

IVMSP-P18.10 3D POSE ESTIMATION FROM MONOCULAR VIDEO WITH CAMERA-BONE ANGLE REGULARIZATION ON THE IMAGE FEATURE

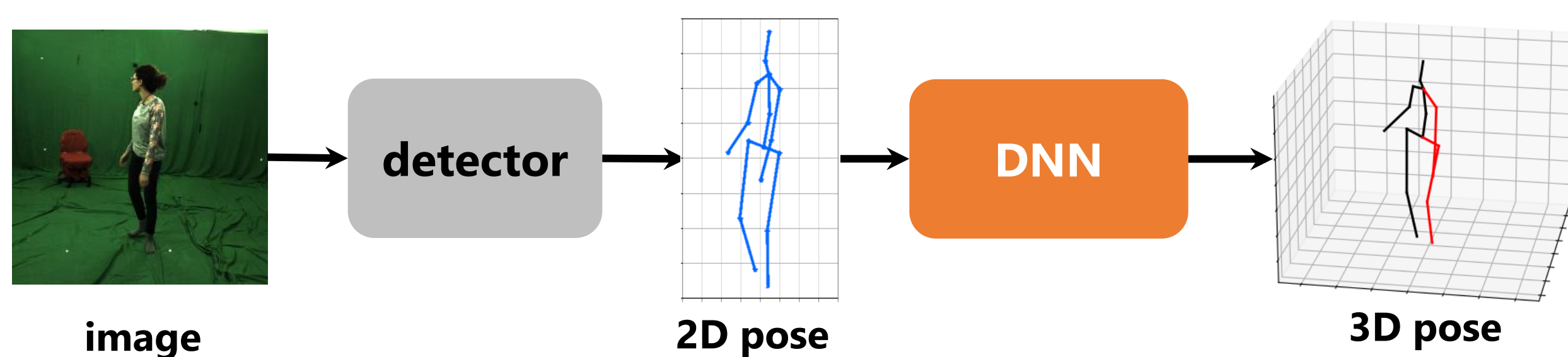
Asuka Ishii, Hiroo Ikeda (NEC)

Appearance information bounds solution space of 2D-to-3D lifting

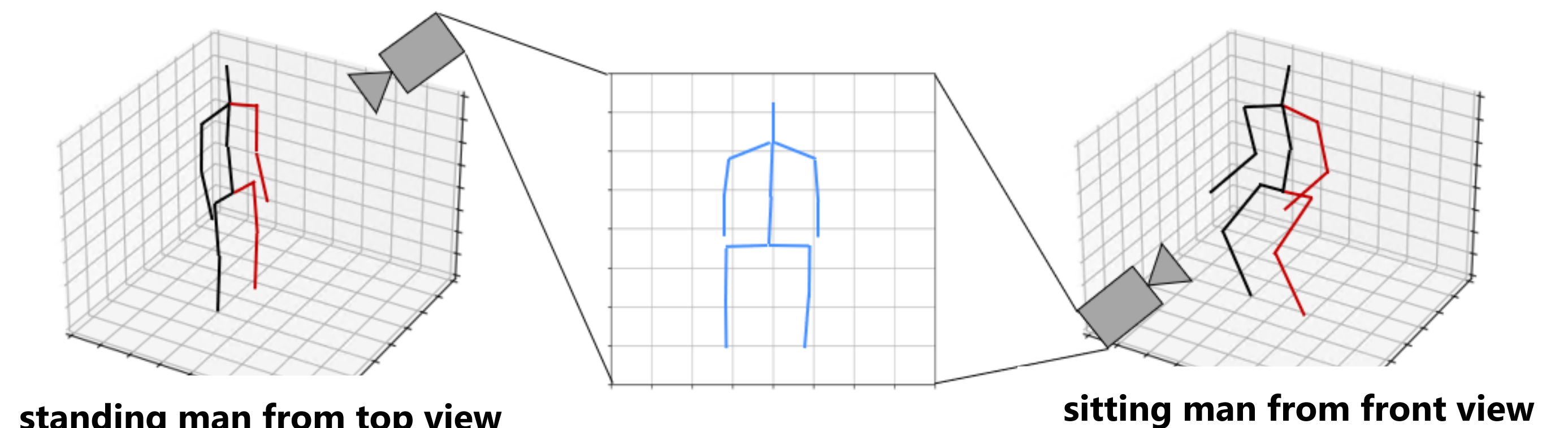
Introduction | 3D human pose estimation from monocular video

2D-to-3D lifting

- estimates 3D pose(s) only from 2D pose(s) detected by another detector
- much lower error than image-based methods [14]
- relative 3D coordinates from root joint (typically center-hip)



ill-posed problem



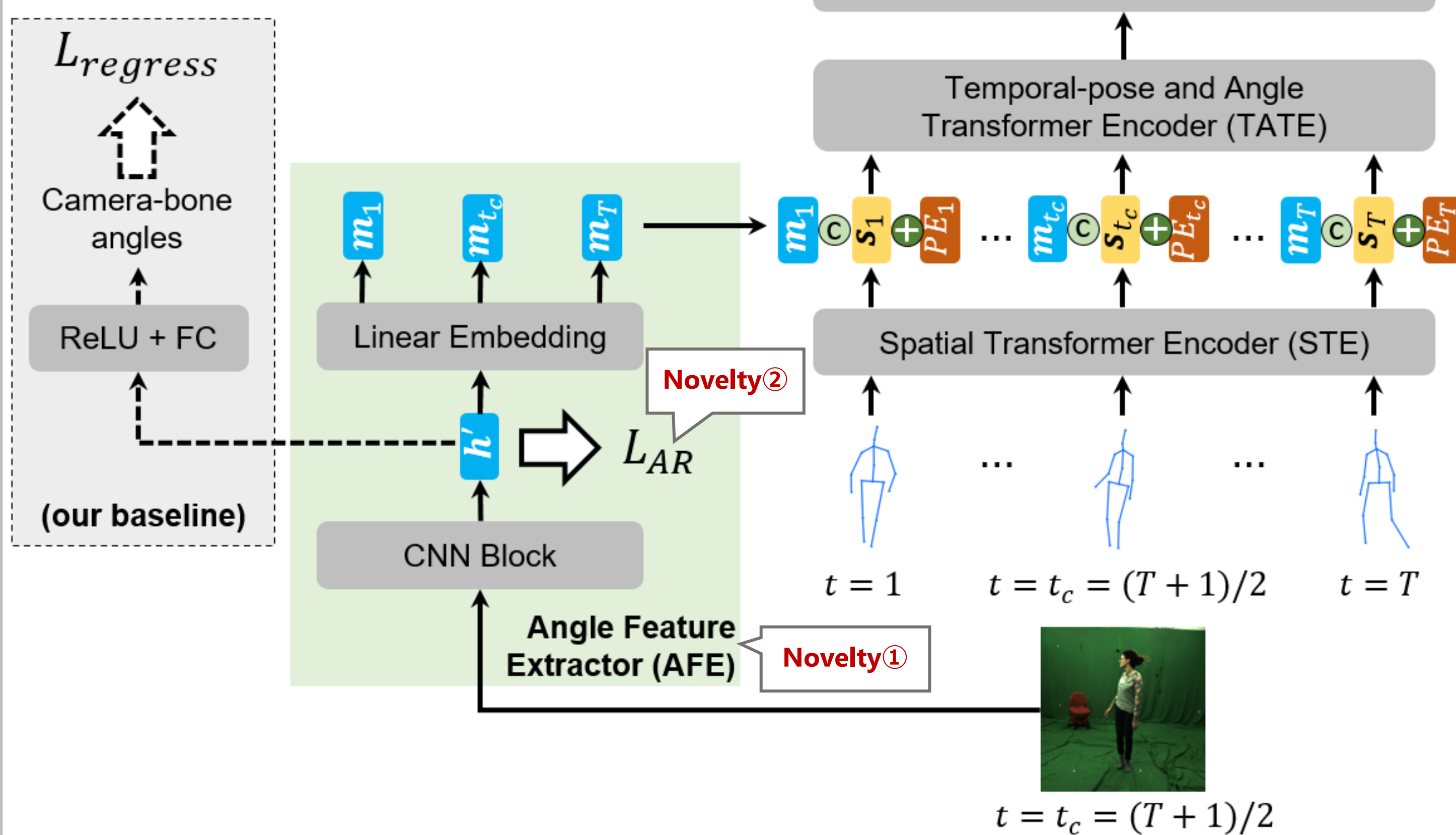
one 2D pose can be mapped to multiple 3D poses

existing methods: consider temporal information only
our method: considers appearance as well

Method | appearance information of subject

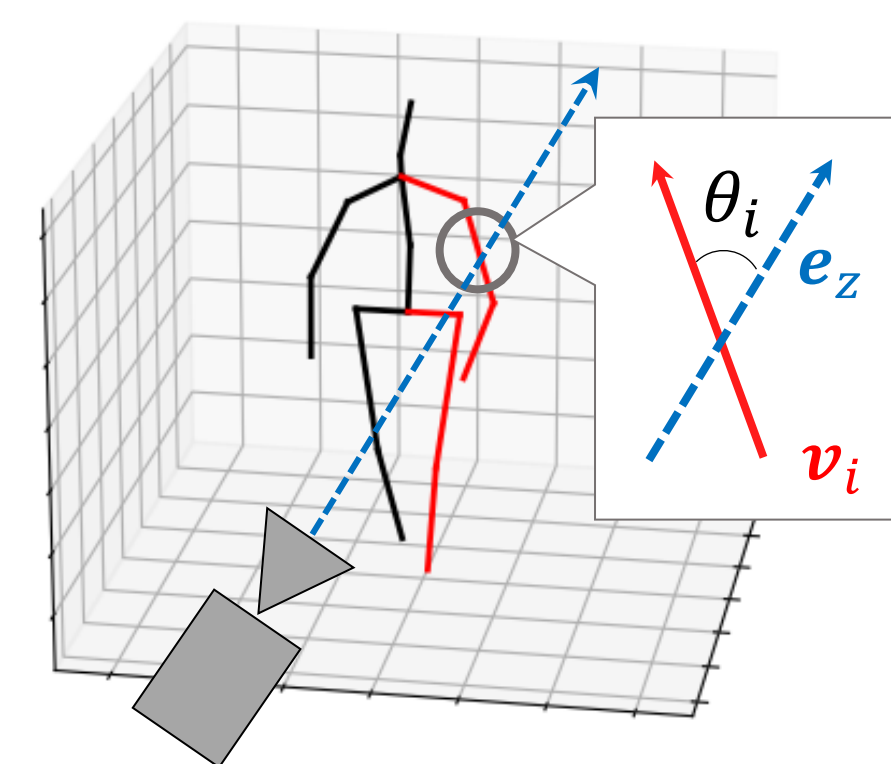
Novelty① CNN (AFE) added on PoseFormer[29], a 2D-to-3D-based network

- concatenation
- addition
- image feature
- 2D pose feature
- position embedding



Novelty② Regularization on image features using camera-bone angles

- Camera-bone angles** between camera optical axis e_z and bones v computed from ground-truth



$c = (\theta_1, \theta_2, \dots, \theta_B)$: camera-bone angle vector

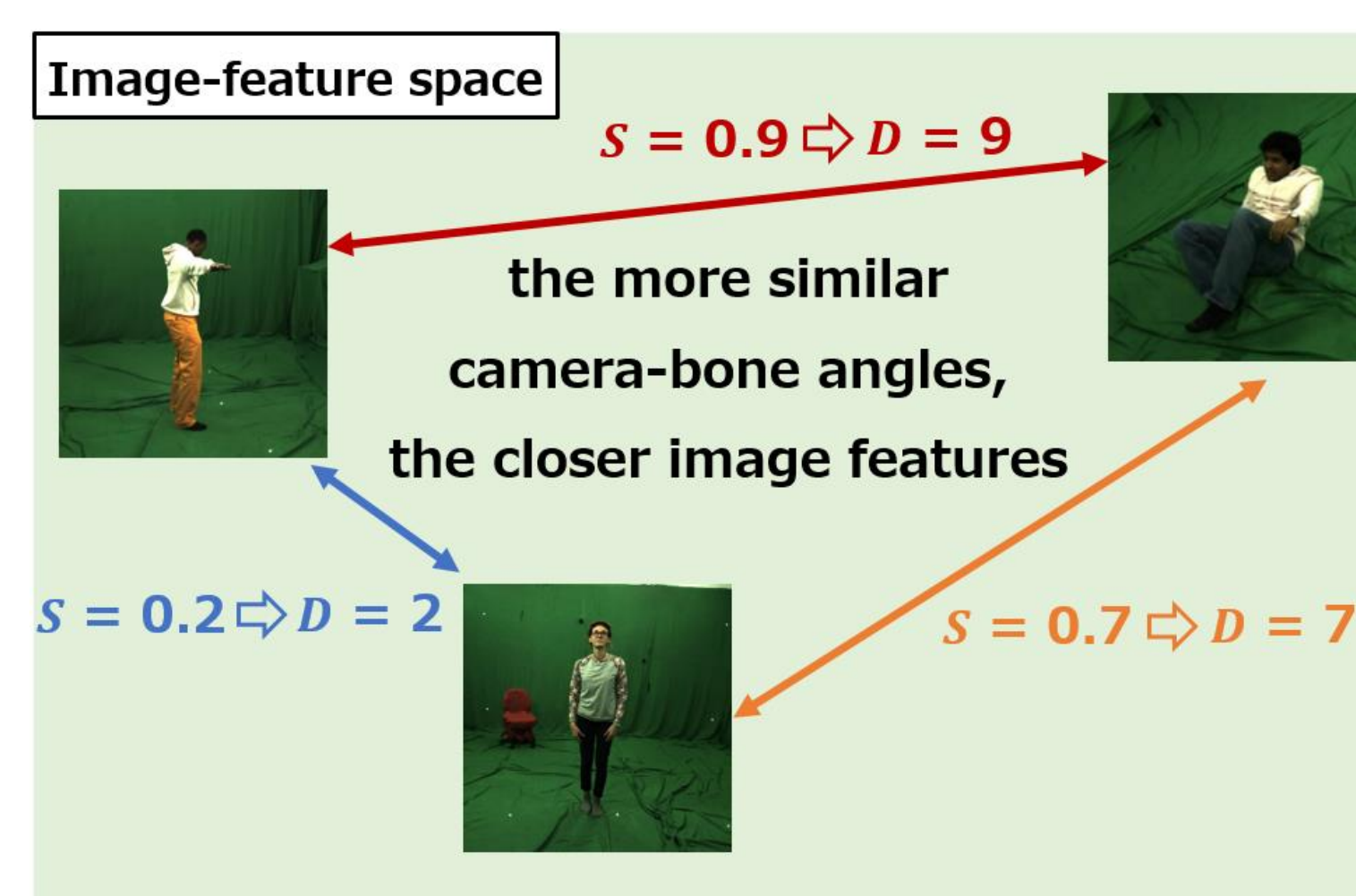
$\theta_i = \arccos\left(\frac{v_i \cdot e_z}{\|v_i\| \|e_z\|}\right)$: angle between i -th bone and camera optical axis

$v_i = \hat{p}_{i1} - \hat{p}_{i2}$: vector of i -th bone

$\hat{p}_{i1}, \hat{p}_{i2}$: ground-truth 3D coordinates of joints at both ends of i -th bone

- Regularization loss L_{AR}**

makes D , distance between image features of minibatch samples, proportional to S , unsimilarity of camera-bone angles of the samples



$$L_{total} = L_p + \lambda L_{AR}$$

$$L_{AR} = \frac{2}{N(N-1)} \sum_{a=1}^N \sum_{b=1}^a l_{a,b}$$

$$l_{a,b} = \|\alpha S(c_a, c_b) - D(h'_a, h'_b)\|_2$$

$$S(c_a, c_b) = 1 - \frac{c_a \cdot c_b}{\|c_a\| \|c_b\|}$$

$$D(h'_a, h'_b) = \|h'_a - h'_b\|_2$$

Why camera-bone angles?

Given $(\|v_i\|, K, \theta_i)$, 2D-to-3D lifting solved analytically (K : camera intrinsic parameters)

However, training model to extract

- $\|v_i\|$ may cause overtraining
- K requires large extra training data [13]

We focus on camera-bone angles

Experiment

* indicates the experiment was conducted in our environment. Otherwise, values were taken from original papers.

Approach	Method	Human3.6M		MPI-INF-3DHP	
		# frames	MPJPE [mm] ↓	# frames	MPJPE [mm] ↓
Image-based	Pavlos+2018 [5]	1	56.2	1	-
	MargiPose [7]	1	55.4	1	85.2
2D-to-3D-based	*PoseFormer [10]	81	49.9	9	50.0
	*PoseFormer + AFE	81	59.8	9	69.6
Ours	*PoseFormer + AFE w/ $L_{regress}$	81	52.4	9	64.8
	*PoseFormer + AFE w/ L_{AR}	81	44.8	9	47.9

Conclusion & future work

Conclusion

- proposed to bound solution space of 2D-to-3D method, an ill-posed problem, by considering appearance information of subject as well.
- proposed regularization loss using camera-bone angles on image features.
- empirically showed the proposed method improves performance.

Future work

- Replacement based 2D-to-3D network with SOTA
- Evaluation on unseen camera angles