

# Supplementary Materials: Joint-Motion Mutual Learning for Pose Estimation in Video

Anonymous Authors

In this supplementary file, we first present additional results and analyses on the effects of different hyperparameters (Section 1). Then, we show more visual comparison results between state-of-the-art method TDMI [1] and our method JM-Pose on challenging scenes (Section 2). Finally, we demonstrate more visual results of our JM-Pose on benchmark datasets (Section 3).

## 1 SUPPLEMENTARY EXPERIMENTS

In this section, we conduct more experiments on the PoseTrack2017 validation set to study the effects of different hyperparameters, including the ratio  $\alpha$  of the trade-off loss (in the main paper Eq. 16) and temporal span  $\delta$  of supporting frames (in the main paper xxx).

**Ablation study on loss ratio  $\alpha$ .** In the Subsection 3.3 of the main paper, we utilize  $\alpha$  to balance the two loss terms in the total loss (in the main paper Eq. 16). We try utilizing different ratios and present the results in Table 1. We observe that performance drops from 86.4 mAP ( $\alpha = 0.25$ ) to 85.5 mAP ( $\alpha = 0.01$ ). This is consistent with our intuition that a smaller ratio hinder the effect of the proposed information orthogonality objective. Additionally, the model performance tends to plateau as  $\alpha$  increases and the best results (86.4 mAP) are obtained at  $\alpha = 0.25$ .

Table 1: Ablation study on the loss ratio  $\alpha$ .

Ratio $\alpha$	$\alpha = 0.01$	$\alpha = 0.25$	$\alpha = 1$
Mean	85.5	<b>86.4</b>	85.9
Declines	↓ 0.9	-	↓ 0.5

**Ablation study on temporal span  $\delta$ .** Furthermore, we investigate the impact of using different temporal spans  $\delta$ , which constrains the number of supporting (neighboring) frames. We experiment with three different temporal span settings: (a)  $\delta = 1$ , (b)  $\delta = 2$ , and (c)  $\delta = 3$ . From the results in Table 2, we observe that mAP improves with increasing  $\delta$ , namely 85.3 for  $\delta = 1$ , 86.4 for  $\delta = 2$ , and 86.4 for  $\delta = 3$ . This is in line with our intuition that more temporal information can be obtained by employing more supporting (neighboring) frames, which provides more cues for pose estimation. When  $\delta = 2$  or  $\delta = 3$ , the performance of pose estimation tends to be similar, we adopt a temporal span of  $\delta = 2$ .

Table 2: Ablation study on the temporal span  $\delta$ .

Span $\delta$	(a) $\delta = 1$	(b) $\delta = 2$	(c) $\delta = 3$
Mean	85.3	<b>86.4</b>	<b>86.4</b>
Declines	↓ 1.1	-	-

## 2 MORE VISUAL COMPARISON RESULTS

As shown in Fig. 1, we visualize the comparison results for challenging scenes, such as *self-occlusions*, *rapid-motion*, and *pose-occlusions*. We try to investigate the robustness of our method. In figure 1, we observe that TDMI [1] fails to estimate some hard-to-detect joints since TDMI ignores the implicit joint information encoded in initial heatmaps. In contrast, our method JM-Pose displays more accurate pose estimation results in the heavily crowded and challenging scenarios.

## 3 MORE VISUAL RESULTS

As shown in Fig. 2, we present more visual results of our method on challenging datasets. From this figure, we observe that our method achieves accurate and robust pose estimation in various scenes.

## REFERENCES

- [1] Runyang Feng, Yixing Gao, Xueqing Ma, Tze Ho Elden Tse, and Hyung Jin Chang. 2023. Mutual information-based temporal difference learning for human pose estimation in video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 17131–17141.

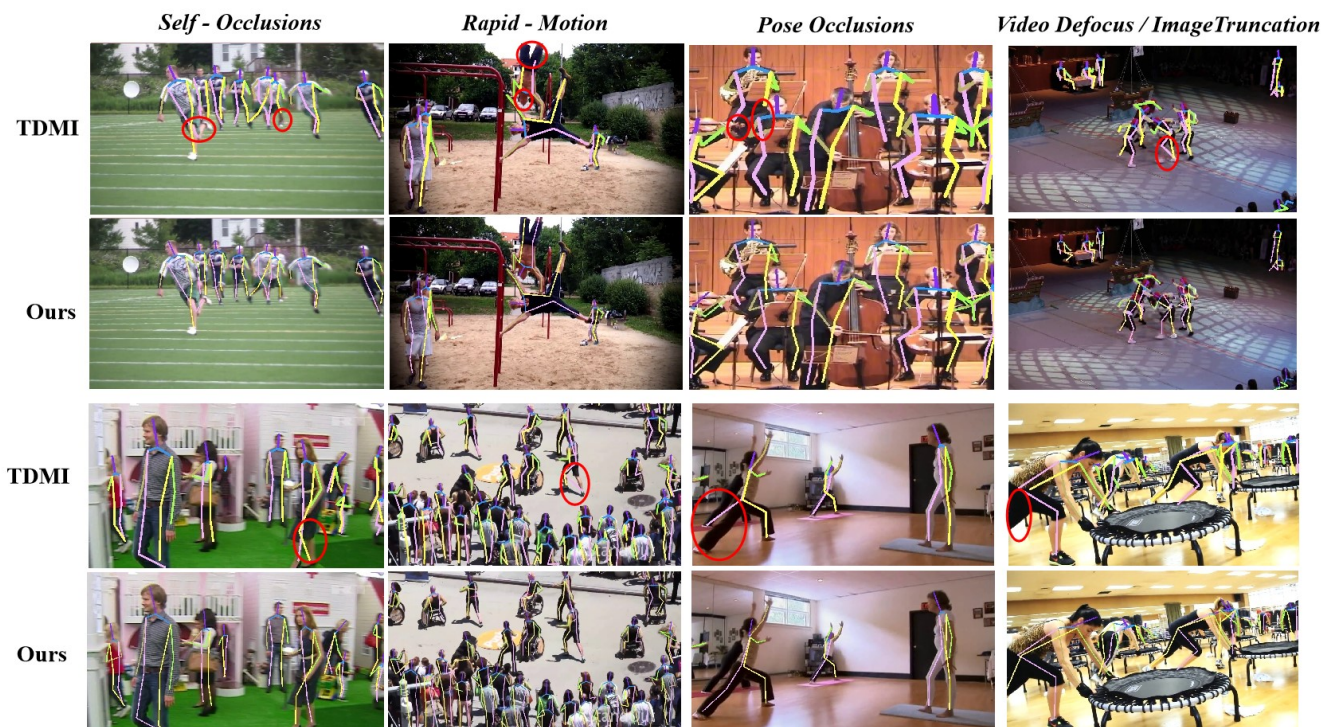


Figure 1: More visual comparisons of detection results obtained from state-of-the-art TDMI and our method JM-Pose on challenging scenes. Inaccurate detections are highlighted with red circles.



Figure 2: More visual results of our JM-Pose on benchmark datasets. Challenging scenes such as fast motion or occlusions are involved.