# Coarse-to-Fine Volumetric Prediction for Single-Image 3D Human Pose

UNIVERSITY of PENNSYLVANIA

Georgios Pavlakos, Xiaowei Zhou, Konstantinos G. Derpanis, Kostas Daniilidis

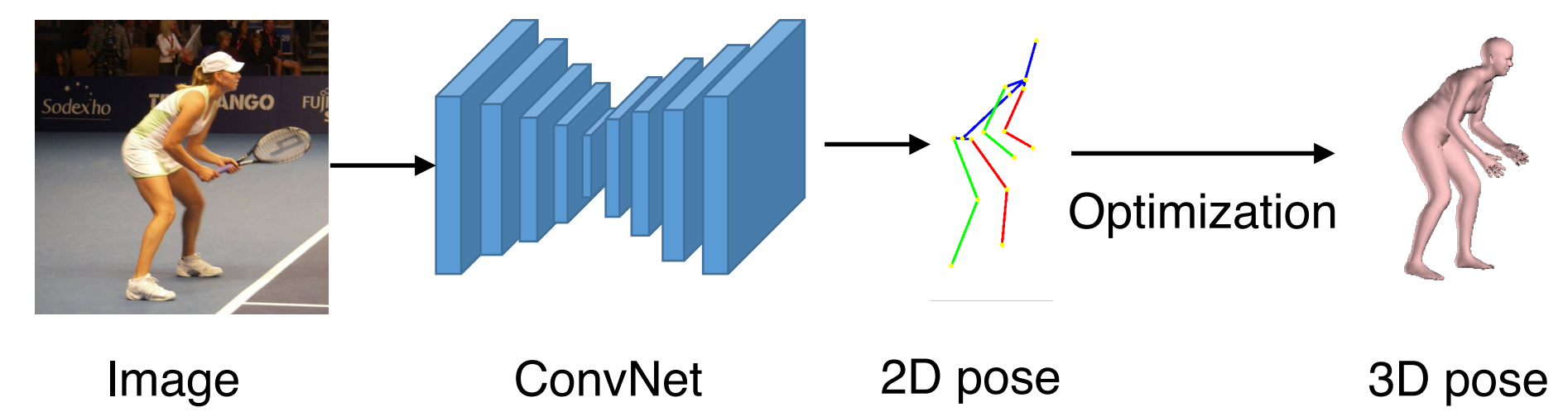CVPR July 21-26 HONOLULU 2017

Training Code
Testing Code
tinyurl.com/PoseVolumetric

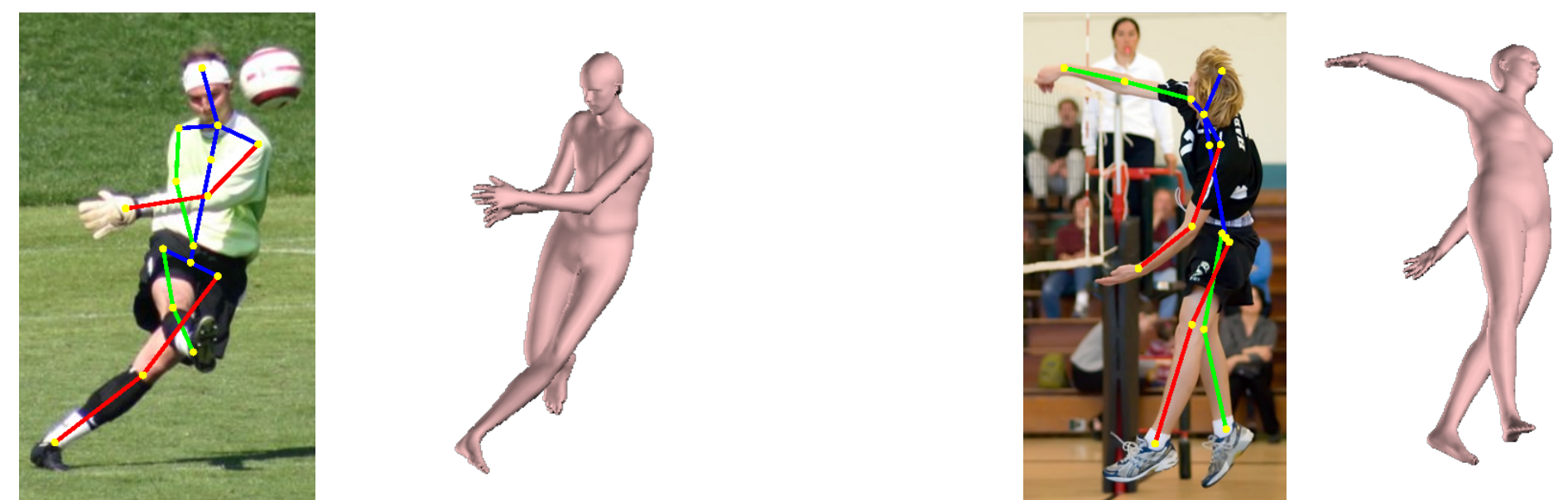## Goal: Estimate 3D human pose from a single color image

Two paradigms dominate this problem.
Reconstruction and discriminative approaches.

### Two-step Reconstruction Approaches

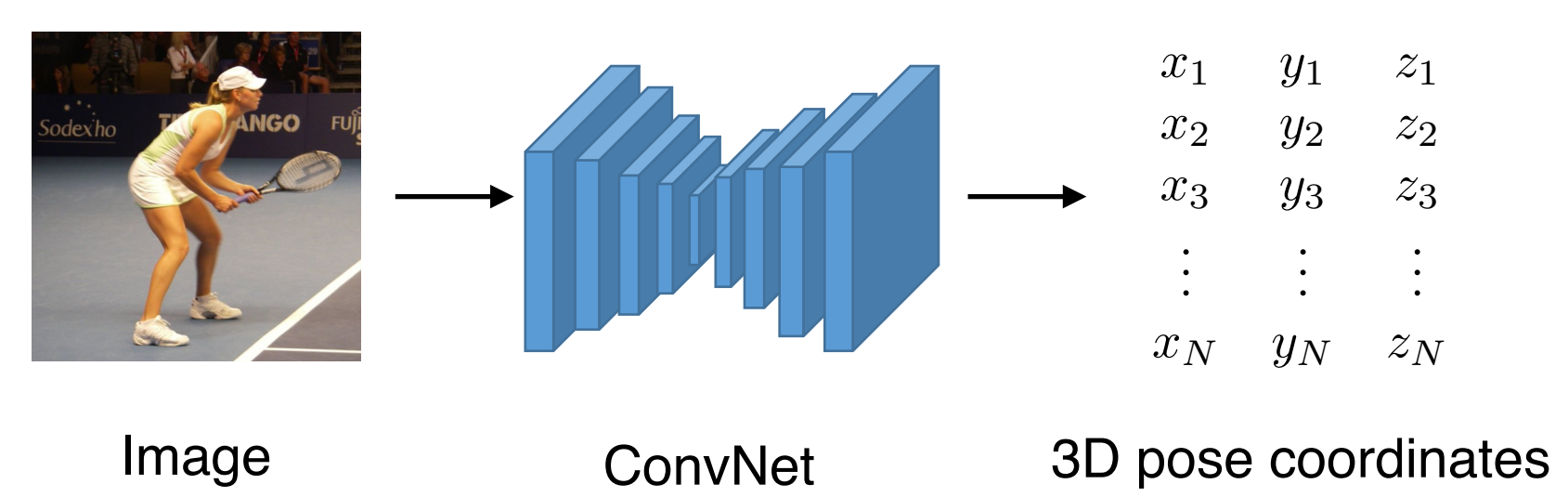**Approach**

2D pose estimation + optimization lifting 2D-to-3D



Image — ConvNet — 2D pose — Optimization — 3D pose

**Limitation**

Reconstruction Ambiguity!



### Discriminative Approaches

e.g.: coordinate regression

**Approach**

Estimate the 3D pose directly from the image.



Image — ConvNet — 3D pose coordinates

$$\begin{matrix} x_1 & y_1 & z_1 \\ x_2 & y_2 & z_2 \\ x_3 & y_3 & z_3 \\ \vdots & \vdots & \vdots \\ x_N & y_N & z_N \end{matrix}$$

**Limitation**

Problem is highly non-linear.

Mapping from images to 3D coordinates is hard to learn.

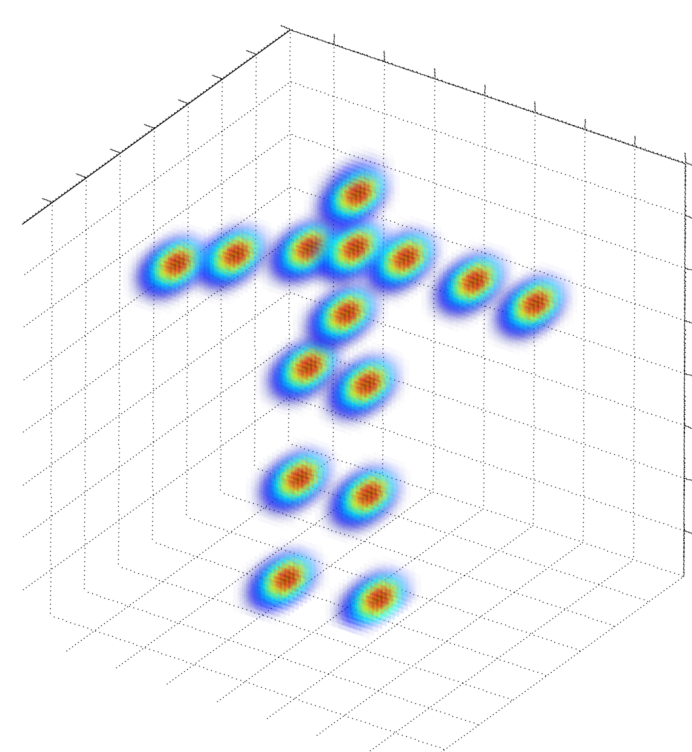Underperforms compared to two-step approaches.

### This work

We attempt to bridge the gap!

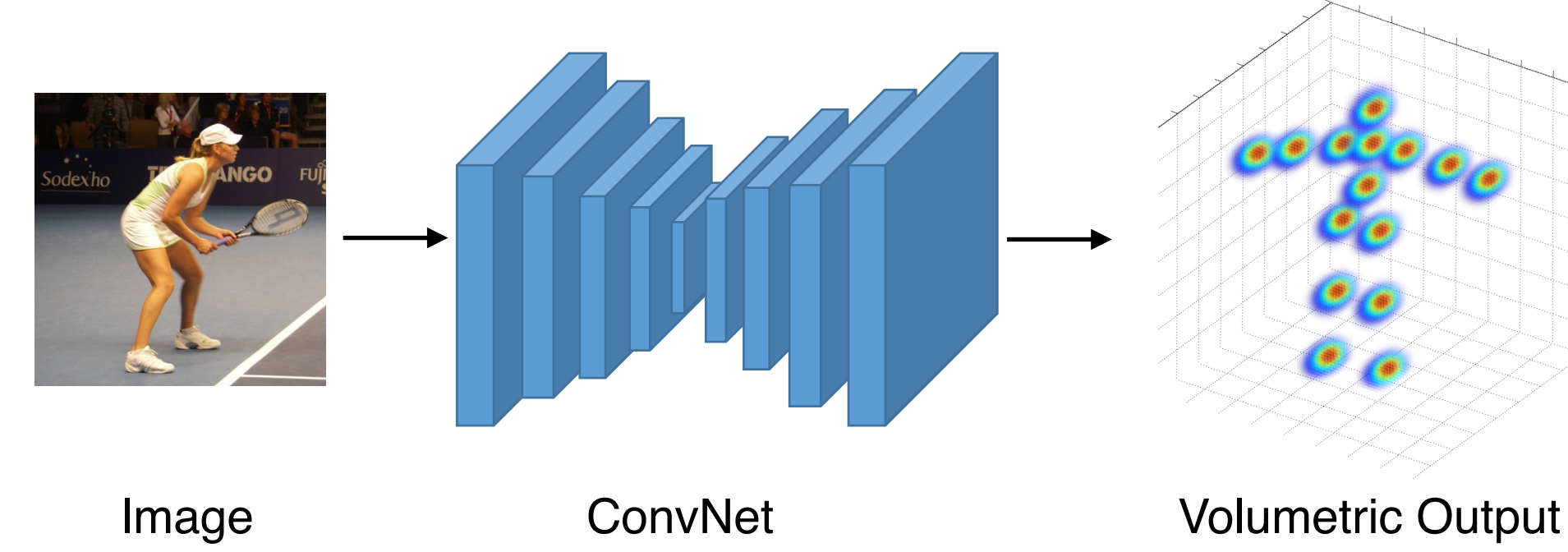We introduce the volumetric representation for 3D human pose.

We use a coarse-to-fine prediction scheme to deal with the excessive dimensionality.

We employ a decoupled architecture for 3D human pose estimation ''in-the-wild''.

We achieve more than **30% relative error reduction** for standard benchmarks!

## Volumetric representation for 3D human pose



Image — ConvNet — Volumetric Output

We cast the problem as 3D keypoint localization.

We regress 3D heat maps of dimensions 64x64x64 for each joint.

### Major advantages

ConvNets can naturally map from 2D images to 3D volumes.

The mapping can be achieved with a Fully Convolutional Network.

Rich output (64x64x64 for each joint). Useful for other tasks/postprocessing.
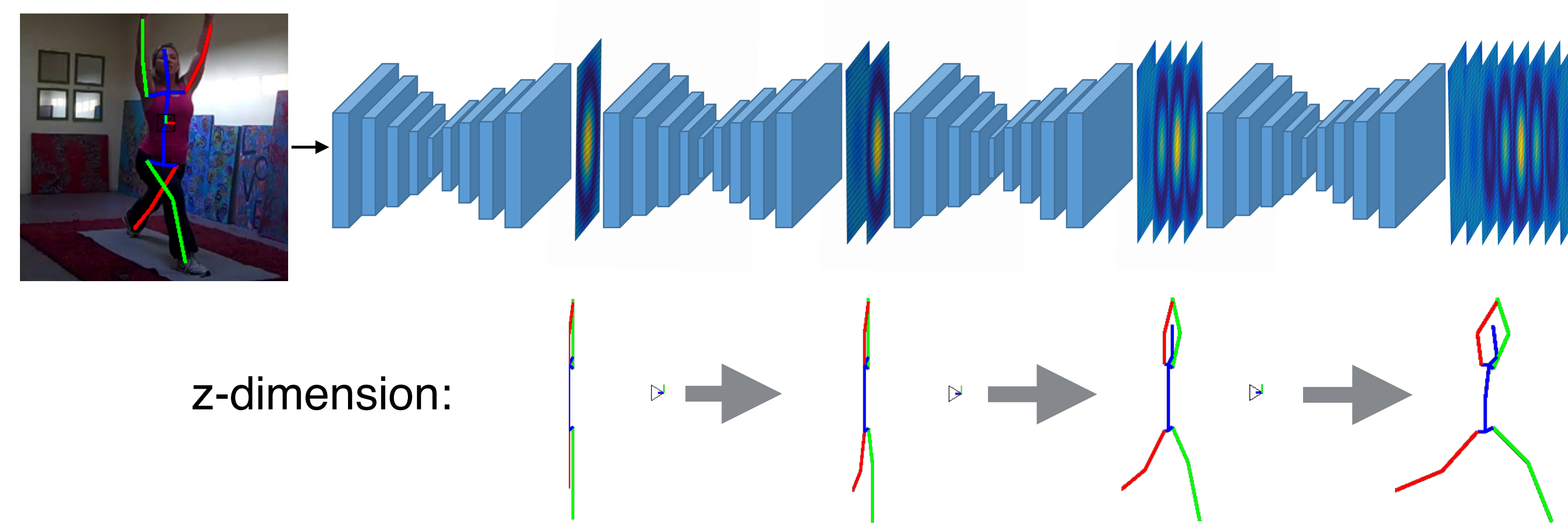
Average Error Human3.6M (mm)

| | |
|---|---|
| Coordinate Regression | 112.41 |
| Volume Regression (depth = 32) | 92.23 |
| Volume Regression (depth = 64) 162 | 85.82 |

## Coarse-to-Fine prediction

Iterative estimation offers diminishing returns because of the excessive dimensionality of our representation.

**We do it in a coarse-to-fine way!**

The resolution of the supervision volume increases gradually for the most challenging z-dimension.
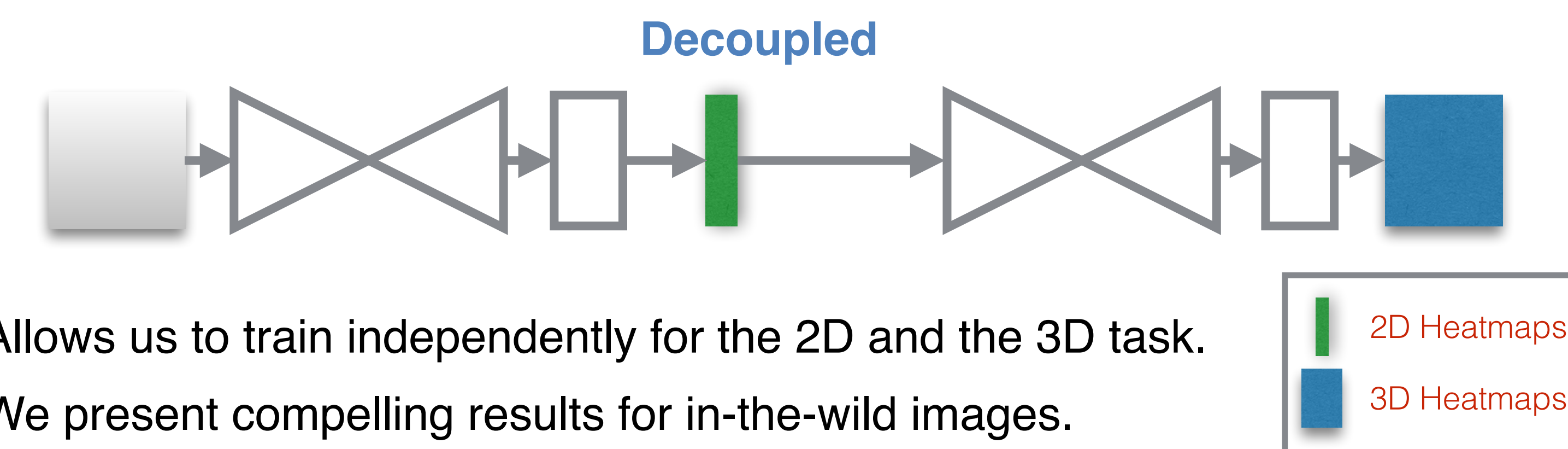


z-dimension:

Average Error Human3.6M (mm)

| | Naive Stacking | Coarse-to-Fine |
|---|---|---|
| Two Hourglasses | 80.14 | 69.77 |
| Three Hourglasses | 78.17 | 68.49 |
| Four Hourglasses | 75.06 | 64.76 |

## Versatility of volumetric representation

Regress 3D heatmaps using 2D heatmaps as input.

### Decoupled



2D Heatmaps
3D Heatmaps

+ Allows us to train independently for the 2D and the 3D task.
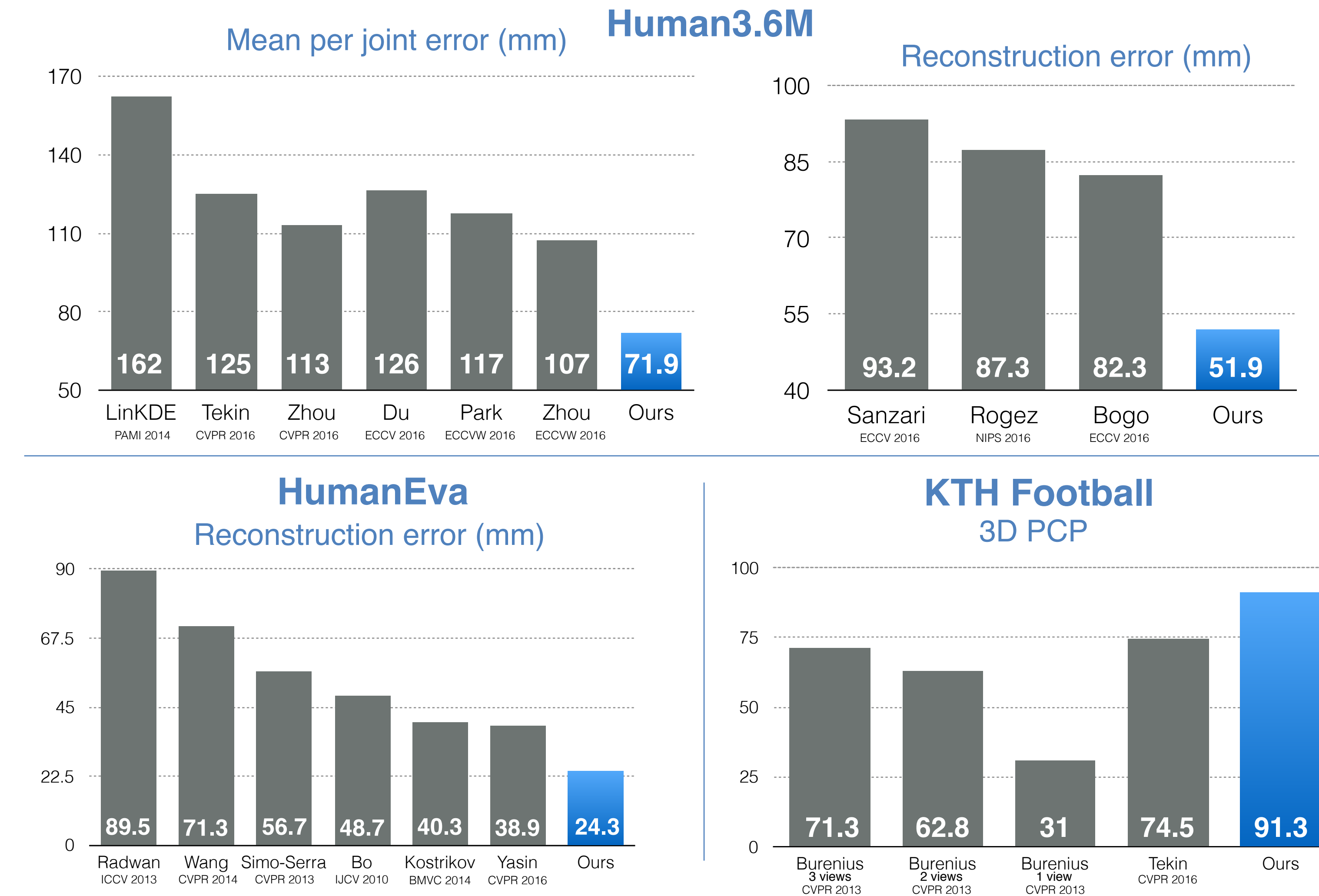+ We present compelling results for in-the-wild images.

− Uses only 2D joint locations and discards additional image evidence.
− When 2D estimates are wrong, 3D prediction can be lead astray.
− Underperforms compared to end-to-end approach.

### Coarse-to-Fine



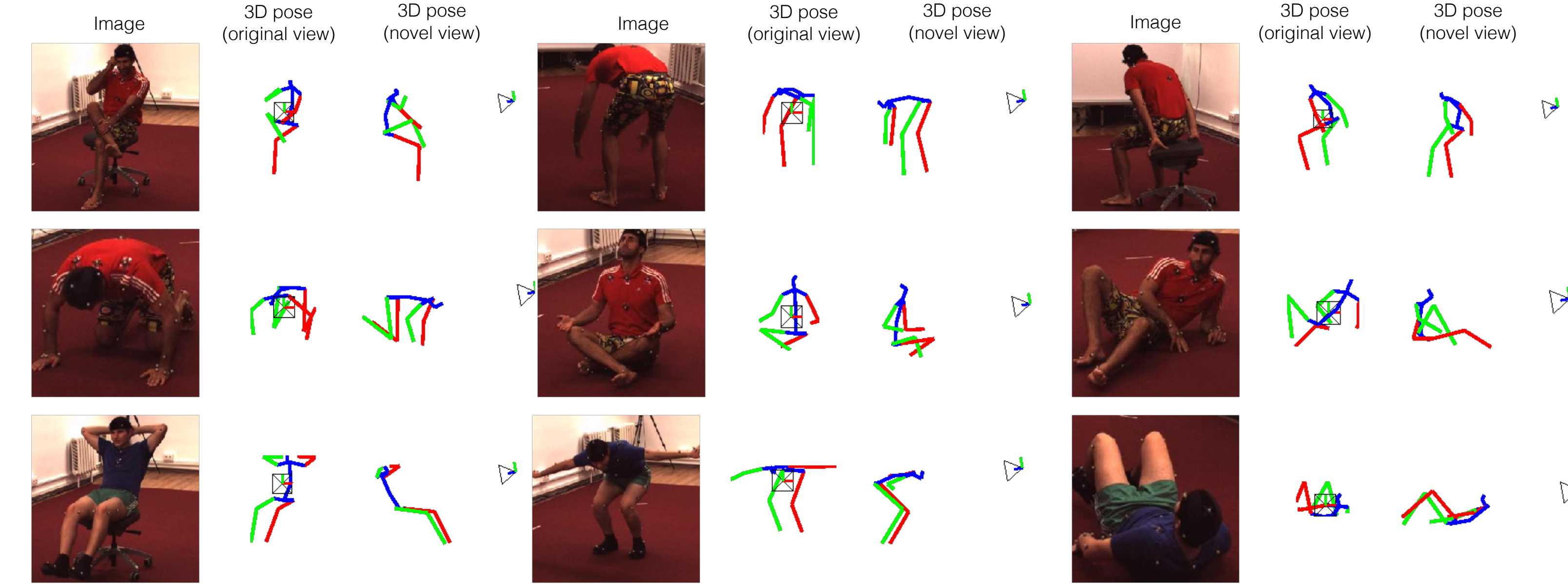Average Error Human3.6M (mm)

| Decoupled | Coarse-to-Fine |
|---|---|
| 78.10 | 69.77 |

## Quantitative results

### Human3.6M

Mean per joint error (mm)



| LinKDE PAMI 2014 | Tekin CVPR 2016 | Zhou CVPR 2016 | Du ECCV 2016 | Park ECCVW 2016 | Zhou ECCVW 2016 | Ours |
|---|---|---|---|---|---|---|
| 162 | 125 | 113 | 126 | 117 | 107 | 71.9 |

Reconstruction error (mm)



| Sanzari ECCV 2016 | Rogez NIPS 2016 | Bogo ECCV 2016 | Ours |
|---|---|---|---|
| 93.2 | 87.3 | 82.3 | 51.9 |

### HumanEva
Reconstruction error (mm)



| Radwan ICCV 2013 | Wang CVPR 2014 | Simo-Serra CVPR 2013 | Bo IJCV 2010 | Kostrikov BMVC 2014 | Yasin CVPR 2016 | Ours |
|---|---|---|---|---|---|---|
| 89.5 | 71.3 | 56.7 | 48.7 | 40.3 | 38.9 | 24.3 |

### KTH Football
3D PCP



| Burenius 3 views CVPR 2013 | Burenius 2 views CVPR 2013 | Burenius 1 view CVPR 2013 | Tekin CVPR 2016 | Ours |
|---|---|---|---|---|
| 71.3 | 62.8 | 31 | 74.5 | 91.3 |

## Qualitative results

### Human3.6M

Image | 3D pose (original view) | 3D pose (novel view)



### HumanEva



### KTH Football



### MPII



### Failure cases



### Decoupled vs Coarse-to-Fine

Image | Decoupled (2D) | Decoupled (3D) | Coarse-to-Fine (3D)