



上海人工智能实验室  
Shanghai Artificial Intelligence Laboratory



# Mask as Supervision: Leveraging Unified Mask Information for Unsupervised 3D Pose Estimation

Yuchen Yang<sup>1,2,\*</sup> Yu Qiao<sup>2</sup> Xiao Sun<sup>2</sup>

<sup>1</sup>Fudan University <sup>2</sup>Shanghai Artificial Intelligence Laboratory

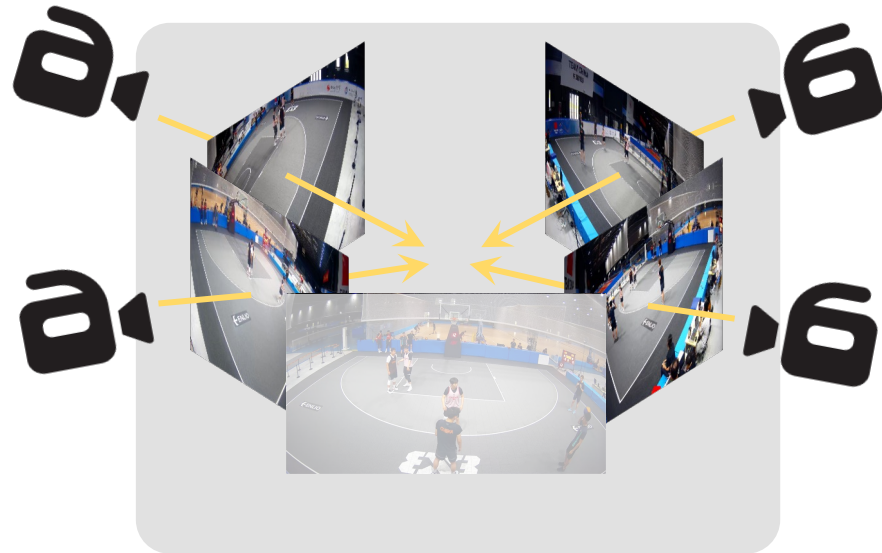
★ Work performed during his internship at Shanghai Artificial Intelligence Laboratory.

## ➤ Why unsupervised pose estimation?

*Problem Definition – 3D HPE*

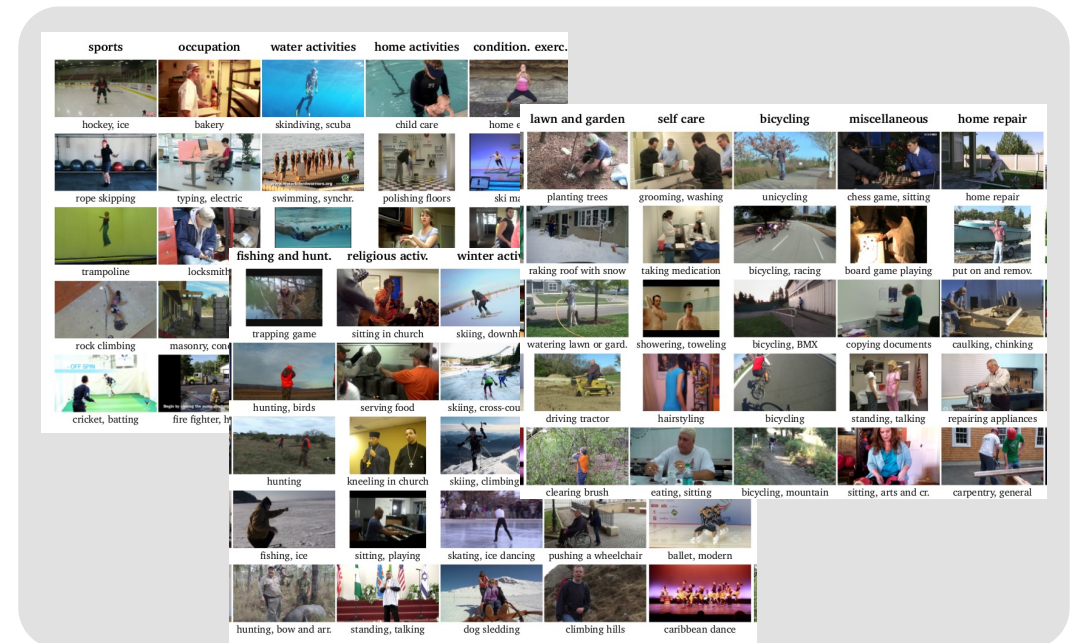
$$X = \phi(I)$$

Given an input image  $I$ , determine a set of joints  $X \in \mathbb{R}^{J \times 3}$ .  
 $J$  represents the number of joints.



Motion Capture Environment

- annotate 3D data remains costly
- vast in-the-wild data for generalization



In-the-wild

## ➤ Why mask as supervision?



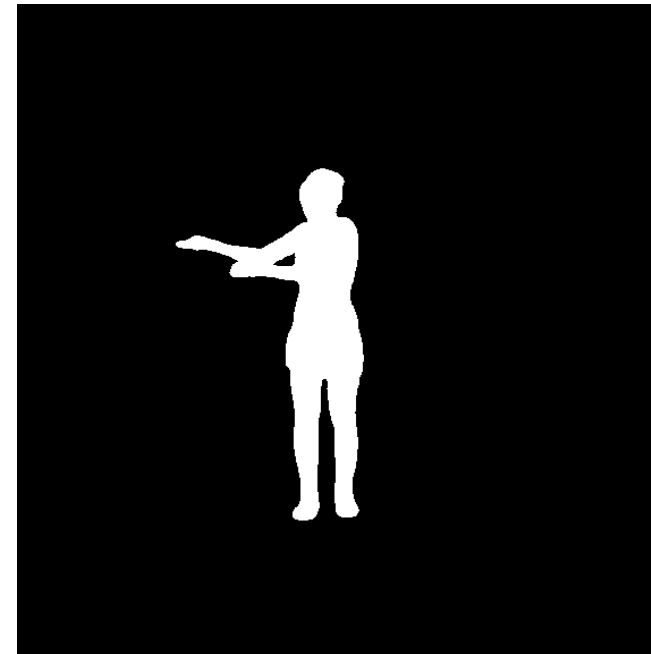
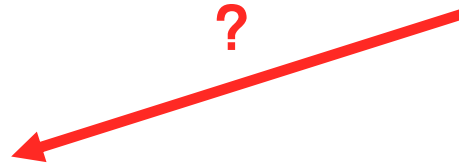
For human, we can easily estimate keypoints from masks



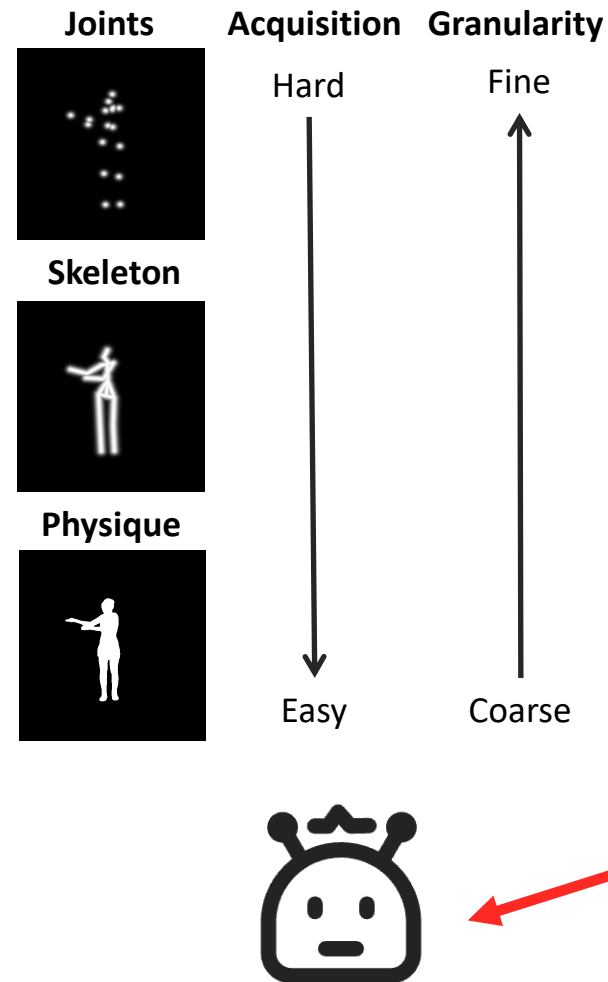
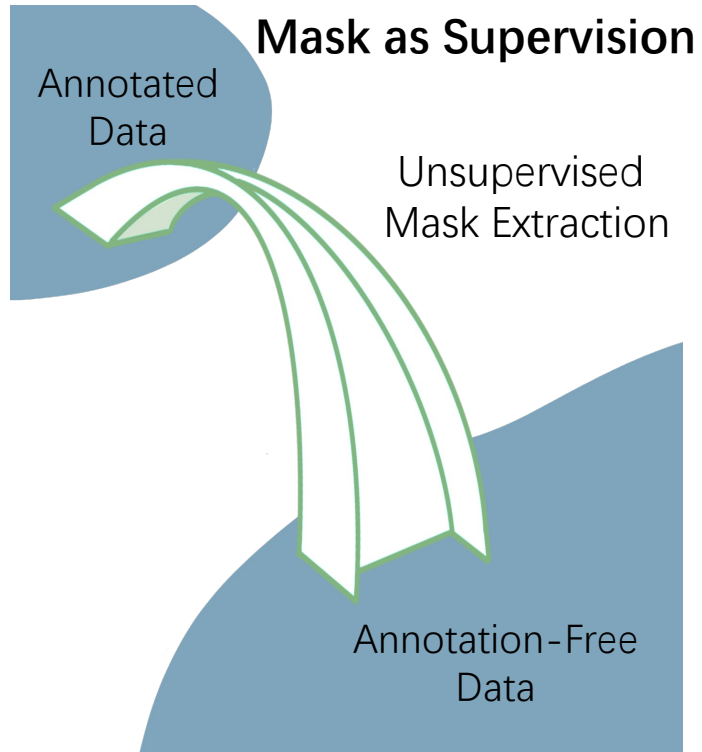
- ◆ Rich priors embedded in mask
- ◆ Human priors already acquired



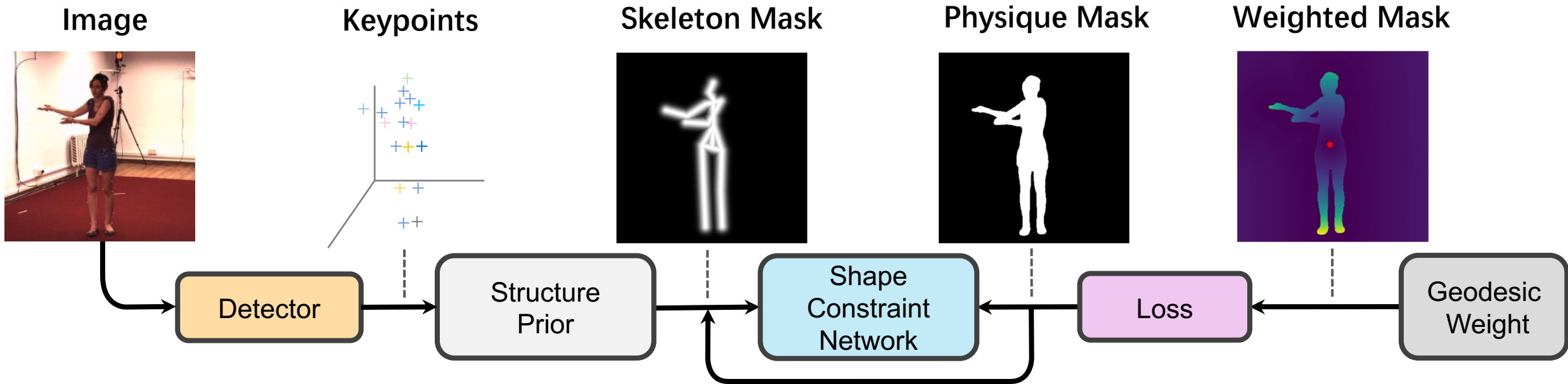
?

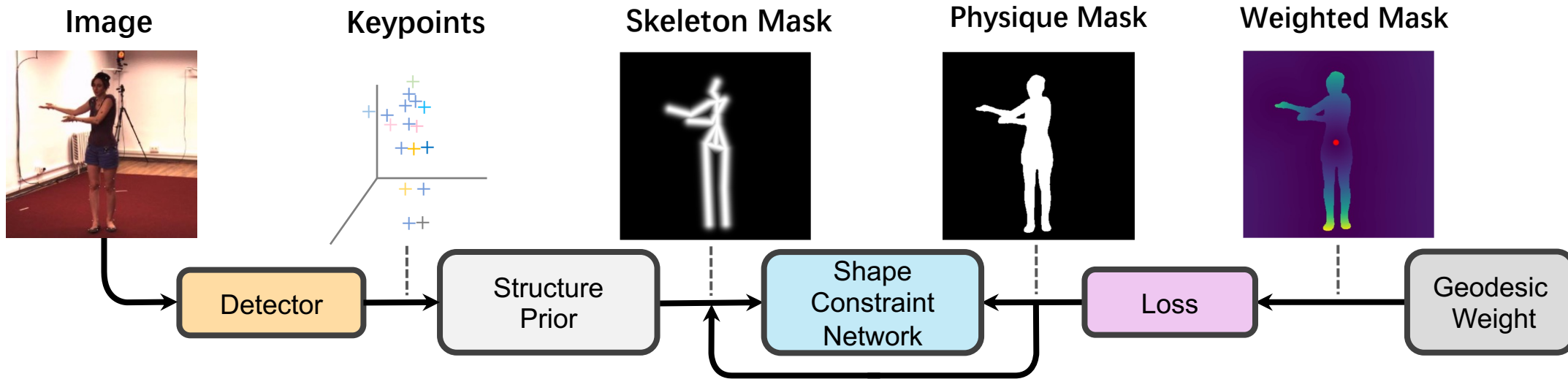


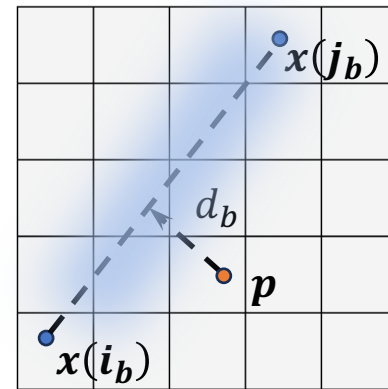
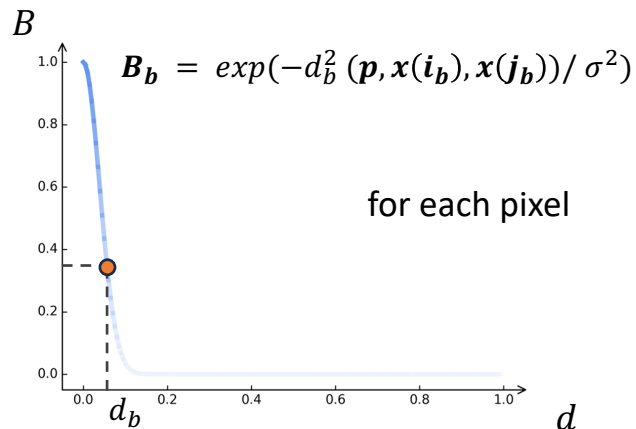
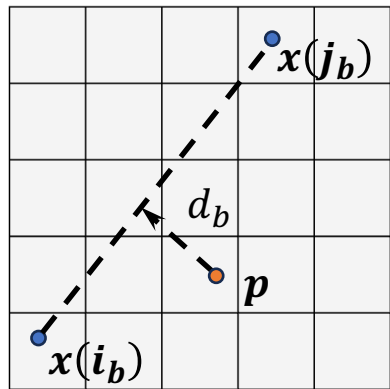
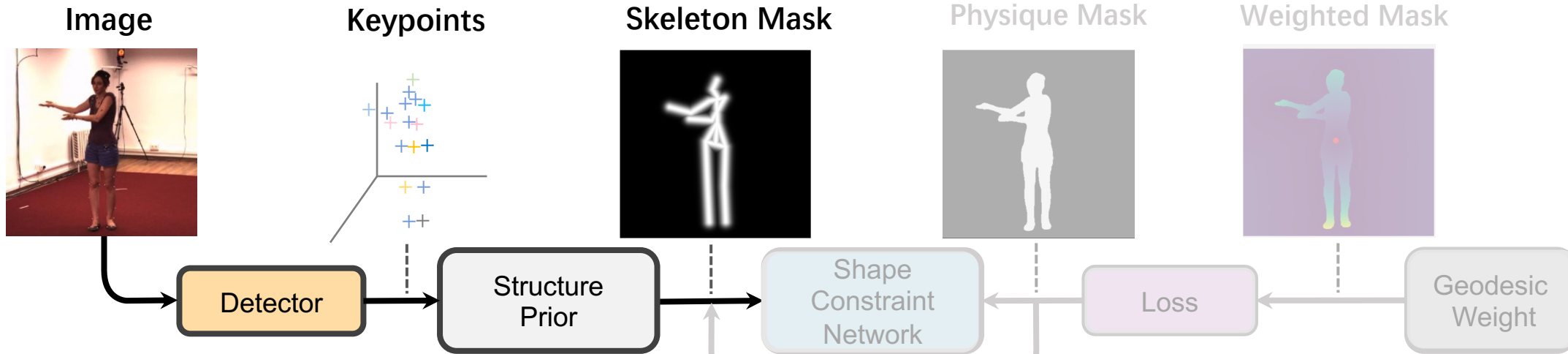
## ➤ Why mask as supervision?



- Coarse-to-fine framework
- Structure prior and shape constraint







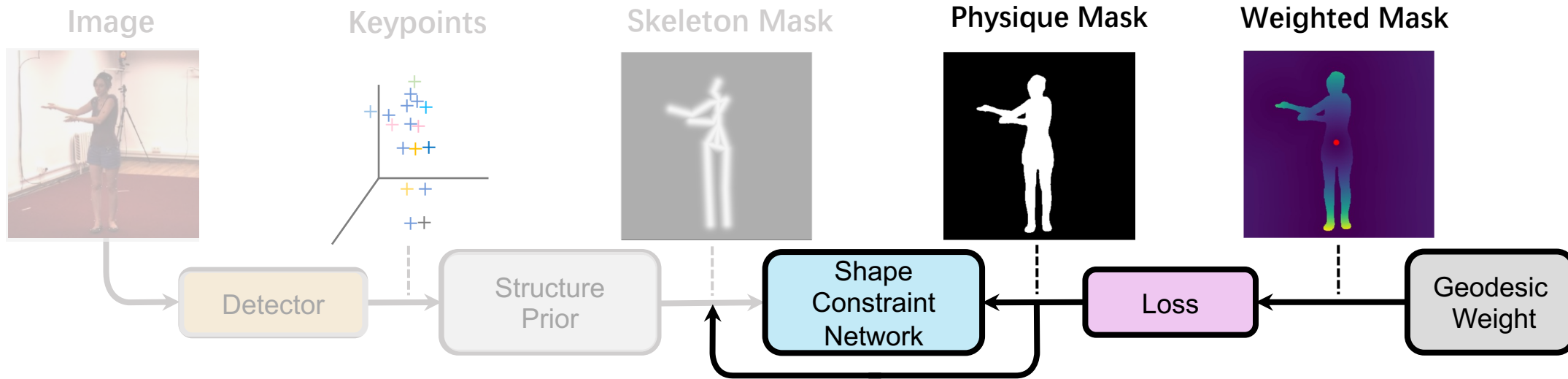
## Skeleton Mask

$$B_b = \exp(-d_b^2(p, x(i_b), x(j_b)) / \sigma^2)$$

$$M_{Skel}(x) = \sum_{b=1}^B B_b(x)$$

$\sigma$  is a hyper-parameter controlling bone width. The line segment defined by bone  $(x_1, x_2)$ . Skeleton Mask  $M_{skel}$  is from all bone maps via pixel-wise summation.





## Physique Mask

$$M_{\text{Physo}}(x) = \psi(M_{\text{Skel}}(x))$$

$\psi$  is implemented by a U-Net ended with a sigmoid function to make  $M_{\text{Physo}} \in (0,1)$

## Loss Function

$$\mathcal{L} = \lambda_s \left\| G \odot (M^{gt} - M_{\text{Skel}}(x)) \right\|_2^2 + \lambda_p \left\| G \odot (M^{gt} - M_{\text{Physo}}(x)) \right\|_2^2$$

Geodesic Distance, denoted as  $G$ , can be computed using the fast marching method

with mask centroid as zero point.

$M^{gt}$  is the given ground truth mask as supervision.

$\lambda$  is the balancing factor for loss.



## SPP-based Method

Step1: Predict landmarks  $\mathbb{R}^{L \times 3}$  ( $L \geq 2 \times J$ )

Step2: Train a mapping network  $\theta: \mathbb{R}^{L \times 3} \rightarrow \mathbb{R}^{J \times 3}$

- Ignore left-right reversal problem
- Involve human annotation

## Quantitative Results

**Table 2:** Comparison with state-of-the-art methods on MPI-INF-3DHP. MPJPE is in *cm*. Note that the first four methods use supervised post-processing and Sosa *et al.* [45] uses unpaired 2D pose to obtain interpretable keypoints.

Method	PCK(†)	AUC(†)	MPJPE(↓)
Denton <i>et al.</i> [6]	-	-	22.28
Rhodin <i>et al.</i> [38]	-	-	20.24
Honari <i>et al.</i> [14]	-	-	20.95
Honari <i>et al.</i> [15]	-	-	14.57
Sosa <i>et al.</i> [45]	69.6	32.8	-
Ours	60.2	24.7	19.36
Ours (SPP)	<b>71.3</b>	<b>42.7</b>	<b>13.67</b>

**Table 1:** Comparison with state-of-the-art methods on Human3.6M. **SPP**: supervised post-processing. **UP**: unpaired ground truth pose or its prior, **T**: manually designed template. **SF**: supervised flip to eliminate left-right ambiguity. † indicates our results do not consider the ambiguity in left-right reversal. †† indicates we do not consider inner skeleton relationships and follow the common SPP settings. The best results in SPP and No SPP groups are marked in red and blue. MSE is in 2D % and MPJPEs are in *mm*.

Method	Settings				Metrics (↓)				
	UP	T	SF	Joint	MSE	MPJPE	N-MPJPE	P-MPJPE	
SPP	Thewlis <i>et al.</i> [50]	×	×	✓	2D	7.51	-	-	-
	Zhang <i>et al.</i> [57]	×	×	✓	2D	4.14	-	-	-
	Lorenz <i>et al.</i> [30]	×	×	✓	2D	2.79	-	-	-
	Suwajanakorn <i>et al.</i> [49]	×	×	×	3D	-	158.7	156.8	112.9
	Sun <i>et al.</i> [47]	×	×	×	3D	-	125.0	-	105.0
	Honari <i>et al.</i> [14]	×	×	✓	3D	-	100.3	99.3	74.9
Honari <i>et al.</i> [15]	×	×	×	3D	<b>2.38</b>	73.8	72.6	63.0	
SPP	<b>Ours</b> <sup>††</sup>	×	×	×	3D	2.52	<b>65.5</b>	<b>66.1</b>	<b>61.9</b>
No SPP	Schmidtke <i>et al.</i> [40]	×	✓	✓	2D	3.31	-	-	-
	Jakab <i>et al.</i> [20]	✓	×	✓	2D	<b>2.73</b>	-	-	-
	Sosa <i>et al.</i> [45]	✓	×	✓	3D	-	-	-	96.4
	Kundu <i>et al.</i> [23]	✓	✓	✓	3D	-	99.2	-	-
	Kundu <i>et al.</i> [24]	✓	×	✓	3D	-	-	-	89.4
No SPP	<b>Ours</b> <sup>†</sup>	×	×	✓	3D	3.17	<b>85.6</b>	<b>85.6</b>	<b>79.3</b>
	<b>Ours</b>	×	×	×	3D	3.63	95.9	96.8	90.4

## Quantitative Results



**Table 3:** Ablation study on shape reconstruction.

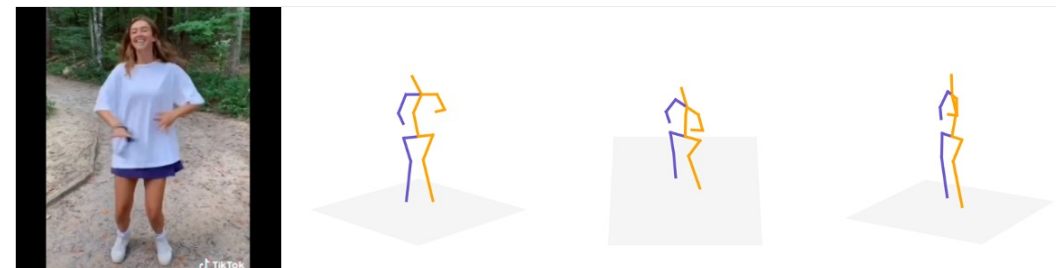
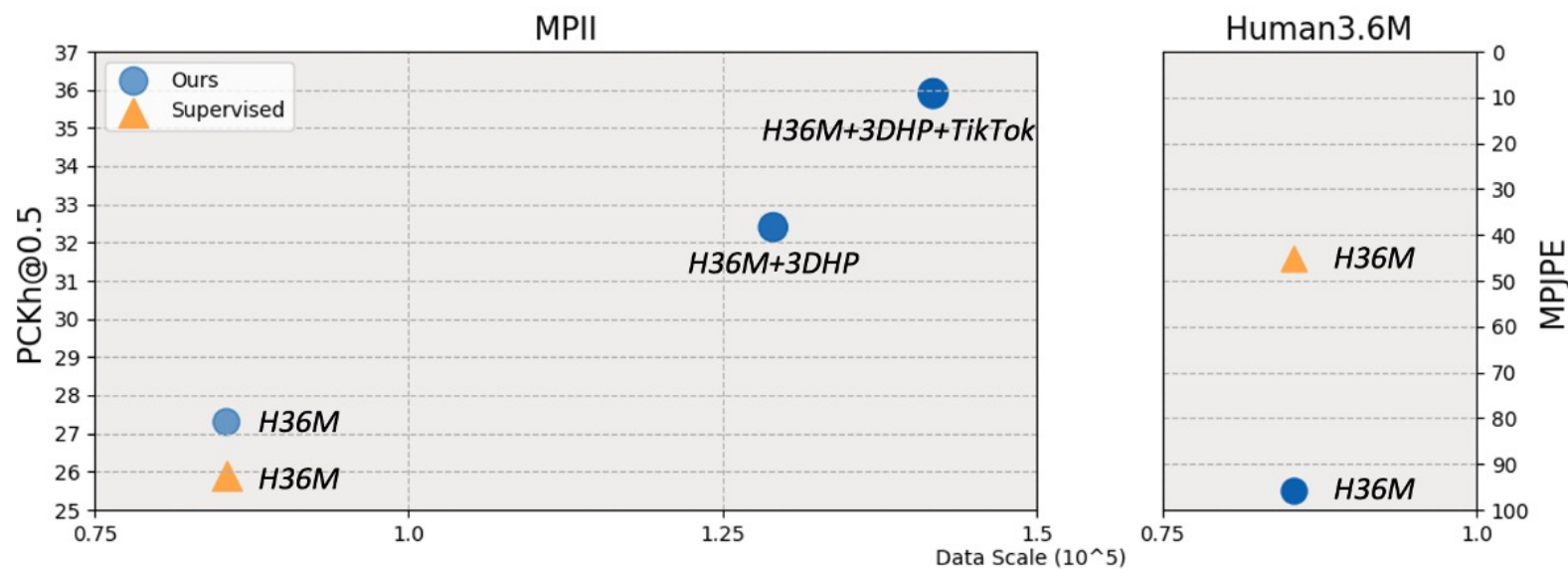
Configurations	MPJPE ( $\downarrow$ )	Ambiguity Ratio ( $\downarrow$ )
<i>wo</i> $\psi$ , $\mathbf{G}$ , $\Delta$	118.1	48.73%
<i>wo</i> $\psi$ , $\mathbf{G}$	127.4	23.34%
<i>wo</i> $\mathbf{G}$	102.6	22.83%
Full	95.9	20.33%

## Ablation Study



- Effectiveness of Skeleton Mask
- Effectiveness of Physique Mask and the rest components

## Leveraging In-the-wild Data





上海人工智能实验室  
Shanghai Artificial Intelligence Laboratory

**Thanks!**