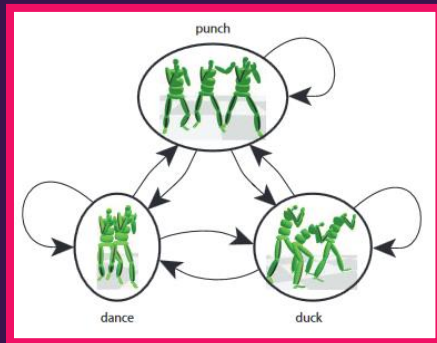# Outline

## Motivation

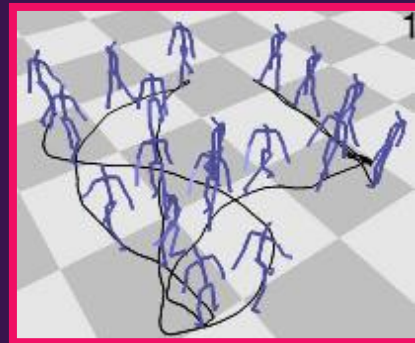## Synthesis

## Editing

## Discussion

# Goal

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

# Previous Work

- Lots of manual processing (Graphs, Trees)
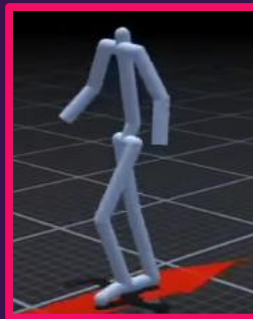
  - Segmentation
  - Alignment
  - Classification

[Heck et al. 2007]    [Kovar et al. 2002]

# Previous Work
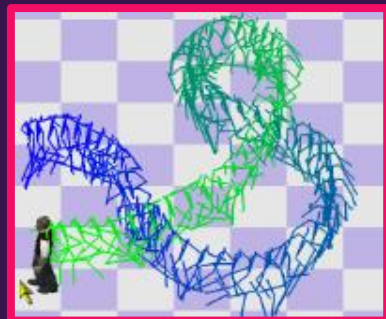
- Scalability Issues (RBF, GP, GPLVM, kNN)
  - Must store whole database in memory
  - Grows O(n²) with number of data points
  - Requires expensive acceleration structures



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



[Mukai and Kuriyama 2005]

# Overview

# Overview

Motion Editing     Motion Synthesis     Disambiguation

Control Parameters

Motion Data     Hidden Units

Footstep Timing

# Overview



Motion Editing
Motion Synthesis
Disambiguation

Motion Data
Hidden Units
Control Parameters
Footstep Timing

# Outline

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

# Convolution

*Filters* convolve over temporal dimension

# Convolution

*Filters* convolve over temporal dimension

# Convolution

*Filters* convolve over temporal dimension

# Convolution

*Filters* convolve over temporal dimension

# Convolution

*Filters* convolve over temporal dimension

# Outline

Motivation

Synthesis

Editing

Discussion

# 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



**Motion Data**        **Hidden Units**

# Constraint Satisfaction

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



Motion Data          Hidden Units

# Constraint Satisfaction

- Local foot velocity must equal global velocity

$$Pos(\mathbf{H}) = \sum_j \|\mathbf{v}_r^{\mathbf{H}} + \omega^{\mathbf{H}} \times \mathbf{p}_j^{\mathbf{H}} + \mathbf{v}_j^{\mathbf{H}} - \mathbf{v}_j'\|_2^2.$$

- Output trajectory must equal input trajectory

$$Traj(\mathbf{H}) \quad = \quad \|\omega^{\mathbf{H}} - \omega'\|_2^2 + \|\mathbf{v}_r^{\mathbf{H}} - \mathbf{v}_r'\|_2^2$$

# A Neural Algorithm of Artistic Style

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

# Style Constraint

- Gram Matrix of *Hidden Units* encode style

- Actual Values of *Hidden Units* encode content

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

$$G(\mathbf{H}) = \frac{\sum_i^n \mathbf{H}_i \mathbf{H}_i^T}{n}$$

- **No correspondence between clips required!**

# Style Constraint

- 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

$$G(\mathbf{H}) = \frac{\sum_i^n \mathbf{H}_i \mathbf{H}_i^T}{n}$$

- **No correspondence between clips required!**

# Style Constraint

- 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

$$G(\mathbf{H}) = \frac{\sum_i^n \mathbf{H}_i \mathbf{H}_i^T}{n}$$

- **No correspondence between clips required!**

# Style Constraint

- Gram Matrix of *Hidden Units* encode style

- Actual Values of *Hidden Units* encode content

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

$$G(\mathbf{H}) = \frac{\sum_i^n \mathbf{H}_i \mathbf{H}_i^T}{n}$$ Gram Matrix

- **No correspondence between clips required!**

# Style Constraint

- Gram Matrix of *Hidden Units* encode style

- Actual Values of *Hidden Units* encode content

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

$$G(\mathbf{H}) = \frac{\sum_i^n \mathbf{H}_i \mathbf{H}_i^T}{n}$$

- **No correspondence between clips required!**

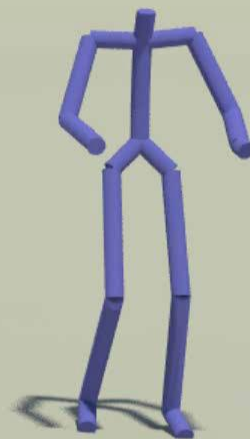# Style

Content                    Transfer

# Outline

Motivation

Synthesis

Editing

Discussion

# 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

# Contribution

- **High quality synthesis** with **no manual preprocessing**

- Motion synthesis and editing in **unified framework**

- **Procedural**, **parallel** technique

# 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?