

# PoseAug: A Differentiable Pose Augmentation Framework for 3D Human Pose Estimation

## (Supplementary Material)

Kehong Gong <sup>\*†</sup>    Jianfeng Zhang <sup>\*</sup>    Jiashi Feng  
 National University of Singapore

gongkehong@u.nus.edu    zhangjianfeng@u.nus.edu    elefjia@nus.edu.sg

### A. Implementation and training details

**Implementation details** We use the single-frame version of VPose [7] as our pose estimator. The pose augmentor consists of four linear layers with Batch Normalization (BN) [3] and leaky ReLU [2]. It first transforms the input 3D pose  $\mathbf{X}$  to 256-D hidden features and then predicts the augmentation parameters  $\gamma_{ba}$ ,  $\gamma_{bl}$ ,  $\mathbf{R}$  and  $\mathbf{t}$  from the hidden features. Both the 3D and 2D pose discriminators consist of 4 linear layers with leaky ReLU and residual connection. The hidden feature dimensions are set as 256 and 100 for 3D and 2D discriminator, respectively. We train our model for 50 epochs on Human3.6M, with batch size of 1024. We adopt Adam optimizer [5] with linear decay and an initial learning rate of 0.001 for all components (*i.e.*, augmentor, discriminator, estimator). The hard ratio  $\beta$  linearly increases from 2 to 20 during the training process. The threshold for regularizing the augmentation parameters  $\gamma_{ba}$  and  $\gamma_{bl}$  is set as 0.1.

**Training details** With the differentiable augmentor design, the pose augmentor  $\mathcal{A}$  with parameters  $\theta_A$ , discriminator  $\mathcal{D}$  with parameters  $\theta_D$  and estimator  $\mathcal{P}$  with parameters  $\theta$  can be jointly trained end-to-end as illustrated in Algorithm 1. In detail, we firstly update the augmentor  $\mathcal{A}$  and discriminator  $\mathcal{D}$  alternatively by minimizing the pose augmentation loss Eqn. (11) and discrimination loss Eqn. (12) with the source pose pairs  $\{\mathbf{x}_i, \mathbf{X}_i\}_{i=1}^N$ , and augmented pose pairs  $\{\mathbf{x}'_i, \mathbf{X}'_i\}_{i=1}^N$ . We then use both the source and augmented pose pairs to train the estimator  $\mathcal{P}$  by minimizing the pose estimation loss Eqn. (8).

### B. More qualitative results

Here we present more qualitative results of the pose estimator (*i.e.*, VPose [7] (1-frame)) trained with and w/o our PoseAug on LSP [4], MPII [1], and MPI-INF-3DHP [6], as

<sup>\*</sup>Equal contribution. <sup>†</sup>Work done during an internship at huawei Singapore research center.

---

#### Algorithm 1: PoseAug training strategy

---

**Input:**  $N$  training samples  $\{\mathbf{x}_i, \mathbf{X}_i\}_{i=1}^N$ , training epochs  $E$  and learning rate  $\alpha$ .

**Output:** pose estimator  $\mathcal{P}$ , augmentor  $\mathcal{A}$  and discriminator  $\mathcal{D}$ .

**for**  $e = 1, \dots, E$  **do**

    // Update augmentor  $\mathcal{A}$

    Generate augmented sample  $\{\mathbf{x}'_i, \mathbf{X}'_i\}_{i=1}^N$  from  $\{\mathbf{x}_i, \mathbf{X}_i\}_{i=1}^N$  through the augmentor  $\mathcal{A}$

    Calculate the augmentation loss  $\mathcal{L}_A$  (Eqn. (11))

    Calculate gradients  $\nabla_{\theta_A} \mathcal{L}_A$  and update the parameters of  $\mathcal{A}$  by  $\theta_A = \theta_A - \alpha \nabla_{\theta_A} \mathcal{L}_A$

    // Update discriminator  $\mathcal{D}$

    Calculate the discrimination loss  $\mathcal{L}_D$  (Eqn. (12))

    Calculate gradients  $\nabla_{\theta_D} \mathcal{L}_D$  and update the parameters of  $\mathcal{D}$  by  $\theta_D = \theta_D - \alpha \nabla_{\theta_D} \mathcal{L}_D$

    // Update estimator  $\mathcal{P}$

    Calculate the estimation loss  $\mathcal{L}_P$  (Eqn. (8)) by feeding  $\{\mathbf{x}_i, \mathbf{X}_i\}_{i=1}^N$  and  $\{\mathbf{x}'_i, \mathbf{X}'_i\}_{i=1}^N$  alternatively to  $\mathcal{P}$

    Calculate gradients  $\nabla_{\theta} \mathcal{L}_P$  and update the parameters of  $\mathcal{P}$  by  $\theta = \theta - \alpha \nabla_{\theta} \mathcal{L}_P$

**end**

---

shown in Fig. S1. The comparison of inference results between with and w/o our PoseAug demonstrate that our auto-augmentation framework helps the pose estimator achieve better performance on challenge in-the-wild scenes. In addition, Fig. S2 shows three examples of source and corresponding augmented pose pairs. These examples demonstrate that our augmentor applies meaningful augmentation operation in generating harder cases, *e.g.*, harder posture (1st row), harder view point (2nd row), and even unseen pose (3rd row), which increases the data diversity and leads

to better generalization for estimator.

### C. Ablation on feedback

Besides, the error feedback strategy brings significant improvement on 3DPW dataset (from 134.1 to 130.3, Table S1). This clearly verifies its effectiveness, especially on challenging scenarios.

Table S1: Ablation study on error feedback strategy on 3DPW. Augmentation denotes the combination of BA, BL, RT operations.

Method	Augmentation	Feedback	PA-MPJPE ( $\downarrow$ )
Baseline			145.2
Variant A	✓		134.1 (-11.1)
PoseAug	✓	✓	<b>130.3</b> (-14.9)

### References

- [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *CVPR*, 2014.
- [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, 2015.
- [3] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv*, 2015.
- [4] Sam Johnson and Mark Everingham. Clustered pose and nonlinear appearance models for human pose estimation. In *BMVC*, 2010.
- [5] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICCV*, 2015.
- [6] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *3DV*, 2017.
- [7] Dario Pavlo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *CVPR*, 2019.

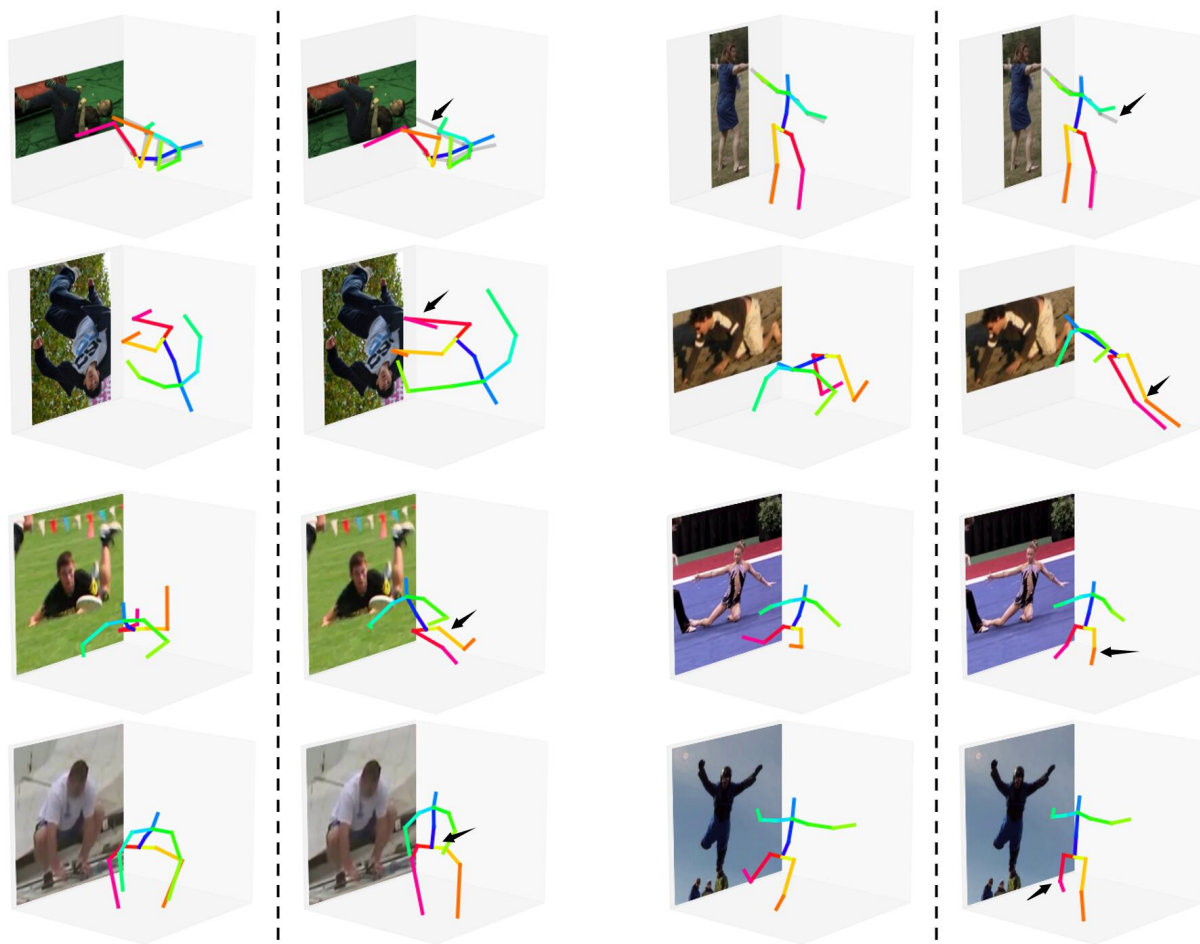


Figure S1: Example 3D pose estimations from 3DHP (1st row), LSP (2nd row), and MPII (3rd and 4th row). PoseAug results are shown in the 1st and 3rd columns. The 2nd and 4th columns show the results of *Baseline*, *i.e.*, VPose [7] (1-frame) trained without PoseAug. The grey skeletons in 3DHP (1st row) are ground truth. Errors are highlighted by black arrows.

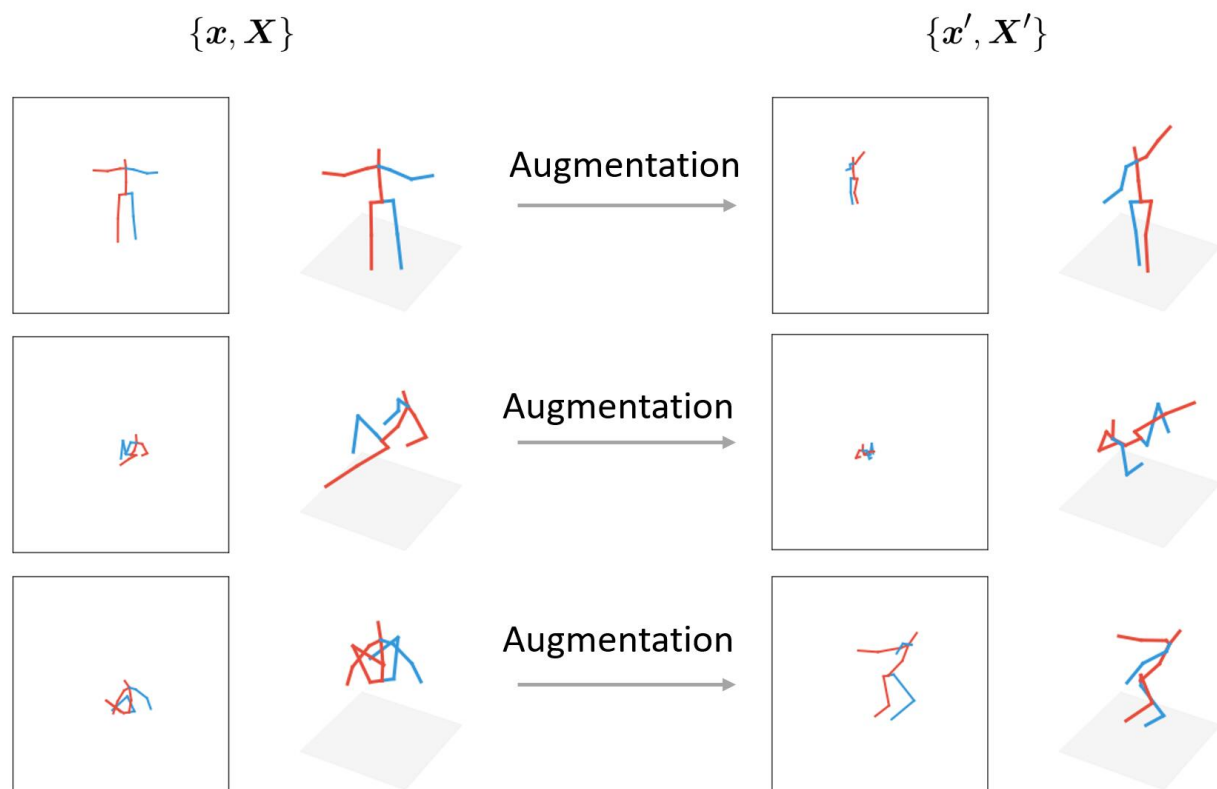


Figure S2: Examples of the source pose pair ( $x$ : 1st column,  $X$ : 2nd column) and its augmented pose pair ( $x'$ : 3rd column,  $X'$ : 4th column). The examples include harder posture (1st row), harder view point (2nd row), and unseen pose (3rd row).