

PAFUSE: Part-based Diffusion for 3D Whole-Body Pose Estimation

Nermin Samet¹, Cédric Rommel¹, David Picard², and Eduardo Valle¹

¹ Valeo.ai, Paris, France

² LIGM, Ecole des Ponts, Univ Gustave Eiffel, CNRS, Marne-la-Vallée, France

Abstract. We introduce a novel approach for 3D whole-body pose estimation, addressing the challenge of scale- and deformability-variance across body parts brought by the challenge of extending the 17 major joints on the human body to fine-grained keypoints on the face and hands. In addition to addressing the challenge of exploiting motion in unevenly sampled data, we combine stable diffusion to a hierarchical part representation which predicts the relative locations of fine-grained keypoints within each part (e.g., face) with respect to the part’s local reference frame. On the H3WB dataset, our method greatly outperforms the current state of the art, which fails to exploit the temporal information. We also show considerable improvements compared to other spatiotemporal 3D human-pose estimation approaches that fail to account for the body part specificities. Code is available at <https://github.com/valeoai/PAFUSE>.

Keywords: 3D Whole-body Pose Estimation · Diffusion for Pose Estimation · 2D to 3D Whole-body Pose Lifting

1 Introduction

As a challenging computer vision task, 3D human pose estimation aims to localize 3D human body keypoints in images and videos. 3D Human pose estimation has an important role in several vision tasks and applications such as action recognition [24, 29, 57, 60], human mesh recovery [16, 23], motion generation [51, 56], sign language [17, 22, 26, 34], augmented/virtual reality [1, 4], and robotics [9–12, 48]. In recent years, the biggest leap forward in 3D human pose estimation was including the temporal aspect and predicting entire sequences of skeletons at once [61, 63]. Indeed, adding the temporal information removes some of the depth-scale ambiguities since predicted poses also have to be compatible with the motion of the human body.

3D whole-body pose estimation further expands its scope by aiming to detect not only standard human body keypoints but also face, hand, and foot keypoints to enable more precise applications. Recently Zhu *et al.* [64] addressed the missing dataset and benchmark issue by introducing a novel dataset for 3D whole-body pose estimation, called Human3.6M 3D WholeBody or H3WB for short. H3WB extends the Human3.6M keypoints dataset [3, 15] by further annotating

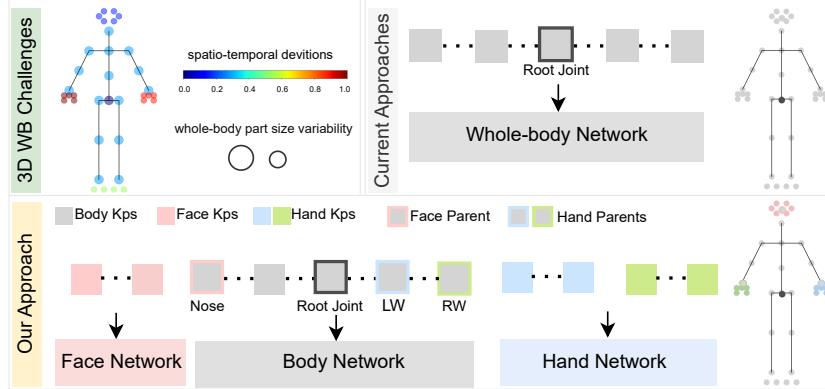


Fig. 1: In a whole-body skeleton, different keypoints have different scales and variations (top left) which presents a challenge for spatio-temporal prediction. Current approaches (top right) process all keypoints in a single network and, as such, have difficulties adapting to the different statistics of each body part. Our approach (bottom) groups keypoints by body parts that share similar behavior and processes them with dedicated networks, allowing better-adapted predictions.

face, hand, and foot keypoints according to the layout of the COCO Whole-Body dataset [18] (see Fig. 2), and thus provides a unified benchmark for 3D whole-body pose estimation compatible with existing 2D whole-body methods. 3D whole-body human pose estimation is much more challenging than regular 3D pose estimation in that the additional keypoints have different scales and diversity in poses. That is exacerbated by considering motion, since some body parts have a much greater motion range (*e.g.*, hands compared to the hip, see Fig. 1 top left).

In this paper we tackle 3D whole-body human pose estimation from such temporal perspective. Inspired by the progress in 2D whole-body pose estimation and the success of stable diffusion for 3D pose estimation, we propose a new part-based diffusion approach, PAFUSE, that explicitly handles the aforementioned challenges. To this end, we build a 3-level 3D pose estimation network for body, hands and face. We condition each body part on their part root joint, leading to a hierarchical system. We jointly train each body part and estimate all whole-body keypoints in their specific coordinate systems (see Fig. 1 bottom). Such hierarchical bottom-up approach brings two benefits. First, the scale variance among different body parts are handled naturally as the relative distances within each part are in a similar range and each part-type is processed by a separate sub-network. Second, it brings flexibility to adjust the capacity of each sub-network according to the difficulty and deformability level of its corresponding body part. Additionally, our method is modular and can be easily added to any 3D pose estimation method to extend them to whole-body [25, 37, 38, 43, 61, 63].

To test our ideas, we extend the H3WB [64] dataset to spatio-temporal prediction. We show that our part-based approach is able to provide significant improvement to a wide variety of spatio-temporal baselines, and that our diffu-

sion based method PAFUSE is able to obtain state of the art results at 41.4mm MPJPE (against 88.3mm for the previous state of the art) while also generalizing to in-the-wild sequences.

To sum up, we are the first to successfully combine a part-based approach with denoising diffusion in 3D whole-body pose estimation, thus handling the variance of both scale and motion among body parts. Our contributions are the following:

- ✓ We introduce a hierarchical part-based approach for 3D whole-body human pose estimation that solves the scale and motion variation issues without any additional computational cost.
- ✓ Based on this part-based approach, we propose PAFUSE, a denoising diffusion model for 3D whole-body human pose estimation that obtains state-of-the-art performance.
- ✓ We extend the H3WB dataset to spatio-temporal prediction and perform an extensive comparison of state-of-the-art 3D human pose estimation methods on this new benchmark by adapting recent method from the literature with our hierarchical part-based approach.

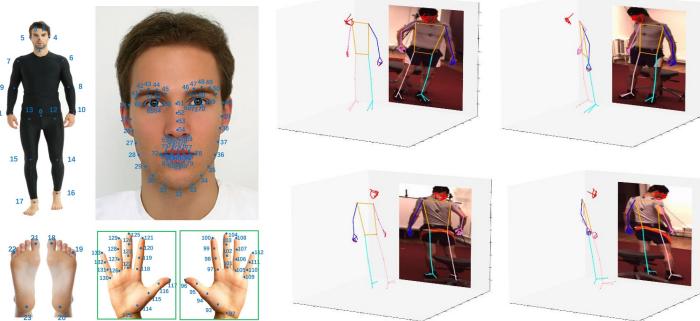


Fig. 2: (Left:) COCO-whole-body [18] layout used in the H3WB [64] dataset, with 133 keypoints. In addition to the standard 17 main-body keypoints, there are 68 face-, 42 hand- (21 keypoints for each), and 6 foot- (3 for each) keypoints. (Right:) Example 2D and 3D whole-body pose pairs from the H3WB dataset. Images taken from [18, 64].

2 Related Work

3D whole-body pose estimation Before the H3WB [64] dataset, 3D whole-body pose estimation methods can be categorized into two groups: parametric models and distillation-based non-parametric approaches. The majority of works belong to the first category, based on parametric human body models such as Adam [19] and SMPL-X [35]. For example, MTC [53] builds upon the Adam model [19] by optimizing its parameters after extracting 2.5D predictions. Another popular parametric model, SMPLify-X, optimizes the parameters of the

SMPL-X model [35] to align with 2D keypoints. Parametric models offer the advantage of sampling nearly infinite keypoints from the mesh [6–8]. However, they have several drawbacks, including being slow and sensitive to parameter initializations. Moreover, their accuracy typically falls behind that of detection-based methods, especially on finer body parts like hands [64].

On the other hand, nonparametric methods [5, 39, 52] employ different strategies to circumvent heavy optimization procedures. Both DOPE [52] and FrankMocap [39] initially train separate models for the body, hands, and face, which are subsequently integrated within a unified learning framework. DOPE [52] obtains pseudo-ground annotations from these individual body models and utilizes them to supervise the distillation model. Similarly, ExPose [5] begins by obtaining a pseudo-ground truth dataset through fitting the SMPL-X model to in-the-wild images, then proceeds to train a joint model to generate whole-body poses. One significant drawback of these methods is their reliance on different datasets for each body part. Consequently, each method produces varying whole-body layouts with differing numbers of keypoints.

Recently, Zhu *et al.* utilized existing image-based 3D pose estimation methods on the H3WB dataset [64] to establish a benchmark. However, as these methods do not directly tackle the challenges of 3D whole-body pose estimation, their performances are suboptimal.

2D Whole-Body Pose Estimation Following the release of the COCO WholeBody dataset [18], significant progress has been made in 2D whole-body pose estimation. Several methods have been proposed, with the primary focus being to overcome the challenge of scale variance in whole-body pose estimation tasks. ZoomNet [18] and ZoomNas [54] are among the pioneering methods in this field. They initially predict body keypoints and subsequently address scale variance by cropping the hand and face areas and transforming them to higher resolutions to further refine face and hand keypoint estimation. HPRNet [41] employs a different bottom-up strategy where keypoints on each body part are separately regressed after being offset according to its part-center. TCFormer [59] introduces dynamic tokenization to handle size variance in body parts by clustering tokens. Keypoint Communities [58] assigns different weights to body parts after constructing a skeleton graph. Seong *et al.* [42] propose a keypoint-wise adaptive method to handle the scale difference between body parts.

3D human pose estimation We can categorize the monocular 3D human pose estimation into two groups: deterministic and generative approaches. Early deterministic approaches, followed end-to-end pipeline to predict 3D keypoints directly from images [31, 32, 36, 47]. Later, two-stage approach become dominant where first 2D keypoints are predicted and then they are lifted into 3D using lightweight networks [30]. Later in order to exploit temporal information video based approaches are proposed based on convolutional neural networks [37] and graph convolutional networks [2, 27, 55, 62, 65]. More recently, attention based spatial-temporal transformer architectures gained attention such as MixSTE [61] and Poseformer [63].

Lifting 2D pose to 3D space is inherently ambiguous as multiple 3D poses can map onto the same 2D input. One drawback of deterministic methods lies in their inability to effectively address this ambiguity. Motivated by this challenge and limitations of deterministic methods, researchers explored multi-hypothesis approaches, including variational autoencoders [45], normalizing flows [20], and, more recently, diffusion models [14, 38, 44]. D3DP [44] adopts the MixSTE backbone as its denoiser and directly integrates raw 2D keypoints as conditions. In contrast to employing a straightforward averaging approach for hypothesis aggregation, it introduces a novel method based on 2D reprojections. DiffPose [14] employs a diffusion based on Gaussian mixture models trained on 2D heatmaps. Additionally, apart from utilizing MixSTE as a denoiser, it employs a pre-trained MixSTE to initialize the reverse diffusion during inference. On the other hand, the denoiser architecture of DiffHPE [38] is based on CSDI [49] and TCD [40], where they use graph-convolutional layers to strike a good balance between computation and accuracy. During both training and inference, they utilize a pre-trained MixSTE for conditioning.

Similar to D3DP [44] and DiffPose [14], our proposed model PAFUSE, also uses MixSTE as its denoiser network. We also follow [38, 44] and incorporate raw 2D part keypoints to condition diffusion process as it is shown to be simple and very effective.

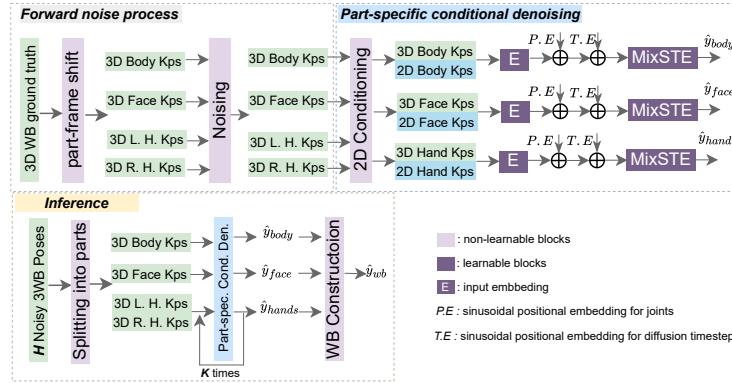


Fig. 3: Overall processing pipeline. During training, we split the input samples into body-part-specific tensors before performing the forward noising process (top-left). Then we train the part-based conditional denoising diffusion models (top-right). During inference (bottom), we start from random Gaussian noise and iterate the part-based conditional denoising diffusion models K times to obtain the skeleton parts, which we reconstruct into a whole-body skeleton. For simplicity, we omitted the temporal aspect, although our method actually processes video sequences consisting of N frames.

3 Method

We assume 2D keypoints are available and focus on lifting them to the 3D space. Thus, our task is to convert an input tensor $x \in \mathbb{R}^{2 \times J \times N}$ into a predicted output

tensor $\hat{y} \in \mathbb{R}^{3 \times J \times N}$, where J is the number of joints (keypoints) in the reference skeleton, and N is the length of the video sequence in frames. We assume that a training dataset of (x, y) pairs is available, y being the ground-truth 3D output. The subtensors \bullet_{body} , \bullet_{hands} , \bullet_{face} refer to the joints on the body parts indicated in the subscript. We indicate joint-wise tensor concatenation by a vertical bar, e.g., $x = x_{\text{body}} | x_{\text{hands}} | x_{\text{face}}$.

Overview Fig. 3 shows the overall pipeline of PAFUSE. From a sequence x of 2D keypoints representing a whole-body pose evolving in time, we predict the corresponding 3D positions. PAFUSE combines a part-based approach with a diffusion-based inference. The former splits the keypoints into main body, face, and hands, anchoring each part into a local frame of reference, as explained below. The latter leverages generative modelling, separately on each part, to propose multiple pose hypotheses. The combined approach allows addressing both the depth-ambiguity inherent to monocular human-pose estimation, and the specific needs of body parts with very different scales, motions, and deformabilities.

Whole-Body-Frame vs. Part-Frame Shift. Usually, y is represented in a way that it sets the body root joint (keypoint 0, at the center of the hips) at the origin of the coordinate system. We propose an alternative representation, with different body parts centered at local coordinate frames: keypoint 0 for the main body, the nose (keypoint 1) for the face, and the corresponding wrists (keypoints 10 and 11) for each hand.

Generative-based prediction. Our learning engine is a denoising diffusion probabilistic models (DDPM) [13], typically employed for generative tasks, but here leveraged to model the distribution $\Pr(y|x)$. DDPMs are trained to reverse the incremental addition of Gaussian noise to data, until they are able to “restore” convincing data samples from pure Gaussian noise.

Forward noise process. Following the usual DDPM procedure, the diffusion process samples a “timestep” $t \sim U(0, T)$ and then corrupts ground-truth 3D tensor by adding a Gaussian noise whose intensity increases with t , i.e., $y^t = y + \epsilon^t$, $\epsilon^t \sim \mathcal{N}(0, f(t)I)$, where f is a monotonically increasing function of t , which is designed according to the desired noising schedule. It should be clear that the idea of “time” expressed by t refers to the noising process and has nothing to do with the pose motion: the entire pose-in-motion tensor (with its N frames) is noised and denoised all at once. In this work, we did not optimize the noise scheduler and use a popular cosine noise scheduler [33].

Part-specific conditional denoising. In PAFUSE, we learn independent denoising models for the main body, the face, and the hands. At the denoising step, we split the noised ground-truth, such that $y^t = y_{\text{body}}^t | y_{\text{face}}^t | y_{\text{L-hand}}^t | y_{\text{R-hand}}^t$. The denoising is conditioned by the timestep t , which we positionally-encode with a family of sinusoidal functions before feeding it to the decoder (similar to how token positions are encoded in transformers). Crucially, the denoising of each part is also conditioned on the corresponding input 2D tensor (e.g., x_{face} for y_{face}^t), by simply concatenating the tensors across the space dimension (i.e., conditioned inputs $\in \mathbb{R}^{5 \times J_{\text{part}} \times N}$). Note that input 2D tensors are used as it is, i.e., no off-

setting for each part. Although DDPM denoising is often learned stepwise, by asking the denoiser to reconstruct \hat{y}^{t-1} from y^t , here we follow D3DP [43] and learn to reconstruct the denoised input directly:

$$\hat{y}_{\text{part}} = D_{\text{part}}(t, x_{\text{part}}, y_{\text{part}}^t) \quad (1)$$

where D_{part} is the part-specific denoising network.

Conditioned inference.

Test time inference predicts \hat{y} by concatenating the 2D input tensor x to a pure-noise tensor sampled from a Gaussian distribution $\hat{y}^T \sim \mathcal{N}(0, I)$. Although the training objective is one-step denoising (Eq. (1)), inference still requires incremental denoising ($\hat{y}^T \rightarrow \hat{y}^{T-1} \rightarrow \dots \rightarrow \hat{y}$) to produce high-quality samples from the expected conditional distribution. Following D3DP, we employ DDIM [46] to bridge the gap between the one-step results produced by the learned denoiser and the high-quality incremental denoising needed as final results.

Losses. We experimented with two losses. The *part loss* considers the keypoints in their local (part-based) frames of reference. Because the denoisers are trained in those local frames, that corresponds to computing the loss between the reassembled keypoint tensors, without further processing:

$$\mathcal{L}_{\text{part}} = \ell(\hat{y}_{\text{body}} | \hat{y}_{\text{hands}} | \hat{y}_{\text{face}}, y_{\text{body}} | y_{\text{hands}} | y_{\text{face}}) \quad (2)$$

where ℓ is either the MPJPE metric or the mean-squared error (MSE).

The whole-body or the *WB loss* uses the whole-body frame of reference, and thus, must recover the global coordinates from the estimated root-joint positions of each local coordinates. Let r be the tensor containing the position (relative to keypoint 0) of the local coordinates of each joint, such that $y_{\text{wb-frame}} = y_{\text{part-frame}} + r$ and $\hat{y}_{\text{wb-frame}} = \hat{y}_{\text{part-frame}} + \hat{r}$. Then, the WB loss is:

$$\begin{aligned} \mathcal{L}_{\text{WB}} = \ell(&\hat{y}_{\text{body}} + \hat{r}_{\text{body}} | \hat{y}_{\text{hands}} + \hat{r}_{\text{hands}} | \hat{y}_{\text{face}} + \hat{r}_{\text{face}}, \\ &y_{\text{body}} + r_{\text{body}} | y_{\text{hands}} + r_{\text{hands}} | y_{\text{face}} + r_{\text{face}}) \end{aligned} \quad (3)$$

where ℓ is, again, either the MPJPE or the MSE.

Our experiments show that the part loss works better than the WB loss, and that the variant with $\ell = \text{MPJPE}$ works better than $\ell = \text{MSE}$ (Sec. 4.3).

Novelty. The *part-based design*, associating the use of local frames of reference for each body part and the use of separate part-based models are the main novelty of PAFUSE, in addition to the application of those techniques to a *diffusion-based inference*. Our results show that our part-based design systematically improves techniques and results in state-of-the-art results when associated to diffusion.

4 Results

4.1 Experimental setup

Dataset. Despite its growing importance, 3D whole-body pose estimation has only recently been addressed by the computer-vision community, and there is

only one annotated video whole-body dataset publicly available, the Human 3.6M WholeBody (H3WB) dataset [64], which is a reannotation of the Human 3.6M (H36M) dataset [3, 15]. It contains 10^5 ground-truth triplets of frames, 2D coordinates and 3D coordinates in camera space obtained from subjects S1, S5, S6, S7, S8; with 80% of triplets for training (S1, S5, S6, S7) and 20% for testing (S8). Remark that, while the frames of H36M are evenly sampled, the annotated frames of H3WB are not, posing additional challenges. Our experiments follow the H3WB benchmark protocol and use the ground-truth 2D keypoints provided in the H3WB dataset (Fig. 2). For in-the-wild videos, we extract 2D whole-body keypoints with OpenPifPaf [21].

Since H3WB is unique in literature, we sought ways to mitigate the use of single dataset in our experiments and provided results on challenging in-the-wild videos featuring intense occlusions and motion blur.

Metrics. Our main evaluation metric is the mean-per-joint-position error (MPJPE), which dominates the literature of 3D human pose estimation. We compute it across all whole-body keypoints, following H3WB’s official benchmark, after translating the root joint to its accurate position (“Protocol #1”). That metric appears as **WB** in our tables. We further compute part-specific MPJPE scores for the **Body**, **Hands**, and **Face**, by considering only the keypoints of those parts, after aligning them to their respectively local root joints (keypoints 0, 1, 10 and 11, respectively for body, face, left hand and right hand). Finally, the average of those part-specific MPJPE appears as **PB** in our tables.

Following the usual practice for generative-based models, including our baseline D3DP [43], we mainly evaluate our method by selecting the hypothesis that most closely aligns with the ground truth (**P-Best**). We also employ an averaged hypothesis (**P-Agg**) as an auxiliary additional metric.

Training and inference details We employ the AdamW optimizer, with momentum parameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$, and a weight decay of 0.1. Consistent with MiXSTE [61], we train for 400 epochs, starting with an initial learning rate of 6×10^{-5} . Due to the uneven sampling of H3WB dataset, we limited the experiments to frame windows of $N = 27$ and 81, shorter than the typical maximum of $N = 243$ found in the 3D human pose estimation literature. For $N = 27$, we set the batch size to 36, and for $N = 81$, to 12. For training, we set the number of hypotheses $H = 1$, and the sampling iterations $K = 1$. For inference, we evaluate variations on those settings ($H = 1, 20$, and 300, and

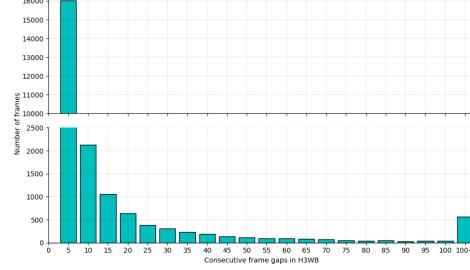


Fig. 4: Distribution of gaps between annotated frames of H3WB dataset, showing a long tail with many gaps of 100 or more frames. Contrast to Human3.6M, which is evenly annotated at every 5 frames. Please remark the discontinuity at the y-axis, to make room for the mode at 5 frames.

$K = 1$ and 5). Unless otherwise specified, we trained PAFUSE on the MPJPE-based part loss (Eq. (2)), sequence length $N = 27$, a single network for both hands, and MixSTE [61] denoiser backbones with 384, 256 and 224 channels, respectively, for the body, hands, and face.

Following D3DP [43], we implement PAFUSE in PyTorch. We train our models on a single NVidia A100 40 GB GPU for 400 epochs, which takes about 18 hours.

Method	WB	PB	Body	Face	Hands
<i>Single-frame</i>					
SMPL-X [35]	188.9	55.3	166.0	23.7	44.4
CanonPose [50]	186.7	58.1	193.7	24.6	48.9
SimpleBaseline [30]	125.4	45.5	125.7	24.6	42.5
CanonPose [50] <i>with 3D supervision</i>	117.7	39.5	117.5	17.9	38.3
Large SimpleBaseline [30]	112.3	34.9	112.6	14.6	31.7
Jointformer [28]	88.3	36.0	84.9	17.8	43.7
<i>Spatio-temporal</i>					
Videopose ($N = 27$)	70.1	27.8	62.2	11.5	34.7
+our part-based design	68.4	22.4	57.2	7.9	26.1
Videopose ($N = 81$)	71.7	28.9	64.9	11.8	36.1
+our part-based design	70.5	22.7	57.2	8.1	26.5
Poseformer ($N = 27$)	59.8	21.6	53.4	7.5	26.4
+our part-based design	58.0	18.2	49.7	5.7	20.6
Poseformer ($N = 81$)	68.7	28.7	62.4	12.4	36.6
+our part-based design	59.0	18.9	47.9	6.3	22.6
MixSTE ($N = 27$)	55.3	21.1	49.5	8.3	25.5
+our part-based design	52.8	20.2	45.2	8.1	24.1
MixSTE ($N = 81$)	54.5	20.7	48.1	8.0	25.7
+our part-based design	50.6	19.8	46.1	7.7	24.3
<i>Spatio-temporal + diffusion</i>					
D3DP($N = 27, H = 1, K = 1$)	50.9	20.1	46.3	8.2	24.3
PAFUSE ($N = 27, H = 1, K = 1$)	45.6	16.7	37.8	6.1	21.9
D3DP($N = 27, H = 20, K = 10$, P-Agg)	49.9	19.6	45.5	7.8	23.8
D3DP($N = 27, H = 300, K = 5$, P-Agg)	50.1	19.4	45.2	7.6	23.6
PAFUSE ($N = 27, H = 20, K = 10$, P-Agg)	45.5	16.6	37.4	5.8	22.3
PAFUSE ($N = 27, H = 300, K = 5$, P-Agg)	45.6	16.6	37.5	5.9	21.9
D3DP($N = 27, H = 20, K = 10$, P-Best)	46.9	19.4	44.6	7.9	23.6
D3DP($N = 27, H = 300, K = 5$, P-Best)	45.3	18.6	43.0	7.5	22.5
PAFUSE ($N = 27, H = 20, K = 10$, P-Best)	43.0	16.4	37.1	5.8	21.7
PAFUSE ($N = 27, H = 300, K = 5$, P-Best)	41.4	15.9	36.3	5.9	20.5

Table 1: State-of-the-art comparison for 2D→3D lifting on the H3WB dataset (MPJPE in mm). Despite H3WB’s irregular frame sampling, spatio-temporal methods consistently improve results, as does our **part-based design**. PAFUSE, combining multiple generative hypothesis and our part-based design has the best results. **The part-based models have the same total number of parameters as their vanilla counterparts.**

4.2 Comparison with the state-of-the-art

Baselines. Single-frame methods currently dominate the state of the art on H3WB [28, 30, 50]. We added to those baselines well-established spatio-temporal 3D human pose estimation methods, selecting those that allow reproducible results [37, 61, 63]. The most important baseline is D3DP [43] which is both spatio-temporal and diffusion-based, but not part-based. We adapted all baselines to accept the 133 keypoints of whole-body poses, and trained them according to their official instructions. The results of the state-of-the-art comparison appear in Tab. 1.

Spatio-temporal consistently outcompete single-frame. Despite the irregular frame-sampling of H3WB (see Fig. 4), which poses challenges for spatio-temporal methods, they still systematically outcompete single-frame ones. However, due to that irregular sampling, all methods favor short windows ($N = 27$), instead of the long windows ($N = 81$) traditionally preferred in literature [14, 37, 38, 43]. Indeed, with such long windows, some of the sequences will have very large gaps between the frames leading which makes them virtually impossible to learn because of their rarity in the training set and the lack of correlation between the frames that are to be predicted.

Our part-based design consistently improves techniques. Our part-based design consistently improves all spatio-temporal techniques, in all metrics, sometimes by several mm, as shown in the [highlighted lines](#) of Tab. 1. Remark that, for a fair comparison, we used the same total number of parameters in the part-based models as for the original models, which means that there is no additional memory cost for our part-based design. The most dramatic improvements appear, as expected, on hands and face, but the body keypoints also benefit from the part separation.

Precise prediction of keypoints outcompete body mesh generation. The results demonstrate that accurately predicting keypoints for the body, face, and hands is more effective than generating body meshes. SMPL-X recorded a MPJPE of 188.9, whereas all keypoint prediction methods outperformed SMPL-X. Our method, in particular, achieved a significantly lower MPJPE of 41.4.

PAFUSE obtains the best results. D3DP showcases that the multiple hypotheses of diffusion allow locating better poses than single-hypothesis techniques. However, the combination of diffusion and part-based of PAFUSE is the one that presents the best results, on all metrics.

The closest method to ours, DOPE [52], based on distillation across different body parts, was omitted from our benchmark due to failing to consistently predict complete skeletons on H3WB, often missing occluded parts, thus precluding meaningful scores.

Qualitative evaluation. We present several qualitative evaluations in Figs. 5 and 6. Fig. 5 shows visual results from the test set of H3WB dataset. We observe

that our method can successfully predict the body, hand and face under challenging deformations. In particular, it tends to better predict body joint that are key for the body parts, like shoulders and knees, which then leads to better aligned extremities like hands and face. Furthermore to show the robustness of our method we provide visual results from in-the-wild scenarios in Fig. 6. Even under challenging conditions such as in row 2 and 3, our method predicts the corresponding 3D poses of the hardly visible left hand of the tennis player or the blurry right hand of ballerina. Notice also that the cameras of these in the wild scenario are very different from the ones of H3WB (different position, focal length, fov, etc) and yet our prediction are well aligned with the whole-body pose of the subject. We provide more in-the-wild qualitative results in the supplementary material.

Method	Shift	Denoisers	WB	PB	Body	Face	Hands
D3DP [43] (baseline)	✗	✗	46.9	19.4	44.6	7.9	23.6
Only part-frame shift	✓	✗	95.5	33.6	123.1	7.7	24.6
Only part denoisers	✗	✓	89.6	21.6	41.2	10.1	29.2
PAFUSE	✓	✓	43.0	16.4	37.1	5.8	21.7

Table 2: Ablation on main components of PAFUSE, all models with the same total number of parameters (MPJPE metric in mm). Both components are crucial to PAFUSE’s performance.

4.3 Ablations

Unless specified otherwise, we set the PAFUSE’s parameters as $N = 27$, $H = 20$, and $K = 10$ in the ablations, and measure the metrics with the P-Best protocol.

Contribution of components. The impact of each key element introduced in PAFUSE is dissected in Table 2, where the baseline is a D3DP [43] minimally adapted to accept 133 whole-body keypoints.

The two main components of the part-based approach proposed in PAFUSE are the part-frame shift of keypoints, which represents their position in a frame local to their local root joints (Sec. 3), and the part-specific denoisers, which employs three separate denoising models for each part (hands, face, and main body). In the ablation with *only part-frame shift*, we employ a single 512-channel MixSTE denoising backbone, with a weighted loss to compensate the imbalance in the amount of keypoints across body parts. In the ablation with *only part denoisers*, all parts are frame-shifted to a single root joint (keypoint 0), but we reintroduce the three part-specific denoisers of 384, 224, and 256 channels, respectively, for the body, face, and hands. Remark that all models in Table 2 have **essentially the same total number of parameters**.

Our ablation shows that both contributions are essential for improving the results over D3DP. This is expected and fairly intuitive as both steps are complementary. Indeed, using only the part-frame shift with a single network introduces a bias that makes, for example, the hands closer to the hip than to the elbow because of this centering. This can create spurious correlation between unrelated

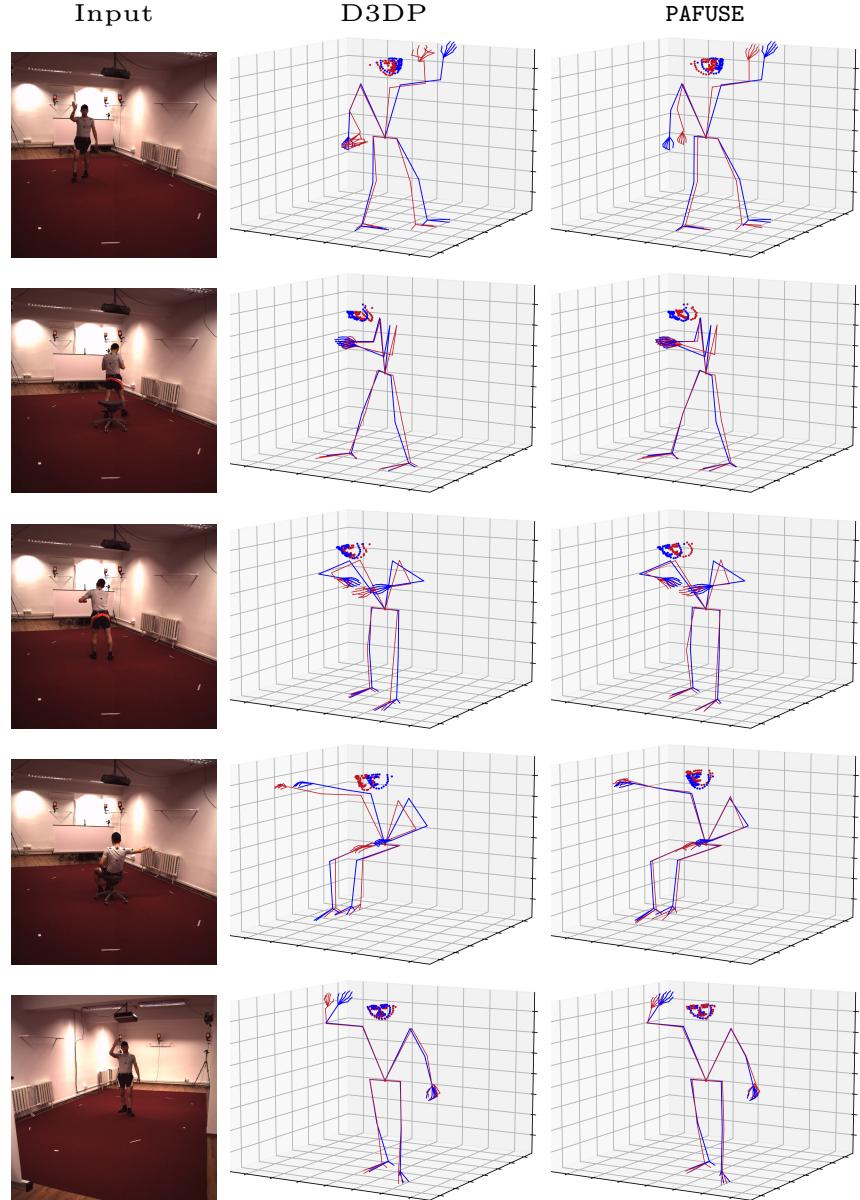


Fig. 5: Qualitative results from the H3WB test set. Blue: ground-truth, Red: best hypothesis. In comparison to D3DP, PAFUSE’s is better-aligned to the body joints (e.g., the shoulders), due to the hierarchical structure of the part-based prediction inducing such alignment. Remark also that PAFUSE’s dedicated networks for hands and face lead to considerably better predictions for those body parts.

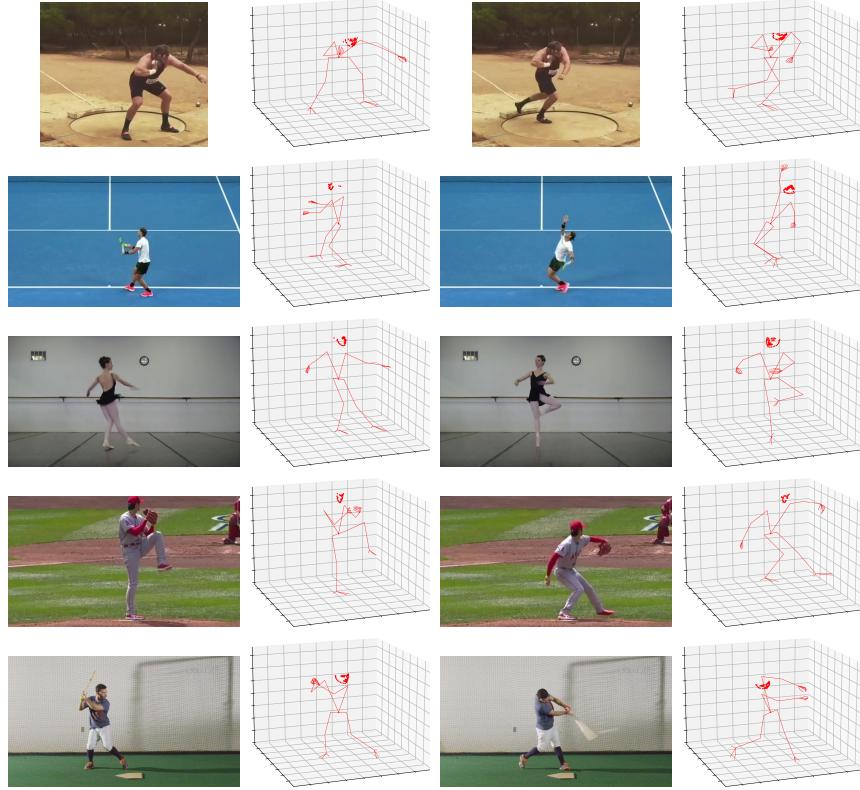


Fig. 6: Qualitative results in the wild. PAFUSE performs well even on challenging scenarios of fast movement, occlusions, and self occlusions. Remark, in particular, the occlusions caused by external objects such as balls (first row), tennis racket (second row), gloves (4th row) and club (last row). Remark also its robustness against motion blur on body parts (ballerina’s hands on 3rd row). 2D whole-body keypoints extracted using OpenPifPaf [21] and 3D poses predicted by PAFUSE.

keypoints as everything is processed by a single network. These spurious correlations are removed when separating the networks since the centered keypoints are now processed independently. Similarly, using separated body-part networks without the part-frame shifts forces the hands and the face networks to predict the absolute pose of these parts without having access to the body pose and in particular the pose of their part root joints (wrists and neck), significantly increasing the difficulty of the task.

Design choices. In addition to the main components, Table 3 showcases results related to finer design choices, each row illustrating the impact of changing one design axis by its baseline alternative. All choices of PAFUSE lead to better results, except for the frame window length (N), in which there seems to be a compromise between balanced performance (with a shorter window of $N = 27$, which is our choice), and emphasizing face and hands performance ($N = 81$).

Model	WB	PB	Body	Face	Hands
PAFUSE (with default choices)	43.0	16.4	37.1	5.8	21.7
Part Loss → WB Loss	47.7	21.5	50.4	7.2	28.1
Balanced → Uniform Denoisers	49.9	17.8	46.0	5.2	22.1
$N = 27 \rightarrow N = 81$	45.1	16.3	38.8	5.6	20.7
MPJPE → MSE	44.1	18.1	42.1	6.5	23.2

Table 3: Ablation experiments on design choices, comparing the effect of substituting each choice from the default base model → an alternative design possibility.

The part-based loss of Eq. (2) brings the most dramatic gains, in contrast to the whole-body loss of Eq. (3), demonstrating that separating the optimization targets for each part results in substantial improvement. Balancing the size of the denoiser models (384, 256 and 224 channels for body, hands, and face respectively) also considerably improved the results in comparison to using an uniform-sized (288 channels) denoiser. Other design choices brought smaller, but still visible improvements. $N=27 \rightarrow N=81$ experiment shows that extending the temporal context adversely affects performance. We attribute this decline to the uneven sampling characteristics of the H3WB dataset. The ablation study comparing *MPJPE* to *MSE* shows that despite the common use of MSE loss in diffusion-based methods for 3D whole-body pose estimation, optimizing with MPJPE yields superior results. This outcome aligns with expectations, given that model performance is evaluated using the MPJPE metric.

5 Conclusion

We introduced PAFUSE, a denoising diffusion model that predicts entire temporal sequences of 3D whole-body skeleton from their 2D counterparts.

While diffusion models are SOTA in 3D pose estimation, making them part-based is well-motivated, as the body, hands, and face present widely different scales and motion ranges. However, designing the best combination of diffusion and part-based is far from trivial, as our ablation experiments showcase.

PAFUSE employs a hierarchical part-based model, where body parts that have different scales and motion are predicted by part-specific networks, conditioned on their parent body parts. We were the first to exploit H3WB [64] temporal information, obtaining outstanding improvements over the previous frame-only state-of-the art (41.4mm MJPJE against 88mm).

Limitations. Performing all quantitative experiments on a single dataset is a necessary limitation, as H3WB is the only publicly annotated dataset for whole-body with sequential frames, allowing us to reconstruct the temporal data. We worked to mitigate that limitation by complementing our method with qualitative results in-the-wild videos, showing that our methods generalizes well to unseen context. We expect that future developments of this exciting research are may bring a growing array of data.

References

1. Bagautdinov, T., Wu, C., Simon, T., Prada, F., Shiratori, T., Wei, S.E., Xu, W., Sheikh, Y., Saragih, J.: Driving-signal aware full-body avatars. *ACM Transactions on Graphics (TOG)* **40**(4), 1–17 (2021)
2. Cai, Y., Ge, L., Liu, J., Cai, J., Cham, T.J., Yuan, J., Thalmann, N.M.: Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 2272–2281 (2019)
3. Catalin Ionescu, Fuxin Li, C.S.: Latent structured models for human pose estimation. In: ICCV (2011)
4. Chessa, M., Maiello, G., Klein, L.K., Paulun, V.C., Solari, F.: Grasping objects in immersive virtual reality. In: 2019 IEEE Conference on Virtual Reality and 3D User Interfaces (VR). pp. 1749–1754. IEEE (2019)
5. Choutas, V., Pavlakos, G., Bolkart, T., Tzionas, D., Black, M.J.: Monocular expressive body regression through body-driven attention. In: ECCV (2020)
6. Fieraru, M., Zanfir, M., Oneata, E., Popa, A.I., Olaru, V., Sminchisescu, C.: Three-dimensional reconstruction of human interactions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7214–7223 (2020)
7. Fieraru, M., Zanfir, M., Oneata, E., Popa, A.I., Olaru, V., Sminchisescu, C.: Learning complex 3d human self-contact. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 1343–1351 (2021)
8. Fieraru, M., Zanfir, M., Pirlea, S.C., Olaru, V., Sminchisescu, C.: Aifit: Automatic 3d human-interpretable feedback models for fitness training. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9919–9928 (2021)
9. Garcia-Salguero, M., Gonzalez-Jimenez, J., Moreno, F.A.: Human 3d pose estimation with a tilting camera for social mobile robot interaction. *Sensors* **19**(22), 4943 (2019)
10. Gong, J., Fan, Z., Ke, Q., Rahmani, H., Liu, J.: Meta agent teaming active learning for pose estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11079–11089 (2022)
11. Gu, Y., Pandit, S., Saraee, E., Nordahl, T., Ellis, T., Betke, M.: Home-based physical therapy with an interactive computer vision system. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops. pp. 0–0 (2019)
12. Gui, L.Y., Zhang, K., Wang, Y.X., Liang, X., Moura, J.M., Veloso, M.: Teaching robots to predict human motion. In: 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 562–567. IEEE (2018)
13. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. *Advances in neural information processing systems* **33**, 6840–6851 (2020)
14. Holmquist, K., Wandt, B.: Diffpose: Multi-hypothesis human pose estimation using diffusion models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15977–15987 (2023)
15. Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *TPMAI* (2014)
16. Iqbal, U., Xie, K., Guo, Y., Kautz, J., Molchanov, P.: Kama: 3d keypoint aware body mesh articulation. In: 2021 International Conference on 3D Vision (3DV). pp. 689–699. IEEE (2021)

17. Ivashechkin, M., Mendez, O., Bowden, R.: Improving 3d pose estimation for sign language. In: 2023 IEEE International Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW). pp. 1–5. IEEE (2023)
18. Jin, S., Xu, L., Xu, J., Wang, C., Liu, W., Qian, C., Ouyang, W., Luo, P.: Whole-body human pose estimation in the wild. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16. pp. 196–214. Springer (2020)
19. Joo, H., Simon, T., Sheikh, Y.: Total capture: A 3d deformation model for tracking faces, hands, and bodies. CVPR (2018)
20. Kolotouros, N., Pavlakos, G., Jayaraman, D., Daniilidis, K.: Probabilistic modeling for human mesh recovery. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 11605–11614 (2021)
21. Kreiss, S., Bertoni, L., Alahi, A.: Pipaf: Composite fields for human pose estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
22. Krishna, S., et al.: Signpose: Sign language animation through 3d pose lifting. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2640–2649 (2021)
23. Kundu, J.N., Rakesh, M., Jampani, V., Venkatesh, R.M., Venkatesh Babu, R.: Appearance consensus driven self-supervised human mesh recovery. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16. pp. 794–812. Springer (2020)
24. Li, M., Chen, S., Chen, X., Zhang, Y., Wang, Y., Tian, Q.: Actional-structural graph convolutional networks for skeleton-based action recognition. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 3595–3603 (2019)
25. Li, W., Liu, H., Tang, H., Wang, P., Van Gool, L.: Mhformer: Multi-hypothesis transformer for 3d human pose estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13147–13156 (2022)
26. Liang, X., Angelopoulou, A., Kapetanios, E., Woll, B., Al Batat, R., Woolfe, T.: A multi-modal machine learning approach and toolkit to automate recognition of early stages of dementia among british sign language users. In: Computer Vision–ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16. pp. 278–293. Springer (2020)
27. Liu, K., Ding, R., Zou, Z., Wang, L., Tang, W.: A comprehensive study of weight sharing in graph networks for 3d human pose estimation. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16. pp. 318–334. Springer (2020)
28. Lutz, S., Blythman, R., Ghosal, K., Moynihan, M., Simms, C., Smolic, A.: Jointformer: Single-frame lifting transformer with error prediction and refinement for 3d human pose estimation. ArXiv (2022)
29. Luvizon, D.C., Picard, D., Tabia, H.: 2d/3d pose estimation and action recognition using multitask deep learning. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5137–5146 (2018)
30. Martinez, J., Hossain, R., Romero, J., Little, J.J.: A simple yet effective baseline for 3d human pose estimation. In: ICCV (2017)
31. Mehta, D., Sridhar, S., Sotnychenko, O., Rhodin, H., Shafiei, M., Seidel, H.P., Xu, W., Casas, D., Theobalt, C.: Vnect: Real-time 3d human pose estimation with a single rgb camera. Acm transactions on graphics (tog) **36**(4), 1–14 (2017)

32. Moreno-Noguer, F.: 3d human pose estimation from a single image via distance matrix regression. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2823–2832 (2017)
33. Nichol, A.Q., Dhariwal, P.: Improved denoising diffusion probabilistic models. In: International Conference on Machine Learning. pp. 8162–8171. PMLR (2021)
34. Parelli, M., Papadimitriou, K., Potamianos, G., Pavlakos, G., Maragos, P.: Exploiting 3d hand pose estimation in deep learning-based sign language recognition from rgb videos. In: Computer Vision–ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16. pp. 249–263. Springer (2020)
35. Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A., Tzionas, D., Black, M.: Expressive body capture: 3D hands, face, and body from a single image. CVPR (2019)
36. Pavlakos, G., Zhou, X., Derpanis, K.G., Daniilidis, K.: Coarse-to-fine volumetric prediction for single-image 3d human pose. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7025–7034 (2017)
37. Pavllo, D., Feichtenhofer, C., Grangier, D., Auli, M.: 3d human pose estimation in video with temporal convolutions and semi-supervised training. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 7753–7762 (2019)
38. Rommel, C., Valle, E., Chen, M., Khalfaoui, S., Marlet, R., Cord, M., Pérez, P.: Diffhpe: Robust, coherent 3d human pose lifting with diffusion. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3220–3229 (2023)
39. Rong, Y., Shiratori, T., Joo, H.: Frankmocap: A monocular 3d whole-body pose estimation system via regression and integration. ICCV (2021)
40. Saadatnejad, S., Rasekh, A., Mofayzei, M., Medghalchi, Y., Rajabzadeh, S., Moridan, T., Alahi, A.: A generic diffusion-based approach for 3d human pose prediction in the wild. In: 2023 IEEE International Conference on Robotics and Automation (ICRA). pp. 8246–8253. IEEE (2023)
41. Samet, N., Akbas, E.: Hprnet: Hierarchical point regression for whole-body human pose estimation. Image and Vision Computing **115**, 104285 (2021)
42. Seong, B., Lee, H.J., Lee, R.: Keypoint-wise adaptive loss for whole-body human pose estimation (2023)
43. Shan, W., Liu, Z., Zhang, X., Wang, Z., Han, K., Wang, S., Ma, S., Gao, W.: Diffusion-based 3d human pose estimation with multi-hypothesis aggregation. arXiv preprint arXiv:2303.11579 (2023)
44. Shan, W., Liu, Z., Zhang, X., Wang, Z., Han, K., Wang, S., Ma, S., Gao, W.: Diffusion-based 3d human pose estimation with multi-hypothesis aggregation. arXiv preprint arXiv:2303.11579 (2023)
45. Sharma, S., Varigonda, P.T., Bindal, P., Sharma, A., Jain, A.: Monocular 3d human pose estimation by generation and ordinal ranking. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 2325–2334 (2019)
46. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 (2020)
47. Sun, X., Shang, J., Liang, S., Wei, Y.: Compositional human pose regression. In: Proceedings of the IEEE international conference on computer vision. pp. 2602–2611 (2017)
48. Svenstrup, M., Tranberg, S., Andersen, H.J., Bak, T.: Pose estimation and adaptive robot behaviour for human-robot interaction. In: 2009 IEEE International Conference on Robotics and Automation. pp. 3571–3576. IEEE (2009)

49. Tashiro, Y., Song, J., Song, Y., Ermon, S.: Csdi: Conditional score-based diffusion models for probabilistic time series imputation. *Advances in Neural Information Processing Systems* **34**, 24804–24816 (2021)
50. Wandt, B., Rudolph, M., Zell, P., Rhodin, H., Rosenhahn, B.: Canonpose: Self-supervised monocular 3d human pose estimation in the wild. In: *CVPR* (2021)
51. Wang, J., Yan, S., Dai, B., Lin, D.: Scene-aware generative network for human motion synthesis. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 12206–12215 (2021)
52. Weinzaepfel, P., Brégier, R., Combaluzier, H., Leroy, V., Rogez, G.: DOPE: distillation of part experts for whole-body 3d pose estimation in the wild. *ECCV* (2020)
53. Xiang, D., Joo, H., Sheikh, Y.: Monocular total capture: Posing face, body, and hands in the wild. In: *CVPR* (2019)
54. Xu, L., Jin, S., Liu, W., Qian, C., Ouyang, W., Luo, P., Wang, X.: Zoomnas: searching for whole-body human pose estimation in the wild. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **45**(4), 5296–5313 (2022)
55. Xu, T., Takano, W.: Graph stacked hourglass networks for 3d human pose estimation. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 16105–16114 (2021)
56. Yamane, S., Yamazoe, H., Lee, J.H.: Human motion generation based on gan toward unsupervised 3d human pose estimation. In: *Pattern Recognition: ACPR 2019 Workshops*, Auckland, New Zealand, November 26, 2019, Proceedings 5. pp. 100–109. Springer (2020)
57. Yan, A., Wang, Y., Li, Z., Qiao, Y.: Pa3d: Pose-action 3d machine for video recognition. In: *Proceedings of the ieee/cvf conference on computer vision and pattern recognition*. pp. 7922–7931 (2019)
58. Zauss, D., Kreiss, S., Alahi, A.: Keypoint communities. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 11057–11066 (2021)
59. Zeng, W., Jin, S., Liu, W., Qian, C., Luo, P., Ouyang, W., Wang, X.: Not all tokens are equal: Human-centric visual analysis via token clustering transformer. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 11101–11111 (2022)
60. Zhang, C., Yang, T., Weng, J., Cao, M., Wang, J., Zou, Y.: Unsupervised pre-training for temporal action localization tasks. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 14031–14041 (2022)
61. Zhang, J., Tu, Z., Yang, J., Chen, Y., Yuan, J.: Mixste: Seq2seq mixed spatio-temporal encoder for 3d human pose estimation in video. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 13232–13242 (2022)
62. Zhao, L., Peng, X., Tian, Y., Kapadia, M., Metaxas, D.N.: Semantic graph convolutional networks for 3d human pose regression. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 3425–3435 (2019)
63. Zheng, C., Zhu, S., Mendieta, M., Yang, T., Chen, C., Ding, Z.: 3d human pose estimation with spatial and temporal transformers. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 11656–11665 (2021)
64. Zhu, Y., Samet, N., Picard, D.: H3wb: Human3. 6m 3d wholebody dataset and benchmark. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 20166–20177 (2023)
65. Zou, Z., Tang, W.: Modulated graph convolutional network for 3d human pose estimation. In: *Proceedings of the IEEE/CVF international conference on computer vision*. pp. 11477–11487 (2021)