

# Supplementary Materials: Denoising and Selecting Pseudo-Heatmaps for Semi-Supervised Human Pose Estimation

Zhuoran Yu<sup>\*†</sup>, Manchen Wang<sup>\*‡</sup>, Yanbei Chen, Paolo Favaro, Davide Modolo  
AWS AI Labs

zhuoran.yu@wisc.com, {manchenw, yanbec, pffavaro, dmodolo}@amazon.com

## 1. Additional Ablation Study

### 1.1. Pseudo-heatmaps accuracy

As described in Section 3.2, we employ the winner-take-all strategy to compute the ensemble pseudo-heatmaps by taking the maximum score per pixel over outputs of all augmented views. In Table 1, we compare pseudo-heatmaps accuracy calculated with Percentage of Correct Keypoint (PCK) when using average or maximum operation to ensemble multi-view augmentations. We find that winner-take-all yields a substantially higher performance (66.0 VS 63.9) by considering human pose characteristics: A clean pseudo-heatmap should have one strong peak. However, average operation often results in multiple weak peaks (see Figure 1), leading to a noisier learning target and poorer performance. On the other hand, maximum operation removes spurious noise by selecting the most confident peak.

Aggregate	None	Average	Maximum
PCK	63.5	63.9	<b>66.0</b>

Table 1. **Comparison of pseudo-heatmaps accuracy when using different aggregations.** Taking the maximum response generates the most accurate pseudo-heatmaps.

### 1.2. Teacher-Student VS Dual-Student

As described in Section 3, we employ a dual-student framework [2] which trains two student networks jointly and uses the output from one network as the learning targets to supervise the other. In Table 2, we show the comparison between our method and the adaptation of a teacher-student network Unbiased Teacher [1] to the human pose estimation task. We follow Unbiased Teacher to evolve both Teacher and Student models via mutual learning scheme, where Teacher model generates pseudo-heatmaps to train

\*These authors contributed equally to this work.

†Currently at The University of Wisconsin–Madison. Work conducted during an internship with AWS AI Labs.

‡Corresponding author.

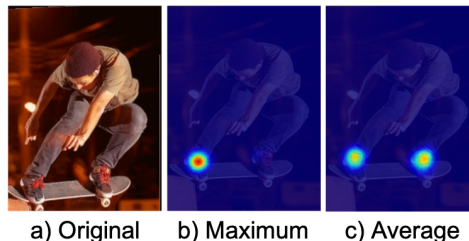


Figure 1. **Pseudo-heatmaps visualization for right ankle by using the maximum and average ensemble.** Maximum operation generate a cleaner and more accurate pseudo-heatmap.

Row	Method	Threshold-and-Refine	mAP
1	Teacher-Student [1]		32.83
2	Teacher-Student [1]	✓	36.53
3	Dual-Student [2]		42.67
4	Dual-Student [2]	✓	<b>44.13</b>

Table 2. **Comparison between Teacher-Student and Dual-Student framework.** Results on COCO-Partial protocol using 1K labeled images.

the Student model, and the knowledge that the Student model learned is then transferred to Teacher model via exponential moving average (EMA) on network weights. Results show that dual-student achieves significantly better results than teacher-student. This justifies why we use dual-student as baseline throughout this work. Furthermore, adding Threshold-and-Refine scheme to denoise the pseudo-heatmaps always improves the model performance in the teacher-student and dual-student frameworks. This suggests our design is agnostic to the underlying semi-supervised learning frameworks.

## 2. Pseudo Code

We provide the pseudo code of pseudo heatmap denoising and uncertainty-guided selection in Algorithm 1.

---

**Algorithm 1** Pseudo Heatmap Denoising and Uncertainty-guided Selection

---

**Input:** Batch of unlabeled examples  $U = (u_b; b \in (1, \dots, B))$ , two student networks  $\theta$  and  $\xi$ , number of multi-view augmentations  $K$ , number of joints  $J$ , uncertainty margin  $\Delta$ , weak augmentation strength  $S_w$ , strong augmentation strength  $S_s$ , multi-view augmentation strength  $S_m$

```
1: for  $b \leftarrow 1$  to  $B$  do
2:    $T_w, I_w \leftarrow \text{Augment}(u_b, S_w)$  ▷ Apply weak augmentation to  $u_b$ 
3:    $T_s, I_s \leftarrow \text{Augment}(u_b, S_s)$  ▷ Apply strong augmentation to  $u_b$ 
4:    $H_{\theta,w} \leftarrow \theta(I_w)$  ▷ Predict heatmaps for weakly and strongly augmented images
5:    $H_{\xi,w} \leftarrow \xi(I_w)$ 
6:    $H_{\theta,s} \leftarrow \theta(I_s)$ 
7:    $H_{\xi,s} \leftarrow \xi(I_s)$ 
8:    $H_{\theta,w} \leftarrow \text{Transform}(H_{\theta,w}, (T_w)^{-1})$  ▷ Transform weakly augmented pseudo heatmap back to original view
9:    $H_{\xi,w} \leftarrow \text{Transform}(H_{\xi,w}, (T_w)^{-1})$ 
10:  for  $k \leftarrow 1$  to  $K$  do
11:     $T_{s,k}, I_{s,k} \leftarrow \text{Augment}(u_b, S_m)$  ▷ Apply  $k^{\text{th}}$  round of multi-view augmentation to  $u_b$ 
12:     $H_{\theta,s}^k \leftarrow \theta(I_{s,k})$  ▷ Predict heatmaps on  $k^{\text{th}}$  multi-view augmented image  $I_s^k$ 
13:     $H_{\xi,s}^k \leftarrow \xi(I_{s,k})$ 
14:     $H_{\theta,s}^k \leftarrow \text{Transform}(H_{\theta,s}^k, (T_{s,k})^{-1})$  ▷ Transform multi-view augmented pseudo heatmap back to original view
15:     $H_{\xi,s}^k \leftarrow \text{Transform}(H_{\xi,s}^k, (T_{s,k})^{-1})$ 
16:  end for
17:   $P_\theta \leftarrow \max\{H_{\theta,w}, \{H_{\theta,s}^k\}_{k=1}^K\}$  ▷ Ensemble pseudo-heatmaps
18:   $P_\xi \leftarrow \max\{H_{\xi,w}, \{H_{\xi,s}^k\}_{k=1}^K\}$ 
19:   $U_\theta \leftarrow \text{stdev}\{H_{\theta,w}, \{H_{\theta,s}^k\}_{k=1}^K\}$  ▷ Estimate uncertainty as pixel-wise standard deviation
20:   $U_\xi \leftarrow \text{stdev}\{H_{\xi,w}, \{H_{\xi,s}^k\}_{k=1}^K\}$ 
21:   $u_\theta \leftarrow \{\max U_\theta^j\}_{j=1}^J$  ▷ Take the maximum value on uncertainty maps for every joint
22:   $u_\xi \leftarrow \{\max U_\xi^j\}_{j=1}^J$ 
23:   $\hat{P}_{\theta,j} \leftarrow P_{\theta,j}$  if  $u_{\theta,j} + \Delta < u_{\xi,j}$  else  $P_{\xi,j}$  ▷ Learning targets selected by uncertainty
24:   $\hat{P}_{\xi,j} \leftarrow P_{\xi,j}$  if  $u_{\xi,j} + \Delta < u_{\theta,j}$  else  $P_{\theta,j}$ 
25:   $\mathcal{L}_\theta^U = \|H_{\theta,s} - \text{Transform}(\hat{P}_\xi, T_s)\|^2$  ▷ Compute unsupervised loss on unlabeled data
26:   $\mathcal{L}_\xi^U = \|H_{\xi,s} - \text{Transform}(\hat{P}_\theta, T_s)\|^2$ 
27: end for
```

---

### 3. Additional qualitative results

We provide more qualitative results of predicted human pose on COCO val2017 dataset in Figure 2 to demonstrate the strong performance of our model.

### References

- [1] Yen-Cheng Liu, Chih-Yao Ma, Zijian He, Chia-Wen Kuo, Kan Chen, Peizhao Zhang, Bichen Wu, Zsolt Kira, and Peter Vajda. Unbiased teacher for semi-supervised object detection. *arXiv preprint arXiv:2102.09480*, 2021. 1
- [2] Rongchang Xie, Chunyu Wang, Wenjun Zeng, and Yizhou Wang. An empirical study of the collapsing problem in semi-supervised 2d human pose estimation. In *ICCV*, 2021. 1



Figure 2. Qualitative results of some example images in COCO dataset.