# Humans in 4D: Reconstructing and Tracking Humans with Transformers

Shubham Goel    Georgios Pavlakos    Jathushan Rajasegaran    Angjoo Kanazawa*    Jitendra Malik*

{shubham-goel, pavlakos, jathushan, kanazawa}@berkeley.edu, malik@eecs.berkeley.edu

University of California, Berkeley

Figure 1: **A "transformerized" view of Human Mesh Recovery**. We describe HMR 2.0, a fully transformer-based approach for 3D human pose and shape reconstruction from a single image. Besides impressive performance across a wide variety of poses and viewpoints, HMR 2.0 also acts as the backbone of an improved system for jointly reconstructing and tracking Humans in 4D (4DHumans). Here, we see output reconstructions from HMR 2.0 for each 2D detection in the left image.

## Abstract

*We present an approach to reconstruct humans and track them over time. At the core of our approach, we propose a fully "transformerized" version of a network for human mesh recovery. This network, HMR 2.0, advances the state of the art and shows the capability to analyze unusual poses that have in the past been difficult to reconstruct from single images. To analyze video, we use 3D reconstructions from HMR 2.0 as input to a tracking system that operates in 3D. This enables us to deal with multiple people and maintain identities through occlusion events. Our complete approach, 4DHumans, achieves state-of-the-art results for tracking people from monocular video. Furthermore, we demonstrate the effectiveness of HMR 2.0 on the downstream task of action recognition, achieving significant improvements over previous pose-based action recognition approaches. Our code and models are available on the project website:* https://shubham-goel.github.io/4dhumans/.

## 1. Introduction

In this paper, we present a fully transformer-based approach for recovering 3D meshes of human bodies from single images, and tracking them over time in video. We obtain unprecedented accuracy in our single-image 3D reconstructions (see Figure 1) even for unusual poses where previous approaches struggle. In video, we link these reconstructions over time by 3D tracking, in the process bridging gaps due to occlusion or detection failures. These 4D reconstructions can be seen on the project webpage.

Our problem formulation and approach can be conceived as the "transformerization" of previous work on human mesh recovery, HMR [30] and 3D tracking, PHALP [65]. Since the pioneering ViT paper [15], the process of "transformerization", *i.e.*, converting models from CNNs or LSTMs to transformer backbones, has advanced rapidly across multiple computer vision tasks, *e.g.*, [8, 16, 24, 40, 61, 77]. Specifically for 2D pose (2D body keypoints) this has already been done by ViTPose [81]. We take that as a starting point and we develop a new version of HMR, which we call HMR 2.0 to acknowledge its antecedent.

We use HMR 2.0 to build a system that can simultaneously reconstruct and track humans from videos. We rely on the recent 3D tracking system, PHALP [65], which we simplify and improve using our pose recovery. This system can reconstruct Humans in 4D, which gives the name to our method, 4DHumans. 4DHumans can be deployed on any video and can jointly track and reconstruct people in video. The functionality of creating a tracking entity for every person is fundamental towards analyzing and understanding humans in video. Besides achieving state-of-the-art results for tracking on the PoseTrack dataset [1], we also apply HMR 2.0 on the downstream application of action recognition. We follow the system design of recent work, [63], and we show that the use of HMR 2.0 can achieve impressive improvements upon the state of the art on action recognition on the AVA v2.2 dataset.

This paper is unabashedly a systems paper. We make design choices that lead to the best systems for 3D human reconstruction and tracking in the wild. Our model is publicly available on the project webpage. There is an emerging trend, in computer vision as in natural language processing, of large pretrained models which find widespread downstream applications and thus justify the scaling effort. HMR 2.0 is such a large pre-trained model which could potentially be useful not just in computer vision, but also in robotics [54, 62, 73], computer graphics [76], biomechanics [60], and other fields where analysis of the human figure and its movement from images or videos is needed.

Our contributions can be summarized as follows:

1. We propose an end-to-end "transformerized" architecture for human mesh recovery, HMR 2.0. Without relying on domain-specific designs, we outperform existing approaches for 3D body pose reconstruction.

2. Building on HMR 2.0, we design 4DHumans that can jointly reconstruct and track humans in video, achieving state-of-the-art results for tracking.

3. We show that better 3D poses from HMR 2.0 result in better performance on the downstream task of action recognition, finally contributing to the state-of-the-art result (42.3 mAP) on the AVA benchmark.

## 2. Related Work

**Human Mesh Recovery from a Single Image.** Although, there have been many approaches that estimate 3D human pose and shape relying on iterative optimization, *e.g.*, SMPLify [7] and variants [22, 38, 56, 66, 72, 85], for this analysis we will focus on approaches that directly regress the body shape from a single image input. In this case, the canonical example is HMR [30], which uses a CNN to regress SMPL [45] parameters. Since its introduction,

many improvements have been proposed for the original method. Notably, many works have proposed alternative methods for pseudo-ground truth generation, including using temporal information [3], multiple views [39], or iterative optimization [35, 29, 57]. SPIN [35] proposed an in-the-loop optimization that incorporated SMPLify [7] in the HMR training. Here, we also rely on pseudo-ground truth fits for training, and we use [37] for the offline fitting.

More recently, there have been works that propose more specialized designs for the HMR architecture. PyMAF [89, 88] incorporates a mesh alignment module for the regression of the SMPL parameters. PARE [34] proposes a body-part-guided attention mechanism for better occlusion handling. HKMR [20] performs a prediction that is informed by the known hierarchical structure of SMPL. HoloPose [23] proposes a pooling strategy that follows the 2D locations of each body joints. Instead, we follow a design without any domain-specific decisions and we show that it outperforms all previous approaches.

Many related approaches are making non-parametric predictions, *i.e.*, instead of estimating the parameters of the SMPL model, they explicitly regress the vertices of the mesh. GraphCMR [36] uses a graph neural network for the prediction, METRO [42] and FastMETRO [10] use a transformer, while Mesh Graphormer [43] adopts a hybrid between the two. Since we regress the SMPL model parameters, instead of the locations of mesh vertices, we are not directly comparable to these. However, we show how we can use a fully "transformerized" design for HMR.

**Human Mesh & Motion Recovery from Video.** To extend Human Mesh Recovery over time, most methods use the basic backbone of HMR [30] and propose designs for the temporal encoder that fuses the per-frame features. HMMR [31] uses a convolutional encoder on features extracted from HMR [30]. VIBE [33], MEVA [48] and TCMR [11] use a recurrent temporal encoder. DSD [71] combines convolutional and self-attention layers, while MAED [75] and t-HMMR [57] employ a transformer-based temporal encoder. Baradel *et al*. [5, 4] also used a transformer for temporal pose prediction, while operating directly on SMPL poses. One key limitation of these approaches is that they often operate in scenarios where tracking is simple [31, 90], *e.g.*, videos with a single person or minimal occlusions. In contrast to that, our complete 4DHumans approach is also solving the tracking problem.

**Tracking People in Video.** Recently, there have been approaches that demonstrate state-of-the-art performance for tracking by relying on 3D human reconstruction from HMR models, *i.e.*, T3DP [64] and PHALP [65]. In these methods, every person detection is lifted to 3D using an HMR network [57] and then tracking is performed using the 3D representations from lifting [64] and prediction [65] to track people in video. Empirical results show that PHALP works
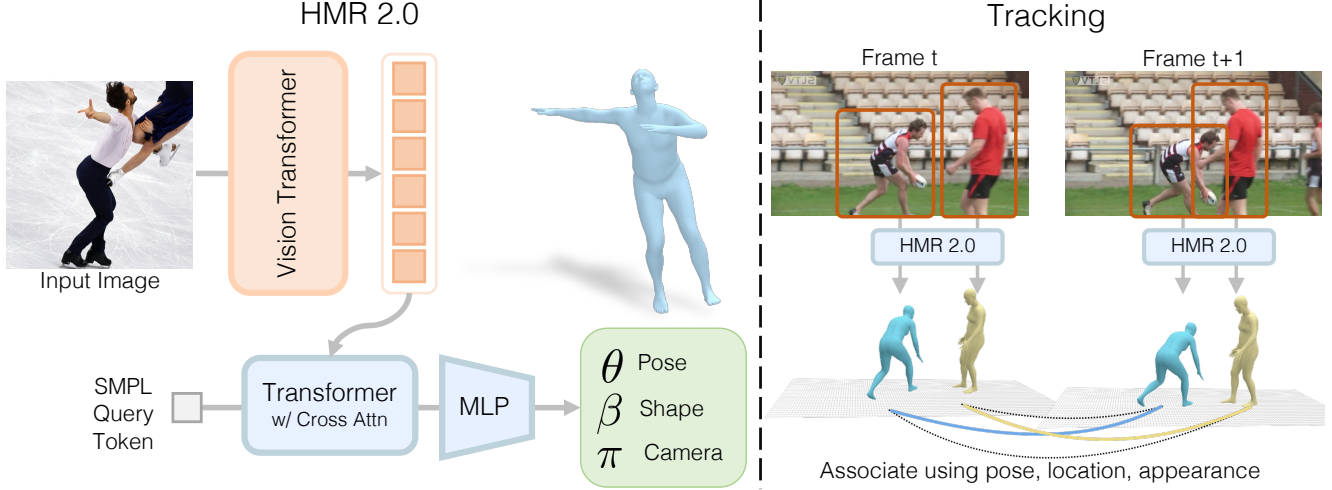
Figure 2: **Overview of our approach.** Left: HMR 2.0 is a fully "transformerized" version of a network for Human Mesh Recovery. Right: We use HMR 2.0 as the backbone of our 4DHumans system, that builds on PHALP [65], to jointly reconstruct and track humans in 4D.

very well on multiple tracking benchmarks (the main requirement is that the images have enough spatial resolution to permit lifting of the people to 3D). We use these tracking pipelines, and particularly PHALP, as a task to evaluate methods for human mesh recovery.

**Action Recognition.** Action recognition is typically performed using appearance features from raw video input. Canonical examples in this category include SlowFast [18] and MViT [16]. Simultaneously, there are approaches that use features extracted from body pose information, *e.g.*, PoTion [12] and JMRN [68]. A recent approach, LART [63], demonstrates state-of-the-art performance for action recognition by fusing video-based features with features from 3D human pose estimates. We use the pipeline of this approach and employ action recognition as a downstream task to evaluate human mesh recovery methods.

## 3. Reconstructing People

### 3.1. Preliminaries

**Body Model.** The SMPL model [46] is a low-dimensional parametric model of the human body. Given input parameters for pose ($\theta \in \mathbb{R}^{24 \times 3 \times 3}$) and shape ($\beta \in \mathbb{R}^{10}$), it outputs a mesh $M \in \mathbb{R}^{3 \times N}$ with $N = 6890$ vertices. The body joints $X \in \mathbb{R}^{3 \times k}$ are defined as a linear combination of the vertices and can be computed as $X = MW$ with fixed weights $W \in \mathbb{R}^{N \times k}$. Note that pose parameters $\theta$ include the body pose parameters $\theta_b \in \mathbb{R}^{23 \times 3 \times 3}$ and the global orientation $\theta_g \in \mathbb{R}^{3 \times 3}$.

**Camera.** We use a perspective camera model with fixed focal length and intrinsics $K$. Each camera $\pi = (R, t)$ consists of a global orientation $R \in \mathbb{R}^{3 \times 3}$ and translation $t \in \mathbb{R}^3$. Given these parameters, points in the SMPL space (*e.g.*, joints $X$) can be projected to the image as $x = \pi(X) = \Pi(K(RX + t))$, where $\Pi$ is a perspective projection with camera intrinsics $K$. Since $\theta$ already includes a global orientation, in practice we assume $R$ as identity and only predict camera translation $t$.

**HMR.** The goal of the human mesh reconstruction (HMR) task is to learn a predictor $f(I)$ that given a single image I, reconstructs the person in the image by predicting their 3D pose and shape parameters. Following the typical parametric approaches [30, 35], we model $f$ to predict $\Theta = [\theta, \beta, \pi] = f(I)$ where $\theta$ and $\beta$ are the SMPL pose and shape parameters and $\pi$ is the camera translation.

### 3.2. Architecture

We re-imagine HMR [30] as an end-to-end transformer architecture that uses no domain specific design choices. Yet, it outperforms all existing approaches that have heavily customized architectures and elaborate design decisions. As shown in Figure 2, we use (i) a ViT [15] to extract image tokens, and (ii) a standard transformer decoder that cross-attends to image tokens to output $\Theta$.

**ViT.** The Vision Transformer, or ViT [15] is a transformer [74] that has been modified to operate on an image. The input image is first patchified into input tokens and passed through the transformer to get output tokens. The output tokens are then passed to the transformer decoder. We use a ViT-H/16, the "Huge" variant with $16 \times 16$ input patch size. Please see SupMat for more details.

**Transformer decoder.** We use a standard transformer decoder [74] with multi-head self-attention. It processes a single (zero) input token by cross-attending to the output image tokens and ends with a linear readout of $\Theta$. We follow [35] and regress 3D rotations using the representation of [91].

## 3.3. Losses

Following best practices in the HMR literature [30, 35], we train our predictor $f$ with a combination of 2D losses, 3D losses, and a discriminator. Since we train with a mixture of datasets, each having different kinds of annotations, we employ a subset of these losses for each image in a minibatch. We use the same losses even with pseudo-ground truth annotations. Given an input image $I$, the model predicts $\Theta = [\theta, \beta, \pi] = f(I)$. Whenever we have access to the ground-truth SMPL pose parameters $\theta^*$ and shape parameters $\beta^*$, we bootstrap the model predictions using an MSE loss:

$$\mathcal{L}_{\mathtt{smpl}} = ||\theta - \theta^*||_2^2 + ||\beta - \beta^*||_2^2.$$

When the image has accurate ground-truth 3D keypoint annotations $X^*$, we additionally supervise the predicted 3D keypoints $X$ with an L1 loss:

$$\mathcal{L}_{\mathtt{kp3D}} = ||X - X^*||_1.$$

When the image has 2D keypoints annotations $x^*$, we supervise projections of predicted 3D keypoints $\pi(X)$ using an L1 loss:

$$\mathcal{L}_{\mathtt{kp2D}} = ||\pi(X) - x^*||_1.$$

Furthermore, we want to ensure that our model predicts valid 3D poses and use the adversarial prior in HMR [30]. It factorizes the model parameters into: (i) body pose parameters $\theta_b$, (ii) shape parameters $\beta$, and (iii) per-part relative rotations $\theta_i$, which is one 3D rotation for each of the 23 joints of the SMPL model. We train a discriminator $D_k$ for each factor of the body model, and the generator loss can be expressed as:

$$\mathcal{L}_{\mathtt{adv}} = \sum_k (D_k(\theta_b, \beta) - 1)^2.$$

## 3.4. Pseudo-Ground Truth fitting

We scale to unlabelled datasets (*i.e.*, InstaVariety [31], AVA [21], AI Challenger [78]) by computing pseudo-ground truth annotations. Given any image, we first use an off-the-shelf detector [40] and a body keypoints estimator [81] to get bounding boxes and corresponding 2D keypoints. We then fit a SMPL mesh to these 2D keypoints using ProHMR [37] to get pseudo-ground truth SMPL parameters $\theta^*$ and $\beta^*$ with camera $\pi^*$.

## 4. Tracking People

In videos with multiple people, we need the ability to associate people across time, *i.e.*, perform tracking. For this we build upon PHALP [65], a state-of-the-art tracker based on features derived from HMR-style 3D reconstructions. The basic idea is to detect people in individual frames, and
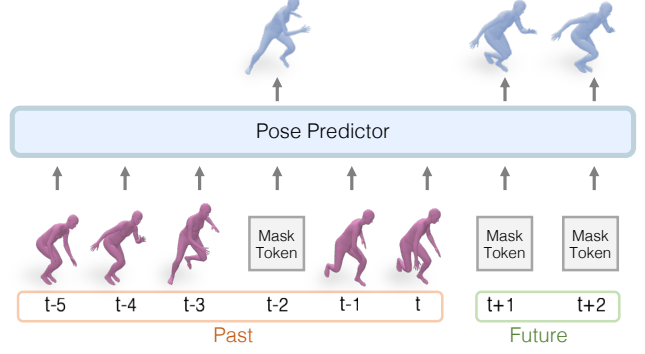


Figure 3: **Pose prediction:** We train a BERT-style [13] transformer model on over 1 million tracks obtained from [63]. This allow us to make future predictions and amodal completion of missing detections using the same model. To predict future poses ($t$+1, $t$+2, ...), we query the model with a mask-token using corresponding positional embeddings. Similarly for amodal completion, we replace missing detections with a masked token.

"lift" them to 3D, extracting their 3D pose, location in 3D space (derived from the estimated camera), and 3D appearance (derived from the texture map). A tracklet representation is incrementally built up for each individual person over time. The recursion step is to predict for each tracklet, the pose, location and appearance of the person in the next frame, all in 3D, and then find best matches between these top-down predictions and the bottom-up detections of people in that frame after lifting them to 3D. The state represented by each tracklet is then updated by the incoming observation, and the process is iterated. It is possible to track through occlusions because the 3D representation of a tracklet continues to be updated based on past history.

We believe that a robust pose predictor should also perform well, when evaluated on this downstream task of tracking, so we use the tracking metrics as a proxy to evaluate the quality of 3D reconstructions. But first we needed to modify the PHALP framework to allow for fair comparison of different pose prediction models. Originally, PHALP used pose features based on the last layer of the HMR network, *i.e.*, a 2048-dimensional embedding space. This limits the ability of PHALP to be used with different pose models (*e.g.*, HMR 2.0, PARE, PyMAF etc.). To create a more generic version of PHALP, we perform the modification of representing pose in terms of SMPL pose parameters, and we accordingly optimize the PHALP cost function to utilize the new pose distance. Similarly, we adapt the pose predictor to operate on the space of SMPL parameters. More specifically, we train a vanilla transformer model [74] by masking random pose tokens as shown in the Fig 3. This allows us to predict future poses in time, as well as amodal completion of missing detections. With these modifications, we can plug in any mesh recovery methods and run them on any videos. We call this modified version PHALP'.

**4DHumans.** Our final tracking system, 4DHumans, uses a sampling-based parameter-free appearance head and a new pose predictor (Figure 3). To model appearance, we texture visible points on the mesh by projecting them onto the input image and sampling color from the corresponding pixels.

To track people in videos, previous approaches relied on off-the-shelf tracking approaches and used their output to reconstruct humans in videos (*e.g.*, take the bounding boxes from tracking output and reconstruct people). For example, PHD [90], HMMR [31] can run on videos with only single person in the scene. In this work, we combine reconstruction and tracking into a single system and show that better pose reconstructions result in better tracking and this combined system can now run on any videos in the wild.

# 5. Experiments

In this section, we evaluate our reconstruction and tracking system qualitatively and quantitatively. First, we show that HMR 2.0 outperforms previous methods on standard 2D and 3D pose accuracy metrics (Section 5.2). Second, we show 4DHumans is a versatile tracker, achieving state-of-the-art performance (Section 5.3). Third, we further demonstrate the robustness and accuracy of our recovered poses via superior performance on the downstream application of action recognition (Section 5.4). Finally, we discuss the experimental investigation when designing HMR 2.0 and ablate a series of design choices (Section 5.5).

## 5.1. Setup

**Datasets.** Following previous work, we use the typical datasets for training, *i.e.*, Human3.6M [27], MPI-INF-3DHP [49], COCO [44] and MPII [2]. Additionally, we use InstaVariety [31], AVA [21] and AI Challenger [78] as extra data where we generate pseudo-ground truth fits.
**Baselines.** We report performance on benchmarks that we can compare with many previous works (Section 5.2), but we also perform a more detailed comparison with recent state-of-the-art methods, *i.e.*, PyMAF [89], CLIFF [41], HMAR [65], PARE [34], and PyMAF-X [88]. For fairness, we only evaluate the body-only performance of PyMAF-X.

## 5.2. Pose Accuracy

**3D Metrics.** For 3D pose accuracy, we follow the typical protocols of prior work, *e.g.*, [35], and we present results on the 3DPW test split and on the Human3.6M val split, reporting MPJPE, and PA-MPJPE in Table 1. Please notice that we only compare with methods that do not use the training set of 3DPW for training, similar to us. We observe that with our HMR 2.0a model, which trains only on the typical datasets, we can outperform all previous baselines across all metrics. However, we believe that these benchmarks are very saturated and these smaller differences in pose metrics tend to not be very significant. In fact, we

|  | Method | 3DPW | | Human3.6M | |
|---|---|---|---|---|---|
|  |  | MPJPE | PA-MPJPE | MPJPE | PA-MPJPE |
| Temporal | Kanazawa *et al.* [31] | 116.5 | 72.6 | - | 56.9 |
| | Doersch *et al.* [14] | - | 74.7 | - | - |
| | Arnab *et al.* [3] | - | 72.2 | 77.8 | 54.3 |
| | DSD [71] | - | 69.5 | 59.1 | 42.4 |
| | VIBE [33] | 93.5 | 56.5 | 65.9 | 41.5 |
| Frame-based | Pavlakos *et al.* [59] | - | - | - | 75.9 |
| | HMR [30] | 130.0 | 76.7 | 88.0 | 56.8 |
| | NBF [53] | - | - | - | 59.9 |
| | GraphCMR [36] | - | 70.2 | - | 50.1 |
| | HoloPose [23] | - | - | 60.3 | 46.5 |
| | DenseRaC [82] | - | - | 76.8 | 48.0 |
| | SPIN [35] | 96.9 | 59.2 | 62.5 | 41.1 |
| | DecoMR [86] | - | 61.7† | - | 39.3† |
| | DaNet [87] | - | 56.9 | 61.5 | 48.6 |
| | Song *et al.* [69] | - | 55.9 | - | 56.4 |
| | I2L-MeshNet [51] | 100.0 | 60.0 | 55.7† | 41.1† |
| | HKMR [20] | - | - | 59.6 | 43.2 |
| | PyMAF [89] | 92.8 | 58.9 | 57.7 | 40.5 |
| | PARE [34] | 82.0 | 50.9 | 76.8 | 50.6 |
| | PyMAF-X [88] | 78.0 | 47.1 | 54.2 | 37.2 |
| | HMR 2.0a | 70.0 | 44.5 | 44.8 | 33.6 |
| | HMR 2.0b | 81.3 | 54.3 | 50.0 | 32.4 |

Table 1: **Reconstructions evaluated in 3D:** Reconstruction errors (in mm) on the 3DPW and Human3.6M datasets. † denotes the numbers evaluated on non-parametric results. Lower ↓ is better. Please see the text for details.

| Method | LSP-Extended | | COCO | | PoseTrack | |
|---|---|---|---|---|---|---|
|  | @0.05 | @0.1 | @0.05 | @0.1 | @0.05 | @0.1 |
| PyMAF [89] | - | - | 0.68 | 0.86 | 0.77 | 0.92 |
| CLIFF [41] | 0.32 | 0.66 | 0.64 | 0.88 | 0.75 | 0.92 |
| PARE [34] | 0.27 | 0.60 | 0.72 | 0.91 | 0.79 | 0.93 |
| PyMAF-X [88] | - | - | 0.79 | 0.93 | 0.85 | 0.95 |
| HMR 2.0a | 0.38 | 0.72 | 0.79 | 0.95 | 0.86 | 0.97 |
| HMR 2.0b | 0.53 | 0.82 | 0.86 | 0.96 | 0.90 | 0.98 |

Table 2: **Reconstructions evaluated in 2D.** PCK scores of projected keypoints at different thresholds on the LSP-Extended, COCO, and PoseTrack datasets. Higher ↑ is better.

observe that by a small compromise of the performance on 3DPW, our HMR 2.0b model, which trains for longer on more data (AVA [21], AI Challenger [78], and InstaVariety [31]), achieves results that perform better on more unusual poses than what can be found in Human3.6M and 3DPW. We observe this qualitatively and from performance evaluated on 2D pose reprojection (Table 2). Furthermore, we observe that HMR 2.0b is a more robust model and use it for evaluation in the rest of the paper.

**2D Metrics.** We evaluate 2D image alignment of the generated poses by reporting PCK of reprojected keypoints at different thresholds on LSP-Extended [28], COCO validation set [44], and Posetrack validation set [1]. Since

PyMAF(-X) [89, 88] were trained using LSP-Extended, we do not report numbers for that part of the table. Notice in Table 2, that HMR 2.0b consistently outperforms all previous approaches. On LSP-Extended, which contains unusual poses, HMR 2.0b achieves PCK@0.05 of 0.53, which is $1.6\times$ better than the second best (CLIFF) with 0.32. For PCK@0.05 on easier datasets like COCO and PoseTrack with less extreme poses, HMR 2.0b still outperforms the second-best approaches but by narrower margins of 9% and 6% respectively. HMR 2.0a also outperforms all baselines, but is worse than HMR 2.0b, especially on harder poses in LSP-Extended.

**Qualitative Results.** We show qualitative results of HMR 2.0 in Figure 4. We are robust to extreme poses and partial occlusions. Our reconstructions are well-aligned with the image and are valid when seen from a novel view. Moreover, we compare with our closest competitors in Figure 5. We observe that PyMAF-X and particularly PARE often struggle with more unusual poses, while HMR 2.0 returns more faithful reconstructions.

## 5.3. Tracking

For tracking, we first demonstrate the versatility of the modifications introduced by PHALP′, which allow us to evaluate 3D pose estimators on the downstream task of tracking. Then, we evaluate our complete system, 4DHumans, with respect to the state of the art.

**Evaluation Setting.** Following previous work [64, 65], we report results based on IDs (ID switches), MOTA [32], IDF1 [67], and HOTA [47] on the Posetrack validation set using the protocol of [65], with detections from Mask R-CNN [25].

**Versatility of PHALP′.** With the modifications of PHALP′, we abandon the model-specific latent space of [65] and instead, we operate in the SMPL space, which is shared across most mesh recovery systems. This makes PHALP′ more versatile and allows us to plug in different 3D pose estimators and compare them based on their performance on the downstream task of tracking. We perform this comparison in Table 3 where we use pose and location cues from state-of-the-art 3D pose estimators (while still using appearance from HMAR [65]). We observe that HMR 2.0 performs the best and PARE [34], HMAR [65], and PyMAF-X [88] closely follow on the Posetrack dataset, with minor differences between them. Note that tracking is often most susceptible to errors in predicted 3D locations with body pose having a smaller effect in performance [65]. This means that good tracking performance can indicate robustness to occlusions, so it is helpful to consider this metric, but it is less helpful to distinguish fine-grained differences in pose. As a result, the competitive results of PARE [34], HMAR [65], and PyMAF-X [88] indicate that they handle occlusions gracefully, but their pose estimation

| Tracker | Pose Engine | Posetrack | | | |
| --- | --- | --- | --- | --- | --- |
| | | HOTA↑ | IDs↓ | MOTA↑ | IDF1↑ |
| PHALP′ | PARE [34] | 53.6 | 510 | 59.4 | 76.8 |
| | PyMAF-X [88] | 53.7 | 472 | 59.2 | 76.9 |
| | CLIFF [41] | 53.5 | 551 | 58.7 | 76.5 |
| | PyMAF [89] | 53.0 | 623 | 58.6 | 76.1 |
| | HMAR [65] | 53.6 | 482 | 59.3 | 77.1 |
| | HMR 2.0 | 54.1 | 456 | 59.4 | 77.4 |

Table 3: **Tracking with different 3D pose estimators.** With the modifications of PHALP′, we have a versatile tracker that allows different 3D pose estimators to be plugged into it. HMR 2.0, PARE, and PyMAF-X perform the best in this setting.

| Method | Posetrack | | | |
| --- | --- | --- | --- | --- |
| | HOTA↑ | IDs↓ | MOTA↑ | IDF1↑ |
| Trackformer [50] | 46.7 | 1263 | 33.7 | 64.0 |
| Tracktor [6] | 38.5 | 702 | 42.4 | 65.2 |
| AlphaPose [17] | 37.6 | 2220 | 36.9 | 66.9 |
| Pose Flow [79] | 38.0 | 1047 | 15.4 | 64.2 |
| T3DP [64] | 50.6 | 655 | 55.8 | 73.4 |
| PHALP [65] | 52.9 | 541 | 58.9 | 76.4 |
| 4DHumans | 54.3 | 421 | 59.8 | 77.9 |
| 4DHumans + ViTDet | 57.8 | 382 | 61.4 | 79.1 |

Table 4: **Comparison of 4DHumans with the state of the art on the Posetrack dataset.** 4DHumans achieve state-of-the-art tracking performance for all metrics. Incorporating a better detection system [40] leads to further performance improvements.

might still be less accurate (as observed from Table 2). See also Figure 5 and SupMat for more qualitative comparisons.

**4DHumans.** Table 4 evaluates tracking performance of our complete system, 4DHumans, on the PoseTrack dataset. Using the same bounding box detector as [64, 65], 4DHumans outperforms existing approaches on all metrics, improving ID Switches by 22%. Using the improved ViTDet detector [40] can improve performance further. As a byproduct of our temporal prediction model (Figure 3), we can perform amodal completion and attribute a pose to missing detections. We show examples of this in the SupMat.

## 5.4. Action Recognition

**Evaluation setting.** The approach of [63] is the state of the art for action recognition in videos. Given a video as input, the authors propose using per-frame 3D pose and location estimates (using off-the-shelf HMR models [65]) as an additional feature for predicting action labels. They also show results for a "pose-only" baseline that predicts action labels
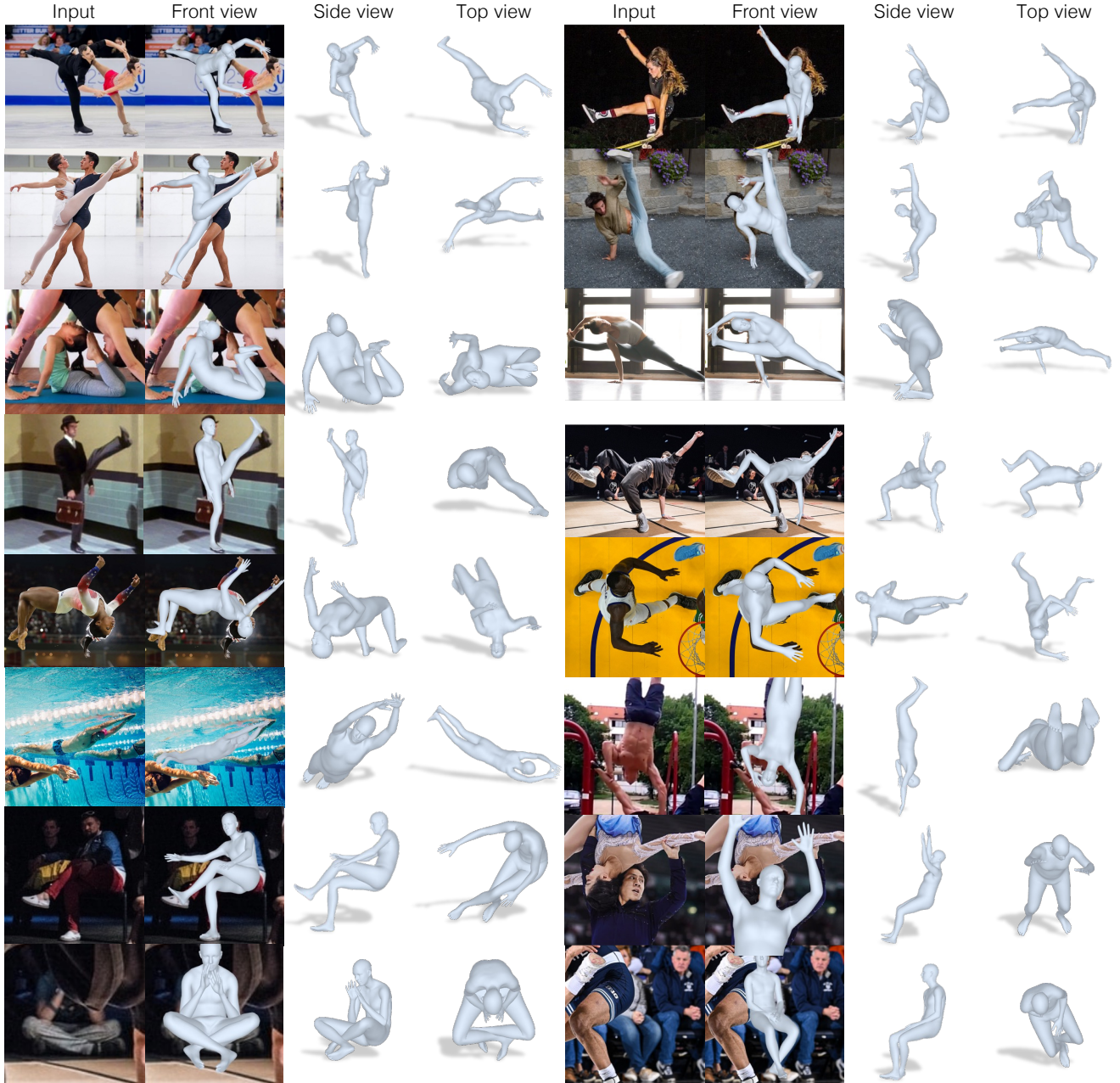
Figure 4: **Qualitative evaluation of HMR 2.0.** For each example we show: a) the input image, b) the reconstruction overlay, c) a side view, d) the top view. To demonstrate the robustness of HMR 2.0, we visualize results for a variety of settings - for unusual poses (rows 1-4), for unusual viewpoints (row 5) and for images with poor visibility, extreme truncations and extreme occlusions (rows 6-8).

using only 3D pose and location estimates. We use this setting to compare our model with baselines on the downstream task of action recognition on the AVA dataset [21]. In [63], the authors train a transformer that takes SMPL poses as input and predicts action labels. Following their setup, we train a separate action classification transformer for each baseline.

**Comparisons.** Comparing results in Table 5, we observe that HMR 2.0 outperforms baselines on the different class categories (OM, PI, PM) and overall. It achieves an mAP of 22.3 on the AVA test set, which is 14% better than the second-best baseline. Since accurate action recognition from poses needs fine-grained pose estimation, this is strong evidence that HMR 2.0 predicts more accurate poses than

Figure 5: **Qualitative comparison of state-of-the-art mesh recovery methods.** HMR 2.0 returns more faithful reconstructions for unusual poses compared to the closest competitors, PyMAF-X [88] and PARE [34].
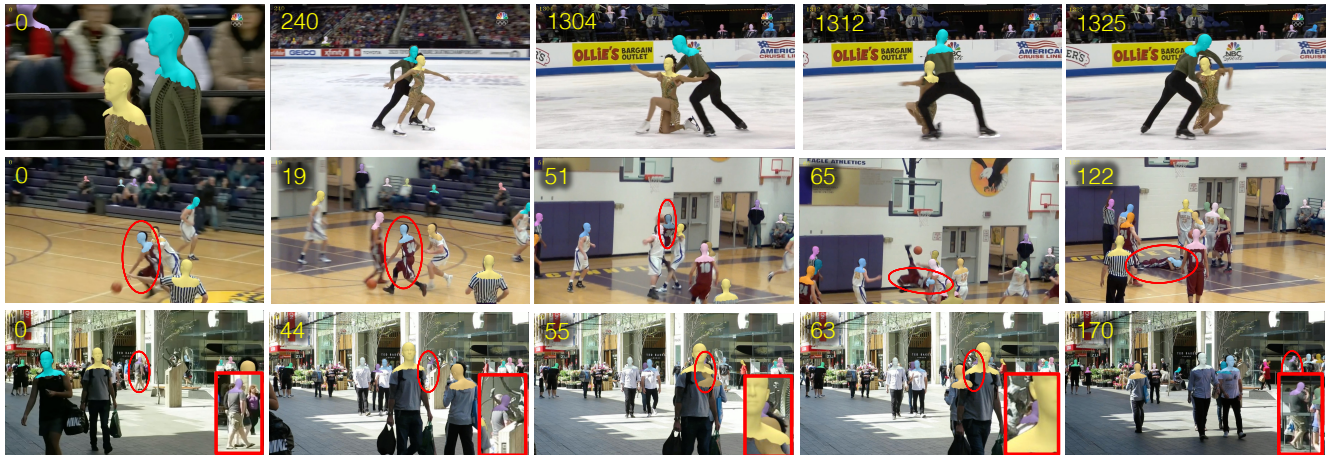


Figure 6: **Qualitative tracking results of 4DHumans**. We use head masks (frame number is on the top left). First row: We track people skating on ice with challenging poses and heavy occlusions, in a minute long video without switching identities. Second row: The main person is tracked through multiple interactions with other players. Third row: The person of interest is tracked through long occlusions.

| Action Model | Pose Engine | OM | PI | PM | mAP |
|---|---|---|---|---|---|
| [63] | PyMAF [89] | 7.3 | 16.9 | 34.7 | 15.4 |
| | CLIFF [41] | 9.2 | 20.0 | 40.3 | 18.6 |
| | HMAR [65] | 8.7 | 20.1 | 40.3 | 18.3 |
| | PARE [34] | 9.2 | 20.7 | 41.5 | 19.1 |
| | PyMAF-X [88] | 10.2 | 21.4 | 40.8 | 19.6 |
| | HMR 2.0 | 11.9 | 24.6 | 45.8 | 22.3 |

Table 5: **Action recognition results on the AVA dataset.** We benchmark different mesh recovery methods on the downstream task of pose-based action recognition. Here, *OM* : Object Manipulation, *PI* : Person Interactions, and *PM* : Person Movement.

existing approaches. In fact, when combined with appearance features, [63] shows that HMR 2.0 achieves the state of the art of 42.3 mAP on AVA action recognition, which is 7% better than the second-best of 39.5 mAP.

## 5.5. HMR 2.0 Model Design

In the process of developing HMR 2.0, we investigated a series of design decisions. Figure 7 briefly illustrates this exploration. We experimented with over 100 settings and we visualize the performance of 100 checkpoints for each run. For the visualization, we use the performance of each checkpoint on the 3DPW and the LSP-Extended dataset.

Our investigation focused on some specific aspects of the model, which we document here as a series of "lessons learnt" for future research. In the following paragraphs, we will regularly refer to Table 6, which evaluates these aspects
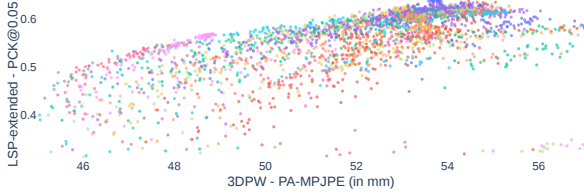
Figure 7: **Extensive model search.** With each dot, we visualize the performance of a checkpoint when evaluated on 3DPW and LSP-Extended. Colors indicate different runs. We explore more than 100 settings, and visualize ∼100 checkpoints from each run.

| Models | 3DPW | | Human3.6M | | LSP-Extended | |
|---|---|---|---|---|---|---|
| | MPJPE | PA-MPJPE | MPJPE | PA-MPJPE | PCK@0.05 | PCK@0.1 |
| HMR 2.0b | 81.3 | 54.3 | 50.0 | 32.4 | 0.53 | 0.82 |
| B1 | 85.2 | 56.8 | 58.9 | 41.4 | 0.35 | 0.66 |
| B2 | 79.7 | 53.4 | 51.4 | 34.4 | 0.48 | 0.81 |
| D1 | 84.1 | 54.8 | 54.5 | 35.1 | 0.45 | 0.79 |
| D2 | 80.2 | 53.3 | 52.4 | 34.9 | 0.46 | 0.79 |
| P1 | 98.9 | 61.7 | 89.9 | 58.7 | 0.24 | 0.52 |
| P2 | 82.7 | 55.6 | 49.3 | 32.4 | 0.52 | 0.81 |

Table 6: **Ablations**: Evaluation for different model designs on the 3DPW, Human3.6M, and LSP-Extended datasets.

on 3D and 2D metrics, using the 3DPW, Human3.6M, and LSP-Extended datasets.

**Effect of backbone.** Unlike the majority of the previous work on Human Mesh Recovery that uses a ResNet backbone, our HMR 2.0 method relies on a ViT backbone. For a direct comparison of the effect of the backbone, Model B1 implements HMR with a ResNet-50 backbone and an MLP-based head implementing IEF (Iterative Error Feedback [9, 30]). In contrast, Model B2 uses a transformer backbone (ViT-H) while keeping the other design decisions the same. By updating the backbone, we observe a significant improvement across the 3D and 2D metrics, which justifies the "transformerization" step.

**Effect of training data.** Besides the architecture, we also investigated the effect of training data. Model D1 trains on the typical datasets (H3.6M, MPII, COCO, MPI-INF) that most of the previous works are leveraging. In comparison, model D2 adds AVA in the training set, following [21]. Eventually, we also train using AI-Challenger and Insta-Variety (model B2), to further expand the training set. As we can see, adding more training data leads to improvements across the board for the reported metrics, but the benefit is smaller compared to the backbone update.

**ViT pretraining.** Another factor that had significant effect on the performance of our model was the pretraining of the ViT backbone. Starting with randomly initialized weights (model P1) results in slow convergence and poor performance. Results improve if our backbone is pretrained with MAE [24] on Imagenet (P2). Eventually, our model of choice (HMR 2.0b), which is first pretrained with MAE on ImageNet and then on the task of 2D keypoint prediction [81], achieves the best performance.

**SMPL head.** We also investigate the effect of the architecture for the head that predicts the SMPL parameters. Our proposed transformer decoder (HMR 2.0b) improves performance when it comes to the image-model alignment (*i.e.*, 2D metrics) compared to the traditional MLP-based head with IEF steps (B2).

**Dataset quality.** Similar to previous work, *e.g.*, [35], it was crucial to keep the quality of the training data high; we filter out low quality pseudo-ground truth fits (high fitting error) and prune images with low-confidence 2D detections.

## 6. Conclusion

We study the problem of reconstructing and tracking humans from images and video. First, we propose HMR 2.0, a fully "transformerized" version of a network for the problem of Human Mesh Recovery [30]. HMR 2.0 achieves strong performance on the usual 2D/3D pose metrics, while also acting as the backbone for our improved video tracker. The full system, 4DHumans, jointly reconstructs and tracks people in video and achieves state-of-the-art results for tracking. To further illustrate the benefit of our 3D pose estimator, HMR 2.0, we apply it to the task of action recognition, where we demonstrate strong improvements upon previous pose-based baselines.

Our work pushes the boundary of the videos that can be analyzed with techniques for 3D human reconstruction. At the same time, the improved results also demonstrate the type of limitations that need to be addressed in the future. For example, the use of the SMPL model [45] creates certain limitations, and leveraging improved models would allow us to model hand pose and facial expressions [56], or even capture greater age variation, *e.g.*, infants [26] and kids [55, 70]. Moreover, since we consider each person independently, our reconstructions are less successful at capturing the fine-grained nature of people in close proximity, *e.g.*, contact [19, 52]. Besides this, our reconstructions "live" in the camera frame, so for proper understanding of the action in a video, we need to consider everyone in a common world coordinate frame, by reasoning about the camera motion too [58, 83, 84]. Finally, lower input resolution can affect the quality of our reconstructions, which could be addressed by resolution augmentations [80].

# References

[1] Mykhaylo Andriluka, Umar Iqbal, Eldar Insafutdinov, Leonid Pishchulin, Anton Milan, Juergen Gall, and Bernt Schiele. PoseTrack: A benchmark for human pose estimation and tracking. In *CVPR*, 2018.

[2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2D human pose estimation: New benchmark and state of the art analysis. In *CVPR*, 2014.

[3] Anurag Arnab, Carl Doersch, and Andrew Zisserman. Exploiting temporal context for 3D human pose estimation in the wild. In *CVPR*, 2019.

[4] Fabien Baradel, Romain Brégier, Thibault Groueix, Philippe Weinzaepfel, Yannis Kalantidis, and Grégory Rogez. Pose-BERT: A generic transformer module for temporal 3D human modeling. *PAMI*, 2022.

[5] Fabien Baradel, Thibault Groueix, Philippe Weinzaepfel, Romain Brégier, Yannis Kalantidis, and Grégory Rogez. Leveraging MoCap data for human mesh recovery. In *3DV*, 2021.

[6] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles. In *ICCV*, 2019.

[7] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *ECCV*, 2016.

[8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.

[9] Joao Carreira, Pulkit Agrawal, Katerina Fragkiadaki, and Jitendra Malik. Human pose estimation with iterative error feedback. In *CVPR*, 2016.

[10] Junhyeong Cho, Kim Youwang, and Tae-Hyun Oh. Cross-attention of disentangled modalities for 3D human mesh recovery with transformers. In *ECCV*, 2022.

[11] Hongsuk Choi, Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. Beyond static features for temporally consistent 3D human pose and shape from a video. In *CVPR*, 2021.

[12] Vasileios Choutas, Philippe Weinzaepfel, Jérôme Revaud, and Cordelia Schmid. PoTion: Pose motion representation for action recognition. In *CVPR*, 2018.

[13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[14] Carl Doersch and Andrew Zisserman. Sim2real transfer learning for 3D human pose estimation: Motion to the rescue. *NeurIPS*, 2019.

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021.

[16] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, 2021.

[17] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. RMPE: Regional multi-person pose estimation. In *ICCV*, 2017.

[18] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, 2019.

[19] Mihai Fieraru, Mihai Zanfir, Elisabeta Oneata, Alin-Ionut Popa, Vlad Olaru, and Cristian Sminchisescu. Three-dimensional reconstruction of human interactions. In *CVPR*, 2020.

[20] Georgios Georgakis, Ren Li, Srikrishna Karanam, Terrence Chen, Jana Košecká, and Ziyan Wu. Hierarchical kinematic human mesh recovery. In *ECCV*, 2020.

[21] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, and Jitendra Malik. AVA: A video dataset of spatio-temporally localized atomic visual actions. In *CVPR*, 2018.

[22] Peng Guan, Alexander Weiss, Alexandru O Bălan, and Michael J Black. Estimating human shape and pose from a single image. In *ICCV*, 2009.

[23] Riza Alp Guler and Iasonas Kokkinos. HoloPose: Holistic 3D human reconstruction in-the-wild. In *CVPR*, 2019.

[24] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022.

[25] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017.

[26] Nikolas Hesse, Sergi Pujades, Michael J Black, Michael Arens, Ulrich G Hofmann, and A Sebastian Schroeder. Learning and tracking the 3D body shape of freely moving infants from RGB-D sequences. *PAMI*, 2019.

[27] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments. *PAMI*, 2013.

[28] Sam Johnson and Mark Everingham. Learning effective human pose estimation from inaccurate annotation. In *CVPR*, 2011.

[29] Hanbyul Joo, Natalia Neverova, and Andrea Vedaldi. Exemplar fine-tuning for 3D human model fitting towards in-the-wild 3D human pose estimation. In *3DV*, 2021.

[30] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018.

[31] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3D human dynamics from video. In *CVPR*, 2019.

[32] Rangachar Kasturi, Dmitry Goldgof, Padmanabhan Soundararajan, Vasant Manohar, John Garofolo, Rachel Bowers, Matthew Boonstra, Valentina Korzhova, and Jing Zhang. Framework for performance evaluation of face, text, and vehicle detection and tracking in video: Data, metrics, and protocol. *PAMI*, 2008.

[33] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. VIBE: Video inference for human body pose and shape estimation. In *CVPR*, 2020.

[34] Muhammed Kocabas, Chun-Hao P Huang, Otmar Hilliges, and Michael J Black. PARE: Part attention regressor for 3D human body estimation. In *ICCV*, 2021.

[35] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3D human pose and shape via model-fitting in the loop. In *ICCV*, 2019.

[36] Nikos Kolotouros, Georgios Pavlakos, and Kostas Daniilidis. Convolutional mesh regression for single-image human shape reconstruction. In *CVPR*, 2019.

[37] Nikos Kolotouros, Georgios Pavlakos, Dinesh Jayaraman, and Kostas Daniilidis. Probabilistic modeling for human mesh recovery. In *ICCV*, 2021.

[38] Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J Black, and Peter V Gehler. Unite the people: Closing the loop between 3D and 2D human representations. In *CVPR*, 2017.

[39] Vincent Leroy, Philippe Weinzaepfel, Romain Brégier, Hadrien Combaluzier, and Grégory Rogez. SMPLy benchmarking 3D human pose estimation in the wild. In *3DV*, 2020.

[40] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *ECCV*, 2022.

[41] Zhihao Li, Jianzhuang Liu, Zhensong Zhang, Songcen Xu, and Youliang Yan. CLIFF: Carrying location information in full frames into human pose and shape estimation. In *ECCV*, 2022.

[42] Kevin Lin, Lijuan Wang, and Zicheng Liu. End-to-end human pose and mesh reconstruction with transformers. In *CVPR*, 2021.

[43] Kevin Lin, Lijuan Wang, and Zicheng Liu. Mesh graphormer. In *ICCV*, 2021.

[44] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014.

[45] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015.

[46] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015.

[47] Jonathon Luiten, Aljosa Osep, Patrick Dendorfer, Philip Torr, Andreas Geiger, Laura Leal-Taixé, and Bastian Leibe. HOTA: A higher order metric for evaluating multi-object tracking. *IJCV*, 2021.

[48] Zhengyi Luo, S Alireza Golestaneh, and Kris M Kitani. 3D human motion estimation via motion compression and refinement. In *ACCV*, 2020.

[49] 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.

[50] Tim Meinhardt, Alexander Kirillov, Laura Leal-Taixe, and Christoph Feichtenhofer. TrackFormer: Multi-object tracking with transformers. In *CVPR*, 2022.

[51] Gyeongsik Moon and Kyoung Mu Lee. I2L-MeshNet: Image-to-lixel prediction network for accurate 3D human pose and mesh estimation from a single RGB image. In *ECCV*, 2020.

[52] Lea Müller, Vickie Ye, Georgios Pavlakos, Michael Black, and Angjoo Kanazawa. Generative proxemics: A prior for 3D social interaction from images. *arXiv preprint arXiv:2306.09337*, 2023.

[53] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Peter Gehler, and Bernt Schiele. Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In *3DV*, 2018.

[54] Austin Patel, Andrew Wang, Ilija Radosavovic, and Jitendra Malik. Learning to imitate object interactions from internet videos. *arXiv preprint arXiv:2211.13225*, 2022.

[55] Priyanka Patel, Chun-Hao P Huang, Joachim Tesch, David T Hoffmann, Shashank Tripathi, and Michael J Black. AGORA: Avatars in geography optimized for regression analysis. In *CVPR*, 2021.

[56] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3D hands, face, and body from a single image. In *CVPR*, 2019.

[57] Georgios Pavlakos, Jitendra Malik, and Angjoo Kanazawa. Human mesh recovery from multiple shots. In *CVPR*, 2022.

[58] Georgios Pavlakos, Ethan Weber, Matthew Tancik, and Angjoo Kanazawa. The one where they reconstructed 3D humans and environments in TV shows. In *ECCV*, 2022.

[59] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3D human pose and shape from a single color image. In *CVPR*, 2018.

[60] Owen Pearl, Soyong Shin, Ashwin Godura, Sarah Bergbreiter, and Eni Halilaj. Fusion of video and inertial sensing data via dynamic optimization of a biomechanical model. *Journal of Biomechanics*, 155:111617, 2023.

[61] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022.

[62] Xue Bin Peng, Angjoo Kanazawa, Jitendra Malik, Pieter Abbeel, and Sergey Levine. Sfv: Reinforcement learning of physical skills from videos. *ACM Transactions On Graphics (TOG)*, 37(6):1–14, 2018.

[63] Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, Christoph Feichtenhofer, and Jitendra Malik. On the benefits of 3D tracking and pose for human action recognition. In *CVPR*, 2023.

[64] Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, and Jitendra Malik. Tracking people with 3D representations. In *NeurIPS*, 2021.

[65] Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, and Jitendra Malik. Tracking people by predicting 3D appearance, location and pose. In *CVPR*, 2022.

[66] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J Guibas. HuMoR: 3D human motion model for robust pose estimation. In *ICCV*, 2021.

[67] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In *ECCV*, 2016.

[68] Anshul Shah, Shlok Mishra, Ankan Bansal, Jun-Cheng Chen, Rama Chellappa, and Abhinav Shrivastava. Pose and joint-aware action recognition. In *WACV*, 2022.

[69] Jie Song, Xu Chen, and Otmar Hilliges. Human body model fitting by learned gradient descent. In *ECCV*, 2020.

[70] Yu Sun, Wu Liu, Qian Bao, Yili Fu, Tao Mei, and Michael J Black. Putting people in their place: Monocular regression of 3D people in depth. In *CVPR*, 2022.

[71] Yu Sun, Yun Ye, Wu Liu, Wenpeng Gao, Yili Fu, and Tao Mei. Human mesh recovery from monocular images via a skeleton-disentangled representation. In *ICCV*, 2019.

[72] Garvita Tiwari, Dimitrije Antić, Jan Eric Lenssen, Nikolaos Sarafianos, Tony Tung, and Gerard Pons-Moll. Pose-NDF: Modeling human pose manifolds with neural distance fields. In *ECCV*, 2022.

[73] Vasileios Vasilopoulos, Georgios Pavlakos, Sean L Bowman, J Diego Caporale, Kostas Daniilidis, George J Pappas, and Daniel E Koditschek. Reactive semantic planning in unexplored semantic environments using deep perceptual feedback. *IEEE Robotics and Automation Letters*, 5(3):4455–4462, 2020.

[74] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.

[75] Ziniu Wan, Zhengjia Li, Maoqing Tian, Jianbo Liu, Shuai Yi, and Hongsheng Li. Encoder-decoder with multi-level attention for 3D human shape and pose estimation. In *ICCV*, 2021.

[76] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. HumanNeRF: Free-viewpoint rendering of moving people from monocular video. In *CVPR*, 2022.

[77] Chao-Yuan Wu, Justin Johnson, Jitendra Malik, Christoph Feichtenhofer, and Georgia Gkioxari. Multiview compressive coding for 3D reconstruction. In *CVPR*, 2023.

[78] Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipei Zhou, Guosen Lin, Yanwei Fu, Yizhou Wang, and Yonggang Wang. AI Challenger: A large-scale dataset for going deeper in image understanding. *arXiv preprint arXiv:1711.06475*, 2017.

[79] Yuliang Xiu, Jiefeng Li, Haoyu Wang, Yinghong Fang, and Cewu Lu. Pose Flow: Efficient online pose tracking. In *BMVC*, 2018.

[80] Xiangyu Xu, Hao Chen, Francesc Moreno-Noguer, Laszlo A Jeni, and Fernando De la Torre. 3D human pose, shape and texture from low-resolution images and videos. *PAMI*, 2021.

[81] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. ViTPose: Simple vision transformer baselines for human pose estimation. In *NeurIPS*, 2022.

[82] Yuanlu Xu, Song-Chun Zhu, and Tony Tung. DenseRaC: Joint 3D pose and shape estimation by dense render-and-compare. In *ICCV*, 2019.

[83] Vickie Ye, Georgios Pavlakos, Jitendra Malik, and Angjoo Kanazawa. Decoupling human and camera motion from videos in the wild. In *CVPR*, 2023.

[84] Ye Yuan, Umar Iqbal, Pavlo Molchanov, Kris Kitani, and Jan Kautz. GLAMR: Global occlusion-aware human mesh recovery with dynamic cameras. In *CVPR*, 2022.

[85] Andrei Zanfir, Elisabeta Marinoiu, and Cristian Sminchisescu. Monocular 3D pose and shape estimation of multiple people in natural scenes - The importance of multiple scene constraints. In *CVPR*, 2018.

[86] Wang Zeng, Wanli Ouyang, Ping Luo, Wentao Liu, and Xiaogang Wang. 3D human mesh regression with dense correspondence. In *CVPR*, 2020.

[87] Hongwen Zhang, Jie Cao, Guo Lu, Wanli Ouyang, and Zhenan Sun. DaNnet: Decompose-and-aggregate network for 3D human shape and pose estimation. In *Proceedings of the 27th ACM International Conference on Multimedia*, 2019.

[88] Hongwen Zhang, Yating Tian, Yuxiang Zhang, Mengcheng Li, Liang An, Zhenan Sun, and Yebin Liu. PyMAF-X: Towards well-aligned full-body model regression from monocular images. *PAMI*, 2023.

[89] Hongwen Zhang, Yating Tian, Xinchi Zhou, Wanli Ouyang, Yebin Liu, Limin Wang, and Zhenan Sun. PyMAF: 3D human pose and shape regression with pyramidal mesh alignment feedback loop. In *ICCV*, 2021.

[90] Jason Y Zhang, Panna Felsen, Angjoo Kanazawa, and Jitendra Malik. Predicting 3D human dynamics from video. In *ICCV*, 2019.

[91] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *CVPR*, 2019.

# Supplementary Material for:
# "Humans in 4D: Reconstructing and Tracking Humans with Transformers"

Shubham Goel    Georgios Pavlakos    Jathushan Rajasegaran    Angjoo Kanazawa*    Jitendra Malik*

{shubham-goel, pavlakos, jathushan, kanazawa}@berkeley.edu, malik@eecs.berkeley.edu
University of California, Berkeley

We provide more details about HMR 2.0, *i.e.*, the architecture we use (Section S.1), the data (Section S.2) and the training pipeline (Section S.3). Furthermore, we describe the aspect of pose prediction (Section S.4) and we discuss the metrics we use for evaluation (Section S.5). Then, we discuss the experimental settings for tracking (Section S.6), and action recognition (Section S.7). Finally, we provide additional qualitative results (Section S.8).

## S.1. HMR 2.0 architecture details

The architecture of our HMR 2.0 model is based on a ViT image encoder and a transformer decoder. We use a ViT-H/16 ("huge") pre-trained on the task of 2D keypoint localization [25]. It has 50 transformer layers, takes a $256 \times 192$ sized image as input, and outputs $16 \times 12$ image tokens, each of dimension 1280. Our transformer decoder is a standard transformer decoder architecture [23] with 6 layers, each containing multi-head self-attention, multi-head cross-attention, and feed-forward blocks, with layer normalization [2]. It has a 2048 hidden dimension, 8 (64-dim) heads for self- and cross-attention, and a hidden dimension of 1024 in the feed-forward MLP block. It operates on a single learnable 2048-dimensional SMPL query token as input and cross-attends to the $16 \times 12$ image tokens. Finally, a linear readout on the output token from the transformer decoder gives pose $\theta$, shape $\beta$, and camera $\pi$.

## S.2. Data details

In our training, we adopt the training data conventions of previous works [10], using images from Human3.6M [4], COCO [13], MPII [1] and MPI-INF-3DHP [18]. This forms the training set for the version we refer to as HMR 2.0a in the main manuscript. For the eventual HMR 2.0b version, we additionally generate pseudo-ground truth SMPL [14] fits for images from AVA [3], InstaVariety [6] and AI Challenger [24]. Since AVA and InstaVariety include videos, we collect frames by sampling at 1fps and 5fps respectively. For pseudo-ground truth generation, we use ViTDet [11] for bounding box detection and ViTPose [25] for key-

point detection, while fitting happens using ProHMR [10]. We discard detections with very few 2D detected keypoints (less than five) and low detection confidence (threshold 0.5). We also discard fits with unnatural body shapes (*i.e.*, body shape parameters outside $[-3, 3]$), unnatural body poses (computed using a per-joint histogram of poses on AMASS [17]), and large fitting errors (*i.e.*, which indicates that the reconstruction was not successful). For training our HMR 2.0b model, we sample with different probabilities from each dataset, *i.e.*, Human3.6M: 0.1, MPII: 0.1, MPI-INF-3DHP: 0.1, AVA: 0.15, AI Challenger: 0.15, InstaVariety: 0.2, COCO: 0.2.

## S.3. Training details

We train our main model using 8 A100 GPUs with an effective batch size of $8 \times 48 = 384$. We use an AdamW optimizer [15] with a learning rate of 1e-5, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and a weight decay of 1e-4. Training lasts for 1M iterations, which takes roughly six days. For our main model HMR 2.0b, we train the network end-to-end. However, for the HMR 2.0a variant, the ViT encoder remains frozen, allowing a larger effective batch size of $8 \times 512 = 4096$, learning rate of 1e-4, and fewer training iterations of 100K (*i.e.*, roughly equivalent number of epochs).

While training, we weigh the different losses. $\mathcal{L}_{\texttt{kp3D}}$, $\mathcal{L}_{\texttt{kp2D}}$, and $\mathcal{L}_{\texttt{adv}}$ have weights 0.05, 0.01, and 0.0005 respectively. The terms within $\mathcal{L}_{\texttt{smpl}}$ are also weighed differently, the $\theta$ and $\beta$ terms weigh 0.001 and 0.0005 respectively.

## S.4. Pose prediction

For the pose prediction model, we train a vanilla transformer model [23] from the tracklets obtained by [19]. Each tracklet at every time instance contains 3D pose and 3D location information, where the pose is parameterized by the SMPL model [14] and the location is represented as the translation in the camera frame. The transformer has 6 layers and 8 self-attention heads with a hidden dimension of 256. Each output token regresses the 3D pose and 3D loca-

tion of the person at the specified time-step. We train this model by randomly masking input pose tokens and applying the loss on the masked tokens. During inference, to predict a future 3D pose, we query the model by reading out from a future time-step, using a learned mask-token as input to that time-step. Similarly for amodal completion, we replace the missing detections with the learned mask-token and read out from the output at the corresponding time-step. The model is trained with a batch size of 64 sequences and a sequence length of 128 tokens. We use the AdamW optimizer [15] with a learning rate of 0.001 and $\beta_1 = 0.9, \beta_2 = 0.95$.

## S.5. Metrics

For our evaluation, we use the metrics that are common in the literature:

**3D Pose**: We follow [5] and we use MPJPE and PA-MPJPE. MPJPE refers to Mean Per Joint Position Error and it is the average L2 error across all joint, after aligning with the root node. PA-MPJPE is similar but is computed after aligning the predicted pose with the ground-truth pose using Procrustes Alignment.

**2D Pose**: We use PCK as defined in [26]. This is the Percentage of Correctly localized Keypoints, where a keypoint is considered as correctly localized if its L2 distance from the ground-truth keypoint is less than a threshold $t$. We report results using different thresholds (@0.05 and @0.1 of image size).

**Tracking**: Following [20, 21], we use standard tracking metrics. This includes ID switches (IDs), MOTA [7], IDF1 [22], and HOTA [16].

**Action Recognition**: We report results using mAP metrics as defined in the AVA dataset [3]. We further provided a more fine-grained analysis reporting results on different action categories: actions that involve Object Manipulation (OM), actions that involve Person Interactions (PI), and actions that involve Person Movement (PM). The results in these categories are also reported using mAP.

## S.6. Tracking with PHALP′

In the main manuscript, we compare different human mesh recovery systems on the downstream problem of tracking (Table 3 of the main manuscript). For this, we modify the PHALP approach [21], so that pose distance is computed on the SMPL space that all the models share. To make this comparison fair, we keep other variables similar to the original PHALP (*e.g.*, same appearance embedding). Note that this comparison is generous to baselines that do not model appearance themselves. Eventually, our final 4DHumans system uses a sampling-based appearance head and our new pose prediction, which lead to the state-of-the-art performance for tracking on PoseTrack (Table 4

of the main manuscript). To model appearance, we texture visible points on the mesh by projecting them onto the input image and sampling color from the corresponding pixels.

## S.7. Action recognition

As an alternative way to assess the quality of 3D human reconstruction, we evaluate various human mesh recovery systems on the downstream task of action recognition on AVA (please refer to [19] for more details on the task definition). More specifically, we take the tracklets from [19], which were generated by running PHALP [21] on the Kinetics [8] and AVA [3] datasets. Then, we replace the poses from various human mesh recovery models (*i.e.*, PyMAF [28], PyMAF-X [27], PARE [9], CLIFF [12], HMAR [21], HMR 2.0) and evaluate their performance on the action recognition task. In this pose-only setting, the action recognition model has access only to the 3D poses (in the SMPL format) and 3D location and is trained to predict the action of each person. For a fair comparison and to achieve the best performance for each 3D pose regressor, we retrain the action recognition model specifically for each 3D pose method.

## S.8. Additional qualitative results

We have already provided a lot of qualitative results of HMR 2.0, both in the main manuscript and in videos on the project webpage. Here, we provide additional results, including comparisons with our closest competitors (Figure S.1), and a demonstration of our results in a variety of challenging cases, including successes (Figure S.2) and failure cases (Figure S.3).

## 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] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

[3] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, and Jitendra Malik. AVA: A video dataset of spatio-temporally localized atomic visual actions. In *CVPR*, 2018.

[4] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments. *PAMI*, 2013.

[5] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018.
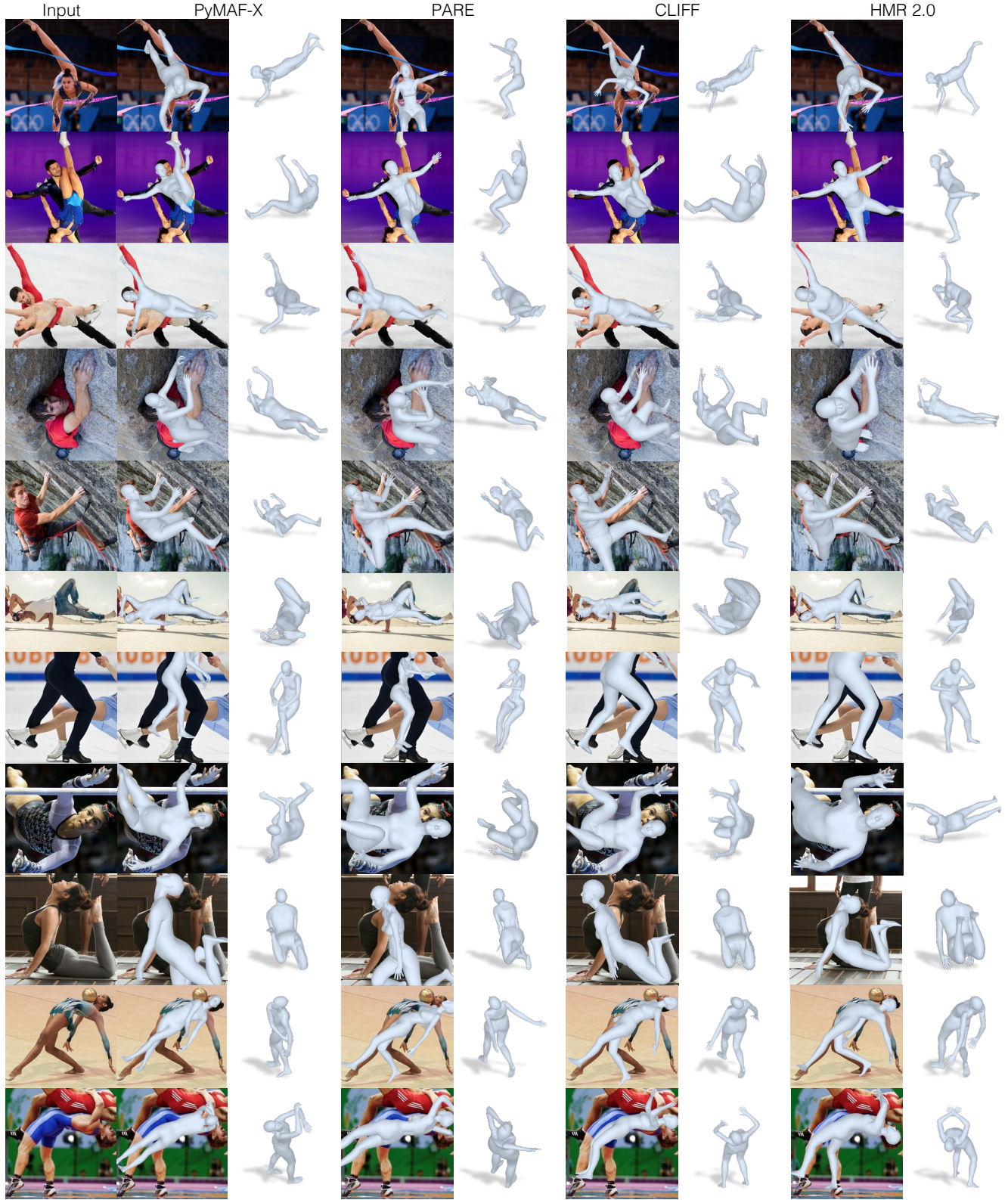
Figure S.1: **Qualitative comparison of our approach with state-of-the-art methods.** We compare HMR 2.0 with our closest competitors, PyMAF-X [27], PARE [9] and CLIFF [12]. For each example, we show the input image, and results from each method (including the frontal and a side view). HMR 2.0 is significantly more robust in a variety of settings, including images with unusual poses, unusual viewpoints and heavy person-person overlap.

Figure S.2: **Qualitative results of our approach on challenging examples.** For each example we show the input image, the reconstruction overlay, a side view and the top view. The examples include unusual poses, unusual viewpoints, people in close interaction, extreme truncations and occlusions, as well as blurry images.
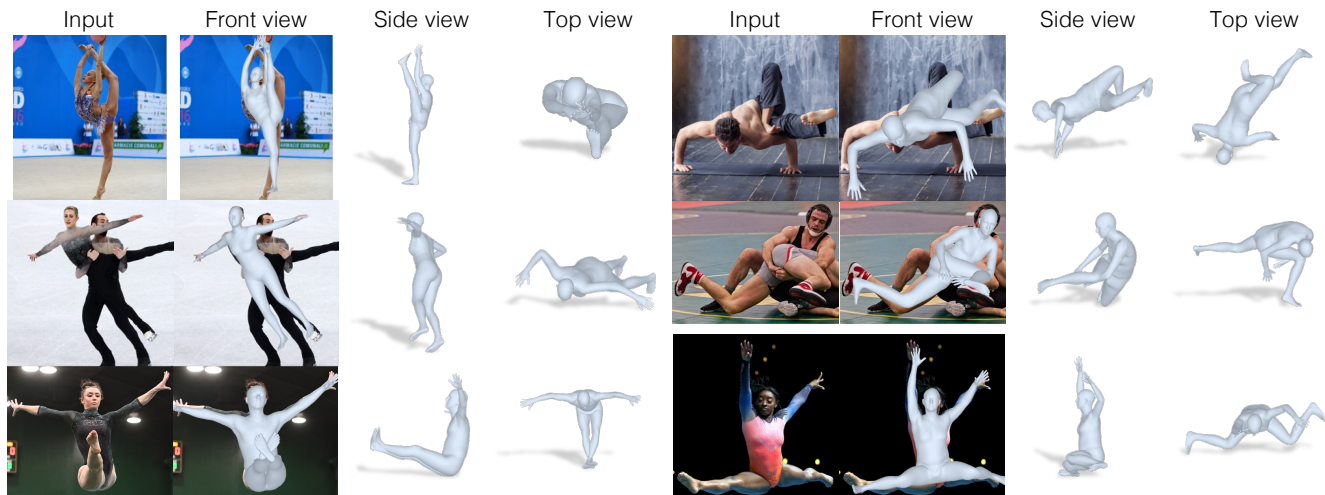
Figure S.3: **Failures of single frame 3D human reconstruction with HMR 2.0.** Despite the increased robustness of our method, we observe that HMR 2.0 occasionally recovers erroneous reconstructions in cases with very unusual articulation (first row), heavy person-person interaction (second row), and very challenging depth ordering for the different body parts (third row).

[6] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3D human dynamics from video. In *CVPR*, 2019.

[7] Rangachar Kasturi, Dmitry Goldgof, Padmanabhan Soundararajan, Vasant Manohar, John Garofolo, Rachel Bowers, Matthew Boonstra, Valentina Korzhova, and Jing Zhang. Framework for performance evaluation of face, text, and vehicle detection and tracking in video: Data, metrics, and protocol. *PAMI*, 2008.

[8] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.

[9] Muhammed Kocabas, Chun-Hao P Huang, Otmar Hilliges, and Michael J Black. PARE: Part attention regressor for 3D human body estimation. In *ICCV*, 2021.

[10] Nikos Kolotouros, Georgios Pavlakos, Dinesh Jayaraman, and Kostas Daniilidis. Probabilistic modeling for human mesh recovery. In *ICCV*, 2021.

[11] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *ECCV*, 2022.

[12] Zhihao Li, Jianzhuang Liu, Zhensong Zhang, Songcen Xu, and Youliang Yan. CLIFF: Carrying location information in full frames into human pose and shape estimation. In *ECCV*, 2022.

[13] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014.

[14] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015.

[15] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

[16] Jonathon Luiten, Aljosa Osep, Patrick Dendorfer, Philip Torr, Andreas Geiger, Laura Leal-Taixé, and Bastian Leibe. HOTA: A higher order metric for evaluating multi-object tracking. *IJCV*, 2021.

[17] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. AMASS: Archive of motion capture as surface shapes. In *ICCV*, 2019.

[18] 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.

[19] Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, Christoph Feichtenhofer, and Jitendra Malik. On the benefits of 3D tracking and pose for human action recognition. In *CVPR*, 2023.

[20] Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, and Jitendra Malik. Tracking people with 3D representations. In *NeurIPS*, 2021.

[21] Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, and Jitendra Malik. Tracking people by predicting 3D appearance, location and pose. In *CVPR*, 2022.

[22] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In *ECCV*, 2016.

[23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.

[24] Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipei Zhou, Guosen Lin, Yanwei Fu, Yizhou Wang, and Yonggang Wang. AI Challenger: A large-scale dataset for going deeper in image understanding. *arXiv preprint arXiv:1711.06475*, 2017.

[25] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. ViTPose: Simple vision transformer baselines for human pose estimation. In *NeurIPS*, 2022.

[26] Yi Yang and Deva Ramanan. Articulated human detection with flexible mixtures of parts. *PAMI*, 2012.

[27] Hongwen Zhang, Yating Tian, Yuxiang Zhang, Mengcheng Li, Liang An, Zhenan Sun, and Yebin Liu. PyMAF-X: Towards well-aligned full-body model regression from monocular images. *PAMI*, 2023.

[28] Hongwen Zhang, Yating Tian, Xinchi Zhou, Wanli Ouyang, Yebin Liu, Limin Wang, and Zhenan Sun. PyMAF: 3D human pose and shape regression with pyramidal mesh alignment feedback loop. In *ICCV*, 2021.