

Method	PCK	AUC	MPJPE	Time[s]
Vanilla-S	82.5	47.6	94.1	0.244
Vanilla-B	83.7	48.4	91.3	0.083
Online-S	83.6	48.2	92.2	0.027
Online-B	84.1	48.7	90.5	0.007

Method	MPJPE
Baseline	45.0
Joint	42.6
ISO	41.8

Method	PCK	AUC	MPJPE
[7]	65.3	33.4	135.7
[9]	75.9	36.3	-
[35]	62.4	30.9	147.4
Ours	83.6	48.2	92.2

Method	PCK	AUC	MPJPE
[14]	60.9	30.2	150.1
[7]	65.3	33.4	135.7
Ours	83.6	48.2	92.2

1 **To R1 Q1. Batch adaptation.** ISO can perform model adaptation with a batch of instances, if available. See in
 2 Table R1, where **B** and **S** denote batch and single-instance ISO. Batch ISO outperforms single-instance ISO in both
 3 accuracy (from 94.1 to 91.3 in MPJPE for *Vanilla*) and efficiency ($\sim 3\times$ and $4\times$ faster for *Vanilla* and *Online*).

4 **Q2. Comparison based on faster version of ISO** will be added in the revision.

5 **Q3. The pose discriminator** is pretrained on the source data. During inference, it is updated using Eqn. (3) for each
 6 new target instance. To alleviate potential training bias from a single instance, we applied horizontal flip augmentation
 7 and the augmentation strategy (see Line 159-161) to update the discriminator. We will add more details in the revision.

8 **Q4. Ablation studies on w/ and w/o ISO over H36M** are given in Table R2. We observe *Joint* improves *Baseline* by
 9 a large margin. ISO makes marginal improvement over *Joint* (from 42.6 to 41.8 in MPJPE) since training and testing
 10 distribution are similar. However, even though there is no significant distribution shift, ISO still makes positive effect.

11 **To R2 Q1. Differences between ISO and model-fitting (MF) methods:** 1) The motivation and methodology are
 12 different. The MF methods (e.g., Kolotouros et al., CVPR 2019) aim to improve the model training by iterating
 13 regression (using a parametric human body model) and optimization. However, they do not consider how to generalize
 14 the model to new testing distributions that are different from the training ones. Our ISO focuses on mining distributional
 15 knowledge about the testing distribution from unlabeled instances via SSL and adapt the models accordingly to gain
 16 better generalizability. 2) The MF methods usually require a parametric model (e.g., SMPL) and a 3D body mesh
 17 initialization for model fitting; whereas ISO does not require these and thus is more general.

18 **Q2. “Main improvements come from Joint”** is only true for *Vanilla* as *Vanilla* is always re-initialized using *Joint* for
 19 each new instance. Thus, *Vanilla* cannot benefit the model too much as distributional information mined from a single
 20 instance is limited. When more distributional information is mined, ISO can make significant improvement. This can
 21 be observed from: 1) *Online* improves *Joint* by -4 in MPJPE (see Table S1). 2) Batch version of *Vanilla* improves the
 22 performance upon single-instance version of *Vanilla* by a large margin (see Q1 in responses to R1 and Table R1).

23 **Q3. Differences from un-/weakly-supervised methods:** un-/weakly-supervised methods are usually used for better
 24 training the model with more unlabeled *training data* (see Line 72-78) and rarely used in a transductive manner for
 25 testing; while ISO aims to improve the model’s generalizability during inference via SSL. Among the suggested methods
 26 [7, 9, 35], [7] uses cycle-supervision and can be used in a transductive manner; while [9,35] requires multi-view data
 27 which are not available on 3DHP and 3DPW during testing. We compare ISO with them on 3DHP in Table R3. ISO
 28 significantly outperforms them. Comparison with them on 3DPW will be added in revision.

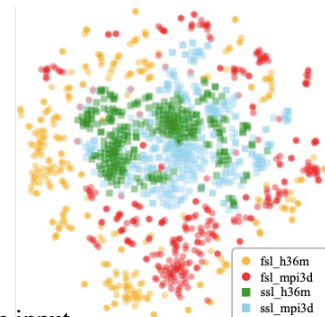
29 **Q4. Qualitative results.** We will visualize GT in figure 4 and show some results with noisy inputs in the revision.

30 **To R3 Q1. 2D pose input.** For fair comparison with [11, 46, 50], in Table 1, we use GT poses as inputs. However,
 31 ISO is robust to 2D pose noise as shown in Table 4 and thus can perform similarly well when using detected poses.

32 **Q2. SSL methods under transductive learning (TL) setup.** [A] cannot be evaluated under TL setup on 3DHP since
 33 it requires multi-view data while the test set only has single-view data. We re-implement other SSL methods [7, 14] and
 34 evaluate them on 3DHP under TL setup (see in Table R4). We can observe ISO outperforms them by a large margin.

35 **Q3. Relation between ISO and TL methods** will be added in the revision.

36 **To R4 Q1. Questions about model components.** To study these, we visualize
 37 the hidden features of different samples (random samples from H36M test set and
 38 samples from 3DHP) extracted by both heads using t-SNE (see right figure). The
 39 samples from 3DHP are the ones showing improvement after ISO, which present
 40 novel viewpoints and body sizes. We observe features from SSL head of both
 41 datasets (green & blue squares) are closer than the ones from FSL head (orange
 42 & red circles). Thus, SSL head learns features more robust to distribution change
 43 whereas FSL head learns more discriminative features from source data. The
 44 features from these two heads are thus disentangled. Shared feature extractor keeps
 45 these two kinds of features. FSL head would fail when facing unusual poses and
 46 viewpoints (e.g., top views) which are very ambiguous with only monocular 2D pose input.



47 **Q2. Will online ISO overfit to the samples come at first?** To answer this, we randomly shuffle the test set of 3DHP
 48 before performing online ISO. We conduct experiments for 8 times and obtain the statistics: PCK: 83.2 ± 0.43 , AUC:
 49 48.0 ± 0.30 and MPJPE: 92.9 ± 1.3 . The small variance implies online ISO does not overfit to the sequential data. As
 50 discussed in Q1, the model overfitting to the SSL task would extract less discriminative features and thus hamper the
 51 performance. We will clarify this in the revision.