

# PedX: Benchmark Dataset for Metric 3D Pose Estimation of Pedestrians in Complex Urban Intersections

Wonhui Kim<sup>1</sup>, Manikandasriram Srinivasan Ramanagopal<sup>1</sup>, Charles Barto<sup>1</sup>, Ming-Yuan Yu<sup>1</sup>, Karl Rosaen<sup>1</sup>, Nick Goumas<sup>1</sup>, Ram Vasudevan<sup>2</sup> and Matthew Johnson-Roberson<sup>3</sup>

**Abstract**—This paper presents a novel dataset titled *PedX*, a large-scale multimodal collection of pedestrians at complex urban intersections. *PedX* consists of more than 5,000 pairs of high-resolution (12MP) stereo images and LiDAR data along with providing 2D and 3D labels of pedestrians. We also present a novel 3D model fitting algorithm for automatic 3D labeling harnessing constraints across different modalities and novel shape and temporal priors. All annotated 3D pedestrians are localized into the real-world metric space, and the generated 3D models are validated using a mocap system configured in a controlled outdoor environment to simulate pedestrians in urban intersections. We also show that the manual 2D labels can be replaced by state-of-the-art automated labeling approaches, thereby facilitating automatic generation of large scale datasets.

## I. INTRODUCTION

Driving in complex urban environments is one of the major challenges for autonomous vehicles (AVs). For AVs to operate in an environment crowded with people, understanding pedestrian pose, motion, behavior, and intention will greatly increase our ability to function safely and efficiently.

In computer vision, estimating human pose has been a long standing problem. The recent application of deep neural networks has generated state-of-the-art results for 2D body pose estimation [1], which has inspired extensions to the 3D pose estimation [2], [3], [4], [5], [6]. However, gathering ground truth 3D pose data is challenging. Motion capture (mocap) systems have been the primary generator of ground truth 3D pose data, but have restricted the variety and complexity of the 3D scenes that can be captured [7], [8]. For example, with mocap systems it is difficult to capture naturalistic in-the-wild scenes with groups of people who are moving and interacting. To overcome those technical limitations, this paper develops both a dataset and a ground truth generation approach to facilitate generating 3D poses on in-the-wild images without relying on mocap.

Most AVs have cameras installed, so this data can serve as a primary source for human pose estimation using computer



Fig. 1: An example frame from the PedX dataset. *Left*: A camera image with bounding boxes around the pedestrians within 5-45m to the camera. *Right*: A rendered image with the 3D human mesh models in metric scale.

vision algorithms. In addition to cameras, LiDAR (Light Detection and Ranging) has become an essential component for AVs due to its precise depth measurements. This motivated the importance of capturing both modalities for this benchmark set of complex urban intersections.

Our dataset has the following unique properties. First, the data are gathered outdoors with real challenges such as varied lighting and weather conditions and the presence of occlusions. Second, the pedestrian data are collected at intersection length scales of up to 45m range, which is relevant for deployment of pose estimation systems at application relevant scales. The captured scenes are also naturalistic. The pedestrians in our dataset are not actors, so they move and interact in a myriad of realistic ways. In addition, our dataset includes multi-person images capturing crowds of people. Lastly, we release multimodal pedestrian data including high resolution images and point clouds that are synchronously captured from stereo cameras and LiDAR sensors.

Annotations of our dataset also have distinctive features. All the annotated 3D pedestrians lie in a global metric-space coordinate frame, as opposed to many existing datasets that operate in hip joint or camera center relative coordinate frames. We stress the importance of determining where a person is in the 3D world so one can plan actions around them. In addition, our multimodal data frames are captured in minutes long sequences with unique tracking IDs for each pedestrian, enabling temporal reasoning.

The contributions of this paper are summarized as follows:

\*This work was supported by a grant from Ford Motor Company via the Ford-UM Alliance under award N022884

<sup>1</sup> Wonhui Kim, Manikandasriram S. R., Charles Barto, Ming-Yuan Yu, Karl Rosaen, and Nick Goumas are research assistants and engineers with UM and Ford Center for Autonomous Vehicles Laboratory, University of Michigan, Ann Arbor, MI 48105, USA <https://fcav.engin.umich.edu/>

<sup>2</sup> Ram Vasudevan is with Faculty of Mechanical Engineering, University of Michigan, Ann Arbor, MI 48105, USA [ramv@umich.edu](mailto:ramv@umich.edu)

<sup>3</sup> Matthew Johnson-Roberson is with Faculty of Naval Architecture & Marine Engineering, Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, MI 48105, USA [mattjr@umich.edu](mailto:mattjr@umich.edu)

1. We release a publicly available large-scale multimodal pedestrian dataset with a rich set of 2D and 3D annotations. The dataset captures the real world challenges of urban intersections.
2. We present an automatic method to obtain full 3D labels from 2D data, enabling labeling of in-the-wild images without mocap.
3. Our automatic 3D labeling method is validated using a mocap system in a controlled outdoor environment that simulates pedestrians in urban intersections.
4. Using state of the art algorithms for 2D annotations, our proposed approach allows generating 3D data in a completely unsupervised manner.

We present this dataset to enable the study of 3D pose estimation while reasoning about pedestrian behavior around vehicles, particularly in crowded urban areas as depicted in Fig. 1. We see this as one of the first areas where human pose estimation can have a tremendous impact on safety and intelligence of mobile robot systems. Understanding the pose of road users affords information about activity, attention, and predictions of future position which are critical to safely navigate around humans.

## II. RELATED WORK

### A. 3D human pose estimation

In many papers, 3D human pose estimation has been formulated as a problem of regressing 3D joint locations by directly extracting visual features from an image [9], [10], [2], [3], [4], [11], [5], [12], or by lifting 2D joint detector outputs to 3D joints in a camera relative frame [6], [13], [14]. To reduce the inherent ambiguity of 3D human pose estimation from a 2D image, the constraints on feasible human pose have additionally been considered [13], [15], or a deformable 3D human model such as Skinned Multi-Person Linear (SMPL) [16] has been used to be fit to known 2D joint locations [14].

Recently, pose representations other than sparse 3D joints have been used to formulate 3D human pose estimation. SMPL model parameters are directly estimated given dense 2D keypoints [17], or the parameters are predicted using end-to-end deep networks with the adversarial loss [18]. UV parameterization of the 3D human body surface is estimated from the dense image-to-surface correspondence regression networks [19]. The per-pixel body depth map for each soccer player in a YouTube video is estimated using a neural network trained on synthetic data [20].

Most approaches take as input an image with a single person or a cropped image patch centered around a person, and return as output 3D pose in a root-relative coordinate frame where the camera is facing toward the person. While the joint estimation of 3D pose and virtual camera parameters have been proposed [13], [5], the outputs are still not in real metric scale. Without knowledge of exactly how far away a person is, controlling a mobile robot safely will be challenging. Most approaches are also unable to handle multi-person images with a single pass of an algorithm. While DensePose [19]

handles in-the-wild images with multiple people, dense pose outputs from a multi-person image cannot be converted into a common unified coordinate frame.

A major impediment to predicting metric space pose for multiple people in a scene is the lack of a suitable dataset with reliable 3D annotations. Addressing this limitation is the focus of this paper.

### B. 3D human pose datasets

Large scale datasets have played an essential role in fueling recent progress in a variety of computer vision tasks. However, building a ground truth 3D human pose dataset in metric space is challenging since annotation in 3D is far more time consuming than the same task in 2D. To avoid manual labeling in 3D, mocap systems are typically used to obtain the ground truth 3D human pose [8], [7], [21], [22].

The data captured with mocap has many limitations. For instance, markers must be attached to subjects which makes images look unnatural. Moreover since mocap is typically restricted to constrained, mostly indoor areas, image backgrounds for mocap datasets have limited variability. In addition the number of makers that can be tracked by mocap systems is limited, which restricts the number of subjects that can appear in any scene.

To counter the limited variability in subject appearance, camera viewpoints, lighting conditions and image backgrounds, synthetic 3D datasets have been proposed [23], [24]. While there have been attempts to improve the photo-realism of these synthetic data generation pipelines [25], [26], [27], state-of-the-art synthetic images are still easy to distinguish from real images. Others have explored techniques to automatically do 3D labeling with limited human intervention [28], [17], [29]. For instance, some have taken advantage of a multiview camera setup to obtain reliable 3D annotations using optimization [28]. Others explore fitting parameterized mesh models [14] to monocular images or multiview images to collect the ground truth 3D labels [17], [29]. However, these methods provide 3D scale models for a virtual camera relative frame and typically for images from various 2D pose datasets cropped around a single detection. This limits the potential utility of these methods while performing mobile robotic tasks safely. In contrast, we construct 3D human pose models in metric space for crowded urban street intersections that include as many as 15 pedestrians in a single image at distances as far as 45m from the camera.

TABLE I: Statistics and characteristics of related datasets.

	num. images	num. instances	num. pixels	3D points?	video?	multi- person?	real?	scene type
H36M [8]	3,600,000	900,000	≈1M	Y	Y	N	Y	mocap indoor
HumanEva [7]	≈80,000	≈80,000	1-1.3M	Y	N	N	Y	mocap indoor
SURREAL [25]	653,6752	653,6752	≈80K	N	Y	N	N	indoor
UP-3D [17]	8,515	8,515	≈300K	N	N	N	Y	indoor+outdoor
DensePose-COCO [19]	33,929	131,129	≈300K	Y	N	Y	Y	indoor+outdoor
Ours	10,152	14,091	≈10M	Y	Y	Y	Y	outdoor

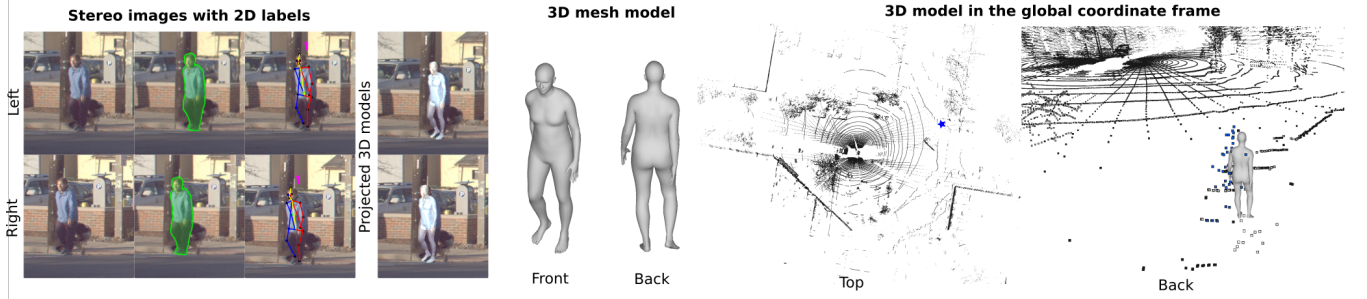


Fig. 2: Visualization of our annotations. For each pedestrian, we provide 2D segmentation, 2D joint locations with visibility of 18 body joints, tracking ID, time-synced LiDAR points, and 3D mesh model localized into the global coordinate frame.

### III. PEDX DATASET

The PedX dataset contains more than 5,000 pairs of high-resolution stereo images with 2,500 frames of 3D LiDAR pointclouds. The cameras and LiDAR sensors are calibrated and time synchronized. We selected three four-way stop intersections with heavy pedestrian-vehicle interaction. Cameras are installed on the roof of the car to obtain driver-perspective images. To cover all four crosswalks at an intersection, the images were captured by two pairs of stereo cameras – one pair facing forward and another facing the incoming road from the left. Our dataset includes more than 14,000 pedestrian models with a distance of 5-45m from the cameras and we provide reliable 2D and 3D labels for each instance. Table I presents statistics of our dataset in comparison to other publicly available 3D human pose datasets.

Instance-level 2D segmentations and body joint locations are manually labeled for all images by instructed annotators. We also provide the unique tracking ID for each instance across the frames if the pedestrian appears in consecutive frames. A 2D segmentation is labeled by outlining a single connected polygon to cover the entire visible area of an object. To label the keypoints, 18 body joints are selected including 4 facial components. In cases where some keypoints are invisible due to occlusion by other objects in the scene or self-occlusion, annotators made a reasonable guess about the joint location, but also indicated the degree of occlusion.

We use the SMPL [16] parameterization to represent our 3D annotations which consist of shape and pose parameters and additionally estimate the global location and orientation of the instance in the global coordinate frame. The best fit 3D SMPL models were computed using the 2D annotations from a pair of stereo images and LiDAR points. The obtained 3D models encode pose, shape, and global location in 3D metric space without scale ambiguities. Fig. 2 illustrates 2D and 3D annotations from our dataset. The automatic algorithm to determine the optimal SMPL parameters is discussed in the next section.

### IV. MULTI-MODAL 3D MODEL FITTING

We perform 3D model fitting on pedestrians in a stereo-LiDAR sequence. In contrast to the previous work [14] that fits a 3D model to a single frame at a time, our approach optimizes over a sequence of stereo images and LiDAR

points. We begin by per-instance model fitting which is then extended to optimize over a sequence of instances using an iterative method. We also propose multi-modal and temporal priors. Note that we use a gender-neutral model for model fitting. Before initiating the model fitting pipeline, we preprocess the LiDAR data to identify regions containing potential pedestrians in 3D space. Using 2D segmentation labels for stereo images with known transformation between LiDAR and camera coordinate frames, we perform the point cloud labeling of each pedestrian instance.

#### A. Fitting to a single instance (at a single time step)

We begin by performing 3D model fitting to a single instance at a single time step. For each pedestrian instance, we are given 2D joint locations  $x_l$  and  $x_r$ , and 2D segmentations  $S_l$  and  $S_r$  for each stereo image. We also have sparse 3D points corresponding to the instance. To find the pose  $\theta$ , shape  $\beta$ , and 3D global position  $t$  that best fit to the instance, we formulate the problem as:

$$\underset{\theta, \beta, t}{\text{minimize}} E_I(\theta, \beta, t) \quad (1)$$

where  $E_I = E_J + E_{3d} + E_P + E_T + E_D$  represents the sum of multiple energy terms. We verify the effectiveness of each energy term through ablative experiments described in Sec. V-B.  $E_J$  is the sum of robust 2D reprojection error [14] for both left and right images,  $E_P$  is the prior term,  $E_{3d}$  is the 3D Euclidean distance term between visible SMPL vertices and the LiDAR points,  $E_T$  is the translation term to constrain the 3D model location, and  $E_D$  is the heading direction term to constrain the body orientation:

$$E_T(t) = \|t - t_0\|_2^2, \quad (2)$$

$$E_D(\theta) = \|f(\theta) - d\|_2^2, \quad (3)$$

$$E_{3d} = \frac{1}{N_v} \sum_i^{N_v} \min_j \|X_i - V_j\|_2^2, \quad (4)$$

where  $t_0$  is the mean of 3D points,  $f$  is a function to convert the axis-angle representation of body orientation to xyz-directional vector,  $d$  is a known heading direction vector,  $X_i$  is the  $i$ -th LiDAR point,  $N_v$  is the total number of 3D points that belong to the instance, and  $V_j$  is the  $j$ -th point of SMPL model vertices.

### B. Fitting to a sequence of single instances

To fit 3D models to a sequence of detections, we develop *shape* and *temporal consistency* constraints across frames in addition to the per-frame constraints.

1) *Global shape consistency*: Suppose one pedestrian appears in  $N$  consecutive frames with full 2D labels. In this instance, while pose parameters and translations change, the shape parameters should remain unchanged across the sequence. To find the pose, shape parameters and translations, we formulate the problem as:

$$\underset{\Theta, \beta, T}{\text{minimize}} E_{seq}(\Theta, \beta, T) \quad (5)$$

where  $\Theta = \{\theta_1, \dots, \theta_N\}$  and  $T = \{t_1, \dots, t_N\}$  are the set of pose parameters and the set of translations for all  $N$  frames.  $\beta$  is the shape parameters shared by all frames where  $\beta = \beta_1 = \dots = \beta_N$ . Optimizing over the entire sequence is challenging due to the high dimension of the decision variables. Since the objective  $E_{seq}$  is separable in terms of a per-frame objective, we decompose this large problem into a set of smaller problems. We rewrite the unconstrained minimization problem over the entire sequence by introducing the consensus variable  $\beta$ :

$$\begin{aligned} & \underset{\theta_{1:N}, \beta_{1:N}, t_{1:N}, \beta}{\text{minimize}} \sum_{k=1}^N E_{I,k}(\theta_k, \beta_k, t_k) \\ & \text{subject to } \beta_k - \beta = 0, k \in \{1, \dots, N\} \end{aligned} \quad (6)$$

where  $k$  denotes the frame ranging from 1 to  $N$  frames in a sequence. This optimization is a constrained minimization problem with a separable objective function and multiple constraints that require that each per-frame shape parameter  $\beta_k$  is equal. The advantage of introducing a consensus variable  $\beta$  is that we can enforce all the frames to have common shape parameters while exploiting parallelism. We solve the problem by using the alternating direction method of multipliers (ADMM) [30]. The augmented Lagrangian to be minimized is:

$$\begin{aligned} \mathcal{L}_\rho(\beta, \beta_k; \Theta, T) = & \sum_{k=1}^N E_{I,k}(\beta_k; \Theta, T) \\ & + \mathbf{y}_k^T(\beta_k - \beta) + \frac{\rho}{2} \|\beta_k - \beta\|_2^2 \end{aligned} \quad (7)$$

$\mathbf{y}_k$  is the dual variable for  $\beta_k$  and  $\rho$  is a positive constant that is experimentally selected.  $\rho=2$  was used for the results reported in this paper. The objective  $\mathcal{L}_\rho$  is optimized using an alternating optimization for the local and global shape parameters with variables  $\{\mathbf{u}_k\}_{k=1}^N$  where  $\mathbf{u}_k = \frac{\mathbf{y}_k}{\rho}$ . The update equations at each iteration are as follows:

$$\beta_k^{t+1} := \underset{\beta_k}{\text{argmin}} E_{I,k}(\beta_k; \Theta, T) + \frac{\rho}{2} \|\beta_k - \beta^t + \mathbf{u}_k^t\|_2^2 \quad (8)$$

$$\beta^{t+1} := \frac{1}{N} \sum_{k=1}^N (\beta_k^{t+1} + \mathbf{u}_k^t) \quad (9)$$

$$\mathbf{u}_k^{t+1} := \mathbf{u}_k^t + \beta_k^{t+1} - \beta^{t+1} \quad (10)$$

We perform the synchronous update for the global shape parameters. The iteration is stopped when  $\|\beta_k^{t+1} - \beta^{t+1}\|_2 < 0.05$  and  $\rho \|\beta^t - \beta^{t+1}\|_2 < 0.05$  or the maximum iteration is reached. (8) is similar to per-frame minimization in (1) with an additional term in the objective function.

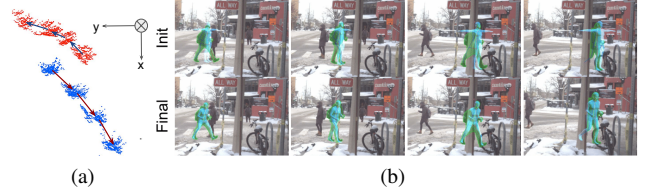


Fig. 3: (a) Overhead view of the pointcloud trajectories of pedestrians in the global frame. (b) Initialized body orientation based on heading direction (top) and the final 3D models after the iterations (bottom).

2) *Temporal pose prior*: In addition to enforcing per-frame shape parameters to share the common values across the sequential frames, we propose a temporal pose prior to give a penalty to unlikely sequences of poses. A 72-dimensional pose vector consists of xyz rotation angles of 23 joints relative to each of their parent nodes, plus the orientation of the root hip in 3D angle-axis representation. We observe that the difference between the pose vectors from two consecutive frames is small and has some patterns common to individual joints. Especially for pedestrians at the intersection, many of them are involved in walking or other actions at slow speed. Since the difference in rotation angles can also be affected by the translation velocity of pedestrians, we define a 75-dimensional vector with the first 3 components consisting of the difference in body translation:  $\Delta \mathbf{x} = (\Delta t, \Delta \theta)$ .

We fit a Gaussian mixture models with 10 distinct distributions for the pose difference vectors using the CMU mocap dataset [22]. Since the frame rate of the mocap and our data capture is different, we use the mocap frame rate when estimating this pose difference vector. We include the negative log of this multivariate normal probability distribution as part of the prior in the objective term:

$$\Delta \mathbf{x}^t = \mathbf{x}^t - \mathbf{x}^{t-1} \sim \sum_i^N w_i \mathcal{N}(\Delta \mathbf{x}^t; \mu_i, \Sigma_i) \quad (11)$$

$$E_{tp}(\mathbf{t}_k, \theta_k; \mathbf{t}_{k-1}, \theta_{k-1}) = -\log \sum_i^N w_i \mathcal{N}(\Delta \mathbf{x}^t; \mu_i, \Sigma_i) \quad (12)$$

where  $\mu_i, \Sigma_i$  are the mean and covariance of the pose difference vector  $\Delta \mathbf{x}$  and  $w_i$  is the weight for the  $i$ -th Gaussian mixture component.

### C. Initialization

One of the challenges while fitting a 3D model is estimating a body orientation with the incorrect sign [29]. To avoid getting a flipped model, we compute the heading direction of each pedestrian from a sequence of LiDAR points and use it to initialize the body orientation via a 3-dimensional angle-axis representation. We assume that pedestrians never move backwards (which held true in our data capture).

Fig. 3a shows some example trajectories from LiDAR points, and Fig. 3b shows projected 3D model onto the left images at four sequential time frames after the initialization. When a sequence is only a single frame, we find the initial body orientation with the template pose by minimizing the stereo reprojection error  $E_J$  only using torso body joints and the translation error  $E_T$  around the mean LiDAR points.



## V. EXPERIMENTS

The data in our dataset is captured at complex outdoor urban intersections where pedestrian-to-camera distances are large (5-45m), with multiple subjects who are often heavily occluded. In contrast, publicly available 3D datasets rely on mocap systems [21], [8], [7] guarantee less than a few mm accuracy for a *single* subject appearing within a controlled *indoor* capture volume of a *few meters* in radius. Given these distinctions, the accuracy of our dataset needs to be evaluated under these realistic conditions. While comparing the accuracy of our proposed approach with a mocap system would reliably validate our dataset, mocap systems cannot be practically setup at urban intersections.

We address this challenge by leveraging the fact that our proposed approach only requires 2D labels and LiDAR data without considering image features. The key factors affecting our approach include the density of the LiDAR returns due to the large distance from the capture vehicle, occlusions by other vehicles and other pedestrians which affect the LiDAR segmentation, the precision of the manual annotation when the pedestrian occupies a small portion of the image and the calibration between LiDAR and camera. The lighting conditions, background appearance, clothing and weather conditions do not affect our approach as we operate on manually labeled joint locations. Therefore, we collect and annotate an evaluation dataset in a controlled outdoor environment with a mocap system and get manual 2D labels while replicating the vehicle to target distances and the clutter and occlusion of the intersection data. By showing that the 3D labels generated by our method, using the same hand annotation process, is comparable to the mocap ground truth, we independently verify the fitting approach and the hand annotation process against a traditional mocap source.

### A. Data Verification

1) *Evaluation dataset*: We use the PhaseSpace mocap system with active LED markers which can be used in outdoor environments. The subject wears a suit with markers placed around body parts and repeats actions such as walking, jogging, and waving that are common for pedestrians. The capture vehicle was parked about 20m away from the mocap setup. To replicate typical occlusions we parked another car between the capture vehicle and the mocap setup as well as having groups of pedestrians walking. We selected 626 frames and obtained manual 2D labels for the images. Since the visual appearance does not affect the evaluation, we restrict the evaluation to a single subject with a single background and focus more on variation in poses and occlusions.

2) *Evaluation metric*: The 3D mean per joint position error (MPJPE) is a standard metric to evaluate pose estimation algorithms which is a mean over all joints of the euclidean distance between ground truth and prediction. In cases where the prediction is not in metric space, the error is computed for a root-relative coordinate frame after allowing a similarity transform to register the prediction to the ground truth. In cases where the prediction is in metric space, we compute

TABLE II: 3D MPJPE in root-relative coordinate frames.

Relative (mm)	rknee	lknee	rankl	lankl	rsho	lsho	relb	lelb	rwri	lwri	head	neck	hip	mean
[3]	113	153	203	184	141	120	130	132	203	194	134	96	88	147
[6]	107	116	172	159	159	160	142	127	197	178	84	88	87	137
[14]	<b>103</b>	89	139	158	77	59	<b>74</b>	71	144	145	85	35	<b>87</b>	97
Ours	104	<b>71</b>	<b>126</b>	<b>136</b>	<b>62</b>	<b>50</b>	83	<b>67</b>	<b>118</b>	<b>130</b>	<b>66</b>	<b>29</b>	106	<b>88</b>

TABLE III: 3D MPJPE in global coordinate frames.

Global (mm)	rknee	lknee	rankl	lankl	rsho	lsho	relb	lelb	rwri	lwri	head	neck	hip	mean
Triangulation	1178	1307	1617	1489	1205	1164	1130	1168	1206	1111	842	1029	1111	1194
Left+disp	766	977	1229	1018	702	648	756	766	825	789	482	437	972	794
[14]+disp	535	576	627	591	620	574	561	570	652	594	588	587	598	593
Ours	<b>205</b>	<b>179</b>	<b>250</b>	<b>255</b>	<b>183</b>	<b>177</b>	<b>187</b>	<b>182</b>	<b>221</b>	<b>204</b>	<b>169</b>	<b>155</b>	<b>161</b>	<b>194</b>

MPJPE in global coordinate frame without any registration. We further report the per joint position errors for both frames. Note that, given the geometry of capture, markers can get completely occluded from the mocap system. Consequently, not all joints are visible in all frames. Moreover, some methods may not predict invisible joints. Therefore, we take the weighted mean while computing the MPJPE where the weight is equal to the number of frames in which the joint was visible in the ground truth and was predicted.

3) *Baseline Methods*: We consider three different families of baseline methods. First, we consider a method that predicts 3D joint coordinates (up to scale) directly from 2D images [3]. Second, we consider methods that take manual 2D joint annotations as inputs [6], [14]. We evaluate these methods in the root relative frame alone. Third, we consider three naive baselines that use stereo geometry information. As we have 2D joint locations for a calibrated rectified stereo pair of images, we directly triangulate these 2D joint locations for visible joints. We refer to this method as *Triangulation*. For the second naive baseline, we use disparity values and 2D joint locations in the left image for the visible joints and the previous triangulation result for invisible joints. We refer to this as *Left+disp*. Finally, we consider a baseline that modifies an existing technique [14], this approach uses the calibrated camera parameters and the estimated skeletons which are scaled to metric space by using the average disparity values at the visible joint locations. We refer to this as *SMPLify [14]+disp*.

4) *Accuracy*: Table II and Table III summarize the results for root-relative and global frames respectively. Although fair comparisons can only be done between methods separated by horizontal lines, the objective of these tables is to highlight that the current state-of-the-art still has room for improvement and the utility of the proposed dataset in closing this gap. The proposed approach achieves lower MPJPE for majority of joints in the root-relative frame. This is expected as our approach leverages additional information including LiDAR returns, temporal priors as well as stereo annotations. The gains while using this additional data are most prominent in the global coordinate frame as no registration is involved and consequently global translation, orientation and scale errors can be seen in the global MPJPE. While naive baselines such as triangulation and left+disp perform poorly as they

do not leverage any prior about proportions of a typical human skeleton, SMPLify+disp which leverages priors about human skeletons still suffers from large errors. In contrast, our proposed approach achieves an MPJPE of 194mm for an average camera-to-pedestrian distance of  $20 \times 10^3$ mm.

### B. Ablation Study

Table IV summarizes the results for using different subsets of energy terms in the optimization. Note that  $E_{J,l}$  and  $E_{J,r}$  represents the reprojection error on left and right images.  $E_T$ ,  $E_{3D}$  and  $E_{tp}$  are defined in Sec. IV. Each column shows per-joint errors in root-relative frame except the last column which shows the MPJPE in global frame.

1) *Effect of stereo*: To see how using stereo reprojection error affects the resulting 3D models, we compute reprojection errors for only the left images (i.e. row 4 of Table IV) and for both stereo images (i.e. row 5 of Table IV). Stereo imagery reduced both global translation error and root-relative pose error as it reduces the depth ambiguities that exist for the monocular approach. The second row of Fig. 4 shows the results from monocular approach which estimates 3D models from the left images. Notice that in several frames, the legs are swapped or the body orientation is estimated incorrectly. Fig. 5 presents additional results with occlusions, and illustrates similar limitations of the existing approach. We can see that the projection of the estimated 3D models onto the right image does not align exactly. Estimating 3D pose from 2D joints is an inherently ill-posed problem because there exist many feasible body configurations. Using stereo information provides more constraints during the optimization and overall, reduces such ambiguities. Furthermore, when some joints are occluded in a single image, when using stereo pairs the second image of the pair may observe these joints which can reduce the uncertainty and produce better 3D models.

2) *Effect of using LiDAR points*: Although stereo images provide reasonable depth estimation in a global coordinate frame, the translation error may be too large since an error of a few pixels during labeling may create a large resulting error in fit. To place 3D models at the correct location in metric space, we include LiDAR information as a form of the translation prior  $E_T$ . The translation prior term localizes the 3D models at the distance observed by the LiDAR. As shown in the first two rows in Table IV, adding this translation prior provides an improvement in estimating global 3D pose.

The translation prior only constrains the location of the root joint (hip) in the 3D metric space. Therefore, depth ambiguities in other parts of the body may still exist, especially when the subject appears sideways with respect to the camera. Adding the 3D distance term  $E_{3D}$  helps to adjust the pose or body orientation to fit to the observed LiDAR points. Consequently the mean column in the root-relative frame significantly reduces from row 2 to row 4.

3) *Temporal prior*: The temporal prior term,  $E_{tp}$ , penalizes unlikely transitions of poses and translations between consecutive frames. It also makes the resulting 3D pose between consecutive frames appear smooth. In Table IV, row

TABLE IV: Ablation study on energy terms. Each column shows per-joint errors in mm in root-relative frames. The last column shows the MPJPE in global frame.

	$E_{J,l}$	$E_{J,r}$	$E_T$	$E_{3D}$	$E_{tp}$	rknee	lknee	rankl	lankl	rsho	lsho	relb	lelb	rwri	lwri	head	neck	hip	mean	global
1	✓	✓				100	79	128	145	66	54	85	66	147	165	72	38	134	98	1050
2	✓	✓	✓			98	76	<b>124</b>	143	68	55	84	66	141	157	73	37	121	95	277
3	✓	✓	✓			117	76	126	138	64	50	81	<b>64</b>	123	137	68	34	116	91	248
4	✓	✓	✓	✓		104	72	125	<b>135</b>	<b>62</b>	<b>49</b>	83	67	122	134	<b>65</b>	<b>29</b>	107	<b>88</b>	250
5	✓	✓	✓	✓	✓	<b>89</b>	84	132	146	68	57	<b>74</b>	65	143	158	80	38	122	97	240
6	✓	✓	✓	✓	✓	104	<b>71</b>	126	136	63	50	83	67	<b>118</b>	<b>130</b>	66	30	<b>106</b>	<b>88</b>	<b>194</b>

TABLE V: Comparison of 3D MPJPE in root-relative coordinate frame for automatic and manual 2D labels. The last column is the MPJPE in global coordinate frame (in mm).

2D labels	rknee	lknee	rankl	lankl	rsho	lsho	relb	lelb	rwri	lwri	head	neck	hip	mean	global
[1]+[31]	<b>104</b>	89	143	156	81	61	89	75	143	134	112	45	<b>59</b>	99	224
GT 2D joints	<b>104</b>	<b>71</b>	<b>126</b>	<b>136</b>	<b>62</b>	<b>50</b>	<b>83</b>	<b>67</b>	<b>118</b>	<b>130</b>	<b>66</b>	<b>29</b>	106	<b>88</b>	<b>194</b>

2 and row 3, row 4 and row 5 show the errors without and with the temporal prior, respectively. By adding this term we obtain similar values for the root-relative errors, and achieve lower global error. Another advantage of the temporal prior is that it makes the model robust to 2D labeling noise for the occluded joints. Fig. 5 illustrates examples with severe occlusions. It is likely that the 2D body joint labels are inconsistent across the frames under severe occlusion, and that affects the estimated pose. As seen in rows 4 and 5 of Fig. 5, the resulting 3D poses are smoother when this temporal prior term is included. However, if the weight of temporal prior is set too high, the transition of poses between the frames can become too restricted.

4) *Global shape consistency*: Fig. 4 compares the results from SMPLify [14] with those from our method which enforces consistency of the global shape parameters across multiple frames. As expected using the global shape consistency constraint produces more consistent shapes across the frames, while the resulting 3D models from SMPLify, which only uses per-frame information, looks inconsistent. Moreover, the models are too skinny especially when the desired pose is far from the template pose.

### C. Effect of Noisy 2D labels

The size of manually labeled datasets are always limited by the time consuming and costly annotation process. This typically restricts the number of frames that can be annotated to only a few thousand. Luckily, recent state-of-the-art methods have made great progress in 2D tasks such as segmentation, 2D joint detection and tracking. In this section, we use the segmentations from a pre-trained Mask R-CNN [31] and 2D joint detections from OpenPose [1] in place of manual labels on our evaluation dataset. Since there is only one subject wearing the mocap suit, the tracking ID is trivially obtained. Table V shows the 3D MPJPE errors for using the 2D estimates and using ground truth labels. Although the accuracy decreases, the method still achieves greater accuracy than ad-hoc versions of monocular methods shown in Table III. This illustrates that our proposed approach is robust to noisy 2D labels and could potentially be scaled for larger datasets using state-of-the-art multi-object tracking algorithms along with segmentation and 2D joint detection networks in urban scenes.



Fig. 4: Results from monocular SMPLify and from our method.

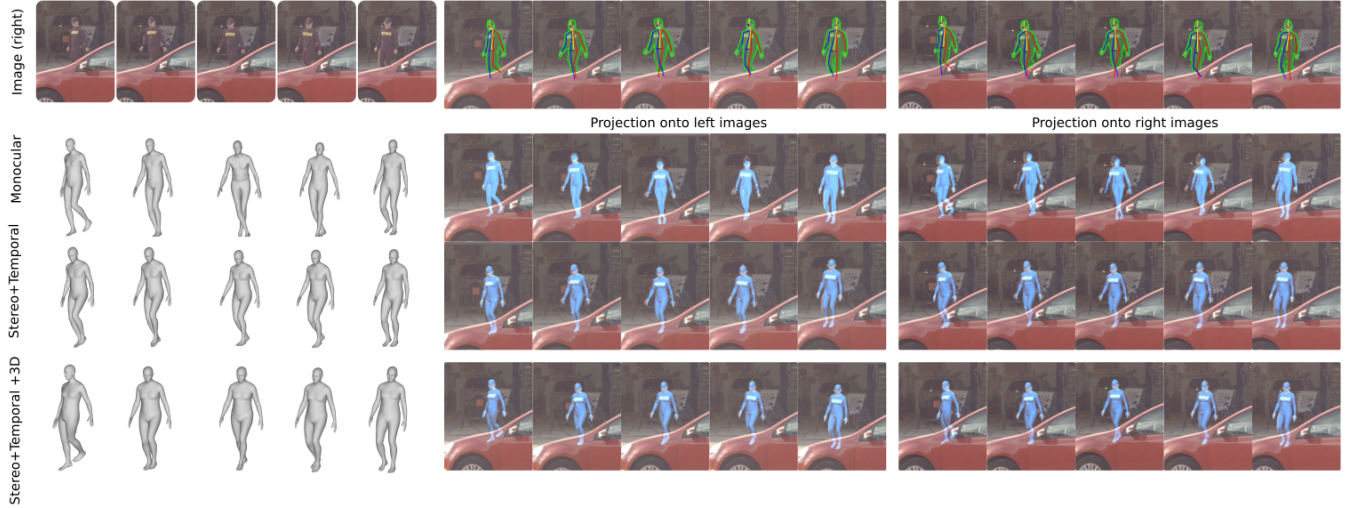


Fig. 5: Results from using different cost terms. The manual annotation for both stereo images is shown and the resulting models in the temporal walking sequence with and without the proposed novel terms. Note the inclusion of the additional cost terms helps to disambiguate complicated and occluded poses where the original fitting approach struggles to deal with the depth ambiguity particularly near the lower limbs that are heavily occluded in this example.



Fig. 6: Representative samples from our dataset illustrates the utility of the dataset. The 3D models from our automatic labeling method are rendered onto the image to show the accuracy of the labels in our dataset even for challenging conditions.

## D. Qualitative Results

Fig. 6 shows some representative examples from *PedX* that illustrate the uniqueness and variety of our dataset. Our dataset covers various actions and poses that are frequently encountered at intersections. Examples include walking, jogging, waving, using a phone, cycling, carrying objects, and talking. Occlusions are another challenge in estimating 3D pose. Our dataset contains many pedestrian instances with severe occlusions by surrounding objects or by other pedestrians. In addition, most frames of the dataset contain more than one pedestrian at a time. Our dataset also contains different weather conditions and rare occurrences such as people in wheelchairs or pedestrians jaywalking.

## VI. CONCLUSION

In this paper, we present a novel large scale multimodal dataset of pedestrians at complex urban intersections with a rich set of 2D/3D annotations. The *PedX* dataset provides a platform for understanding pedestrian behaviors at intersections with real life challenges such as large occlusions, pedestrians walking in group and carrying possessions. This dataset can be used to solve 3D human pose estimation, pedestrian detection and tracking in-the-wild, and to further expand to new problems. Of particular interest is the negotiation of priority at unsignaled intersections and expanding our prior work in forecasting pedestrian trajectories [32] to include body pose as a prior in the motion model.

## REFERENCES

- [1] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *CVPR*, 2017.
- [2] D. Mehta, S. Sridhar, O. Sotnychenko, H. Rhodin, M. Shafiei, H.-P. Seidel, W. Xu, D. Casas, and C. Theobalt, "Vnect: Real-time 3d human pose estimation with a single rgb camera," *arXiv preprint arXiv:1705.01583*, 2017.
- [3] X. Zhou, Q. Huang, X. Sun, X. Xue, and Y. Wei, "Towards 3d human pose estimation in the wild: a weakly-supervised approach," in *IEEE International Conference on Computer Vision*, 2017.
- [4] X. Zhou, M. Zhu, S. Leonardos, K. G. Derpanis, and K. Daniilidis, "Sparseness meets deepness: 3d human pose estimation from monocular video," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4966–4975.
- [5] C.-H. Chen and D. Ramanan, "3d human pose estimation= 2d pose estimation+ matching," *arXiv preprint arXiv:1612.06524*, 2016.
- [6] J. Martinez, R. Hossain, J. Romero, and J. J. Little, "A simple yet effective baseline for 3d human pose estimation," in *ICCV*, 2017.
- [7] L. Sigal, A. O. Balan, and M. J. Black, "Human3.6m: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion," *International journal of computer vision*, vol. 87, no. 1, pp. 4–27, 2010.
- [8] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu, "Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 7, pp. 1325–1339, jul 2014.
- [9] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele, "Deepcrut: A deeper, stronger, and faster multi-person pose estimation model," in *European Conference on Computer Vision (ECCV)*, May 2016.
- [10] A.-I. Popa, M. Zanfir, and C. Sminchisescu, "Deep multitask architecture for integrated 2d and 3d human sensing," *arXiv preprint arXiv:1701.08985*, 2017.
- [11] D. Mehta, H. Rhodin, D. Casas, O. Sotnychenko, W. Xu, and C. Theobalt, "Monocular 3d human pose estimation in the wild using improved cnn supervision."
- [12] D. Tome, C. Russell, and L. Agapito, "Lifting from the deep: Convolutional 3d pose estimation from a single image," *CVPR 2017 Proceedings*, pp. 2500–2509, 2017.
- [13] V. Ramakrishna, T. Kanade, and Y. Sheikh, "Reconstructing 3d human pose from 2d image landmarks," in *European Conference on Computer Vision*. Springer, 2012, pp. 573–586.
- [14] F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black, "Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image," in *Computer Vision – ECCV 2016*, ser. Lecture Notes in Computer Science. Springer International Publishing, Oct. 2016.
- [15] I. Akhter and M. J. Black, "Pose-conditioned joint angle limits for 3d human pose reconstruction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1446–1455.
- [16] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, "SMPL: A skinned multi-person linear model," *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, vol. 34, no. 6, pp. 248:1–248:16, Oct. 2015.
- [17] C. Lassner, J. Romero, M. Kiefel, F. Bogo, M. J. Black, and P. V. Gehler, "Unite the people: Closing the loop between 3d and 2d human representations," *arXiv preprint arXiv:1701.02468*, 2017.
- [18] A. Kanazawa, M. J. Black, D. W. Jacobs, and J. Malik, "End-to-end recovery of human shape and pose," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [19] I. K. Riza Alp Guler, Natalia Neverova, "Densepose: Dense human pose estimation in the wild," *arXiv*, 2018.
- [20] K. Rematas, I. Kemelmacher-Shlizerman, B. Curless, and S. Seitz, "Soccer on your tabletop," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4738–4747.
- [21] F. Ofii, R. Chaudhry, G. Kurillo, R. Vidal, and R. Bajcsy, "Berkeley mhad: A comprehensive multimodal human action database," in *Applications of Computer Vision (WACV), 2013 IEEE Workshop on*. IEEE, 2013, pp. 53–60.
- [22] "CMU Graphics Lab Motion Capture Database," <http://mocap.cs.cmu.edu/>.
- [23] M. F. Ghezghieh, R. Kasturi, and S. Sarkar, "Learning camera viewpoint using cnn to improve 3d body pose estimation," in *3D Vision (3DV), 2016 Fourth International Conference on*. IEEE, 2016, pp. 685–693.
- [24] W. Chen, H. Wang, Y. Li, H. Su, Z. Wang, C. Tu, D. Lischinski, D. Cohen-Or, and B. Chen, "Synthesizing training images for boosting human 3d pose estimation," in *3D Vision (3DV), 2016 Fourth International Conference on*. IEEE, 2016, pp. 479–488.
- [25] G. Varol, J. Romero, X. Martin, N. Mahmood, M. J. Black, I. Laptev, and C. Schmid, "Learning from synthetic humans," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2017)*, 2017.
- [26] G. Rogez and C. Schmid, "Mocap-guided data augmentation for 3d pose estimation in the wild," in *Advances in Neural Information Processing Systems*, 2016, pp. 3108–3116.
- [27] D. Mehta, H. Rhodin, D. Casas, P. Fua, O. Sotnychenko, W. Xu, and C. Theobalt, "Monocular 3d human pose estimation in the wild using improved cnn supervision," in *3D Vision (3DV), 2017 Fifth International Conference on*. IEEE, 2017. [Online]. Available: <http://gvv.mpi-inf.mpg.de/3dhp-dataset>
- [28] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis, "Harvesting multiple views for marker-less 3d human pose annotations," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, 2017, pp. 1253–1262.
- [29] Y. Huang, F. Bogo, C. Lassner, A. Kanazawa, P. V. Gehler, I. Akhter, and M. J. Black, "Towards accurate markerless human shape and pose estimation over time," in *International Conference on 3D Vision (3DV)*, 2017.
- [30] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, et al., "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends® in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [31] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," *arXiv preprint arXiv:1703.06870*, 2017.
- [32] H. O. Jacobs, O. K. Hughes, M. Johnson-Roberson, and R. Vasudevan, "Real-time certified probabilistic pedestrian forecasting," *IEEE Robotics and Automation Letters*, vol. 2, no. 4, pp. 2064–2071, 2017.