# Render the Possibilities SIGGRAPH2016

THE 43RD INTERNATIONAL CONFERENCE AND EXHIBITION ON

Computer Graphics Interactive Techniques
24-28 JULY

ANAHEIM, CALIFORNIA

# Render the Possibilities SIGGRAPH2016

THE 43RD INTERNATIONAL CONFERENCE AND EXHIBITION ON

© Computer Graphics Interactive Techniques





A Deep Learning Framework for Character Motion Synthesis and Editing

Daniel Holden \*, Jun Saito †, Taku Komura \*,

\*The University of Edinburgh †Marza Animation Planet



#### Motivation

Synthesis

Editing

Discussion



Data driven synthesis of motion from high level controls with no manual preprocessing



• Lots of manual processing (Graphs, Trees)

- Segmentation
- Alignment
- Classification



[Heck et al. 2007]



[Kovar et al. 2002]

- Scalability Issues (RBF, GP, GPLVM, kNN)
  - Must store whole database in memory
  - Grows O(n<sup>2</sup>) with number of data points
  - Requires expensive acceleration structures



[Lee et al. 2010] [Park et al. 2002]



[Mukai and Kuriyama 2005]

- Instability issues (GPDM, CRBM, RNN)
  - Limited to some classes of motion
  - Can suffer high frequency noise or "dying out"



#### [Levine et al. 2012] [Wang et al. 2005]



[Taylor et al. 2011]

- Deep Neural Network *Hidden Units* used to represent motion
  - Denoising
  - Retrieval
  - Interpolation



[Holden et al. 2015]

- Deep Learning not always ready for production
- Results can look strange



[Radford et al. 2015]

## Contribution

- High quality synthesis with no manual preprocessing
- Motion synthesis and editing in unified framework
- Procedural, parallel technique















Motivation

Synthesis

Editing

Discussion



# **Convolutional Neural Networks**

- Great success in classification and segmentation for images, video, sound
- We can use CNN on motion data too















# What Happened?

- **Ambiguity**: the same control signal maps to multiple motions
- These motions are averaged in the output



## Foot Contact

- Contact times resolve ambiguity
- Automatically label using foot speed and height
- Learn model that generates contact times from trajectory



### Foot Contact

- Use small neural network to map trajectories to contact durations and frequencies
- Produce timings from durations and frequencies









#### Motivation

Synthesis

Editing

Discussion

# **Motion Editing**

### Once motion is generated it must be edited



# **Motion Editing**

### Post processing may not ensure naturalness



# **Motion Editing**

 We edit using the motion manifold learned by a Convolutional Autoencoding Network [Holden et al. 2015]



**Motion Data** 

**Hidden Units** 



### Autoencoder

• Learns projection operator of motion manifold



# Manifold Surface

- Hidden Unit values parametrise manifold surface
- Adjusting them ensures motion remains natural



# **Constraint Satisfaction**

• Motion editing is a *constraint satisfaction* problem over Hidden Units



## **Constraint Satisfaction**

Local foot velocity must equal global velocity

$$Pos(\mathbf{H}) = \sum_{j} \|\mathbf{v}_{r}^{\mathbf{H}} + \boldsymbol{\omega}^{\mathbf{H}} \times \mathbf{p}_{j}^{\mathbf{H}} + \mathbf{v}_{j}^{\mathbf{H}} - \mathbf{v}_{j}'\|_{2}^{2}.$$

• Output trajectory must equal input trajectory

$$Traj(\mathbf{H}) = \|\omega^{\mathbf{H}} - \omega'\|_{2}^{2} + \|\mathbf{v}_{r}^{\mathbf{H}} - \mathbf{v}_{r}'\|_{2}^{2}$$

















## A Neural Algorithm of Artistic Style

 Combine style of one image with content of another [Gatys et al. 2015]



- Gram Matrix of *Hidden Units* encode style
- Actual Values of Hidden Units encode content

 $Style(\mathbf{H}) = s \|G(\mathbf{\Phi}(\mathbf{S})) - G(\mathbf{H})\|_2^2 + c \|\mathbf{\Phi}(\mathbf{C}) - \mathbf{H}\|_2^2$ 

$$G(\mathbf{H}) = \frac{\sum_{i}^{n} \mathbf{H}_{i} \mathbf{H}_{i}^{T}}{n}$$

- Gram Matrix of *Hidden Units* encode style
- Actual Values of Hidden Units encode content

 $Style(\mathbf{H}) = s \|G(\Phi(\mathbf{S})) - G(\mathbf{H})\|_{2}^{2} + c \|\Phi(\mathbf{C}) - \mathbf{H}\|_{2}^{2}$ Content Term $G(\mathbf{H}) = \frac{\sum_{i}^{n} \mathbf{H}_{i} \mathbf{H}_{i}^{T}}{n}$ 

- Gram Matrix of *Hidden Units* encode style
- Actual Values of Hidden Units encode content

 $Style(\mathbf{H}) = s \|G(\Phi(\mathbf{S})) - G(\mathbf{H})\|_{2}^{2} + c \|\Phi(\mathbf{C}) - \mathbf{H}\|_{2}^{2}$ Style Term  $G(\mathbf{H}) = \frac{\sum_{i}^{n} \mathbf{H}_{i} \mathbf{H}_{i}^{T}}{n}$ 

• Gram Matrix of *Hidden Units* encode style

 $G(\mathbf{H}) = \frac{\sum_{i}^{n} \mathbf{H}_{i} \mathbf{H}_{i}^{T}}{\mathbf{H}_{i}}$ 

• Actual Values of Hidden Units encode content

 $Style(\mathbf{H}) = s \|G(\mathbf{\Phi}(\mathbf{S})) - G(\mathbf{H})\|_2^2 + c \|\mathbf{\Phi}(\mathbf{C}) - \mathbf{H}\|_2^2$ 

- Gram Matrix of *Hidden Units* encode style
- Actual Values of Hidden Units encode content

 $Style(\mathbf{H}) = s \|G(\mathbf{\Phi}(\mathbf{S})) - G(\mathbf{H})\|_2^2 + c \|\mathbf{\Phi}(\mathbf{C}) - \mathbf{H}\|_2^2$ 

$$G(\mathbf{H}) = \frac{\sum_{i}^{n} \mathbf{H}_{i} \mathbf{H}_{i}^{T}}{n}$$





#### Transfer









#### Transfer





#### Motivation

Synthesis

Editing

Discussion

### Representation

• We use joint positions local to hip

 Global rotation / translation removed using hips and shoulders direction

• Root velocity and contact times appended to representation



# Training

- Motion Manifold
  - Several large databases (including whole CMU)
  - Training takes around 6 hours

- Motion Synthesis
  - Task specific data only (e.g. locomotion only)
  - Training takes around 1 hour

# **Procedural vs Simulated**

### Procedural

- Output computed at arbitrary times or in parallel
- Ideal for precise animation

### Simulated

- Output computed frame by frame in series
- Ideal for interactive applications

## Performance

- Using GPU system runs in parallel over frames
- Very fast for long motions or many characters

Task	Duration	Foot Contacts	Synthesis	Editing	Total	FPS
Walking	60s	0.025s	0.067s	1.096s	1.188s	3030
Running	60s	0.031s	0.073s	1.110s	1.214s	2965
Punching	48	-	0.019s	0.259s	0.278s	863
Kicking	4s	-	0.020s	0.302s	0.322s	745
Style Transfer	8s	-	-	2.234s	2.234s	214
Crowd Scene	10s	0.557s	1.335s	2.252s	4.144s	28957

Figure 13: Performance breakdown.

## Performance

- Using GPU system runs in parallel over frames
- Very fast for long motions or many characters

Task	Duration	Foot Contacts	Synthesis	Editing	Total	FPS
Walking	60s	0.025s	0.067s	1.096 s	1.188s	3030
Running	60s	0.031s	0.073s	1.11 <i>)</i> s	1.214s	2965
Punching	48	-	0.019s	0.25 9s	0.278s	863
Kicking	4s	-	0.020s	0.30 2s	0.322s	745
Style Transfer	8s	-	-	2.234 s	2.234s	214
Crowd Scene	10s	0.557s	1.335s	2.252	4.144s	28957

Figure 13: Performance breakdown.



## Future Work

- Need more general solution for ambiguity issue
- Wish to use more high level features with a deeper network
- What changes are required for interactive applications?

