

Temporal Smoothing in 2D Human Pose Estimation for Bouldering

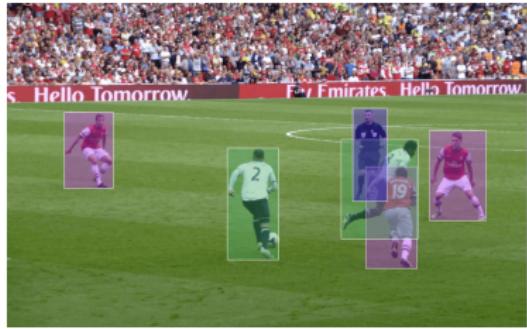
André Oskar Andersen
wpr684

Institution of Computer Science, University of Copenhagen

2023

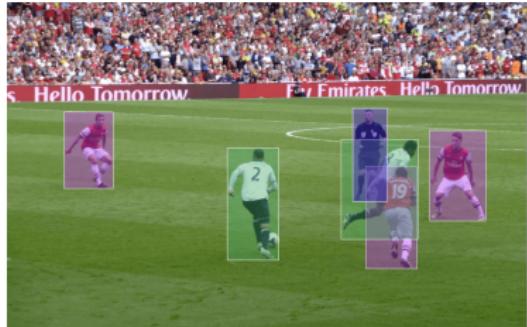
Introduction

- ▶ Increased usage of video analysis in sports.
 - ▶ Help referee
 - ▶ Improve techniques



Introduction

- ▶ Increased usage of video analysis in sports.
- ▶ Often requires the position of the players.
 - ▶ Already developed for popular sports.
 - ▶ Missing for the less popular sports.



Introduction

- ▶ Increased usage of video analysis in sports.
- ▶ Often requires the position of the players.
- ▶ Problems with the data
 - ▶ Methods require large quantities
 - ▶ Unusual poses/movements



Introduction

- ▶ Increased usage of video analysis in sports.
- ▶ Often requires the position of the players.
- ▶ Problems with the data
- ▶ ClimbAlong at NorthTech ApS
 - ▶ Frame-independent pose-detector for bouldering



Introduction

- ▶ Increased usage of video analysis in sports.
- ▶ Often requires the position of the players.
- ▶ Problems with the data
- ▶ ClimbAlong at NorthTech ApS
 - ▶ Frame-independent pose-detector for bouldering
 - ▶ Proposition: Incorporate temporal information



Introduction

- ▶ Aim: extend the ClimbAlong pose-detector to use temporal information.

The Data

ClimbAlong

- ▶ Fully annotated videos of climbers



The Data

ClimbAlong

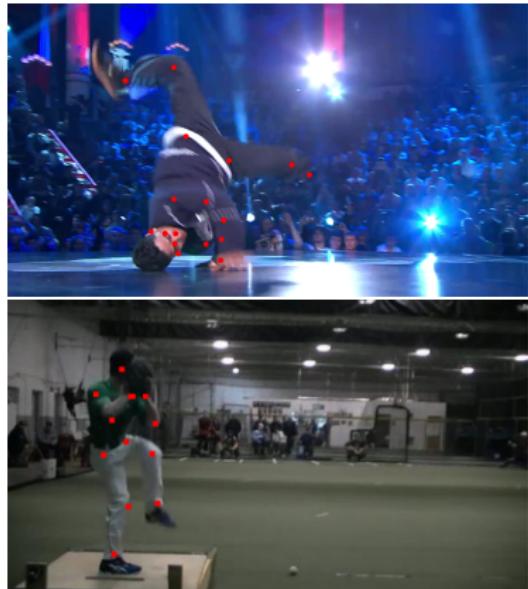
- ▶ Fully annotated videos of climbers
- ▶ Problem: very small dataset
 - ▶ BRACE
 - ▶ Penn Action



The Data

ClimbAlong

- ▶ Fully annotated videos of climbers
- ▶ Problem: Very small dataset
- ▶ Solution: pretrain on related datasets and finetune on ClimbAlong
 - ▶ BRACE
 - ▶ Penn Action



The Models

Generally, three approaches

1. Convolutional layer
2. Recurrent neural network
3. Transformer

The Models

3DConv

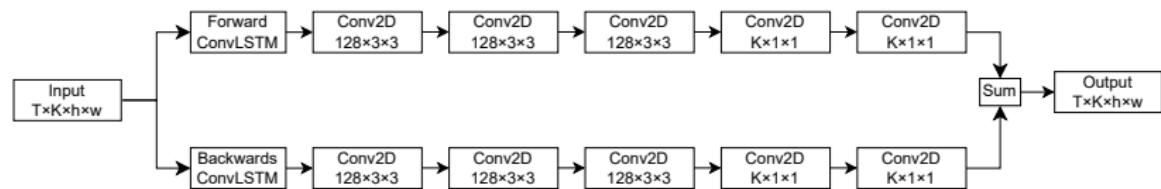
- ▶ 3-dimensional conv. layer + ReLU



The Models

bi-ConvLSTM Model S

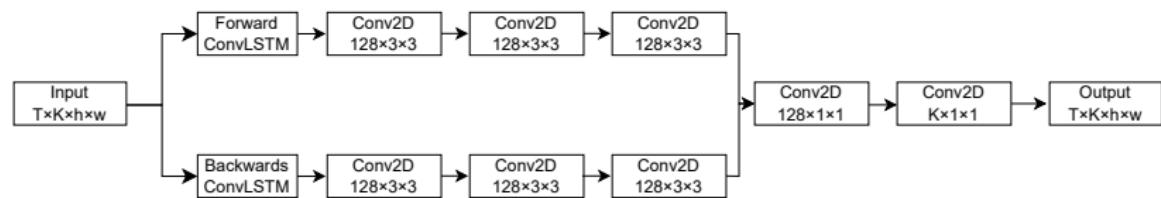
- ▶ Adaptation of Unipose-LSTM by Artacho and Savakis
- ▶ Bidirectional convolutional LSTM + conv. layers and ReLU
- ▶ Processing directions summed together



The Models

bi-ConvLSTM Model C

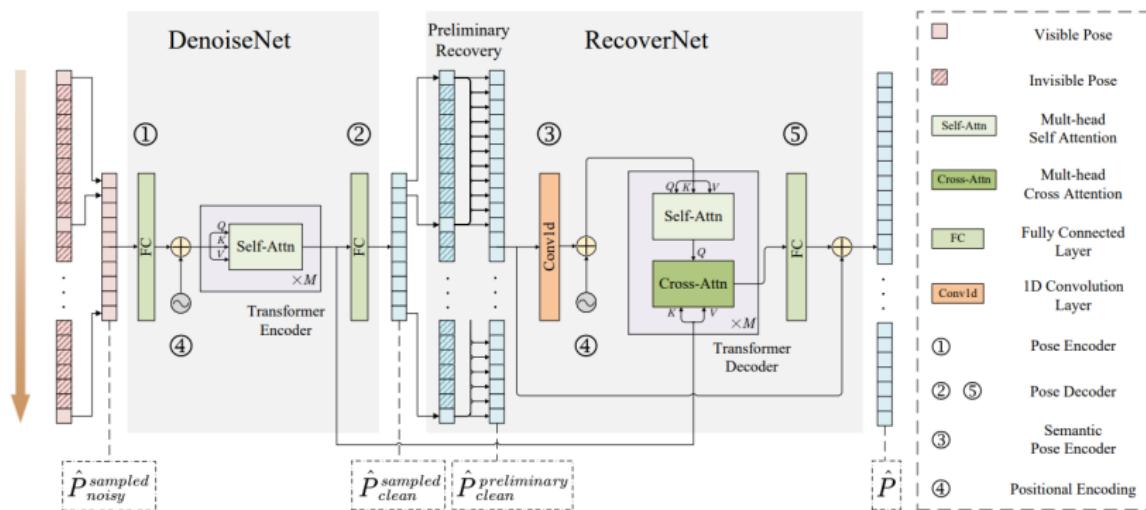
- ▶ Problem: Prioritization of processing direction
- ▶ Solution: Using convolution



The Models

DeciWatch by Zeng *Et al.*

- ▶ Transformer-based
- ▶ Samples every n th frame
- ▶ DenoiseNet + RecoverNet



Pretraining

Procedure

- ▶ Not training already-developed pose-detector
 - ▶ Different input images

Pretraining

Procedure

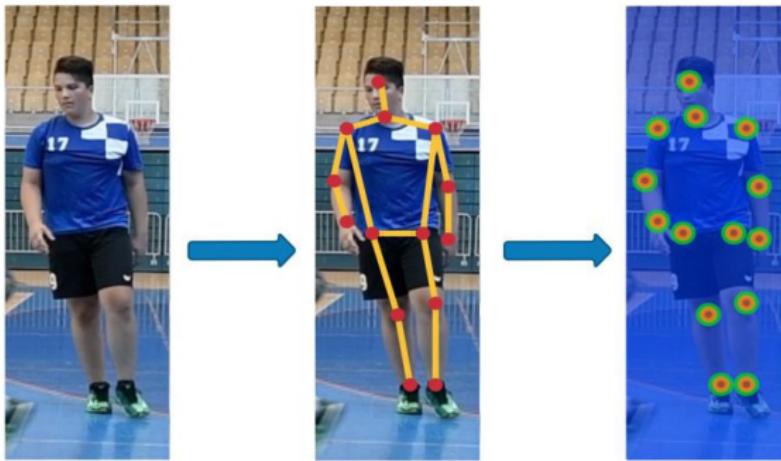
- ▶ Not training already-developed pose-detector
 - ▶ Different input images
- ▶ Instead, simulate pose-detector output by shifting keypoints

Pretraining

Finding optimal setting of models

- ▶ Three experiments

1. Various smearing standard deviations

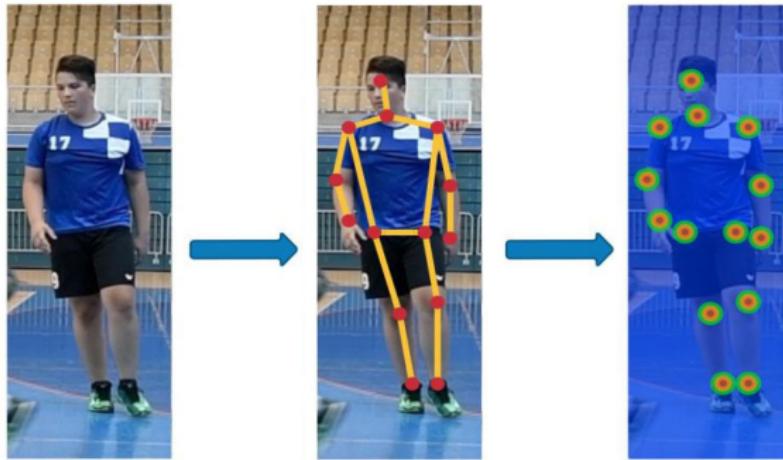


Pretraining

Finding optimal setting of models

► Three experiments

1. Various smearing standard deviations
2. Fixed smearing standard deviation

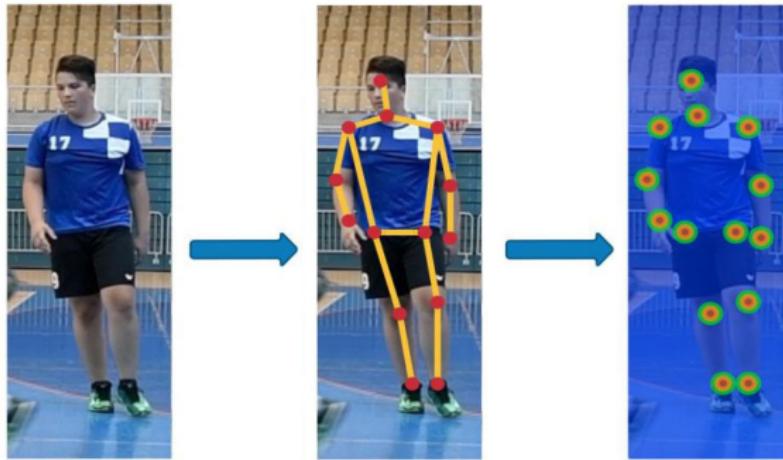


Pretraining

Finding optimal setting of models

► Three experiments

1. Various smearing standard deviations
2. Fixed smearing standard deviation
3. Smearing standard deviations + decreased frame rate



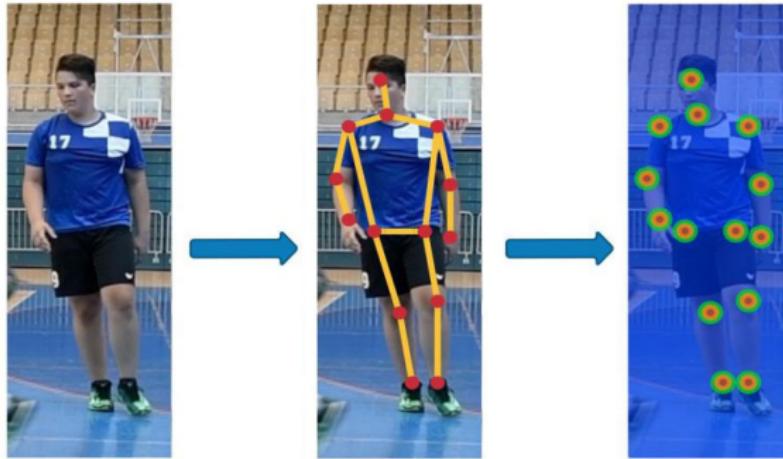
Pretraining

Finding optimal setting of models

- ▶ Three experiments

1. Various smearing standard deviations
2. Fixed smearing standard deviation
3. Smearing standard deviations + decreased frame rate

- ▶ Two different shifting-scalars



Finetuning

- ▶ Using all of the developed models with pose-detector
- ▶ Freezing pose-detector
 - 1. Quicker fitting
 - 2. Greater understanding of results

Finetuning

Test results

- ▶ Shifting-scalar: only minor effect

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
	1.1	1.2	1.3	1.1	1.2	1.3	1.1	1.2	1.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	49.7	52.3	53.1	95.7	95.7	95.8	99.2	99.3	99.3
DeciWatch	76.6	76.7	68.1	94.4	94.3	87.3	99.2	99.2	96.1
bi-ConvLSTM - Model S	37.8	34.9	39.0	91.8	92.1	92.2	99.4	99.7	99.2
bi-ConvLSTM - Model C	35.9	39.0	38.5	93.1	93.6	92.6	99.8	99.7	99.7

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
	2.1	2.2	2.3	2.1	2.2	2.3	2.1	2.2	2.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	46.5	51.6	47.3	95.5	95.5	95.8	99.2	99.3	99.2
DeciWatch	76.0	75.9	36.8	94.2	94.2	74.9	99.2	99.2	92.8
bi-ConvLSTM - Model S	38.8	37.4	35.9	92.7	92.1	91.2	99.4	99.5	99.3
bi-ConvLSTM - Model C	39.2	39.5	37.1	92.5	92.9	92.6	99.6	99.3	99.6

Finetuning

Test results

- ▶ 3DConv: translation + scaling vs only translation

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
	1.1	1.2	1.3	1.1	1.2	1.3	1.1	1.2	1.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	49.7	52.3	53.1	95.7	95.7	95.8	99.2	99.3	99.3
DeciWatch	76.6	76.7	68.1	94.4	94.3	87.3	99.2	99.2	96.1
bi-ConvLSTM - Model S	37.8	34.9	39.0	91.8	92.1	92.2	99.4	99.7	99.2
bi-ConvLSTM - Model C	35.9	39.0	38.5	93.1	93.6	92.6	99.8	99.7	99.7

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
	2.1	2.2	2.3	2.1	2.2	2.3	2.1	2.2	2.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	46.5	51.6	47.3	95.5	95.5	95.8	99.2	99.3	99.2
DeciWatch	76.0	75.9	36.8	94.2	94.2	74.9	99.2	99.2	92.8
bi-ConvLSTM - Model S	38.8	37.4	35.9	92.7	92.1	91.2	99.4	99.5	99.3
bi-ConvLSTM - Model C	39.2	39.5	37.1	92.5	92.9	92.6	99.6	99.3	99.6

Finetuning

Test results

- DeciWatch: effects of decreased frame rate

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
Experiment	1.1	1.2	1.3	1.1	1.2	1.3	1.1	1.2	1.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	49.7	52.3	53.1	95.7	95.7	95.8	99.2	99.3	99.3
DeciWatch	76.6	76.7	68.1	94.4	94.3	87.3	99.2	99.2	96.1
bi-ConvLSTM - Model S	37.8	34.9	39.0	91.8	92.1	92.2	99.4	99.7	99.2
bi-ConvLSTM - Model C	35.9	39.0	38.5	93.1	93.6	92.6	99.8	99.7	99.7

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
Experiment	2.1	2.2	2.3	2.1	2.2	2.3	2.1	2.2	2.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	46.5	51.6	47.3	95.5	95.5	95.8	99.2	99.3	99.2
DeciWatch	76.0	75.9	36.8	94.2	94.2	74.9	99.2	99.2	92.8
bi-ConvLSTM - Model S	38.8	37.4	35.9	92.7	92.1	91.2	99.4	99.5	99.3
bi-ConvLSTM - Model C	39.2	39.5	37.1	92.5	92.9	92.6	99.6	99.3	99.6

Finetuning

Test results

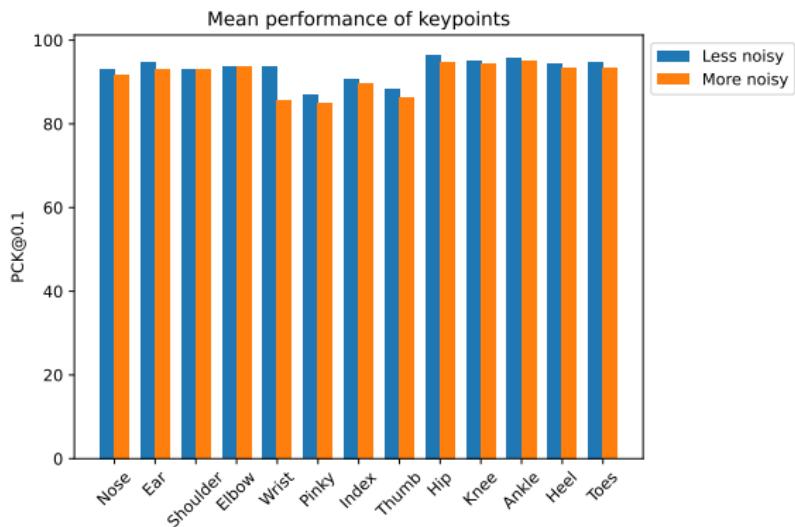
► bi-ConvLSTM: Model S vs Model C

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
	1.1	1.2	1.3	1.1	1.2	1.3	1.1	1.2	1.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	49.7	52.3	53.1	95.7	95.7	95.8	99.2	99.3	99.3
DeciWatch	76.6	76.7	68.1	94.4	94.3	87.3	99.2	99.2	96.1
bi-ConvLSTM - Model S	37.8	34.9	39.0	91.8	92.1	92.2	99.4	99.7	99.2
bi-ConvLSTM - Model C	35.9	39.0	38.5	93.1	93.6	92.6	99.8	99.7	99.7

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
	2.1	2.2	2.3	2.1	2.2	2.3	2.1	2.2	2.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	46.5	51.6	47.3	95.5	95.5	95.8	99.2	99.3	99.2
DeciWatch	76.0	75.9	36.8	94.2	94.2	74.9	99.2	99.2	92.8
bi-ConvLSTM - Model S	38.8	37.4	35.9	92.7	92.1	91.2	99.4	99.5	99.3
bi-ConvLSTM - Model C	39.2	39.5	37.1	92.5	92.9	92.6	99.6	99.3	99.6

Finetuning

Test results



Discussion

Results:

- ▶ Translation vs translation + scaling

Discussion

Results:

- ▶ Translation vs translation + scaling
- ▶ Halving frame rate

Discussion

Results:

- ▶ Translation vs translation + scaling
- ▶ Halving frame rate
- ▶ Easiest vs most difficult joints

Discussion

Results:

- ▶ Translation vs translation + scaling
- ▶ Halving frame rate
- ▶ Easiest vs most difficult joints
- ▶ Experiment 1 vs experiment 2

Discussion

Results:

- ▶ Translation vs translation + scaling
- ▶ Halving frame rate
- ▶ Easiest vs most difficult joints
- ▶ Experiment 1 vs experiment 2
- ▶ Effects of pretraining

Discussion

Results:

- ▶ Translation vs translation + scaling
- ▶ Halving frame rate
- ▶ Easiest vs most difficult joints
- ▶ Experiment 1 vs experiment 2
- ▶ Effects of pretraining
- ▶ Worst performing keypoints

Discussion

All models performed better during finetuning than pretraining

1. More data

Discussion

All models performed better during finetuning than pretraining

1. More data
2. Semantically different videos in pretraining

Discussion

All models performed better during finetuning than pretraining

1. More data
2. Semantically different videos in pretraining
3. Noise in BRACE annotations

Discussion

All models performed better during finetuning than pretraining

1. More data
2. Semantically different videos in pretraining
3. Noise in BRACE annotations
4. Frame rate in Penn Action

Discussion

All models performed better during finetuning than pretraining

1. More data
2. Semantically different videos in pretraining
3. Noise in BRACE annotations
4. Frame rate in Penn Action
5. Performance of identity function

Discussion

Which model is the best?

- ▶ Greatest testing accuracy: DeciWatch 1.1/1.2

Discussion

Which model is the best?

- ▶ Greatest testing accuracy: DeciWatch 1.1/1.2
- ▶ Greatest rough estimation: bi-ConvLSTM Model C 1.1

Discussion

Which model is the best?

- ▶ Greatest testing accuracy: DeciWatch 1.1/1.2
- ▶ Greatest rough estimation: bi-ConvLSTM Model C 1.1
- ▶ Speed and memory: 3DConv

Discussion

General reflections

- ▶ Pretraining
 - ▶ Should have estimated parameters of data

Discussion

General reflections

- ▶ Pretraining
 - ▶ Should have estimated parameters of data
 - ▶ Overlapping video sequences

Discussion

General reflections

- ▶ Pretraining
 - ▶ Should have estimated parameters of data
 - ▶ Overlapping video sequences
- ▶ Finetuning
 - ▶ Groundtruth outside of bbox

Discussion

Future work

1. DeciWatch with all frames

Discussion

Future work

1. DeciWatch with all frames
2. DeciWatch with vision transformer

Discussion

Future work

1. DeciWatch with all frames
2. DeciWatch with vision transformer
3. Avoid overfitting

Discussion

Future work

1. DeciWatch with all frames
2. DeciWatch with vision transformer
3. Avoid overfitting
4. Multiple retraining

Conclusion

Successfully developed and tested the incorporation of temporal smoothing for pose estimation

Extras: Mistakes Were Made!

Misimplemented evaluation-function

- ▶ Generally, same patterns
- ▶ Decreased accuracies
- ▶ Models do not improve the rough estimates
- ▶ bi-ConvLSTM: Model C is now always better than Model S

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
	1.1	1.2	1.3	1.1	1.2	1.3	1.1	1.2	1.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	33.3	33.4	32.8	72.5	72.4	73.1	85.8	85.8	86.0
DeciWatch	32.8	33.8	30.9	68.0	68.1	62.7	85.1	84.9	82.8
bi-ConvLSTM - Model S	31.7	30.1	31.6	71.5	68.3	71.3	86.3	82.5	86.2
bi-ConvLSTM - Model C	32.0	32.2	31.8	72.2	72.2	71.4	86.6	86.5	86.5

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
	2.1	2.2	2.3	2.1	2.2	2.3	2.1	2.2	2.3
Identity function	19.4	19.4	19.4	66.1	66.1	66.1	85.2	85.2	85.2
3DConv	33.1	33.3	32.7	72.1	72.3	72.6	85.6	85.7	85.6
DeciWatch	21.3	26.4	19.4	51.1	58.7	48.1	77.1	81.0	77.0
bi-ConvLSTM - Model S	30.9	31.6	30.6	71.2	71.6	70.7	86.0	86.1	86.2
bi-ConvLSTM - Model C	31.5	32.2	30.8	71.8	72.0	71.9	86.4	86.4	86.5