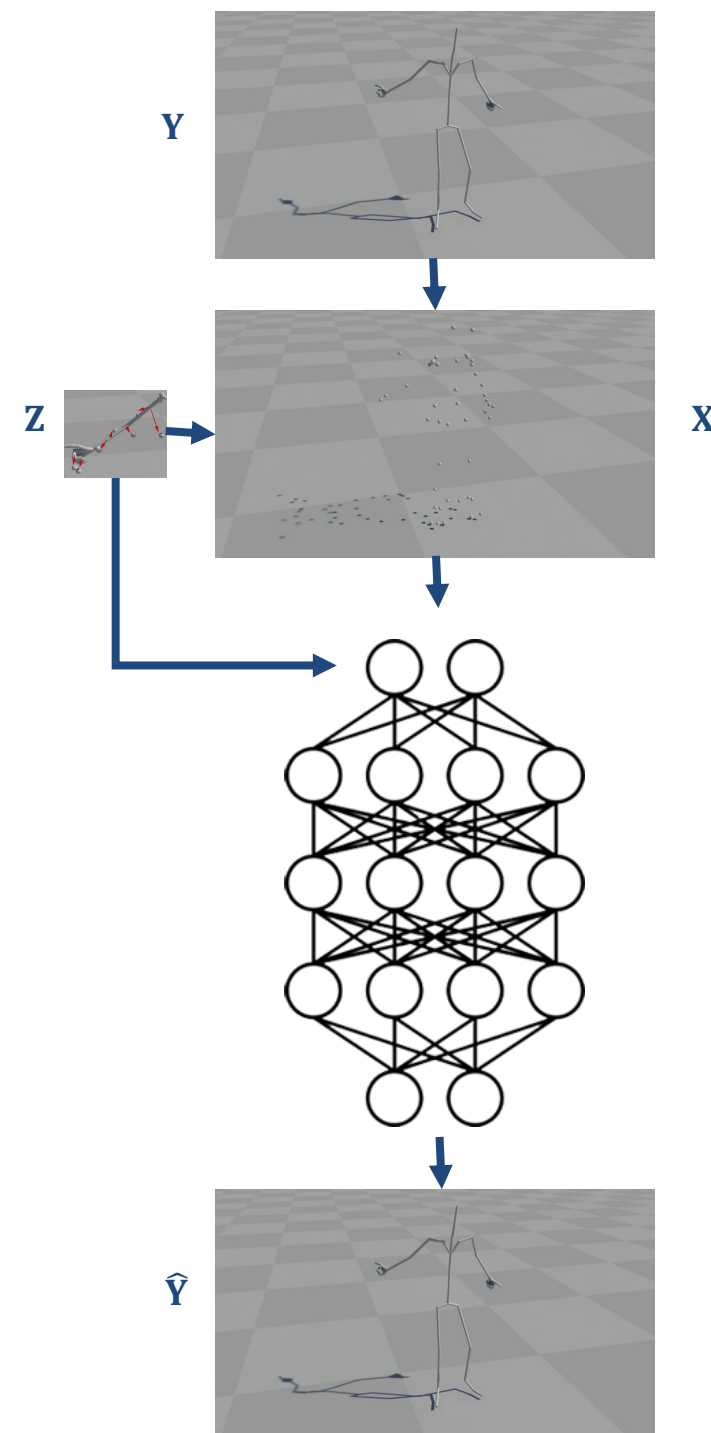
$$\mathbf{X} \leftarrow \text{LBS}(\mathbf{Y}, \mathbf{Z})$$

Corrupt markers.

$$\hat{\mathbf{X}} \leftarrow \text{Corrupt}(\mathbf{X})$$

Normalize data and input into neural network.

$$\hat{\mathbf{Y}} \leftarrow \text{Network}([\hat{\mathbf{X}} \hat{\mathbf{Z}}]; \theta)$$



Training Algorithm

Sample a set of marker configurations.

$$\mathbf{Z} \sim \mathcal{N}(\mathbf{z}^\mu, \mathbf{z}^\Sigma)$$

Compute marker positions via linear blend skinning.

$$\mathbf{X} \leftarrow \text{LBS}(\mathbf{Y}, \mathbf{Z})$$

Corrupt markers.

$$\hat{\mathbf{X}} \leftarrow \text{Corrupt}(\mathbf{X})$$

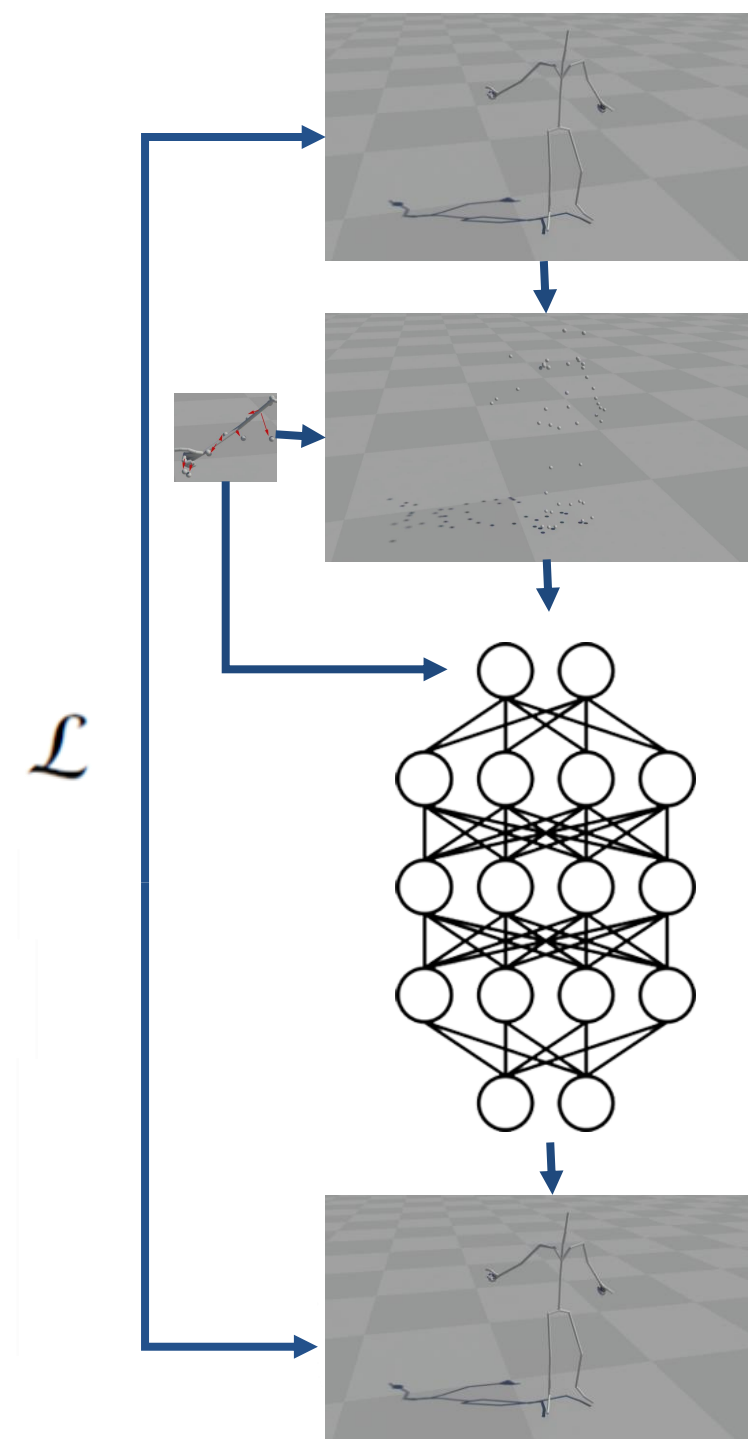
Normalize data and input into neural network.

$$\hat{\mathbf{Y}} \leftarrow \text{Network}([\hat{\mathbf{X}} \ \hat{\mathbf{Z}}]; \theta)$$

Denormalize, calculate loss, and update network parameters.

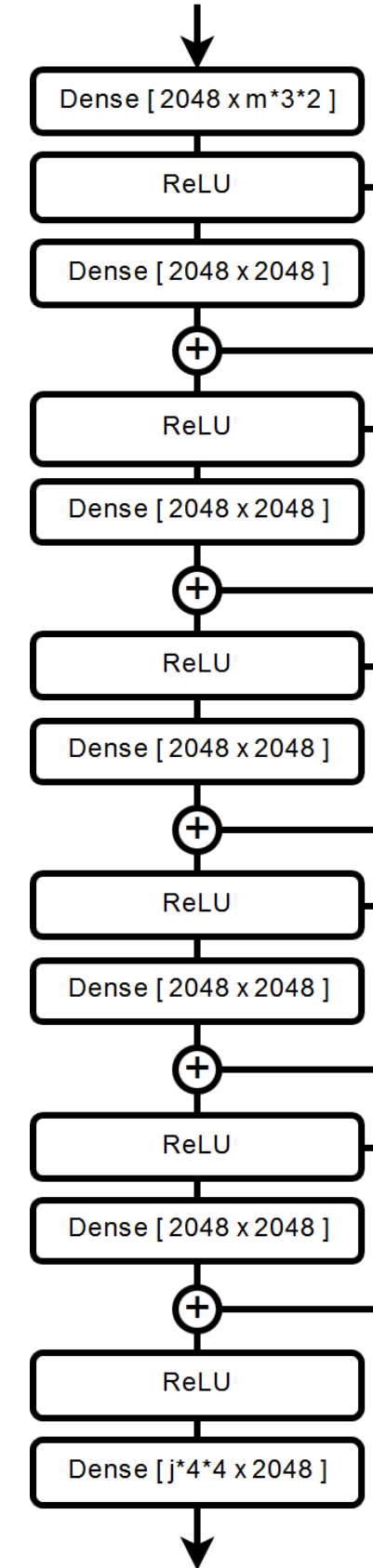
$$\mathcal{L} \leftarrow |\lambda \odot (\hat{\mathbf{Y}} - \mathbf{Y})|_1 + \gamma \|\theta\|_2^2$$

$$\theta \leftarrow \text{Adam}(\theta, \nabla \mathcal{L})$$

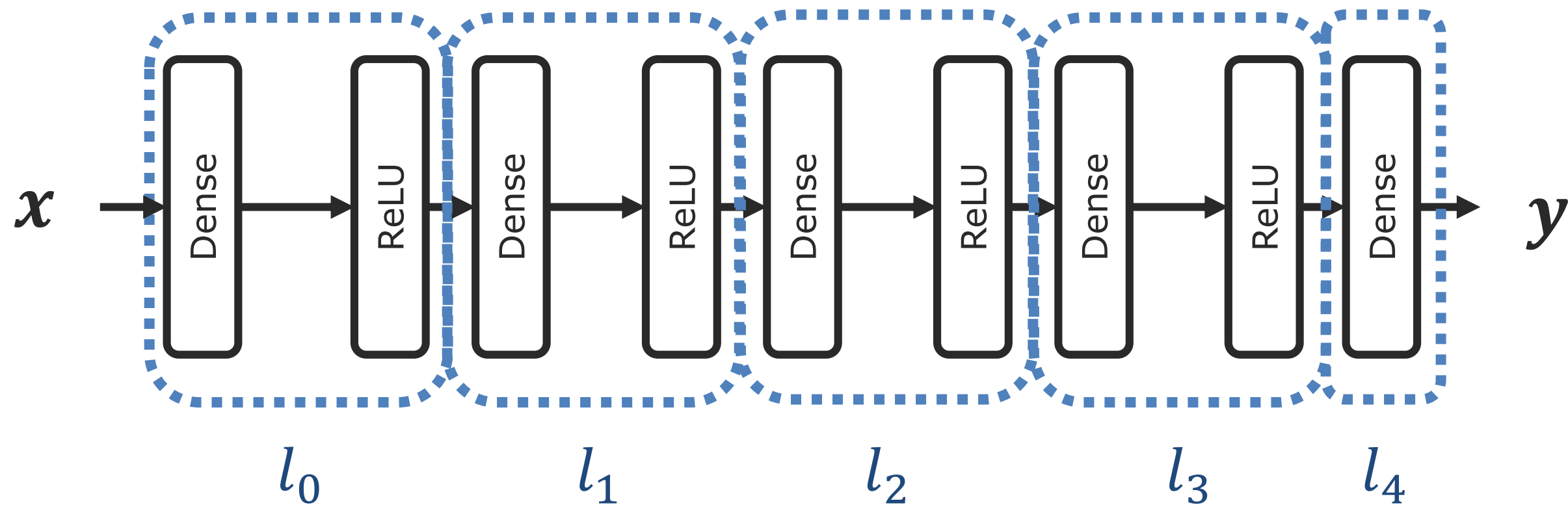


Structure

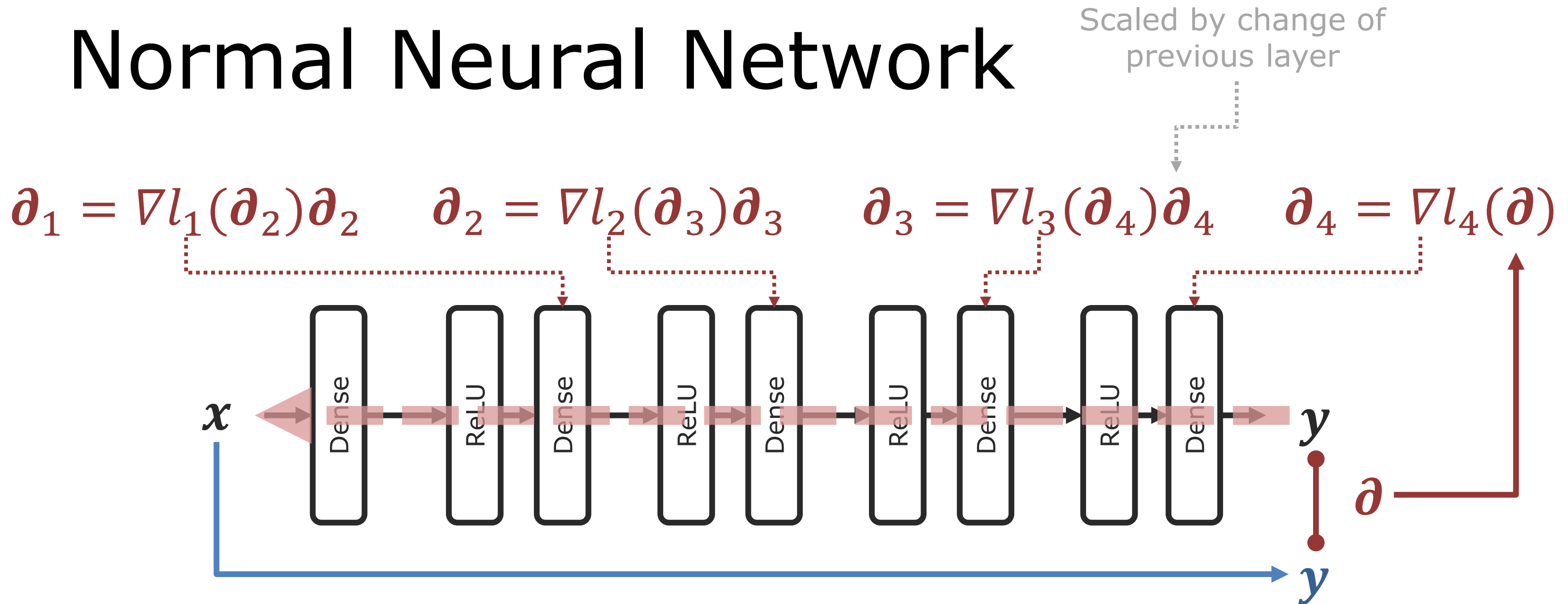
- Feed Forward Residual Neural Network.
- Inputs single pose, outputs single pose.
- More accurate results than normal network.



Normal Neural Network



Normal Neural Network



Normal Neural Network

Problem:

Earlier layers train slower than later layers.

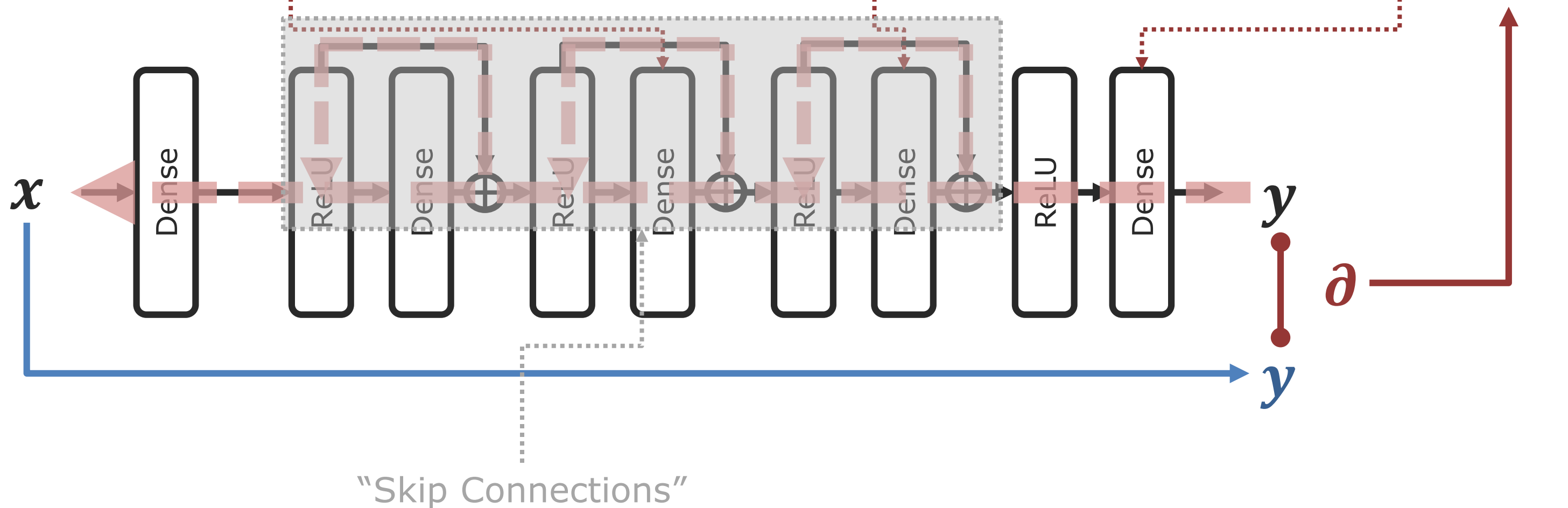
Residual Neural Network

Includes change of
previous layers

$$\partial_2 = \nabla l_2(\partial_3) \partial_3 + \partial_4$$

$$\partial_3 = \nabla l_3(\partial_4) \partial_4$$

$$\partial_4 = \nabla l_4(\partial)$$



Residual Neural Network

Solution:

Earlier layers train faster as error is propagated deeper by *Skip Connections*.

Training Data

- It is ideal if every possible pose is covered in the training data.
- We already have a massive database of motion capture data.
- Even so, we warm-start the data collection process by capturing some **extreme** range of motion takes...

TCG-22:49:31:19

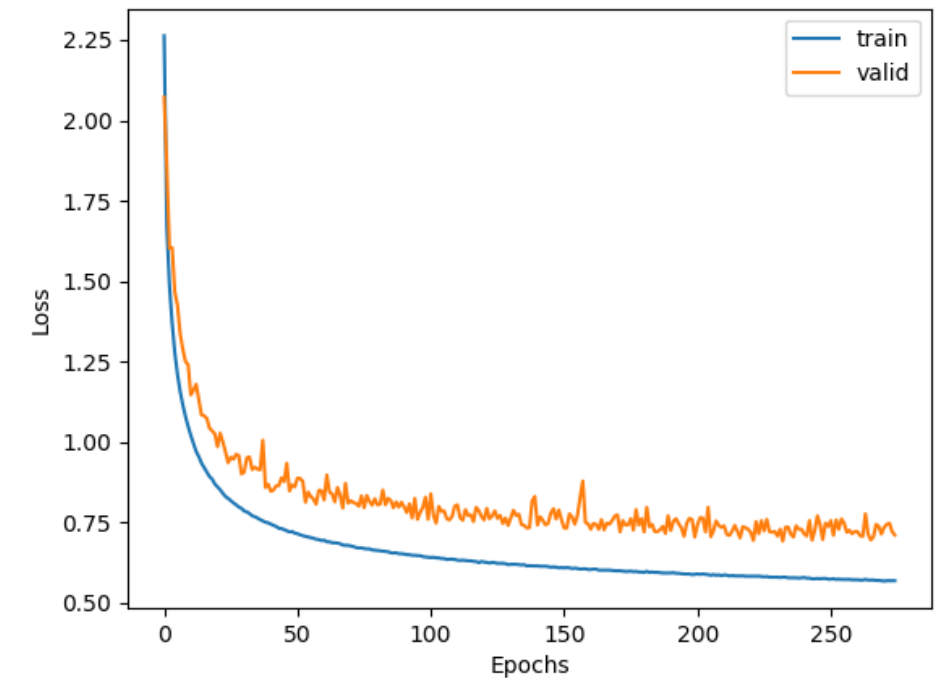


EXT-LK



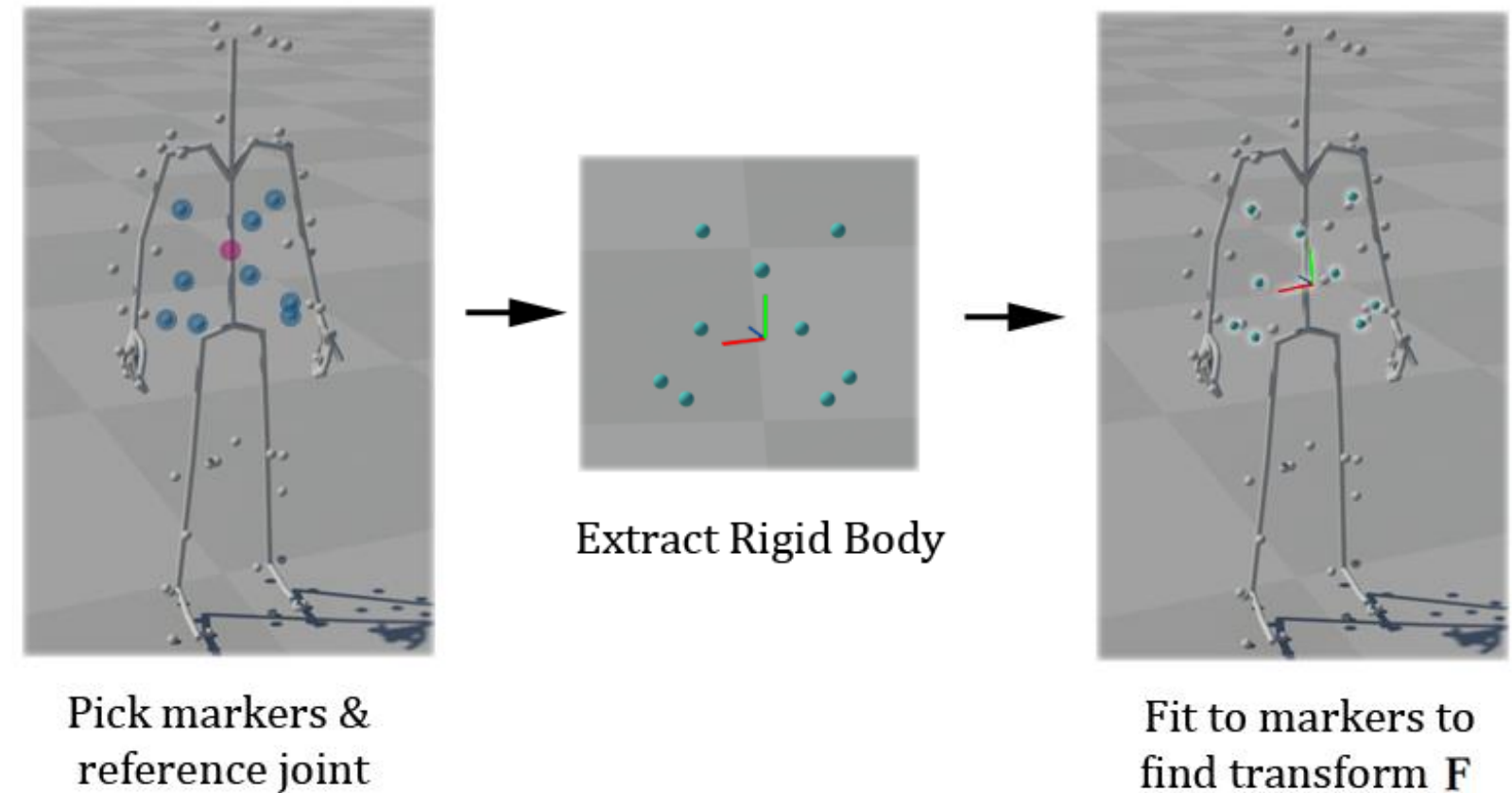
Training

- We train on around 10GB (12 hours) of Motion Capture.
- Train overnight using mid-tier graphics card.
- Perform sampling and corruption dynamically.



Local Reference Frame

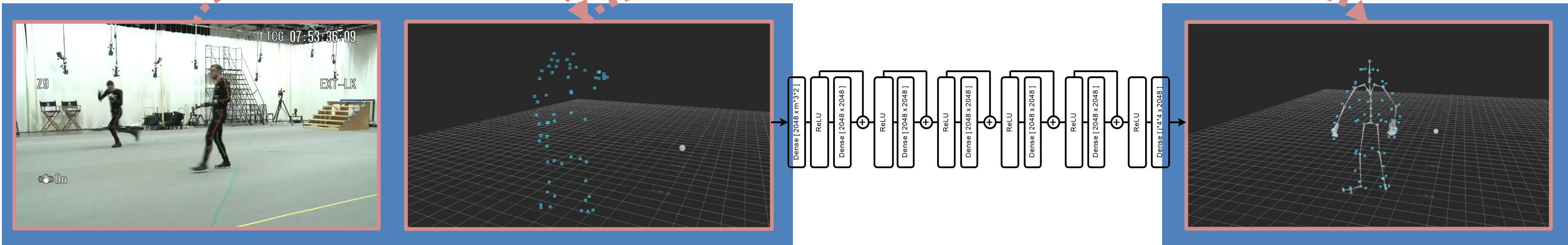
- Markers and joints must be represented local to the character.
- We find this transform using Rigid Body Fitting.
- This process must run before the markers have been cleaned.

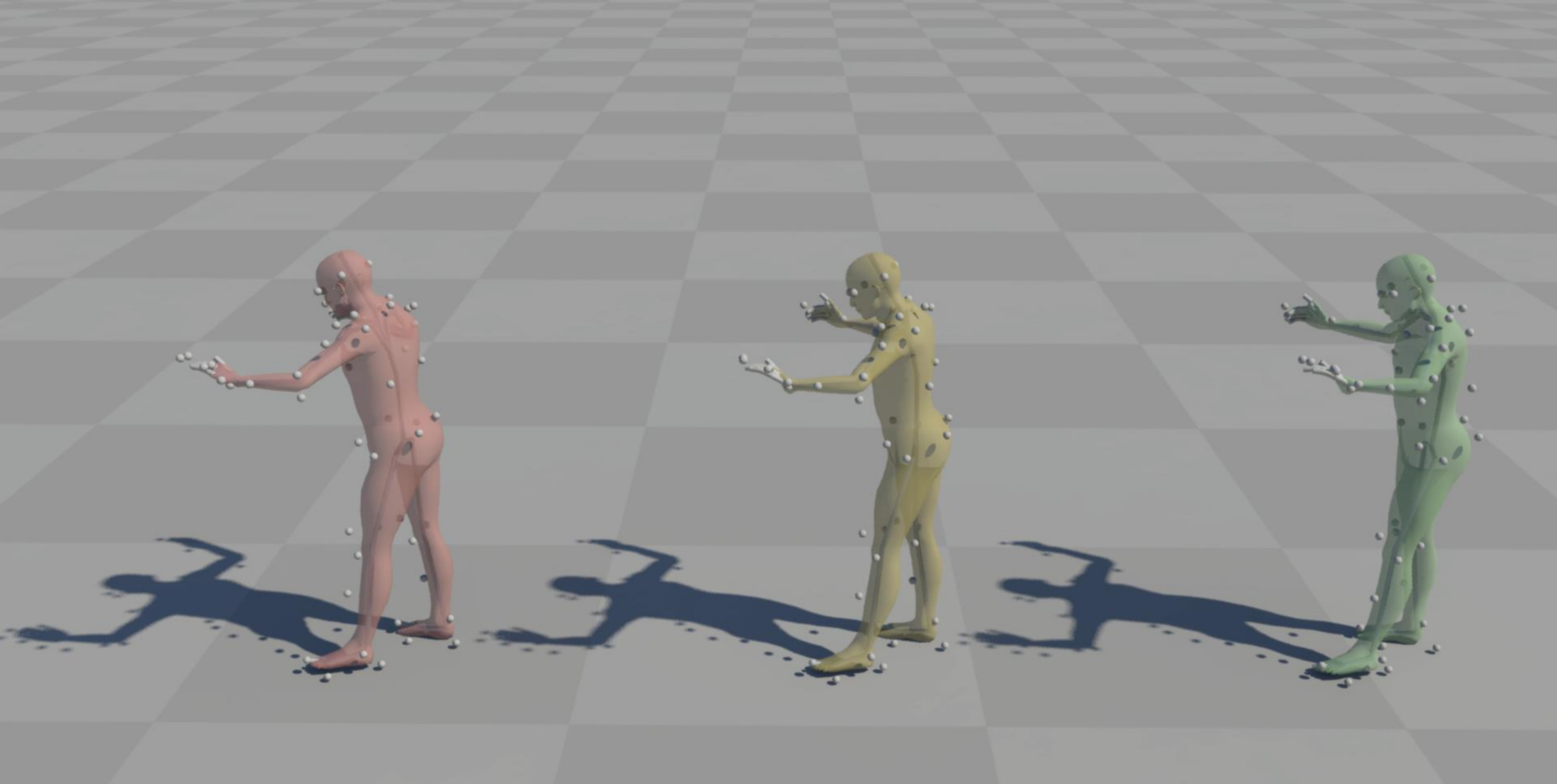


Motion Capture Pipeline

Tracking

Robust Solving

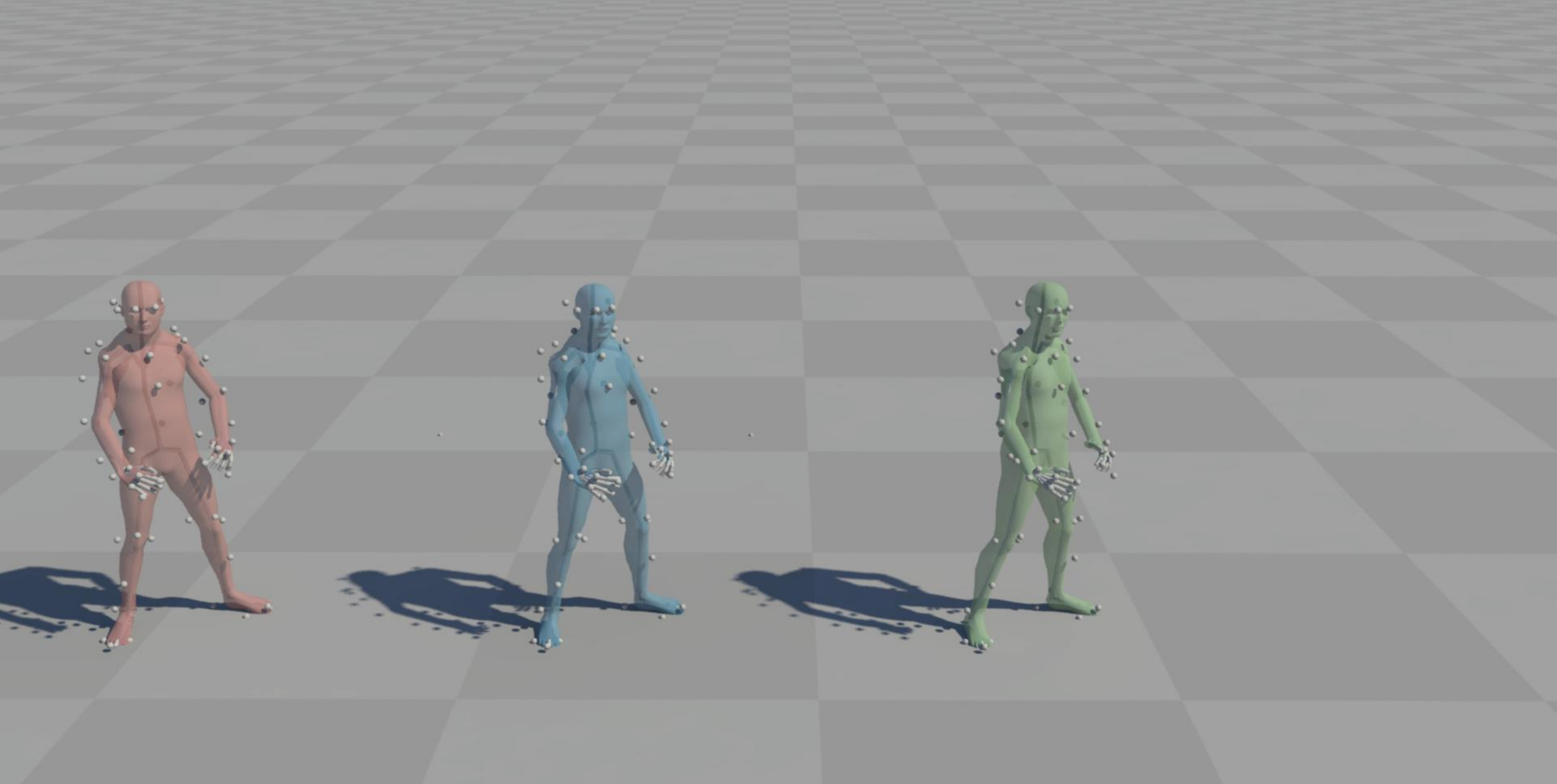




Worst Case

Commercial Software

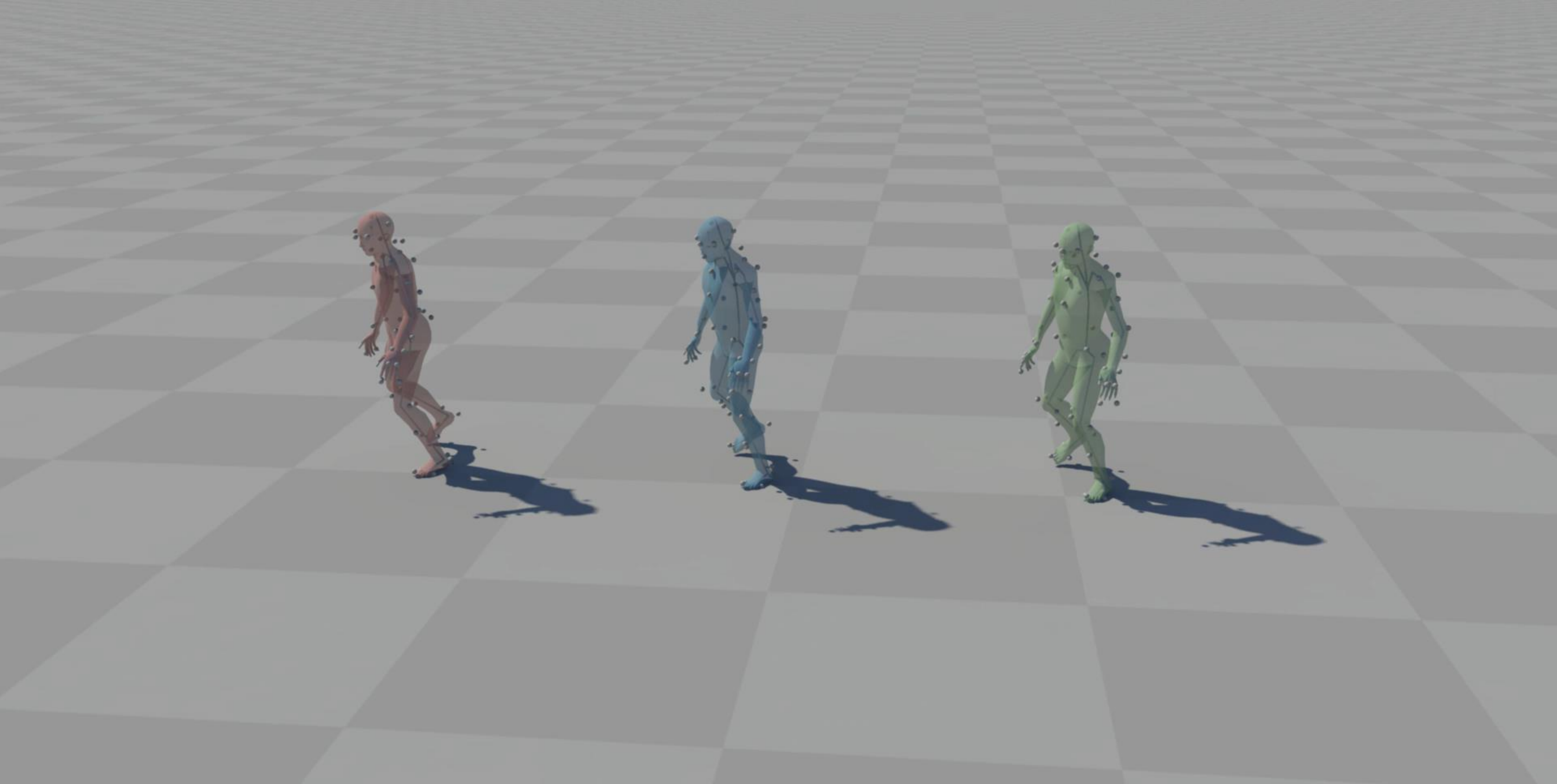
Hand Cleaned



Worst Case

Our Method

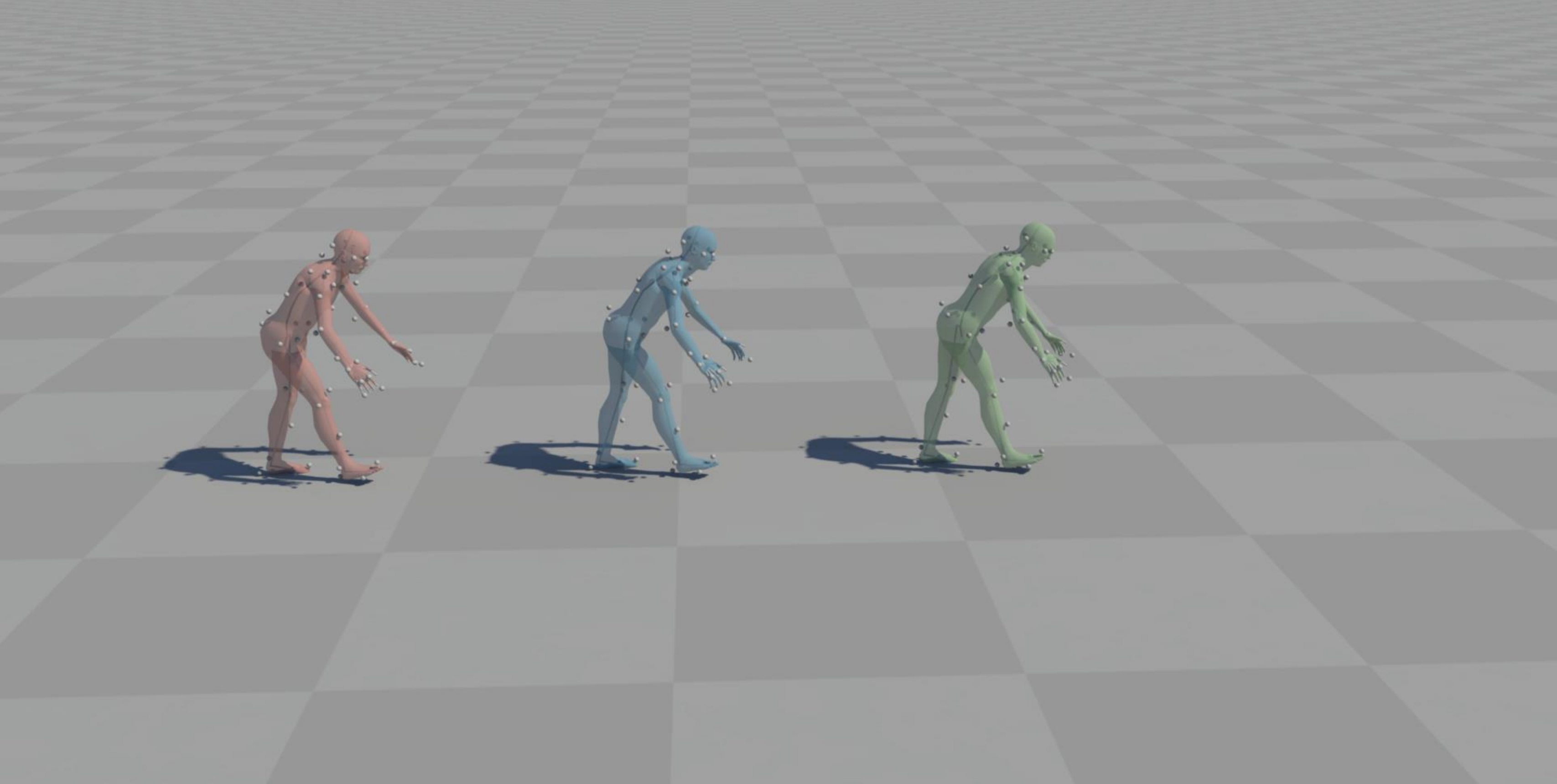
Hand Cleaned



Worst Case

Our Method

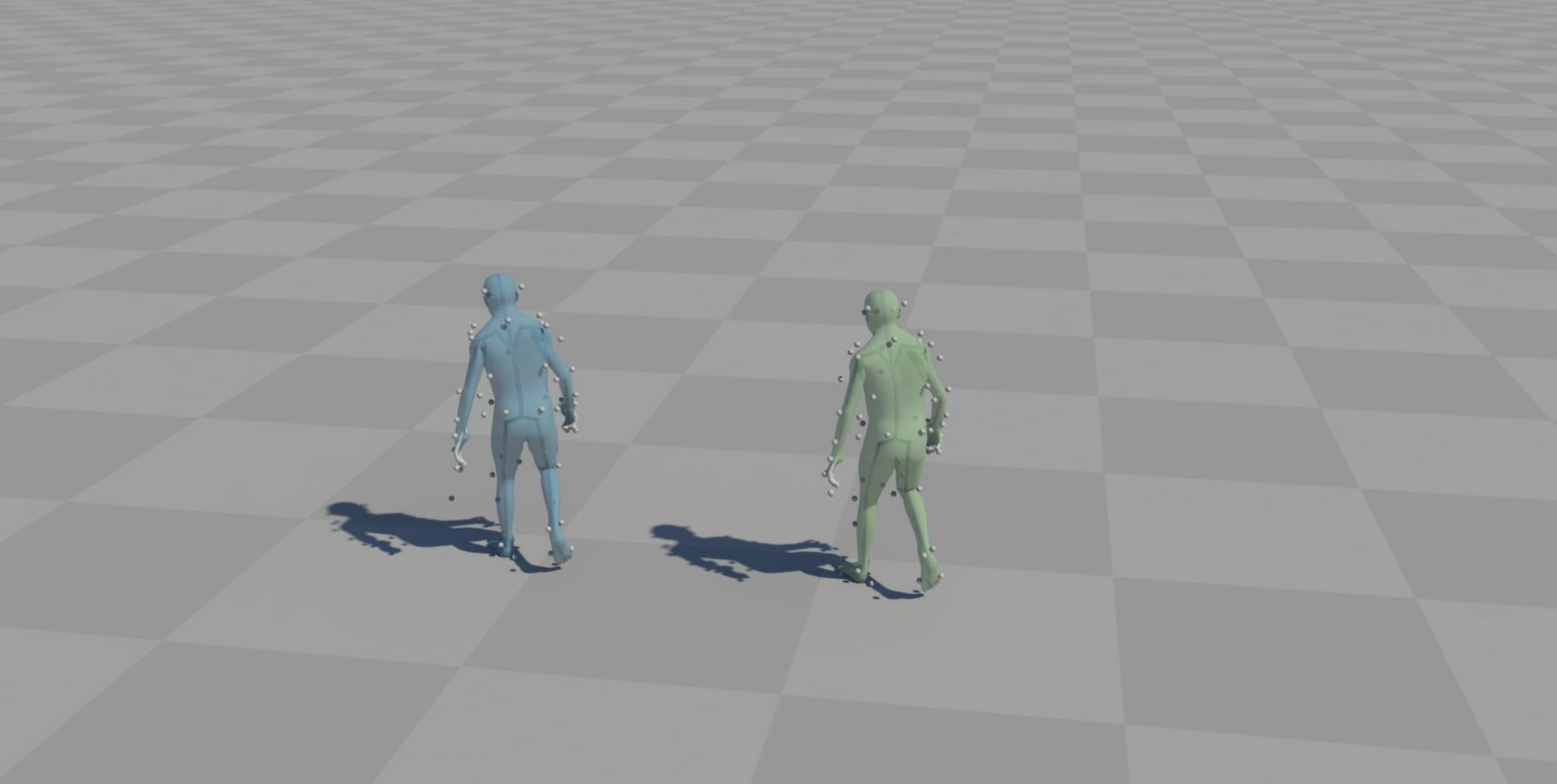
Hand Cleaned



Worst Case

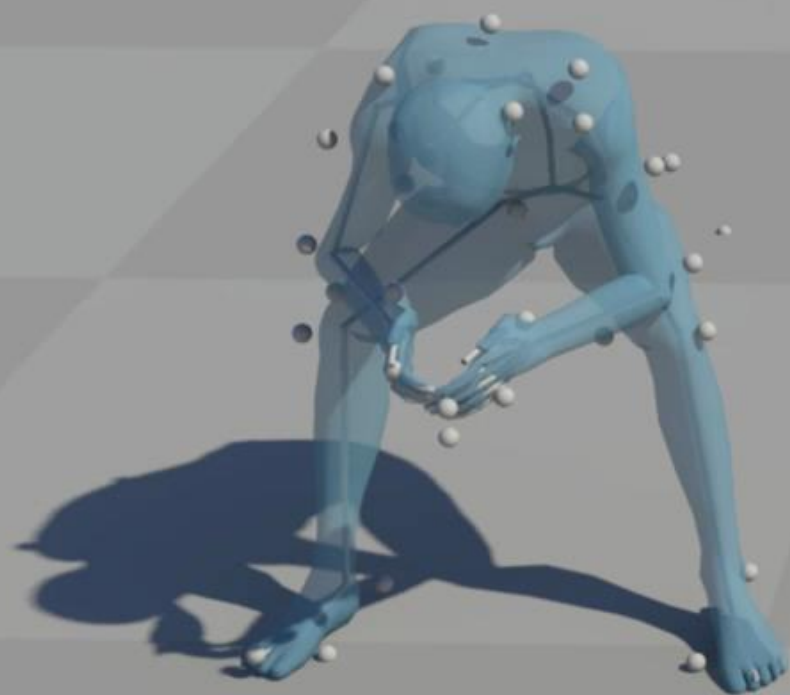
Our Method

Hand Cleaned



Our Method

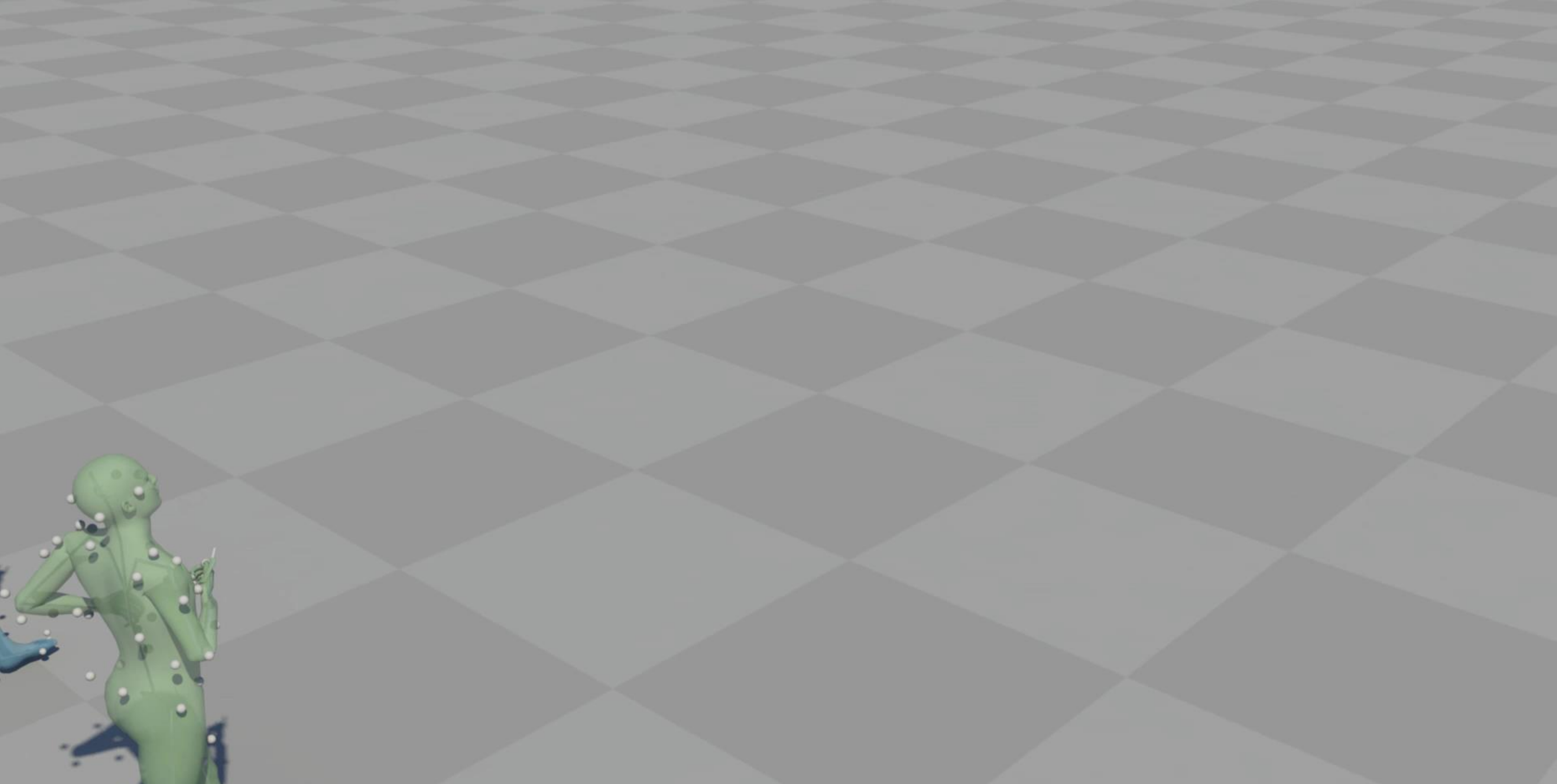
Ground Truth



Our Method



Ground Truth



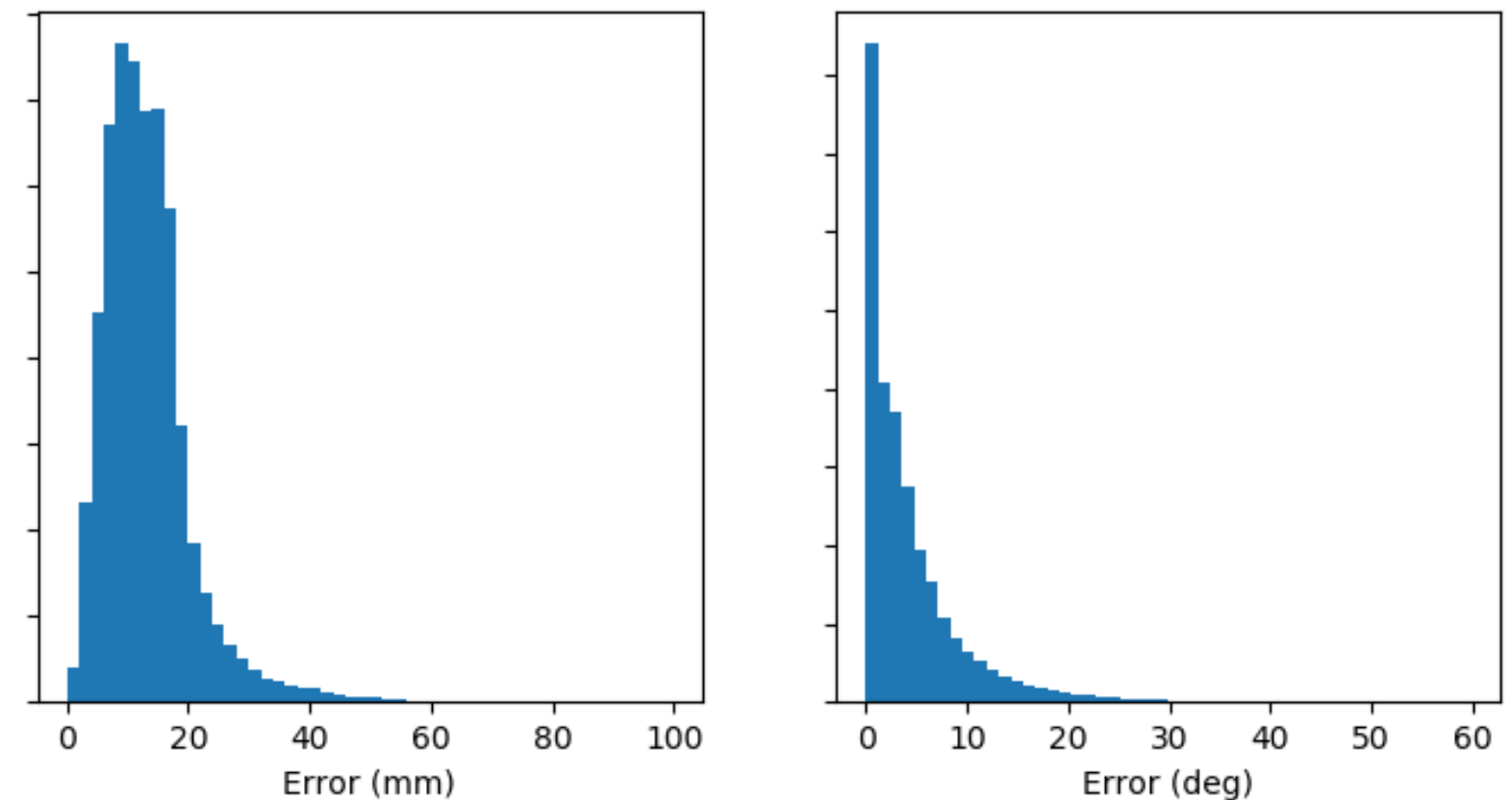
Our Method

Ground Truth

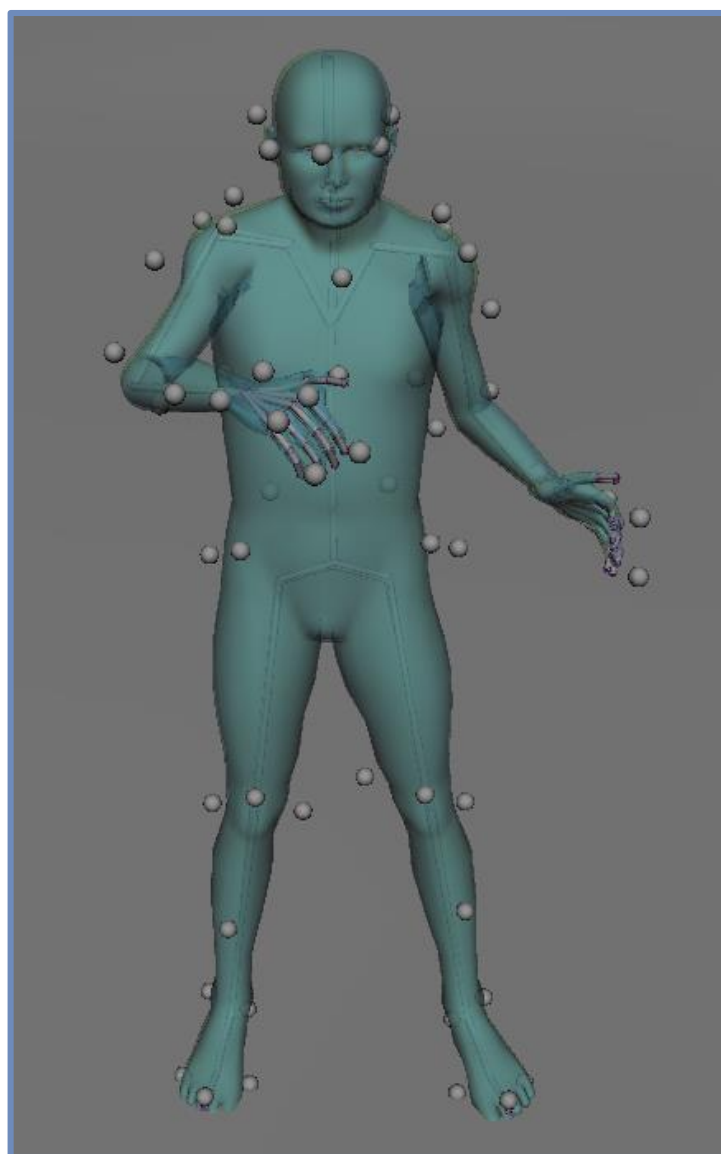
Evaluation

- **90%** of errors are less than **20mm** or **10°**.
- **99.9%** of errors are less than **60mm** or **40°**.

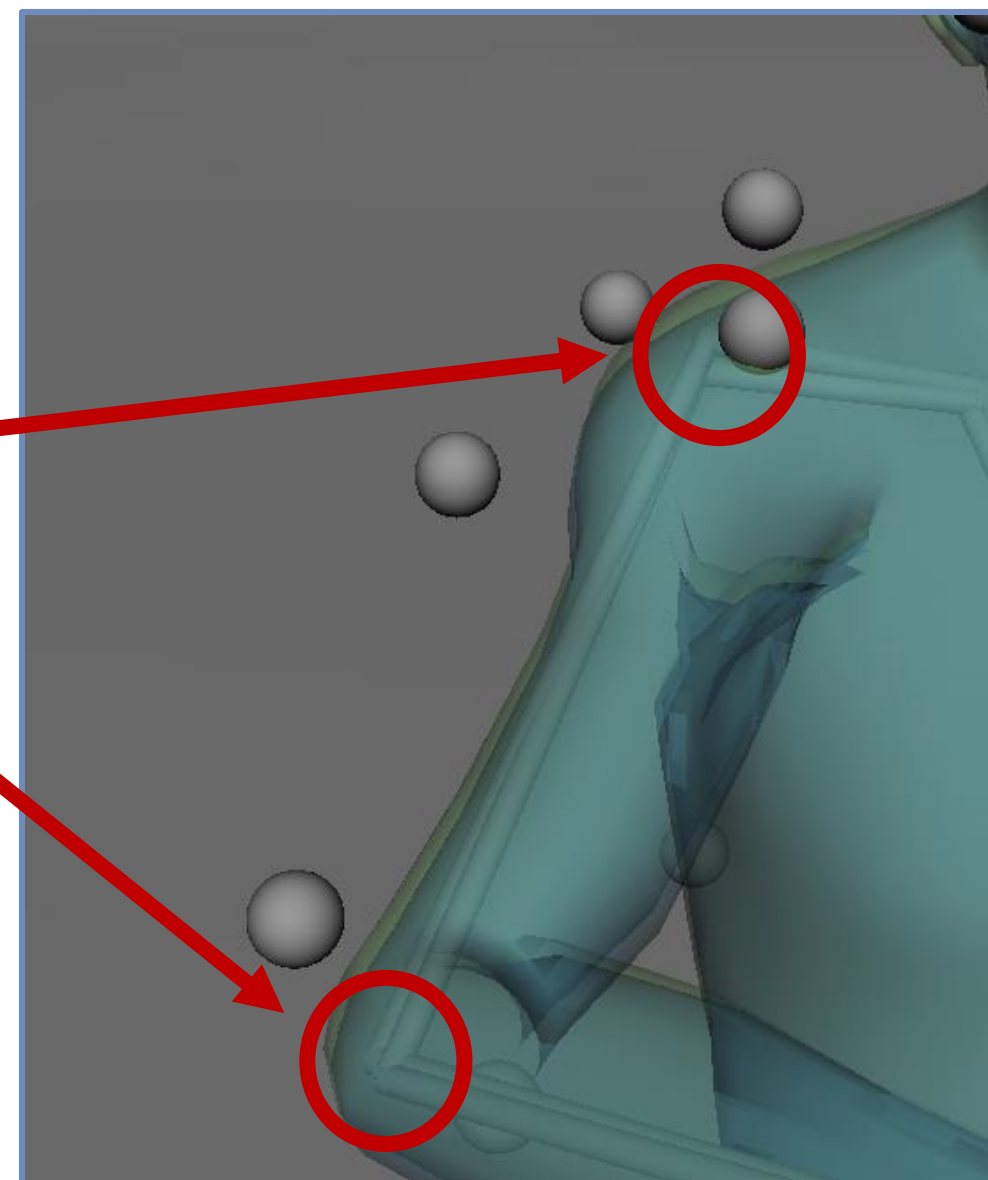
Distribution of Errors in Test Data



Average Case



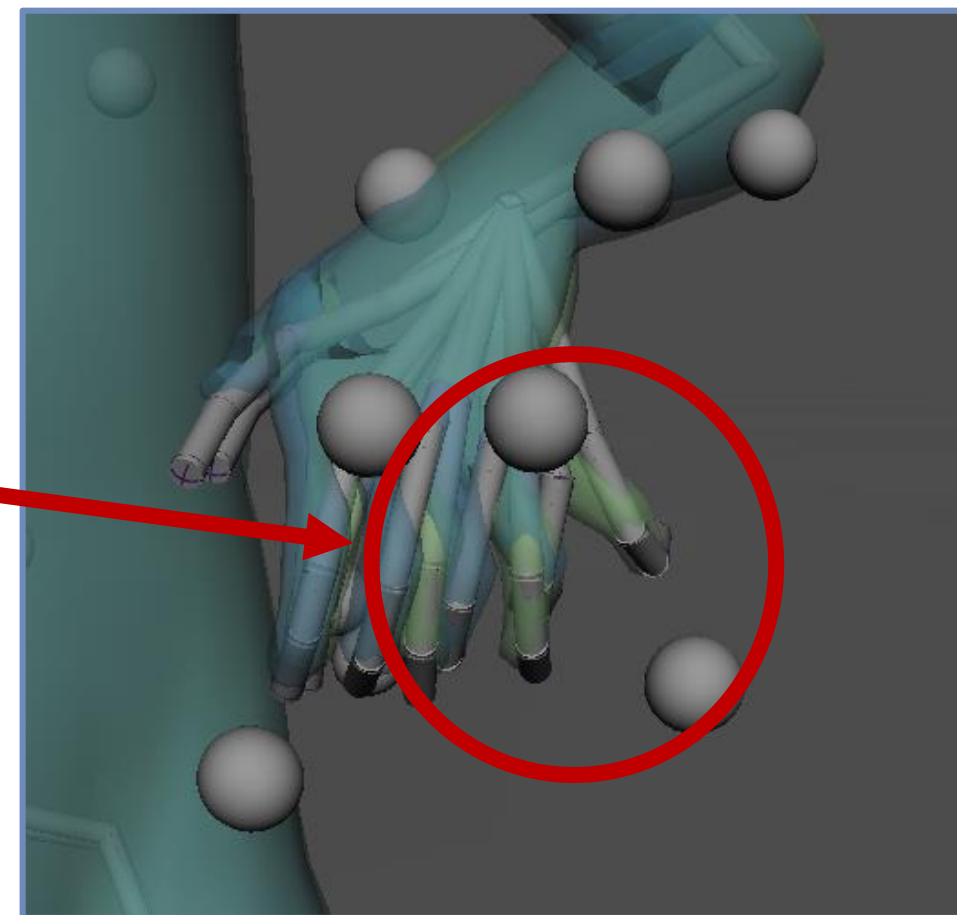
**Around
15mm
difference**



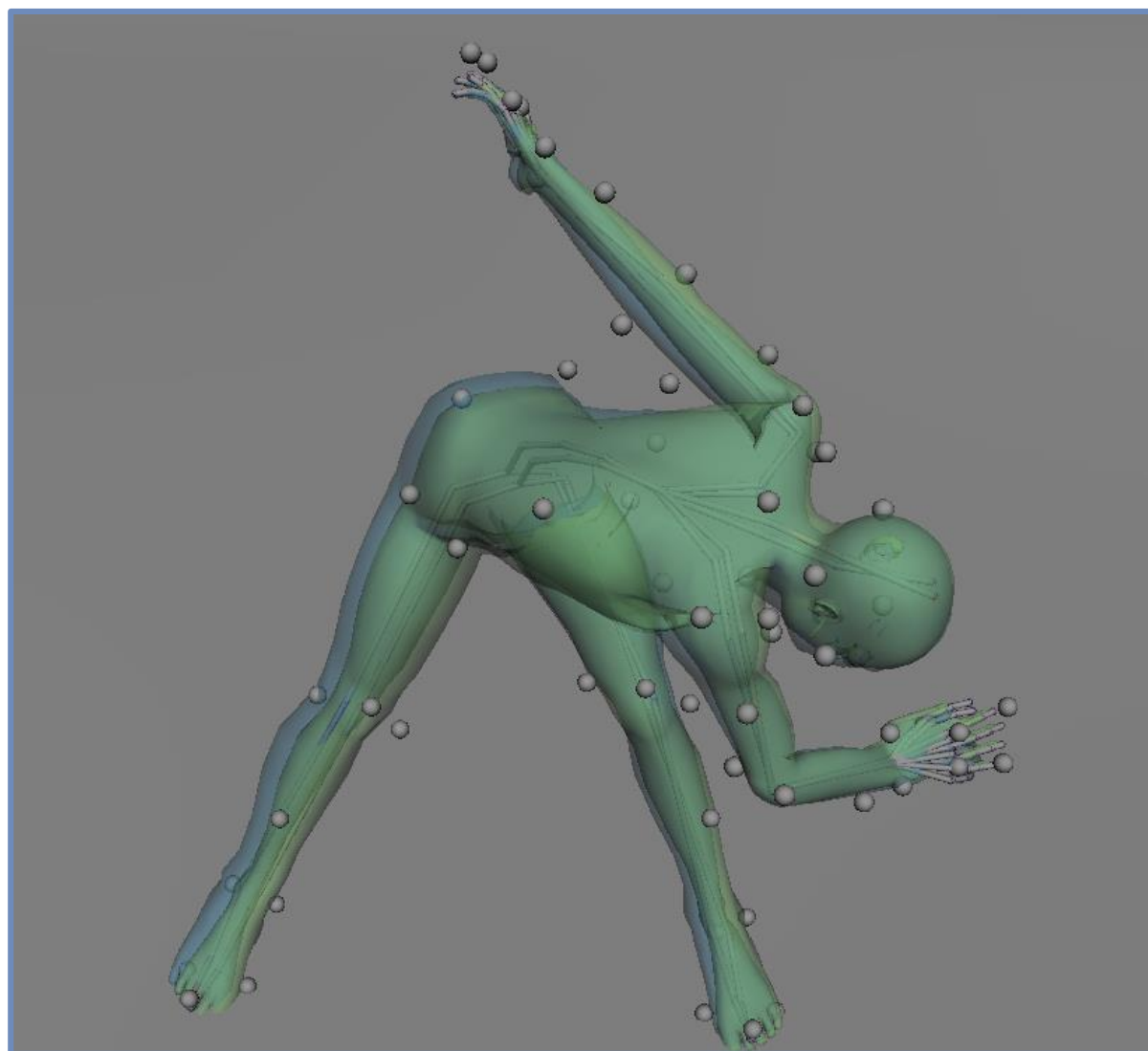
Worst Case



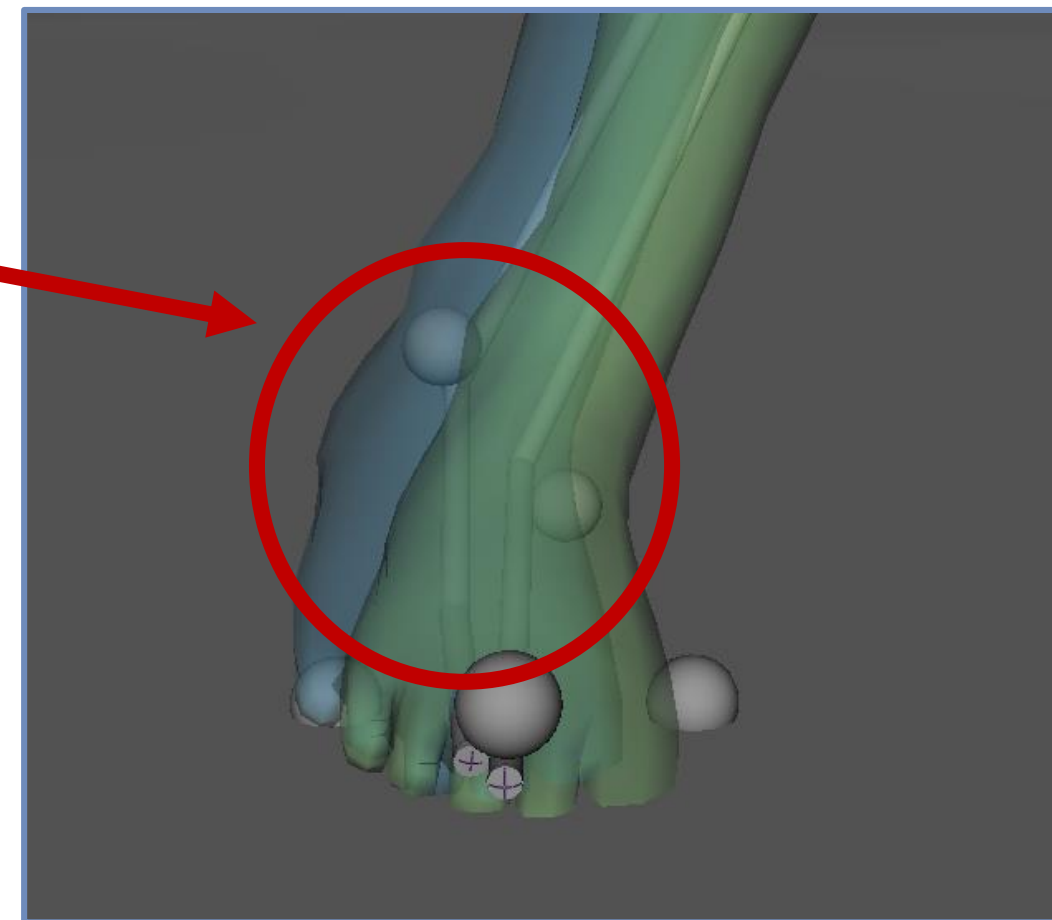
**Around
40mm / 30°
difference**



Worst Case



**Around
50mm
difference**



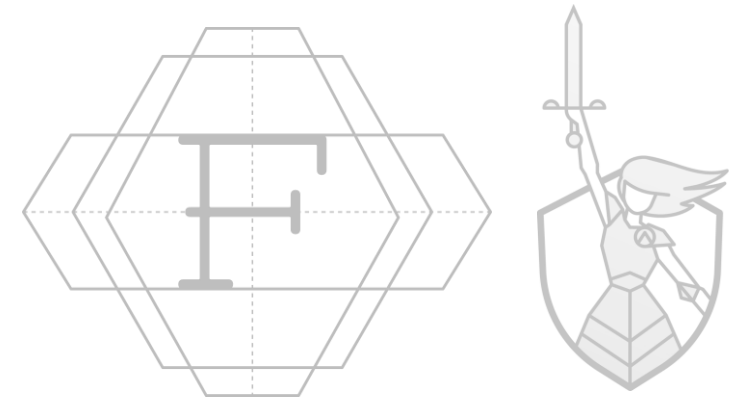
Summary

- We avoid cleaning data by making *solving* robust to errors.
- We dynamically generate data with a custom noise function.
- We train a neural network to perform the solving task.

History

Mocap Cleaning
Facial Tracking
Audio to Facial

The Future



TCG 02:21:06:26

Z21

EXT-LK

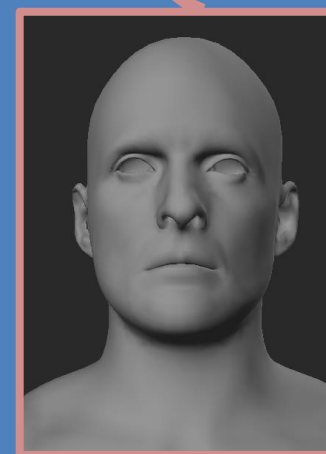
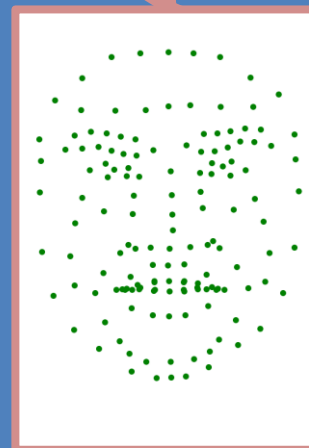
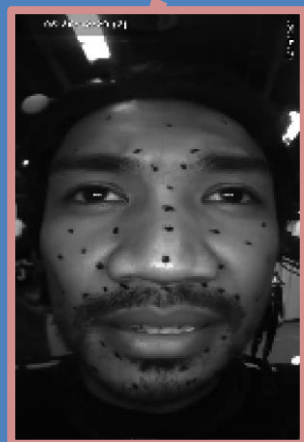
«» On
Full MF

Facial Capture Pipeline

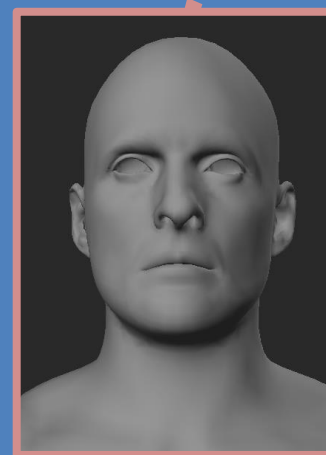
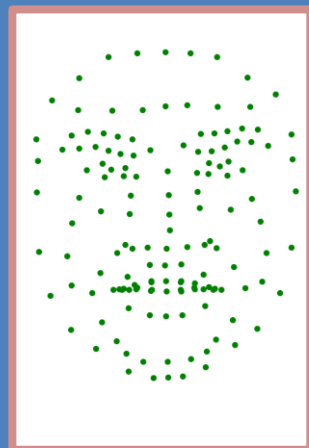
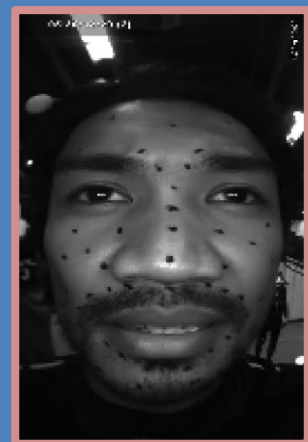
Tracking

Retargeting

Polishing



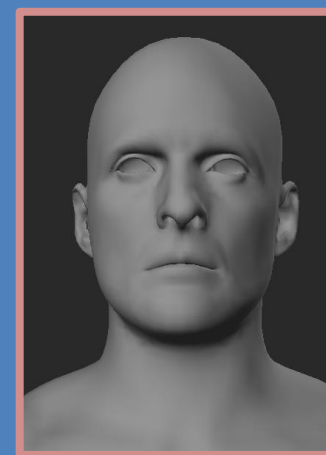
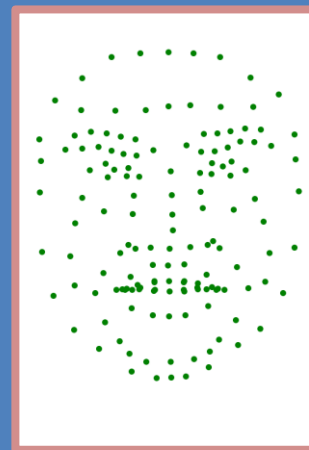
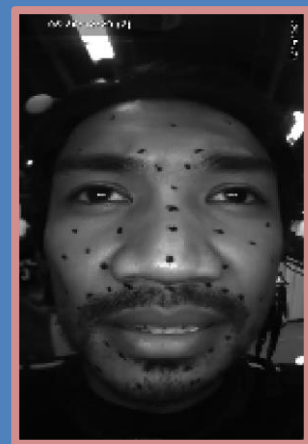
Facial Capture Pipeline



Polishing

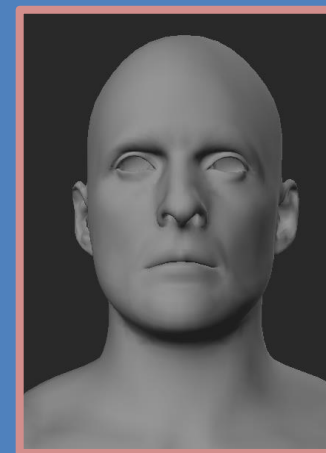
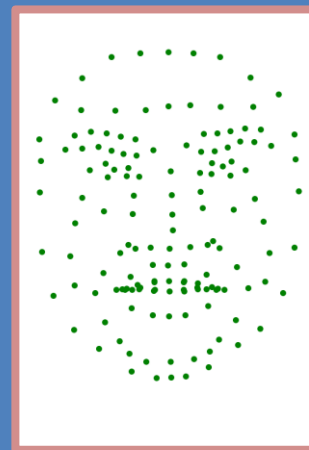
Desired Focus

Facial Capture Pipeline



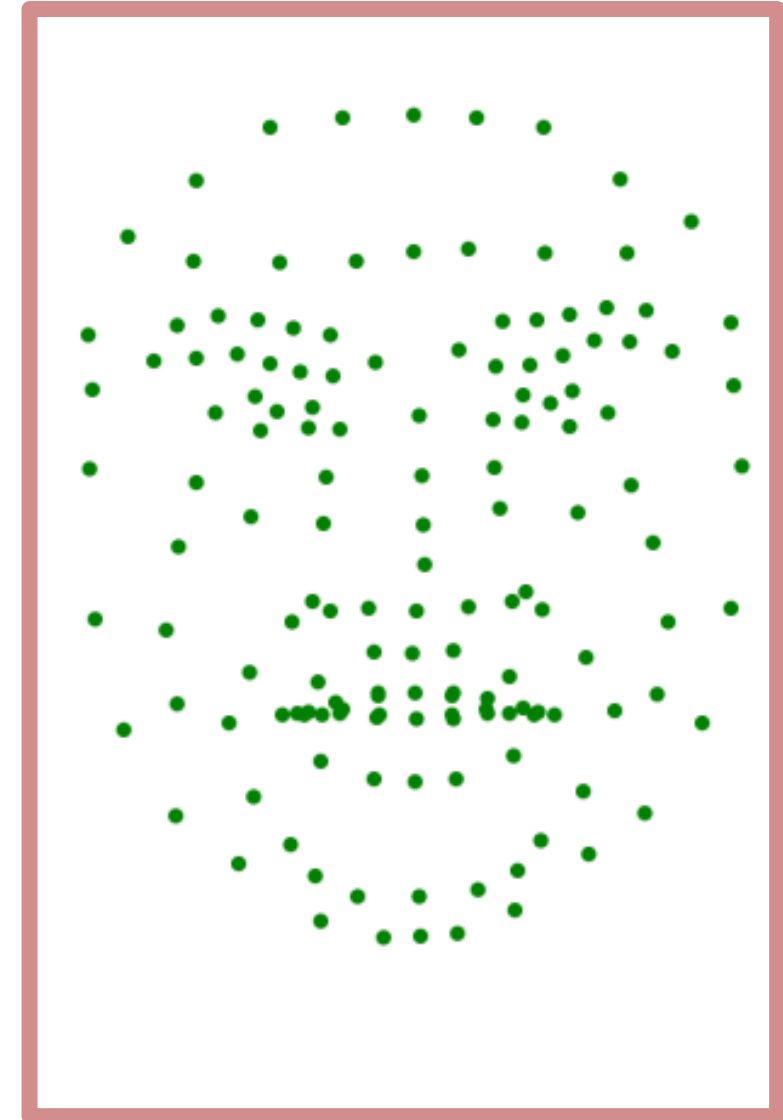
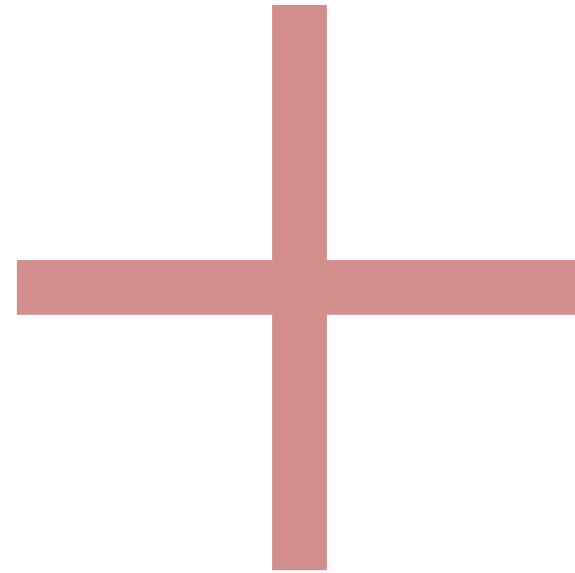
What if we could automate either of these stages of the pipeline?

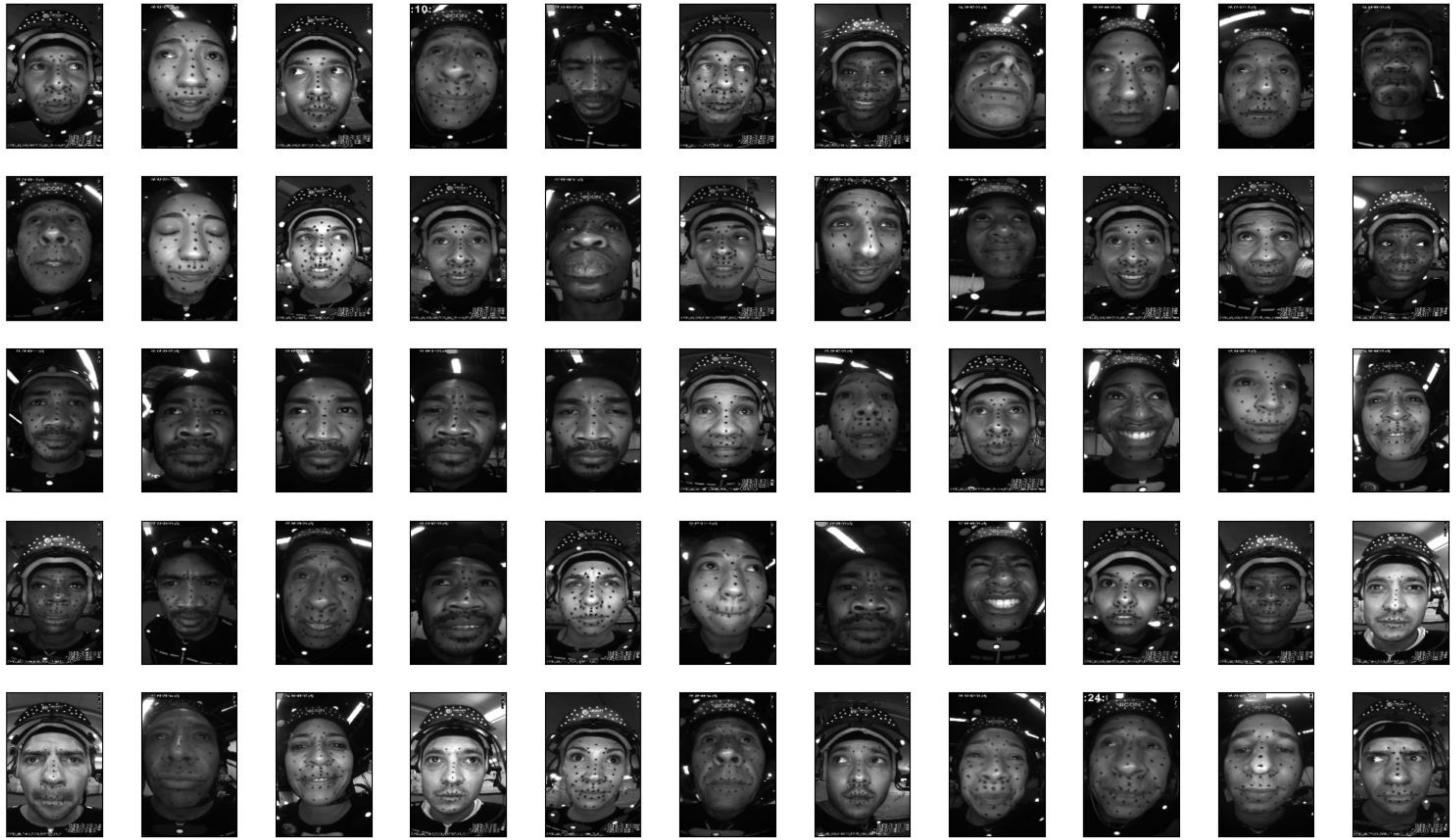
Facial Capture Pipeline



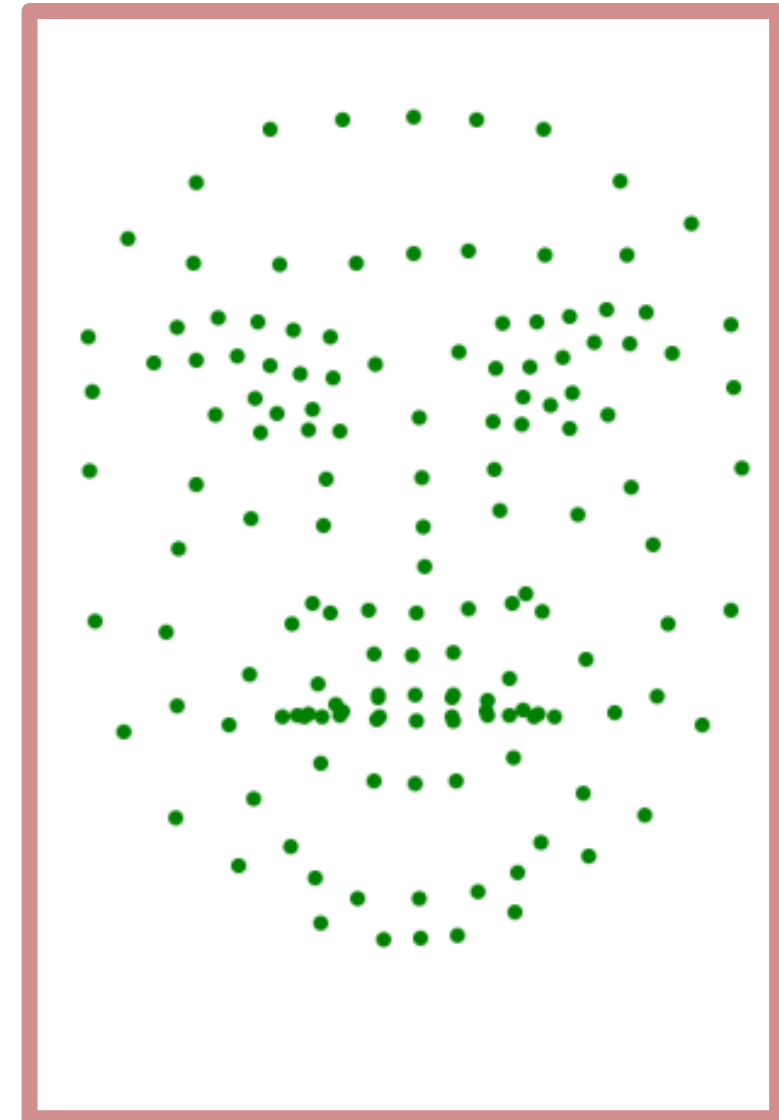
Automatic Facial Tracking

We have many hours of facial animation that is already tracked.

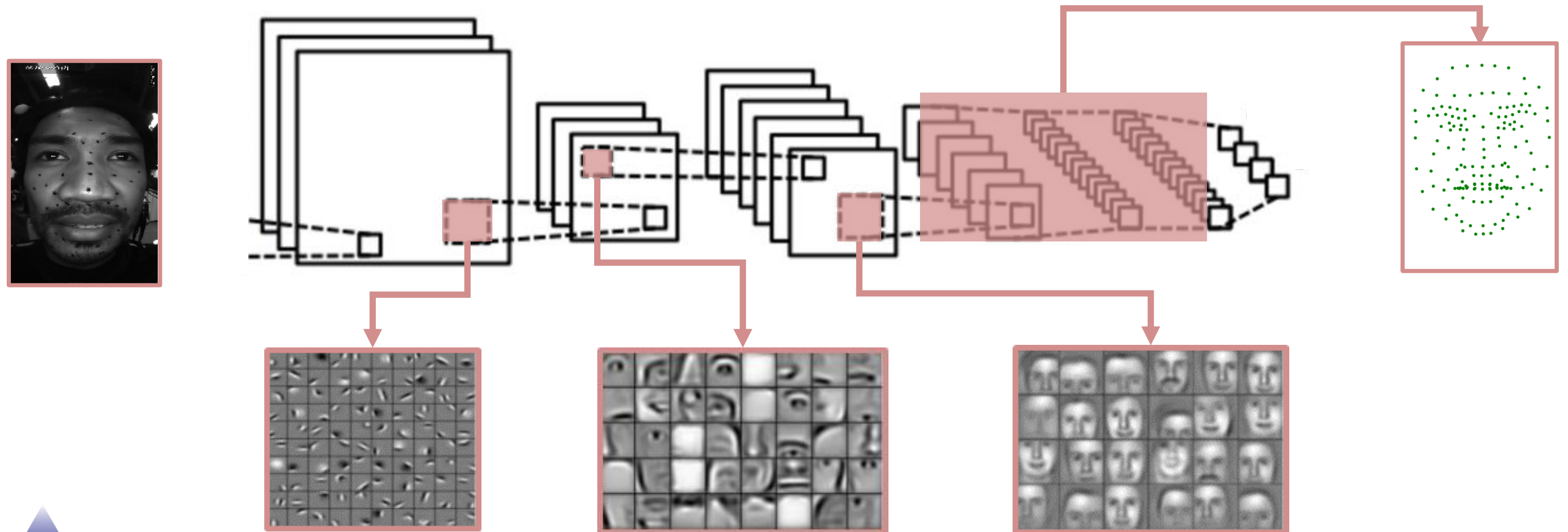


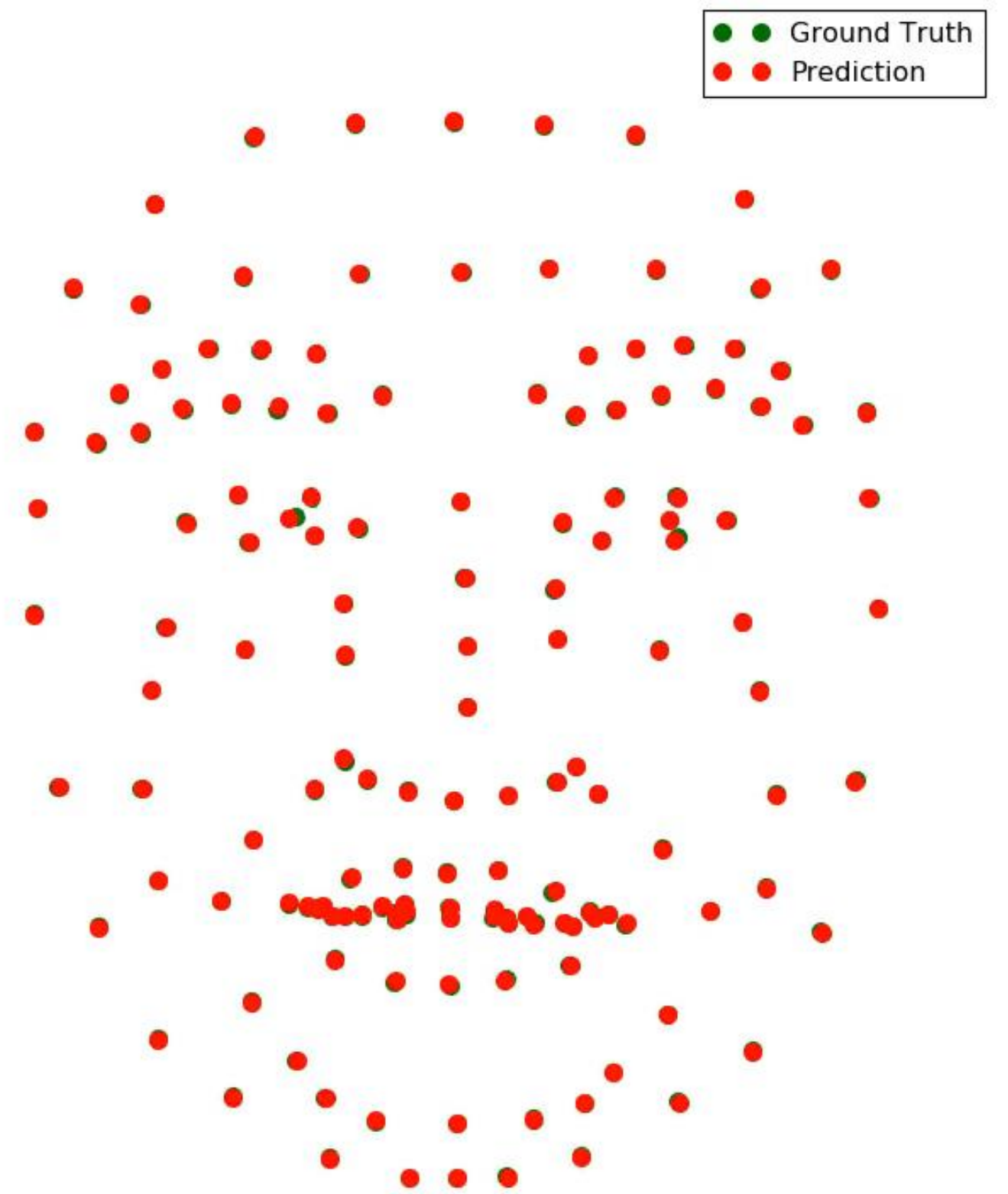


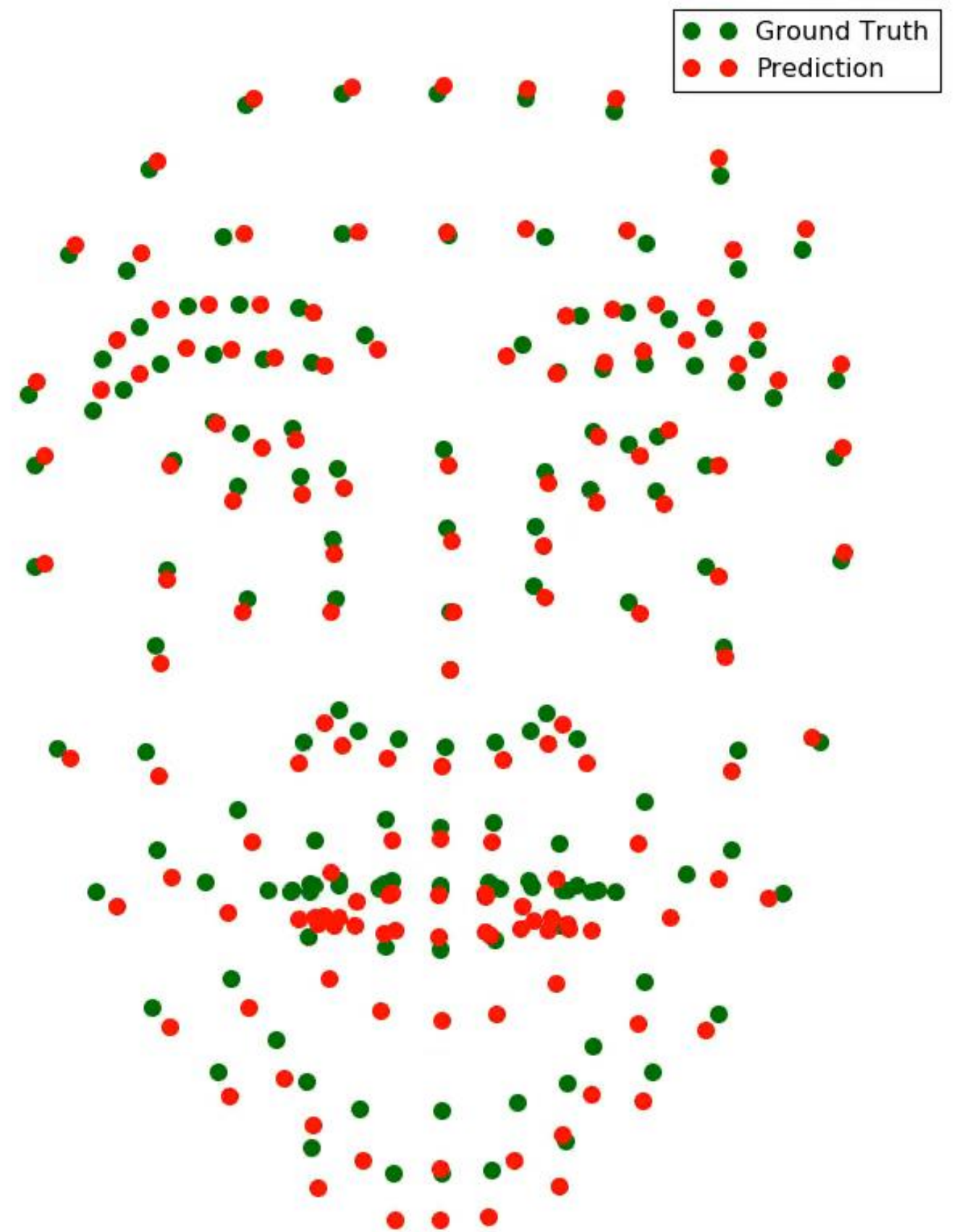
Can we learn this mapping using Machine Learning?



Convolutional Neural Network



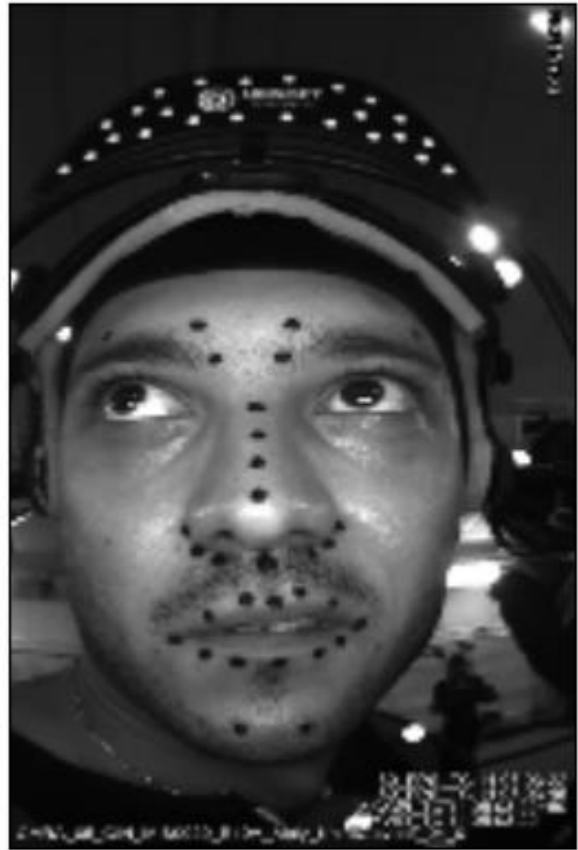




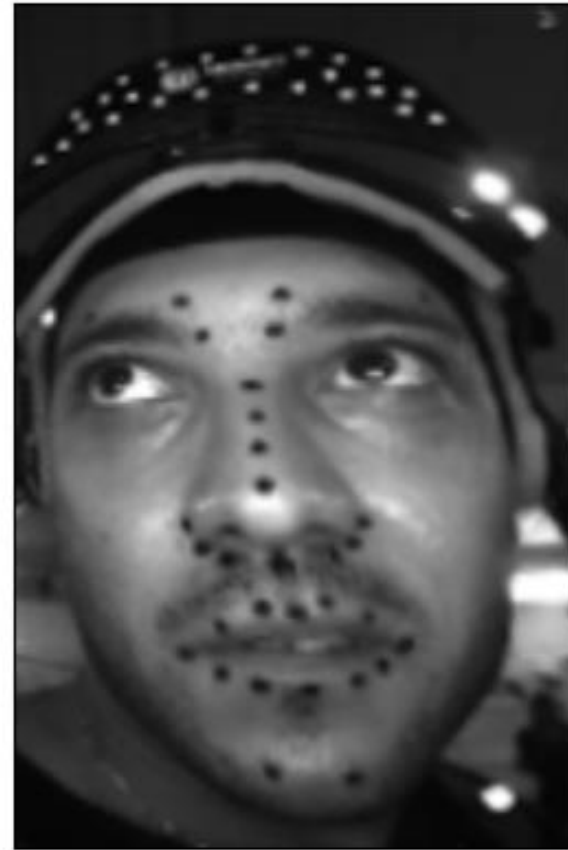
Data Augmentation

- It's impossible to capture every Actor in all lighting conditions.
- Can we do something to increase coverage in our data?

We can try to emulate different lighting conditions and actors.



Noise



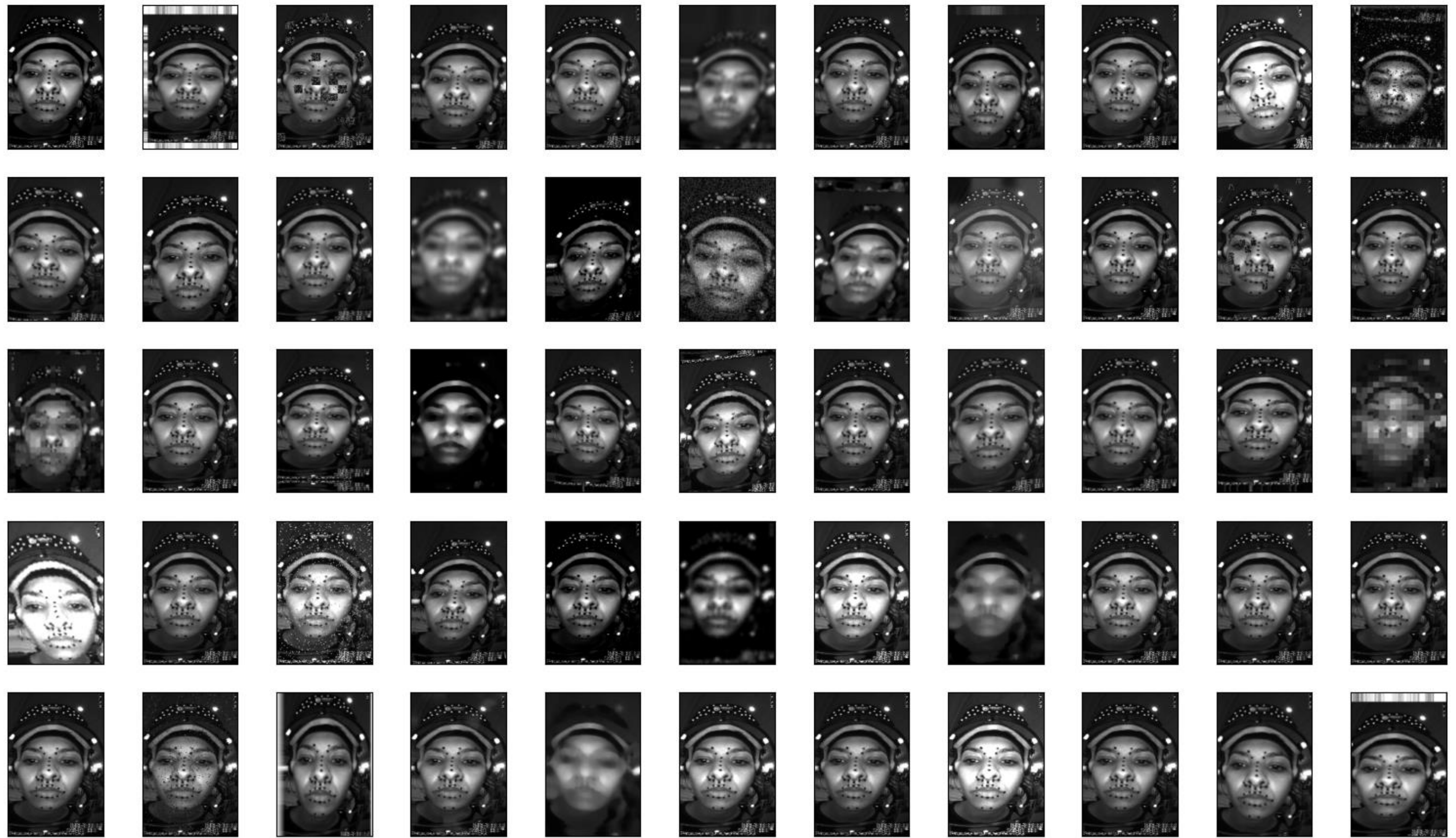
Distortion

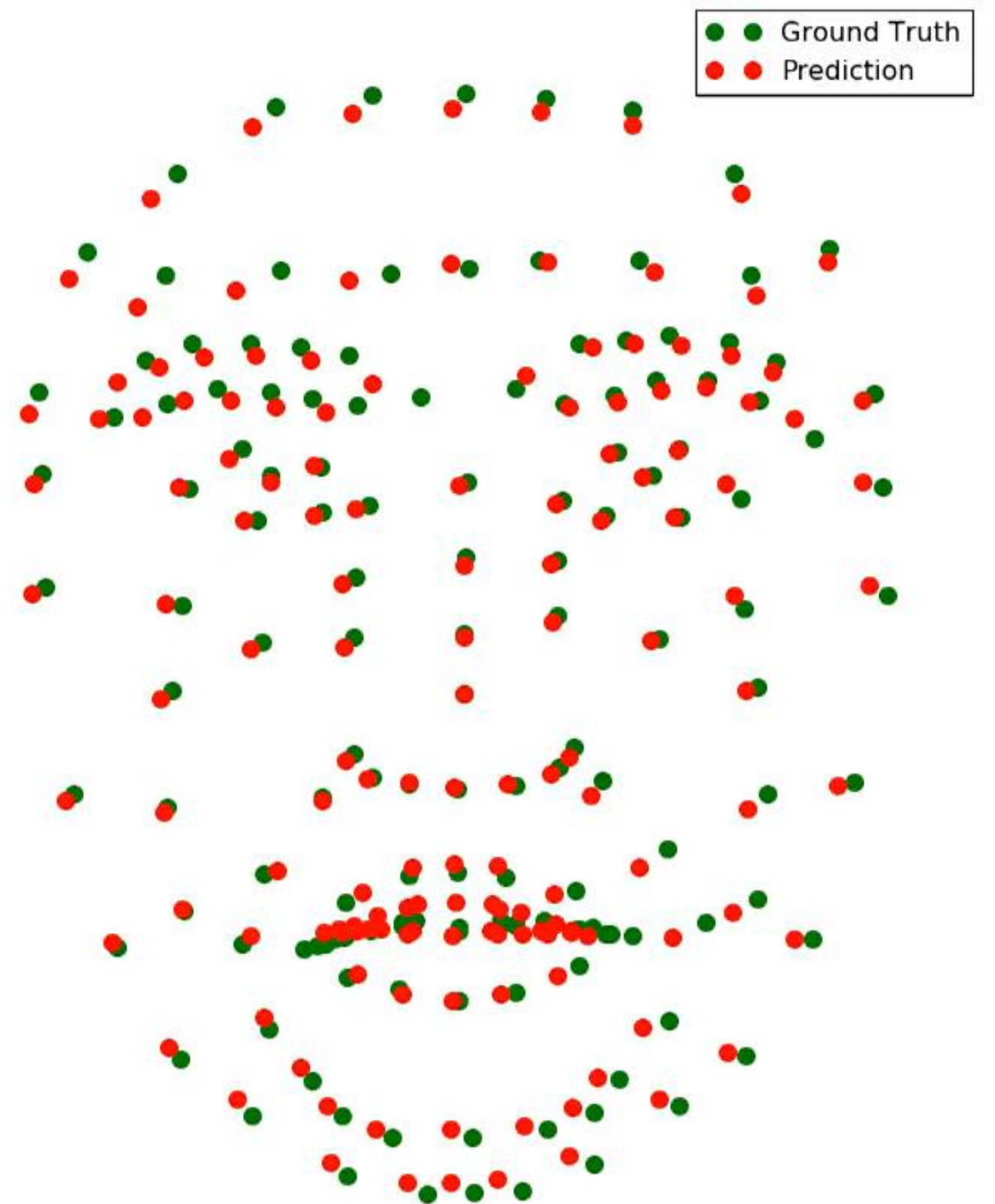


Perspective



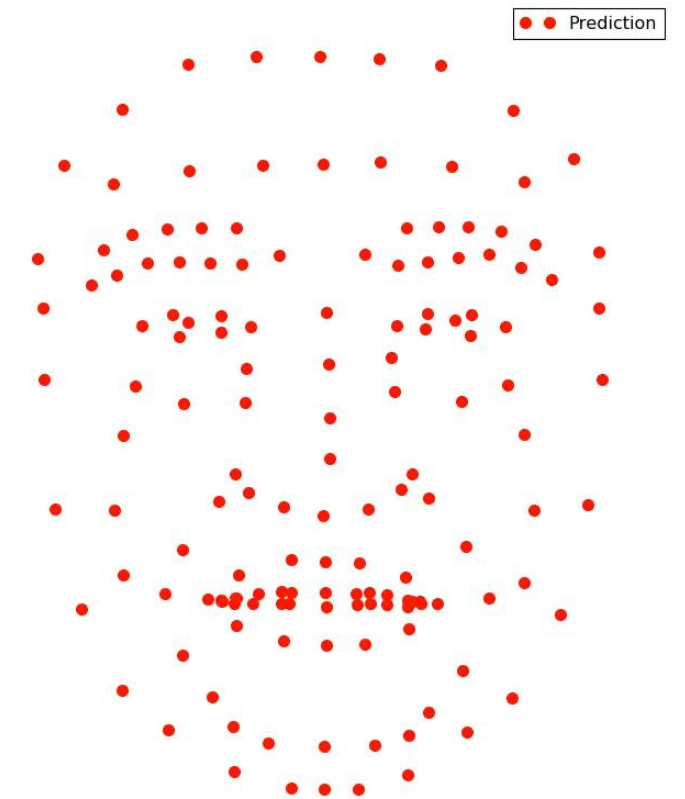
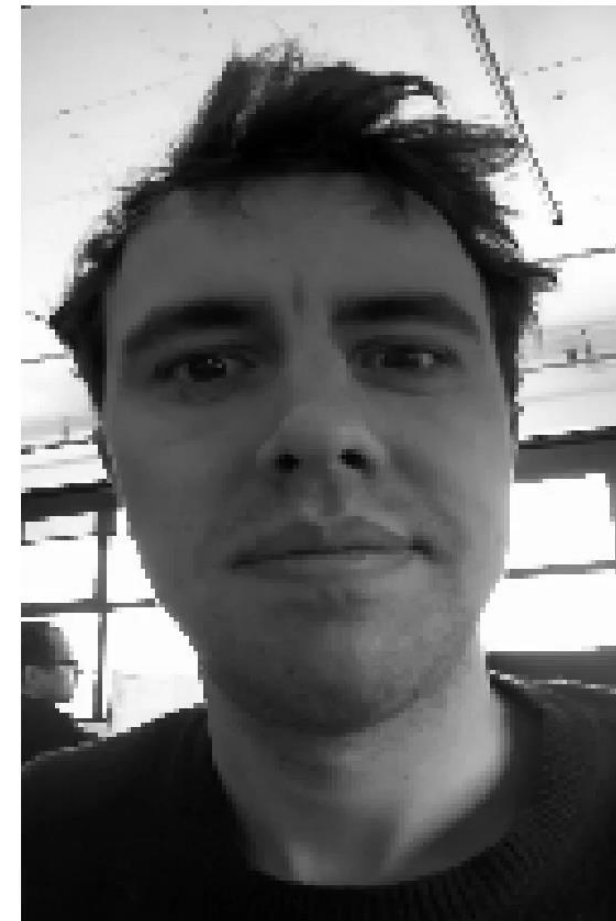
Contrast



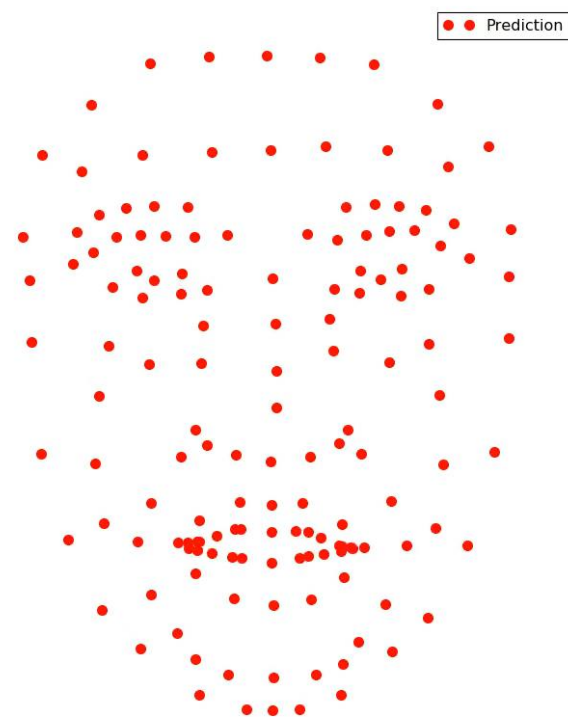


Webcam Capture

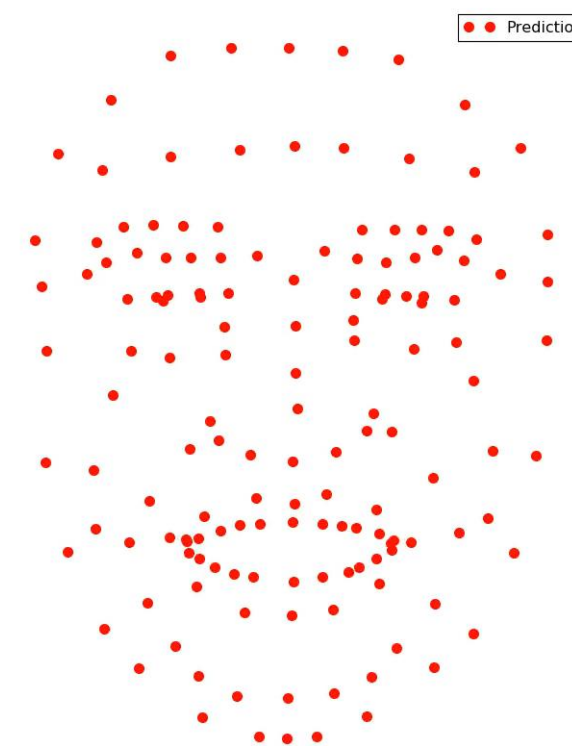
- Even very different capture conditions work to some degree.
- For example using a webcam as input.



Webcam Capture Limitations

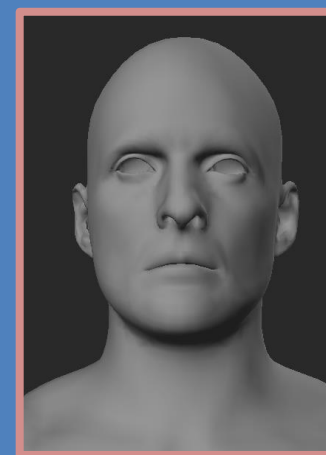
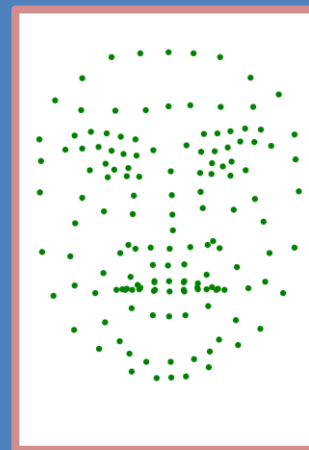
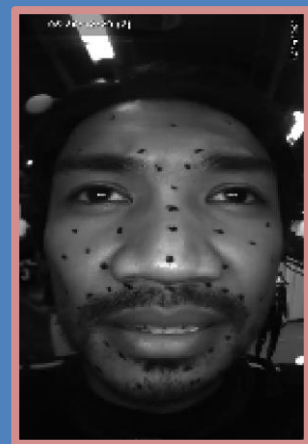


Complicated Expressions



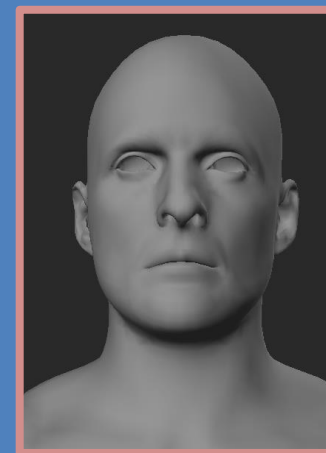
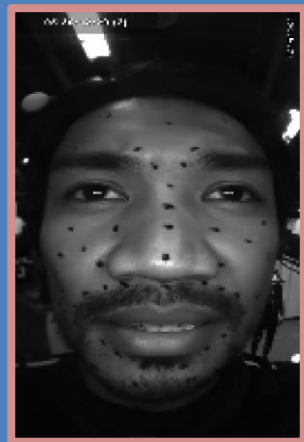
Beards / Glasses

Facial Capture Pipeline



What if we could automate either of these stages of the pipeline?

Facial Capture Pipeline

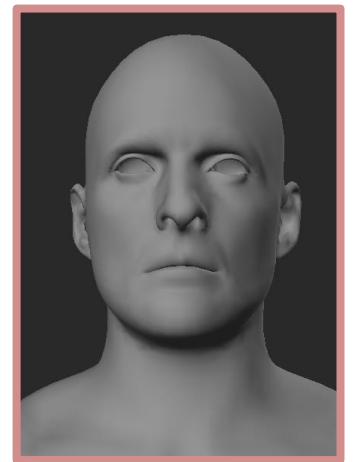
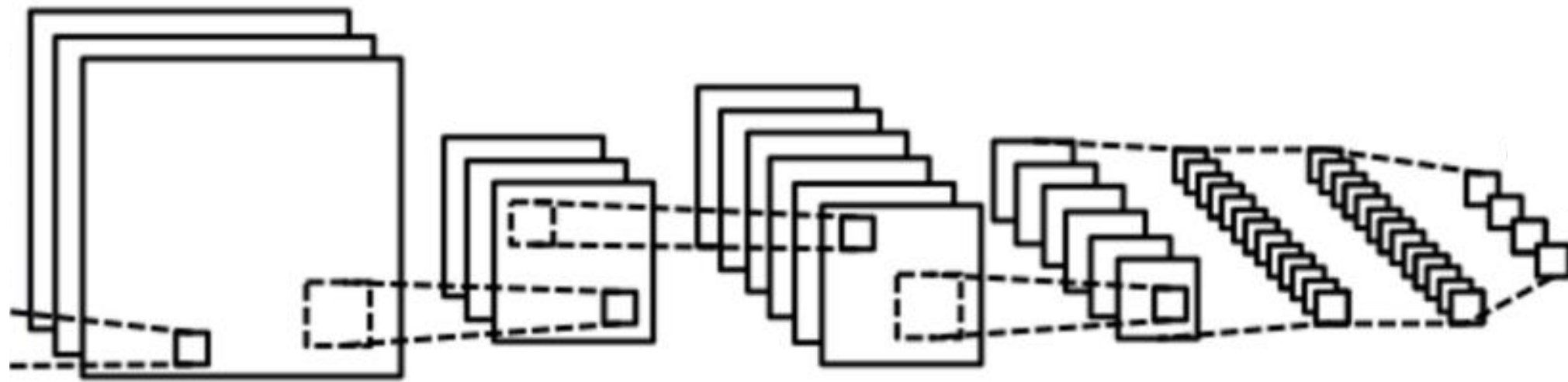


Fully Automatic Capture

Convolutional Neural Network



Frame of
Video



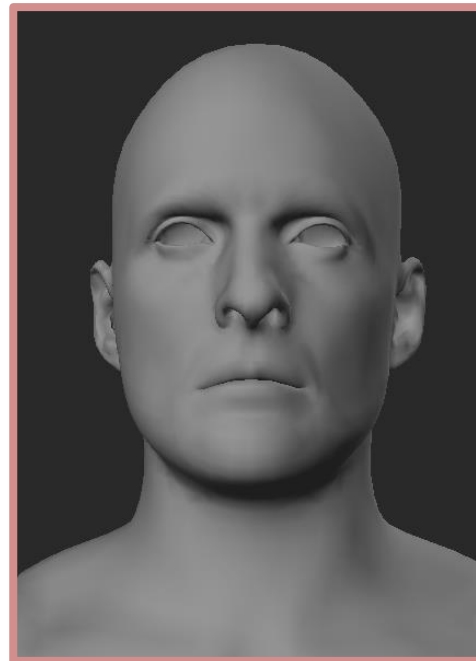
Values of
Facial Rig
Controls

Setup is practically identical...

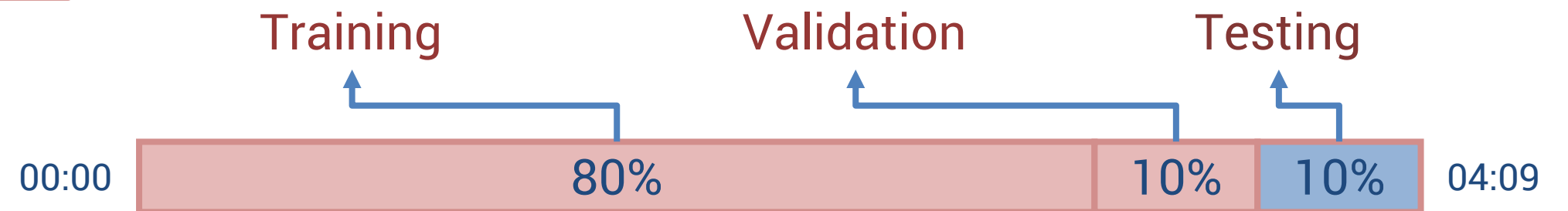
We build an experiment
with one short clip of
training data

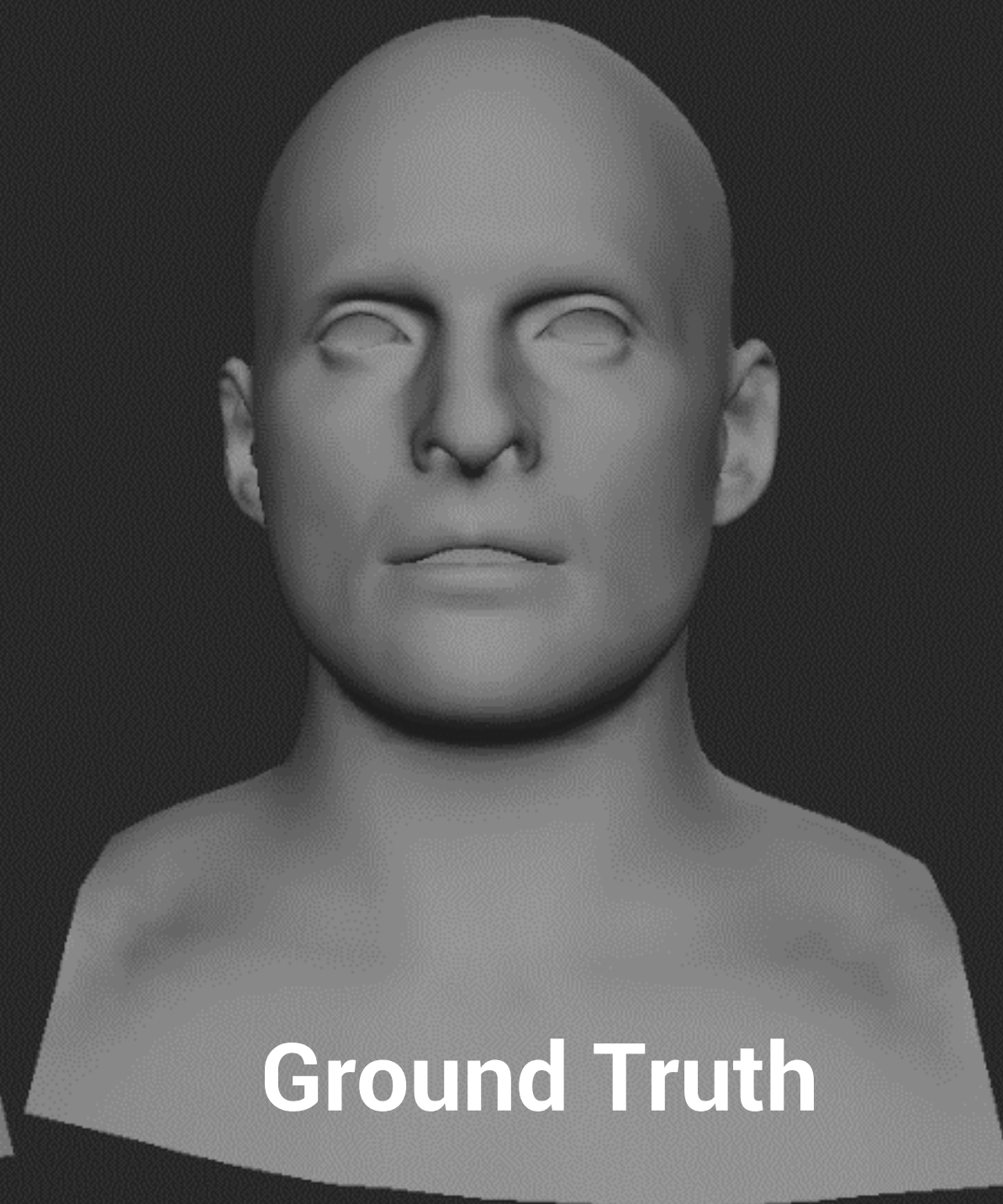
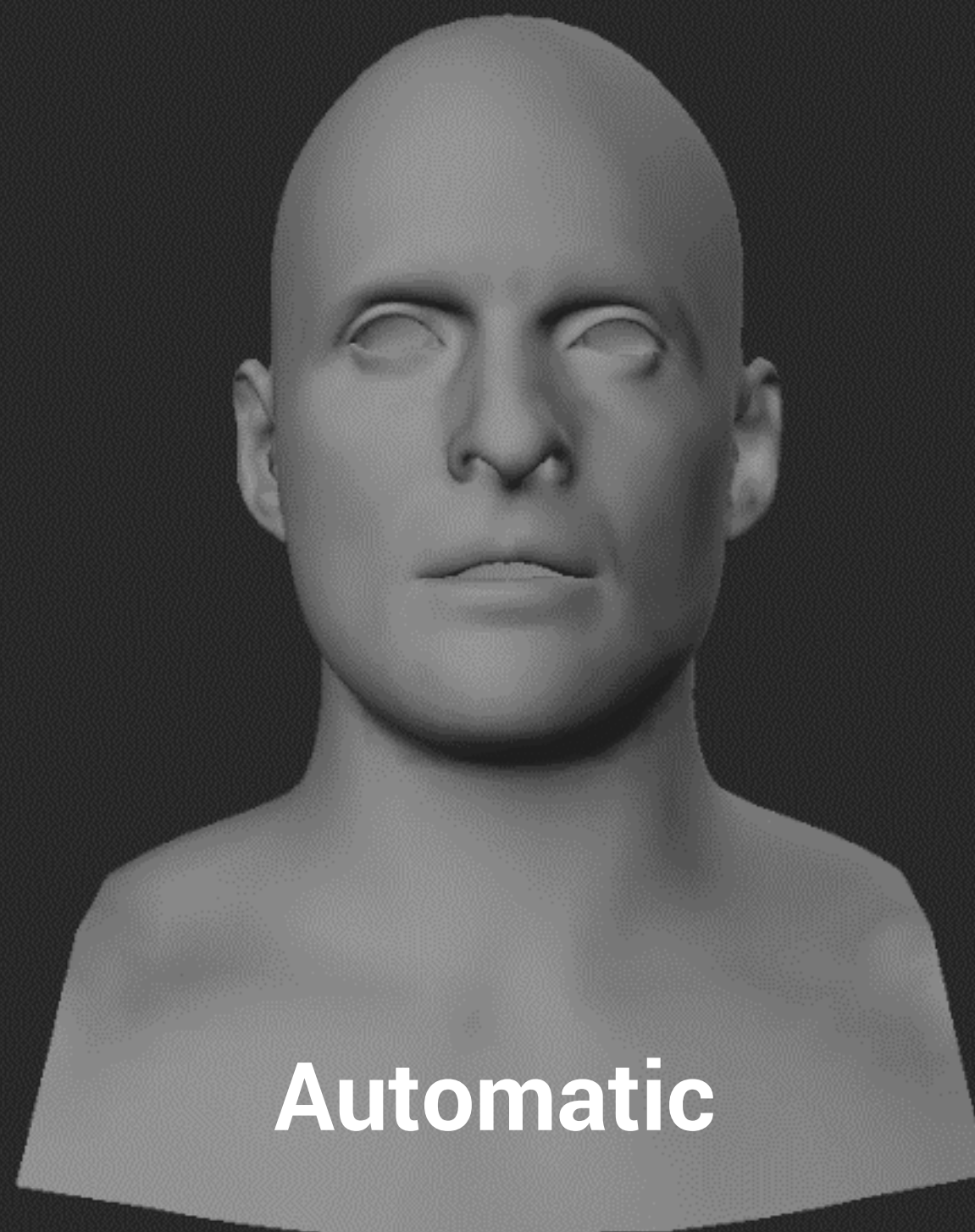


+



Data





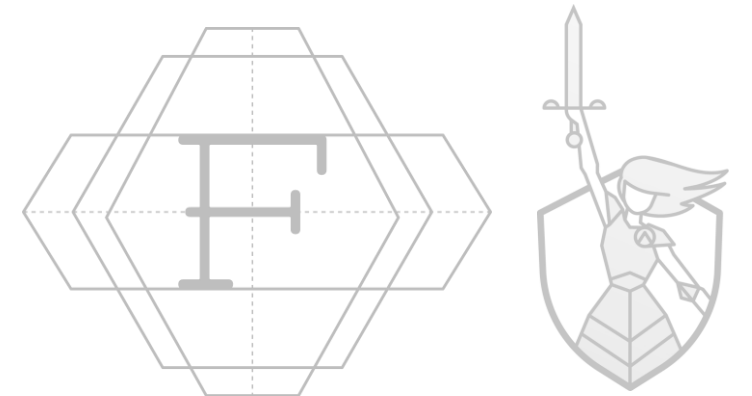
Summary

- CNNs can be applied to parts of the facial capture pipeline.
- Data Augmentation helps us generalize beyond training set.
- Initial promise shown for a fully automatic capture pipeline.

History

Mocap Cleaning
Facial Tracking
Audio to Facial

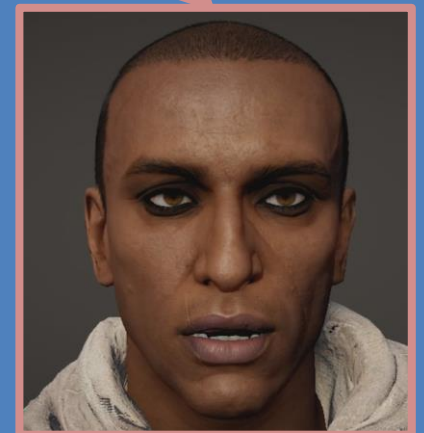
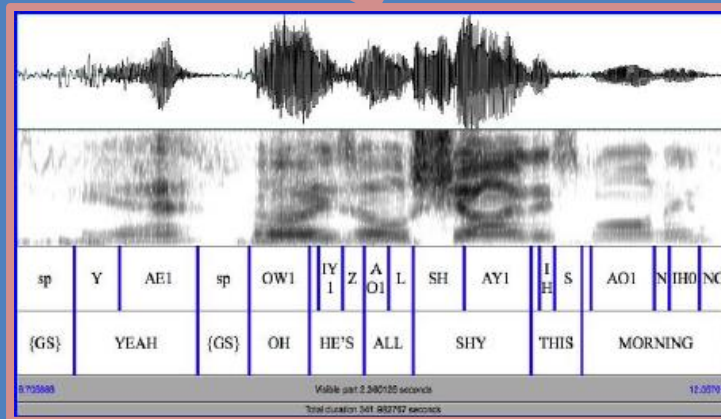
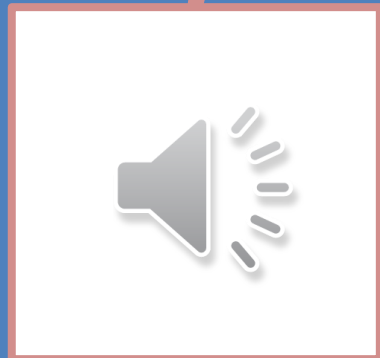
The Future



Audio to Facial Pipeline

Phonetic
Transcription

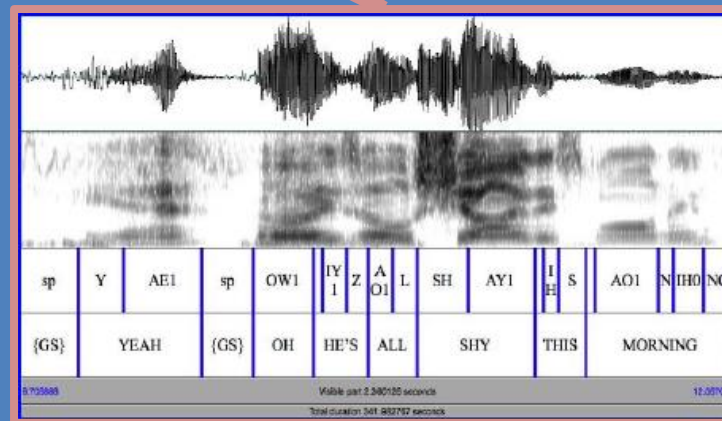
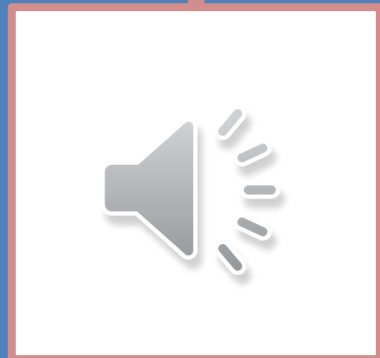
Rules based
Animation



Audio to Facial Pipeline

Phonetic
Transcription

English
Animation

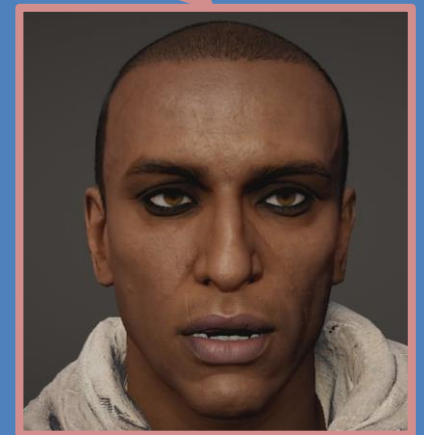
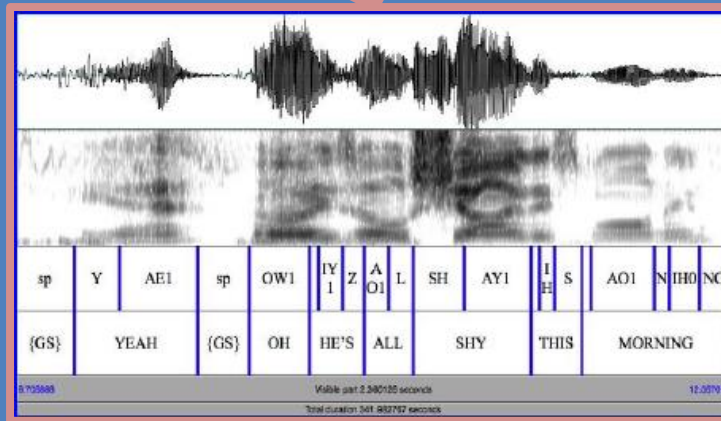


*Spanish
Audio*

Audio to Facial Pipeline

Phonetic
Transcription

Hand Made
Animation

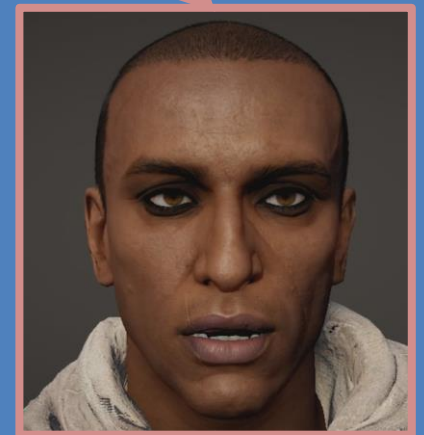
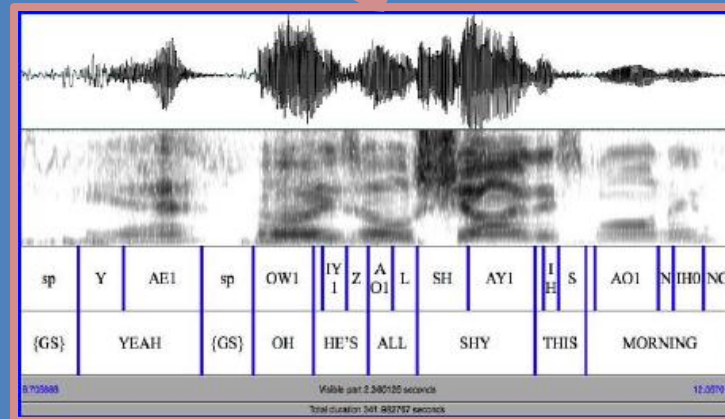
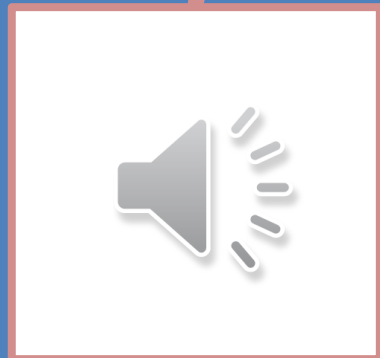


*Grunts and
Screams*

Audio to Facial Pipeline

Phonetic
Transcription

Rules based
Animation

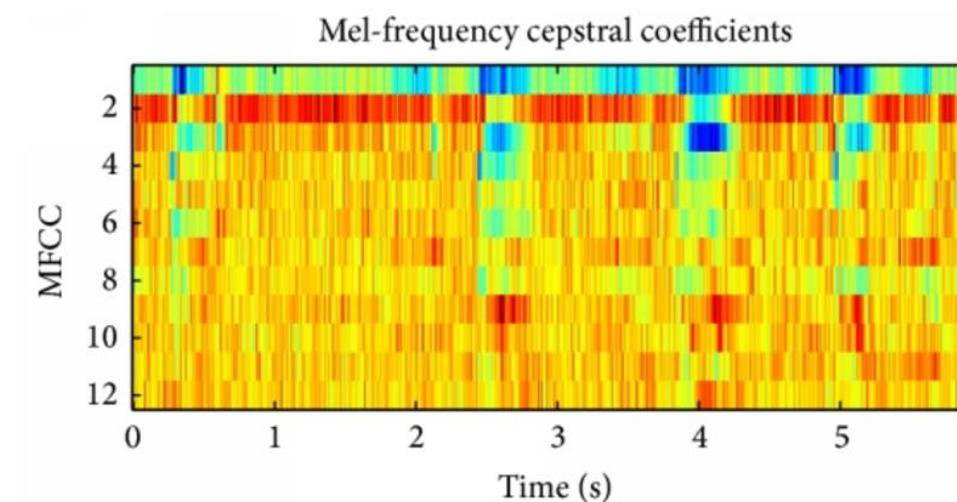


Sound Matching

*What if we could produce
facial animation directly?*

Input / Output

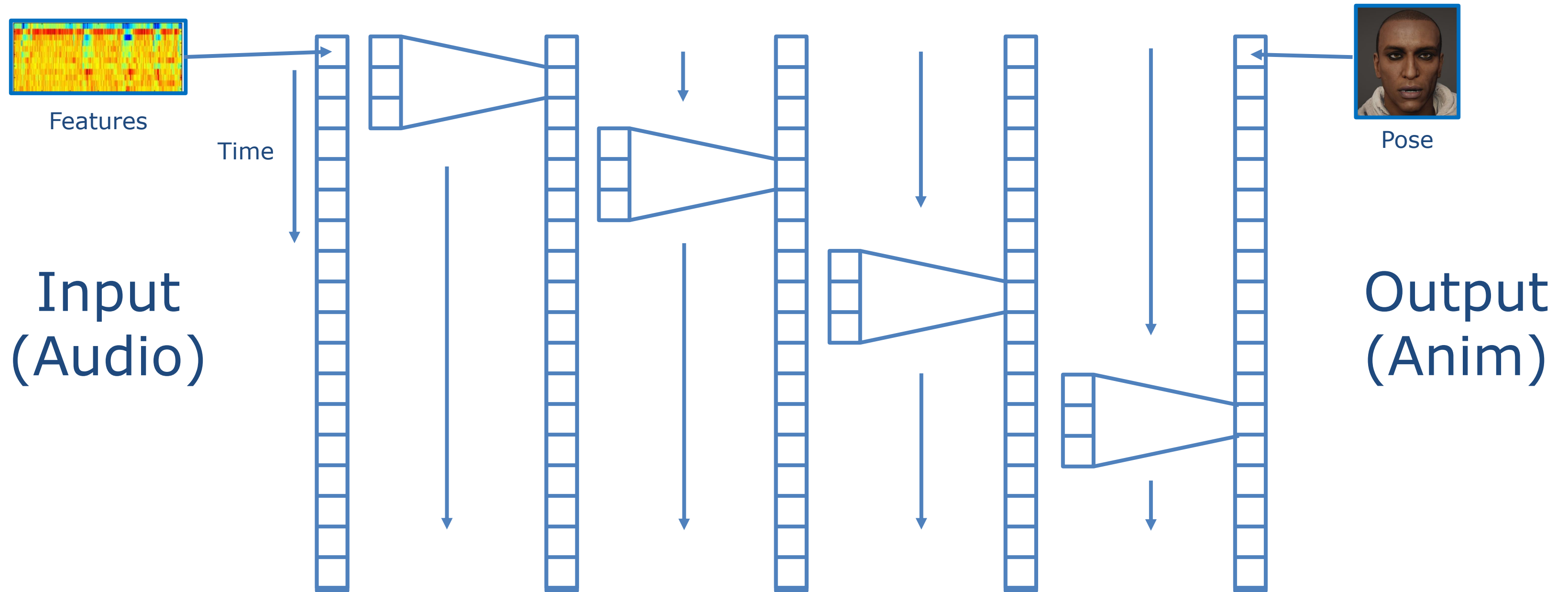
- **Input:** MFCC audio features.



- **Output:** Animation from rules based system.

English Only

Convolutional Neural Networks

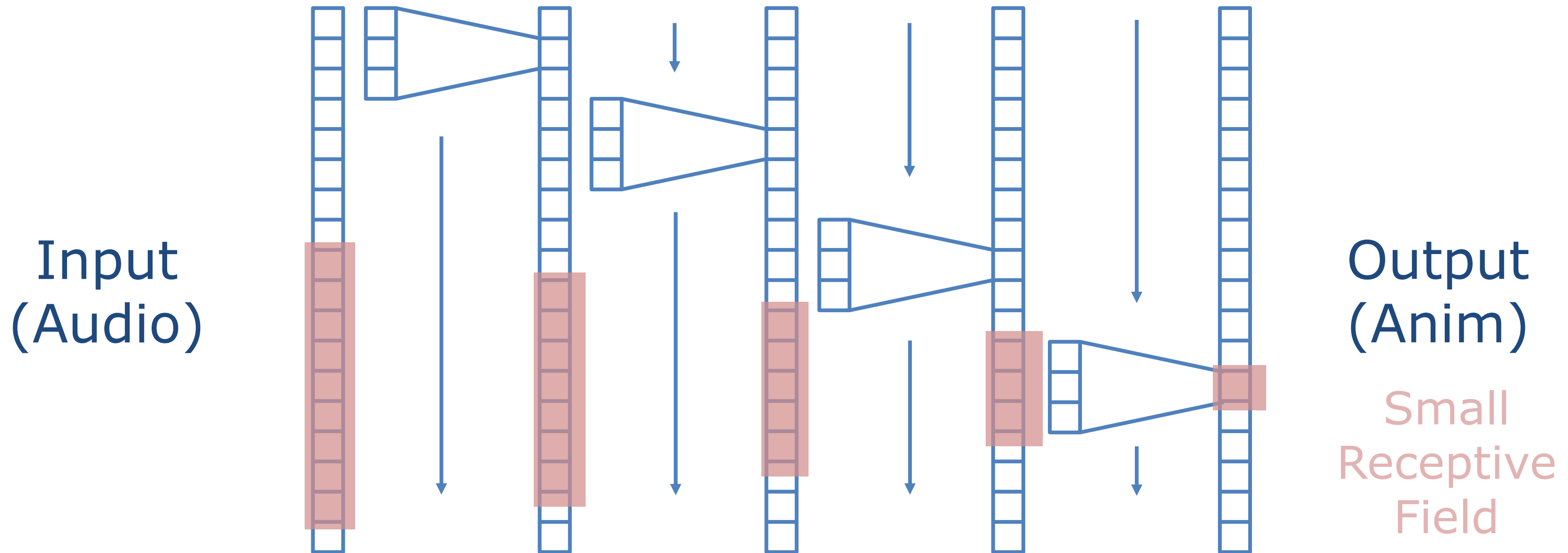


Convolutional Neural Networks

Problem:

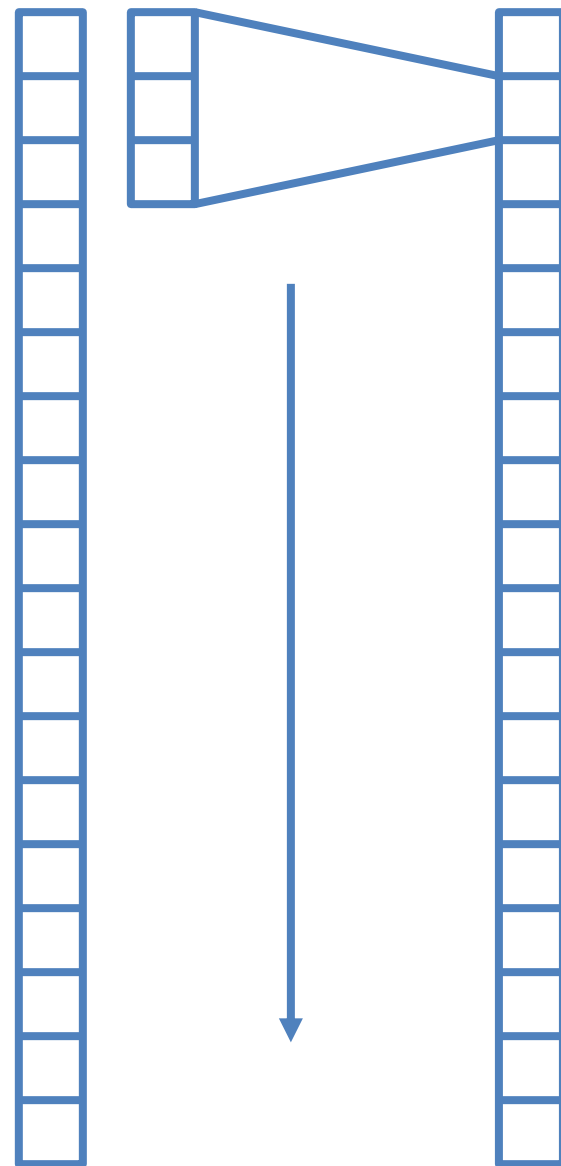
Output pose is only a function of a small window of the input features.

Convolutional Neural Networks

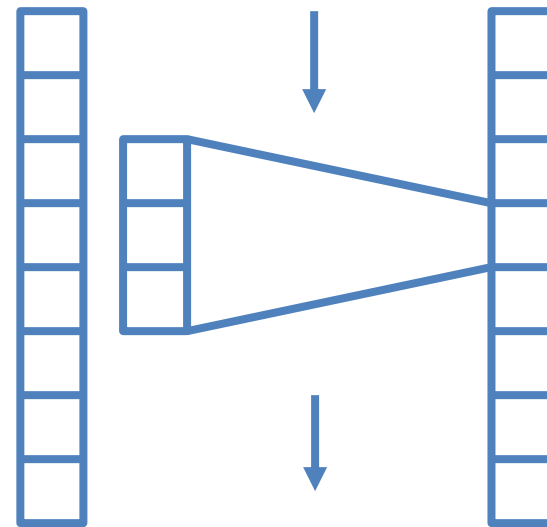


Pooling

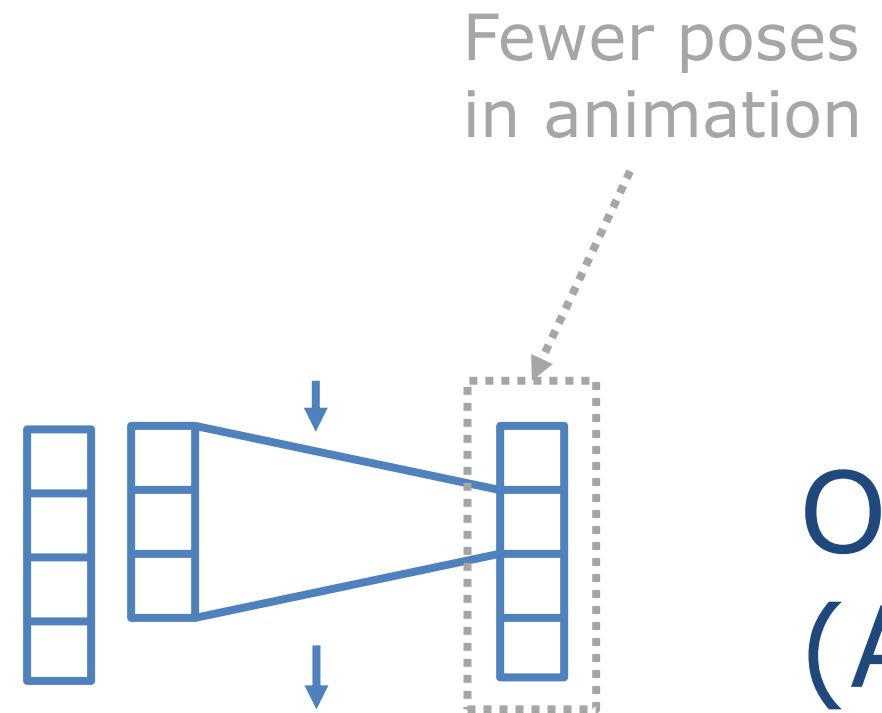
Input
(Audio)



Pool



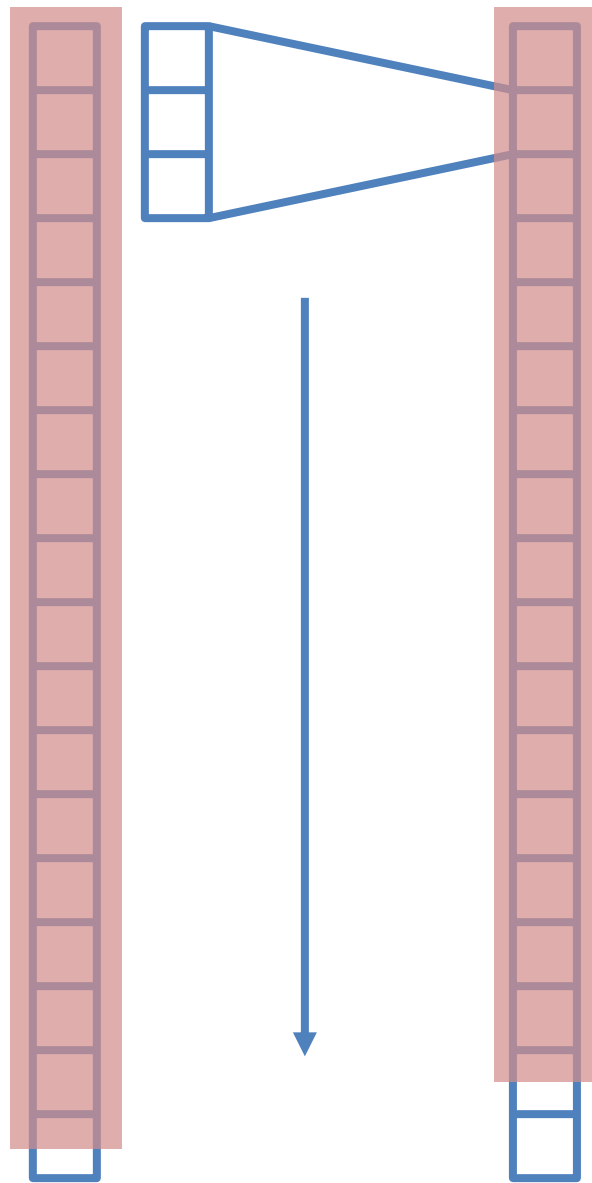
Pool



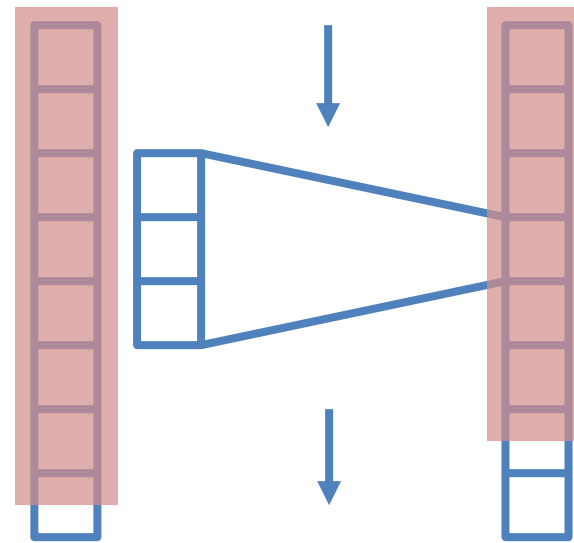
Output
(Anim)

Pooling

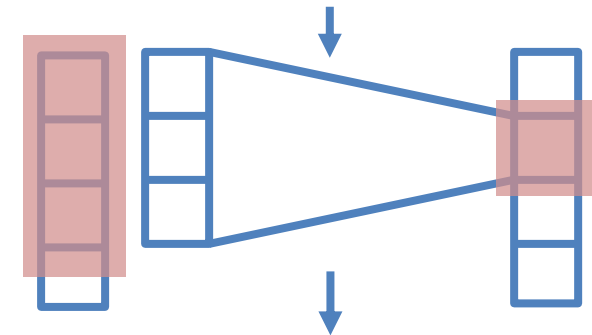
Input
(Audio)



Pool



Pool



Output
(Anim)

Large
Receptive
Field

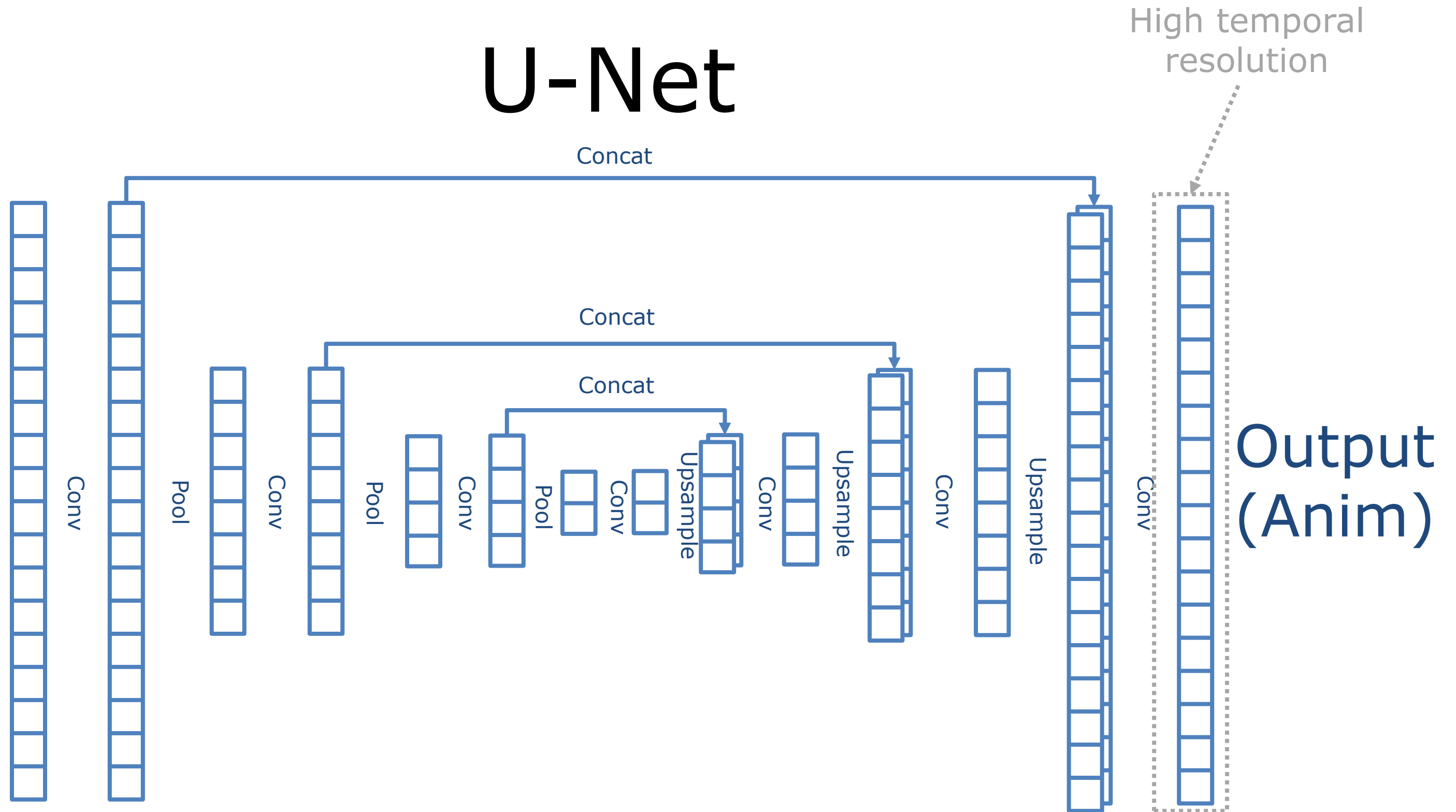
Pooling

Problem:

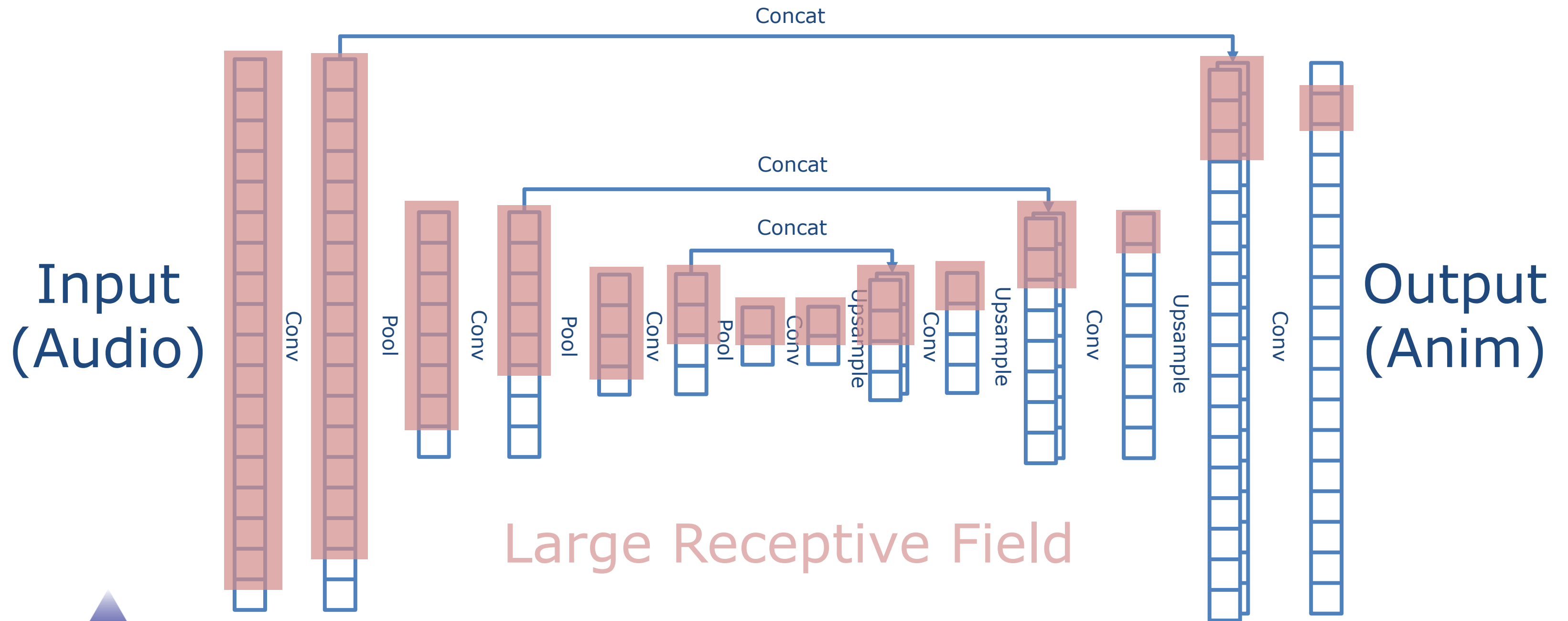
Output is low resolution. High resolution details are not preserved.

U-Net

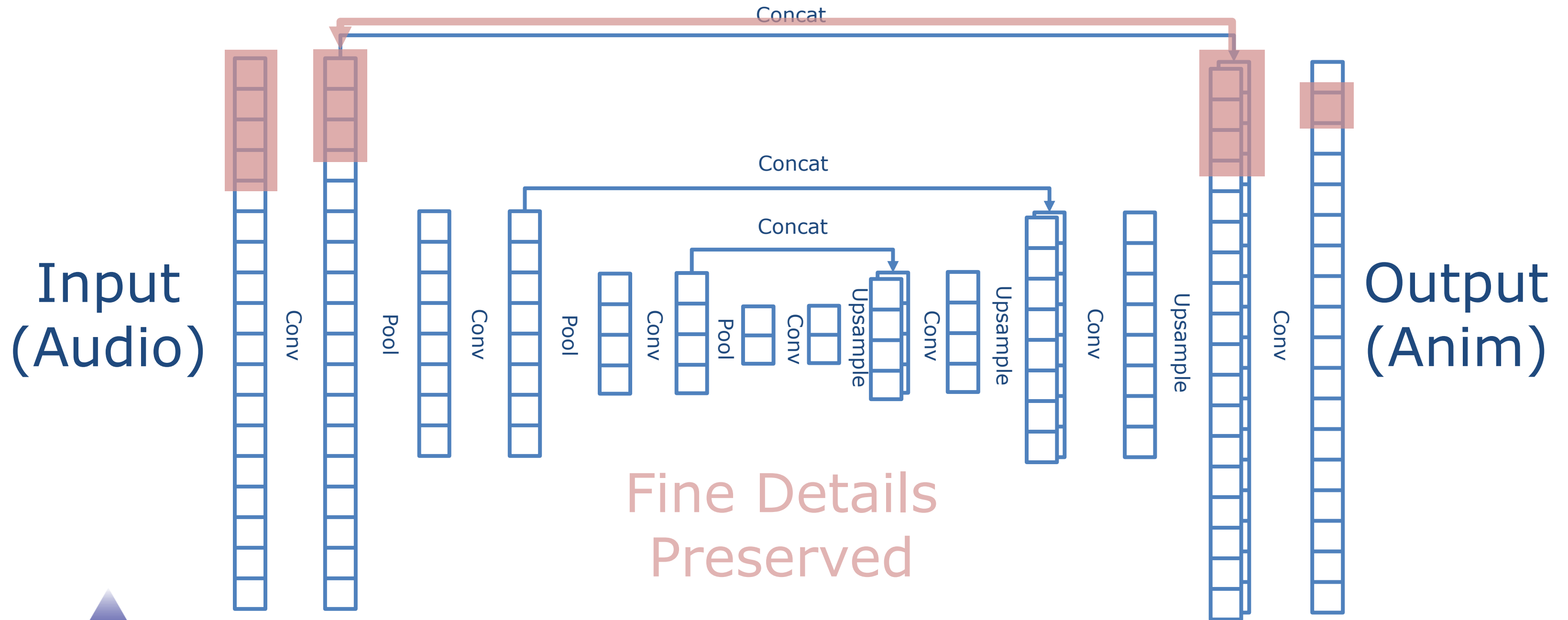
Input
(Audio)



U-Net



U-Net



U-Net

Solution:

Combination of **high resolution input** and **low resolution input** preserves details while maintaining a large receptive field.

Other Tricks

- Augment dataset with sped up and slowed down versions of the training data.
- Post-process output with sharpening filters controlled by the artist.

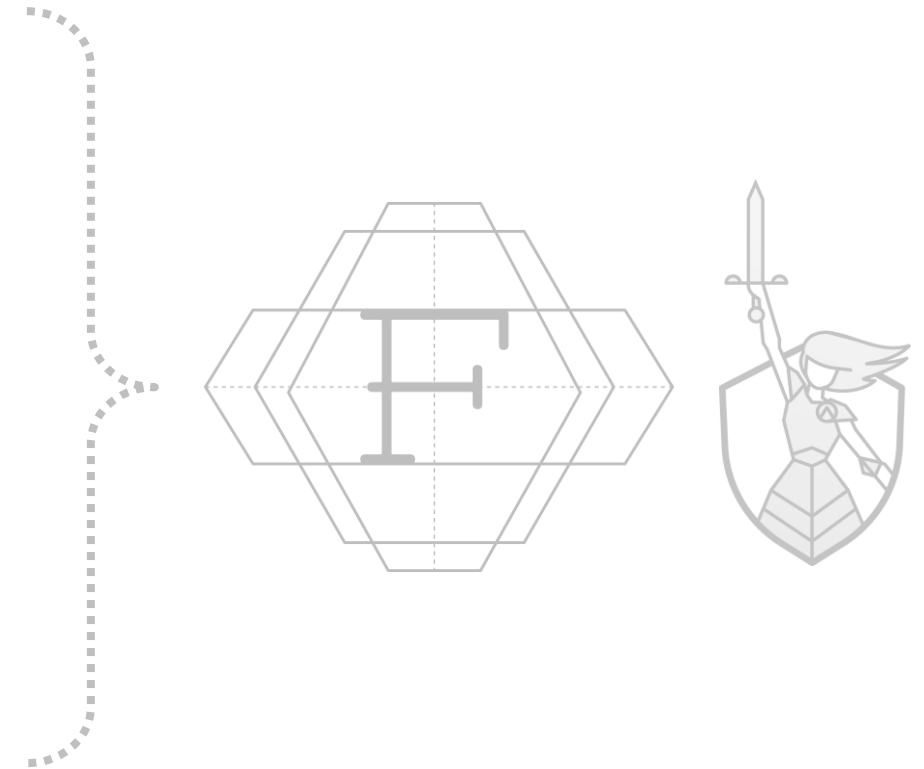
Summary

- CNNs can be used to avoid tricky parts of the pipeline.
- U-Nets are powerful for sequence to sequence tasks.
- Training data can be acquired using existing pipelines.

History

Mocap Cleaning
Facial Tracking
Audio to Facial

The Future





130 km²

The Future

- Bigger worlds means more content and more possibilities.
- Not all of this can be created by artists and designers.
- Performance capture enables artists to direct.

The Future

- There is still a huge way to go with making this tech practical and stable.
- No system can survive as a black box for long – engineers will have to learn the basics of data driven methods.
- There will always be a need to refine the output.

Conclusion

- *The Era of Machine Learning* is almost over!
- Prepare your pipelines for tech you see today.