

Fast On-Board 3D Torso Pose Recovery and Forecasting*

Abhijat Biswas¹, Henny Admoni, and Aaron Steinfeld

Abstract—Anticipatory human intent modeling is important for robots operating alongside humans in dynamic or crowded environments. Humans often telegraph intent through posture cues, such as torso or head cues. In this paper, we describe a computationally lightweight approach to human torso pose recovery and forecasting with a view towards limited sensing for easy on-board deployment. Our end-to-end system combines RGB images and point cloud information to recover 3D human pose, bridging the gap between learning-based 2D pose estimation methods and the 3D nature of the environment that robots and autonomous vehicles must reason about, with minimum overhead. In addition to pose recovery, we use a simple filter-and-fit method to forecast torso pose. We focus on rapidly generating short horizon forecasts, which is the most relevant scenario for autonomous agents that iteratively alternate between data gathering and planning steps in highly dynamic environments. While datasets suited to benchmarking multi-person 3D pose prediction in real-world scenarios are scarce, we describe an easily replicable evaluation method for benchmarking in a near real-world setting. We then assess the pose estimation performance using this evaluation procedure. Lastly, we evaluate the forecasting performance quantitatively on the Human3.6M motion capture dataset. Our simple 3D pose recovery method adds minimum overhead to 2D pose estimators, with comparable performance to 3D pose estimation baselines from a computer vision alternative. Furthermore, our uncomplicated forecasting algorithm outperforms complicated recurrent neural network methods while also being faster on the torso pose forecasting task.

I. INTRODUCTION

Perceiving and anticipating human motion is increasingly important and relevant as mobile autonomous systems are steadily deployed in highly dynamic and cluttered environments with imperfect information about their surroundings. In real-world settings, a crucial aspect of human-robot interaction (HRI) is real-time anticipatory modeling of human motion. Fluid tasks such as collaborative assembly, handovers, and navigating through moving crowds require timely prediction of probable future human motion.

Consider the case where a mobile, convention center robot meets visitors who have requested assistance. First, it must rendezvous with the human. A strong cue that a particular human is ready for interaction is when they turn to face the oncoming robot. Second, the robot must navigate past other humans without crossing their path in a rude manner [1], [2]. Finally, the robot needs to orient itself properly as

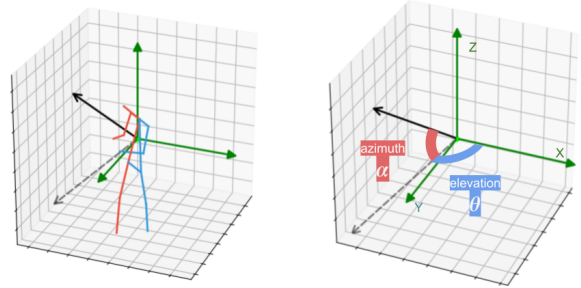


Fig. 1: Torso pose from 3D joints. (Left) Normal to the torso plane is shown in solid black and its projection to the horizontal plane in dashed gray. Torso center is plotted in fluorescent green. (Right) The plane azimuth (α) and plane elevation (θ) are shown in blue and red, respectively.

it approaches the person [3]. Timely perception of human torso pose is important for all of these steps.

More generally, to be accepted by society, mobile robots deployed in public settings need to behave in expected and predictable ways. To meet this goal, robots need to reason not only about individual humans in various trajectories, but about social groups and personal spaces for which, again, body orientation is an important feature [4].

In support of these, and similar interactions, we present a new human torso pose estimation and anticipation model. We focus specifically on the case of mobile robots with limited computational and sensing resources, operating in highly dynamic environments. The typical sensing-perceiving-acting loop in such scenarios involves alternating between data gathering and action or motion re-planning steps in rapid iterations. We show that a simple filter and polynomial fit model outperforms deep neural networks for short-to-medium horizon (under 1 s) predictions, which is the most important case for mobile robots expected to rapidly gather data and re-plan. We also show this method to be much faster, allowing it to be deployed for low-cost on mobile systems since it does not require significant and expensive computation.

Both torso pose recovery and forecasting are challenging problems, so prior approaches have involved computationally expensive solutions. As an illustration, consider that one of the preeminent 2D articulated full-body human pose detectors [5] can perform at around 18Hz using 2x Nvidia 1080 Ti GPUs [6]. Additionally, real world human perception requires the knowledge of 3D pose rather than 2D pose. We attempt to efficiently bridge this 2D to 3D gap with a focus on limited-compute, real-time operation, as well

*This work was funded by a grant from the National Science Foundation (IIS-1734361) and the Mobility21 National University Transportation Center, which is sponsored by the US Department of Transportation.

¹All three authors are with the Robotics Institute, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA-15217, USA. Corresponding author: abhi.jat@cmu.edu

as demonstrate suitability for pose forecasting. Full body articulated human pose forecasting is also challenging due to the associated high dimensional, non-linear dynamics and inherent stochasticity of human motion.

To make the problem more tractable, researchers have approached the forecasting problem by restricting the scope to a particular part of the body relevant to the task, thereby reducing the dimensionality of the problem space. For example, some predictively model human reaching motions for a shared workspace assembly task [7], while others predict future hand locations in egocentric video to allow anticipatory motion planning and assistance [8]. We draw inspiration from this strategy and restrict the problem to modeling the spatio-temporal behaviour of the human torso. Specifically, we aim to detect and forecast the human torso plane position and orientation, the latter being an important cue correlated with motion intent and social engagement [4].

Our algorithm uses multi-modal visual input data, namely RGB with scene depth data, to estimate and forecast a 3D torso plane. This in contrast to most previous body pose forecasting work (e.g., [9]–[12]) that either use 2D or 3D articulated pose, often with initial joint configurations obtained directly from a motion capture system. Such multi-modal sensing not only helps overcome depth ambiguity [13], but also allows us to use monocular 2D body pose estimators (which are more accurate than monocular 3D pose estimators) and project these estimations to 3D easily using an RGB image in conjunction with a registered point cloud. All of our algorithmic design choices are made to prioritize fast running times on generic, portable hardware, such as barebones PCs or embedded systems.

Additionally, we describe a useful evaluation procedure for single-view 3D pose estimation in crowded scenes which can be used by the community for benchmarking. It is difficult to obtain ground truth pose estimates from single-view sensing in real-world scenarios due to occlusions and prior work has tended to use marker-driven motion capture data for these purposes, which inherently contains only clutter and occlusion free scenarios which are artificial. We work around this for evaluation purposes by simulating single viewpoint visual sensing in cluttered scenes using the publicly available Panoptic Studio dataset [14].

Contributions: In this paper, we describe a computationally light-weight end-to-end 3D torso pose estimation and forecasting system combining both depth and color visual data.

Further, we show that a simple filtering and polynomial fitting algorithm outperforms more complicated recurrent neural network based pose forecasting approaches and is $45\times$ faster, trading off speed and accuracy for pose granularity. We evaluate the pose forecasting system quantitatively on the Human 3.6M (H3.6M) dataset [15]. We show superior performance for short-to-medium term forecasts and competitive results for longer term forecasts, especially for predictable activities such as Walking motions in H3.6M.

II. RELATED WORK

Human modeling for robotics has taken various forms including estimating [5], [16]–[18] and forecasting [9], [11], [12] human pose from visual data, modeling human motion trajectories individually [19], [20] and in groups [21], as well as predicting human intent [22].

To this end, previous works have utilized the intrinsic kinematics of the human anatomy [23], eye gaze [24], [25], semantic information of the scene [26], and spatio-temporal structure of the task space [22]. These methods have used graphical models such as Markov Decision Processes or Conditional Random Fields to encode constraints and spatio-temporal relationships. None of these works combine a mobile robot’s viewpoint with realtime forecasting.

A. Pedestrian detection and tracking

Among on-board perception methods, most approaches combine 2D and 3D scene information to build a 3D, human-aware map of the autonomous agent. Some systems fuse handcrafted features from LIDAR and Histogram-of-Gradient features from vision to detect pedestrians [27]. Others align point clusters from two LIDARs and pedestrian detection bounding boxes from three RGB cameras to obtain 3D pedestrian estimates [28]. Such methods prioritize active depth sensors. Our method is agnostic to the source of the depth data and can be used with any source that provides spatial correspondences between an RGB image and points in space, including passive sensors. In our qualitative evaluation we use a stereo camera, the Stereolabs ZED.

B. 3D Human pose estimation

Markerless human pose recovery from visual data is a challenging but useful capability that has recently seen tremendous success in the computer vision community. The focus has mostly been on joint keypoint localization using a single RGB camera in pixel space (2D pose estimation) [16] of a single individual [18], [29], [30] and, more recently, multiple individuals [5], [17], [31]. The most successful models have employed graphical or neural network models trained on large datasets.

Unlike 2D pose, large-scale data is difficult to annotate for 3D pose (predictions in metric space) without dense instrumentation of the environment [14] or of the humans [15], both leading to artificial restrictions on the humans. Despite this, significant interest in 3D pose estimation exists, owing to its numerous applications. Several end-to-end models have been trained on this task that regress the individual skeletal keypoints [32]–[37]. More complex methods may predict 2D and 3D methods jointly. Even though these algorithms have the advantage of being able to work with inexpensive RGB cameras, these are monocular 3D pose estimators and, as such, suffer from depth and scale ambiguities that allow multiple plausible 3D pose hypotheses given a 2D pose estimate [13]. Moreover, the real-time methods among these, such as V-Nect [33], [38], can only produce single-person pose estimates and require tracked bounding boxes for each person, making it unsuitable for use in dynamic crowds.

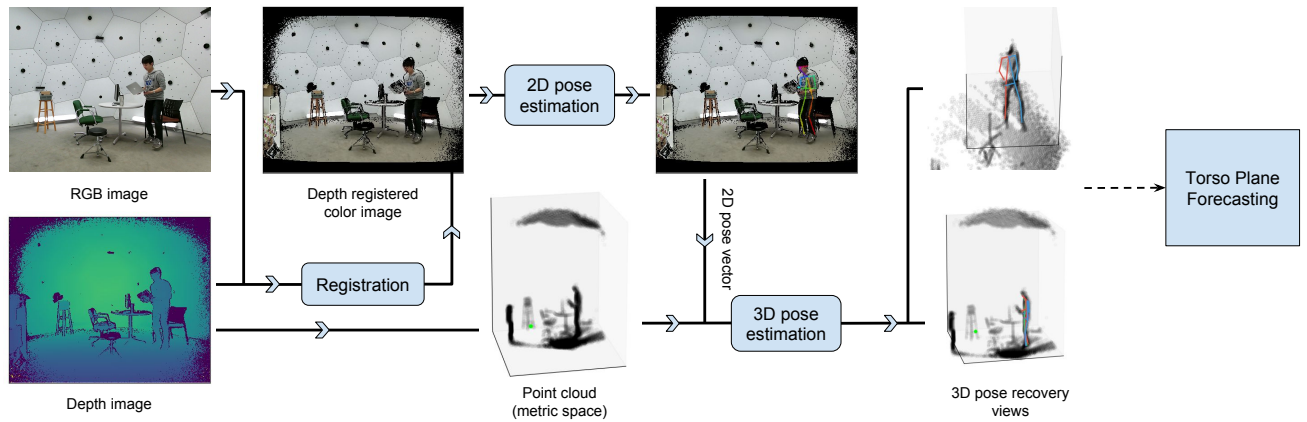


Fig. 2: 3D torso pose estimation algorithm overview. For details see Section III-A.

In this paper, we compare against a strong baseline based on the better-performing real-time method [38]. Note that these methods are not trained to account for global translation and rotation. Our baseline provides global rectification to their method in order to overcome this.

Closest to our work is [39], where the authors use a Kinect V2 sensor and a voxel-based neural network to provide 3D poses in metric space. This method achieves impressive results but is not real time and places an explicit requirement on a single type of depth sensor, which we do not.

C. Human pose forecasting

Much work in human pose forecasting predicts articulated human poses without considering the acquisition of the pose skeletons themselves. For example, [9], [11], [12] all model and forecast human motion using recurrent neural networks (RNNs). However, these methods do not model the global position of their subjects, instead focusing on generating a continuation of observed human motion in a coordinate frame attached to the body. For forecasting, these works are the closest to ours and we compare against the best performing method in this paper [12].

The graphics community also uses deep recurrent neural networks, primarily for character motion synthesis conditioned on human user input. For example, these methods have been used to animate game characters [40].

While each of these methods generate realistic human motions, they fail to match ground truth human poses and suffer from discontinuity artifacts between the ground truth instances and first predicted instance.

Finally, the autonomous and assisted driving community has also investigated the use of pedestrian pose-based features for intent prediction with encouraging results. These works restrict intent to higher-level classes such as “cross/no-cross” [41] or “start/stop/cross/bend” [42] for curbside pedestrians. These methods also use 2D pose instead of 3D which removes the ability of these systems to reason about absolute pose, which is adequate for their application, but we wish to investigate the use of finer-grained pose information.

None of the aforementioned works study real-time human

pose forecasting with the sensing and computation restrictions of a typical mobile system, as we do in this paper.

III. APPROACH

Keeping with example of the perception system of a mobile robot that interacts with humans, such a system would require both real-time performance and an output signal that allows human attention/intent prediction.

For both these reasons, we choose to use human torso pose as the perception output. Acquiring accurate, articulated full-body 3D pose in real-time is challenging given the constraints of on-board sensing, which is prone to occlusions because only a single view-point is available. Hence, we restrict ourselves to 3D torso pose, comprising the global Cartesian coordinates of the torso center-of-mass and the torso plane angles. Torso pose also evolves less rapidly than head pose [43], hence mitigating information loss at lower temporal sampling rates, which in turn allows lower hardware design costs.

This also allows us to incorporate a smooth temporal constraint in our model of human pose, which is a non-trivial consideration since previous learning-based methods such as [9], [11], [12], suffer from discontinuity artifacts at the beginning of the forecast, as shown in [12], which are inconsistent with human anatomical limits. For example, see the relatively large error at the start of each error graph in Fig. 3, corresponding to the HMP method from [12].

Our end-to-end system comprises a torso pose estimation module followed by a forecasting module. The algorithm requires registered RGB and depth inputs with proper calibration. In its most basic form, it is agnostic to the data source. For instance, in the evaluation in Table I we used a Kinect v2 (active sensing) as the input source, while our end-to-end qualitative system demonstration used a ZED camera (stereo). This allows flexibility for system designers to trade-off the requirements for their particular scenario. For example, higher fidelity 3D maps can be obtained with a 3D LIDAR at the cost of higher power consumption and overall expense.

A. Pose Estimation

We use an off-the-shelf 2D human pose detection system [5] in conjunction with registered depth information in a two-step process. The input to the 2D pose detector is an RGB image, which is used to obtain joint locations for humans in the scene. Once 2D joint locations are known, they are projected onto a registered point cloud obtained by triangulation in a separate step, giving us 3D joint locations.

We parameterize torso pose by the position (x, y, z of torso center) and orientation (plane azimuth: α and elevation: θ) of an estimated torso plane. Given a pose skeleton, the torso plane is defined as the plane that minimizes sum of squared distances from each of the 3D torso joint locations. At a given pose skeleton this plane is given by:

$$\mathbf{n}^*, \mathbf{c}^* = \underset{\mathbf{n}, \mathbf{c}}{\operatorname{argmin}} \sum_{i=1}^{|\tau|} |\mathbf{n} \cdot \mathbf{x}_i + \mathbf{c}| \quad (1)$$

where $\mathbf{n} \cdot \mathbf{x}_i + \mathbf{c} = 0$ defines the torso plane (\mathbf{n} is the plane normal and \mathbf{c} is a constant, both in \mathcal{R}^3) and $\mathbf{x}_i \in \tau \subset \mathcal{R}^3$ is the set of all torso joint locations in 3D space.

Hence, the torso center is:

$$\mathbf{C}_{torso} = \frac{\sum_{i=1}^{|\tau|} \mathbf{x}_i}{|\tau|} \quad (2)$$

For pose recovery, we present our method alongside a baseline from a state-of-the-art monocular 3D pose estimate. For each of these methods, once τ is constructed, the plane is calculated using Equation 1, giving the plane azimuth (α) and elevation (θ) directly.

$$\alpha = \arctan \frac{\mathbf{n}_y}{\mathbf{n}_x} ; \theta = \arccos \frac{\mathbf{n}_z}{\|\mathbf{n}\|_2} \quad (3)$$

For our method, we compose the set τ comprising solely the torso points available from the 2D pose detector. For annotated ground-truth skeletons from the datasets used in our evaluation, τ contains the shoulder joints, two hip joints, the mid-spine, and the tip of the tailbone. In this formulation, the registered point-cloud is constructed at every time-step and the points \mathbf{x}_i corresponding to the detected joints in the RGB image are picked from the corresponding point cloud. Hence, \mathbf{x}_i s are in metric 3D space. For each \mathbf{x}_i , temporal consistency is ensured by discarding values that deviate over 10 cm between consecutive time-steps, lending some robustness to temporary occlusion. The discarded values are replaced by the corresponding \mathbf{x}_i from the previous time-step. See Fig. 2 for an overview of this method.

As a competitive baseline, we use the method of [38]. This is a 4-layer shallow neural network trained on the Human 3.6M dataset [15] to “lift” a given 2D pose detection into 3D space. Being a 4-layer neural network with relatively low dimensional inputs and outputs, this only takes about 5 ms per forward pass, making it suitable for real-time deployment, as opposed to other more computationally intensive 3D human pose predictors discussed in Section II-B. Comparing against this method is also fair since it also tries to bridge the gap

between 2D and 3D pose estimation rather than attempting monolithic 3D pose estimation.

In the method from [12], the 3D poses are not guaranteed to be recovered in global-scale. To transform their output pose into the global coordinate frame, we find and apply the best least-squares rigid transform between the 3D predicted pose and 3D ground truth pose. Note that this represents an unattainable gold-standard performance for this method, since ground-truth pose is never available in a real-world setting.

B. Pose Forecasting

For forecasting, we predict elevation, azimuth, and absolute position of the torso plane for a variable time lookahead, from a 2 s history. Average error for multiple forecasting windows (400 ms, 1 s, 2 s) is in Table II. Subsequently, we refer to the elevation and azimuth together as the plane orientation. Together, the three components describe the plane uniquely (see Fig. 1). Once the torso plane is acquired from the pose estimation module, we apply a low pass filter to the two orientation components. A low pass filter was chosen so that we model only the macro-level orientation of a human subject, which is the most relevant signal for many activities. This is followed by fitting an N th order polynomial, which is used to extrapolate a forecast for each individual component. Error analyses of various orders of fitting (N) are presented in Table II.

For the low pass filter, we used a second-order Butterworth filter, with the cutoff frequency empirically set to 5Hz.

IV. EXPERIMENTS

A. Datasets

We evaluated the two modules, pose recovery and forecasting, on two datasets respectively:

- **Panoptic Studio** [14]: We used this for quantitative evaluation of the pose recovery system. It contains RGB-D inputs and multi-person scenarios, representing the closest available data to our target application.
- **Human 3.6M** [15]: We used this for quantitative evaluation of the pose forecasting system. 3D pose in world-coordinates can directly be obtained from their marker-based MoCap system. The lack of RGB-D views prevents us from testing the system end-to-end.

While we would ideally evaluate our work end-to-end, most datasets with grounding for 3D pose are marker-based and do not have associated RGBD data. The Panoptic Studio does have marker-less grounding for 3D pose but the contained activities (e.g. Office:sitting at a desk or Range of Motion: arm movements that keep the torso mostly stationary) are unsuited to evaluation of a torso pose forecasting method, hence we perform modular evaluation.

1) *Pose estimation*: The motivation of this work is to provide a fast method to recover and forecast multi-person 3D pose in the real world, with a focus towards social navigation. Hence, ideal evaluation would be data-driven and with said data collected in the wild. However, instrumenting to recover accurate 3D pose in such scenarios is difficult due

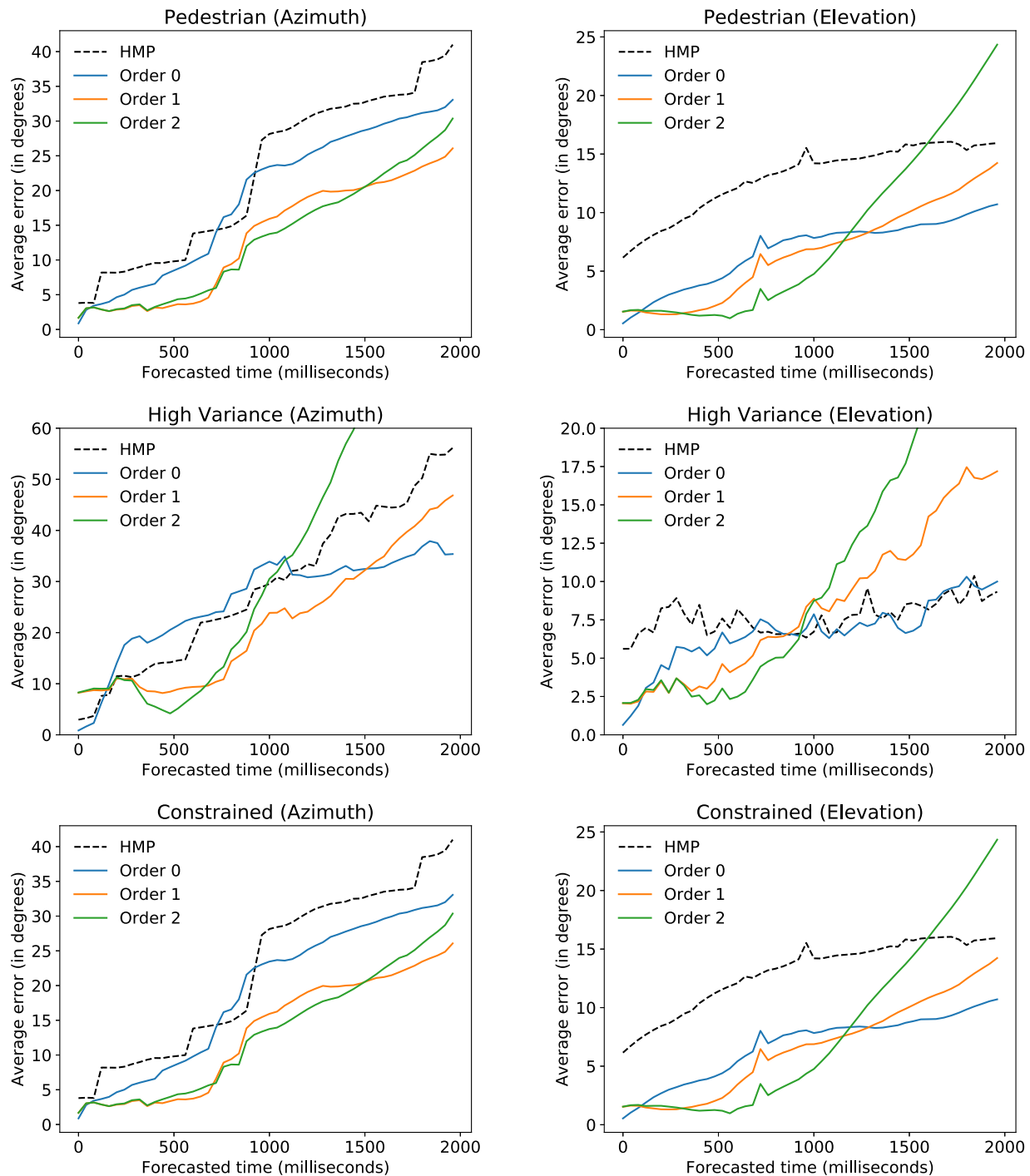


Fig. 3: Average forecasting error (across all test sequences) vs. forecasting time extent for various categories of Human3.6M data (lower is better). The plots show our recommended method (Order 1), two baselines (Order 0 and 2), and the RNN-based method from HMP [12].

to financial, computational, and privacy concerns. Unfortunately, 3D pose datasets with multiple, simultaneous humans and RGBD inputs are rare.

The best effort in this domain is the Panoptic Studio [14], which uses advances in 2D pose recognition with 500 RGB cameras and 10 Kinects to recover the 3D pose of observed humans accurately. While still being an artificial environment, housed in a geodesic sphere of diameter 5.49 m , it

solves the occluded pose recovery problem by dense instrumentation of the environment rather than equipping humans, allowing for more naturalistic movement. This dataset has the added benefit of multi-person capture sequences that present several types of occlusion challenges likely to also be found in dynamic crowds. We use relevant sequences with Kinect inputs and ground-truth 3D human pose present. This amounts to about 100 minutes of data at 30 Hz.

2) *Pose forecasting*: To evaluate the pose forecasting module we chose the Human 3.6M [15] dataset. This is currently the largest publicly available dataset of motion capture data, containing 7 actors performing 15 varied activities such as walking, taking photos or giving directions, with only a single person per task. We group the tasks into three categories: Pedestrian, Constrained, and High Variance. The Pedestrian group contains Walking, Walking Dog, Walking Together activities. The Constrained group comprises of Eating, Smoking, Phoning, Sitting, Sitting Down which all involve the person’s torso being constrained in some fashion (e.g. by being placed in a rotating chair). The High Variance group comprises of all other activities such as Taking Photo, Posing, etc which have mostly stochastic motion where very little intent is telegraphed. In our opinion this is not really relevant to evaluate forecasting models since the premise of motion history based forecasting is that consecutive motions are correlated and motion intent is telegraphed. Nevertheless we perform evaluation for comparative purposes.

Once 3D torso pose is acquired, our analysis is local and does not consider inter-person effects, meaning single-person sequences are equivalent to multi-person scenarios for evaluation purposes. Prior pose forecasting work [11], [12] has also evaluated on this data. Consequently, an evaluation procedure exists for articulated pose forecasting which we adapted to the torso pose scenario.

While this data seems like an ideal candidate for end-to-end system evaluation, lack of registered RGB and depth views render it unusable for that purpose.

B. Evaluation procedures

1) *Pose recovery*: To evaluate this component of our algorithm, we wanted to simulate single view-point visual sensing (e.g. a mobile robot with on-board sensing) by using inputs from a single Kinect v2 RGB-D sensor in the Panoptic Studio [14]. This allows us to test pose recovery in the presence of occlusions, which is important for applications like dynamic pedestrian tracking in busy environments.

For each frame during a sequence, the pose recovery component of the algorithm in Fig. 2 is used. Ground truth articulated pose (which is reconstructed with a combination of over 500 camera views and a 2D pose estimation method [5], [14]) is used to compute the ground truth torso plane, as in Equation 1. The body center-of-mass and plane angle errors are shown in Table I.

2) *Pose forecasting*: For quantitative evaluation on Human 3.6M, we used the same train-test split as [9], [11], [12] and compared against [12] since it is the quantitatively best performing model of the three. In [12], the MoCap data was down-sampled to 25 Hz. During testing, skeletal poses over a 2 s sample (50 frames) were fed to a recurrent neural network (single-layer), which then generated samples over a forecast window of 400 ms (10 frames) sample. The initial 50 frames are referred to as the conditioning ground truth. Their method also has the advantage over previous work [9], [11] in that it trains one-model across all actions in the dataset. We retain this advantage by using the same set of

filter parameters for the entire dataset, eliminating the need to tune for every individual action.

The choice of the 400ms forecasting method follows from previous work [9], [11], [12]. Further, to properly characterize the properties of our method and HMP [12] as well as to enable comparison of the two, we present the analysis for multiple forecast windows.

The aforementioned methods do not estimate the 3D pose of a human from visual data. Rather, they acquire the ground truth 3D poses directly obtained from the MoCap data accompanying Human3.6M. For evaluating our pose forecasting method in this experiment, our pose estimation module was bypassed to keep the quantitative comparison of our forecasting system with [12] fair.

Since our method focuses on torso planes rather than full body articulated pose, we must obtain ground truth planes from the MoCap data. This was done by fitting a least squares plane to hip, shoulder, and neck joints of an articulated pose obtained from the MoCap data, as described in Equation 1.

Additionally, instead of using the Euclidean distance in Euler angle space for all body joints (as in previous work), we computed the angle error of the plane orientation forecast. We chose this measure since it is most indicative of the macro-level expression of torso pose. See Table II for average azimuth and elevation angle error for each of the 15 Human3.6M activities as well as within the subcategories described in IV-A.2.

V. RESULTS AND DISCUSSION

Table I shows the results of our pose estimation method and a baseline using a state-of-the-art, learned, 3D pose predictor [38]. These results show that learning for 3D pose estimation may need more improvement before it can be used for accurate, real-time performance suitable for robot deployment. We see comparable performance for our method and that of [38]. Since, the latter has access to information about ground truth rotation and translation (as a rigid transform), which is not available in real world scenarios, we assert that our method enables 3D pose recovery with far less overhead.

In terms of computational performance, our bottleneck lies in the 2D pose estimation step. We are able to achieve a performance of 10 Hz using an Nvidia 1080Ti. However, our method is not tied to a particular type of 2D pose estimator. A faster pose estimator, such as the recent work in [44] (180 Hz on similar GPU with similar accuracy, real-time performance on CPU) can significantly speed up performance without sacrificing accuracy.

Table II shows the results of pose forecasting methods, including polynomials of degree $N = 1$ and 2, state-of-the-art human motion predictor (HMP, the quantitatively best performing method for on Human3.6M) [12], and a constant prediction baseline ($N = 0$) (where the last ground-truth torso plane orientation is predicted for the entire forecast window). The results also show the importance of the filtering step (see last row, where we omit it and directly fit an N th order polynomial to unfiltered data.)

TABLE I: Torso plane estimation errors on Panoptic studio [14] data (centimetres/degrees)

Activity	Torso Center X (cm)		Torso Center Y (cm)		Torso Center Z (cm)		Plane Azimuth (deg)		Plane Elevation (deg)	
	Ours	[38]	Ours	[38]	Ours	[38]	Ours	[38]	Ours	[38]
Range of Motion	7.64	4.46	3.91	12.26	15.59	4.73	30.25	33.98	14.44	8.24
Office	5.83	12.11	2.85	11.41	12.52	14.85	24.95	41.47	10.60	9.70

TABLE II: Torso plane orientation forecasting errors for various forecasting windows on H3.6M [15] data (degrees)

Forecast time→	400 <i>ms</i>				1 <i>s</i>				2 <i>s</i>			
Activity ↓	Const	<i>N</i> = 1	<i>N</i> = 2	HMP [12]	Const	<i>N</i> = 1	<i>N</i> = 2	HMP [12]	Const	<i>N</i> = 1	<i>N</i> = 2	HMP [12]
Plane Azimuth												
Constrained	4.24	2.87	2.92	7.13	9.83	5.49	5.37	11.48	18.98	13.11	13.11	22.32
HV	10.87	9.38	9.1	8.45	19.13	11.01	11.03	15.72	26.32	21.98	38.92	28.98
Pedestrian	4.4	2.17	2.09	5.02	10.46	4.82	3.22	11.45	18.23	11.32	16.69	22.83
All	7.36	5.77	5.64	7.33	14.3	7.93	7.58	13.45	22.26	16.89	25.87	25.53
All (no filter)	7.36	5.4	13.56	7.33	14.3	14.15	50.81	13.45	22.26	27.2	155.87	25.53
Plane Elevation												
Constrained	2.32	1.45	1.54	8.13	4.62	3.24	1.99	10.83	6.77	6.63	8	13.06
HV	3.59	2.79	2.81	7.21	5.3	4.32	3.53	7.09	6.68	8.48	11.65	7.69
Pedestrian	2.81	1.6	1.68	6.53	4.86	2.6	2.34	8.02	6.25	5.61	8.4	8.55
All	3.01	2.11	2.16	7.38	4.98	3.62	2.78	8.52	6.63	7.29	9.79	9.65
All (no filter)	3.01	3.26	16.35	7.38	4.98	8.84	67.58	8.52	6.63	17.33	202.66	9.65

A visual representation of the error as it evolves with the forecasting extent across time is shown. Fig. 3.

A few trends can be seen in Table II and the error graphs in Fig. 3. First, the plane azimuth is harder to predict than the elevation, given the higher error rates across all 15 activity sequences and various methods. However, the best average error for both torso orientation components is under 5 degrees. This is small enough to not cause ambiguity in most real-world activities.

Second, the filtering step is essential. Without it, we see larger errors in the polynomial fitting and the errors tend to explode in the larger forecasting windows (Table II). This suggests that the forecast becomes unstable for higher order approximations due to susceptibility to high frequency components in the pose variation.

Third, the recurrent neural network model from HMP [12] tends to make much larger errors than our simple 1st degree (linear) polynomial fit, especially over the short-to-medium term (i.e., 400 *ms*-1*s*) and over all windows for the Pedestrian group of activities. The HMP errors also show higher variability across tasks than our method.

This suggests that such models are either over-fitting or that the error they are trained to minimize is unsuitable for our task. That is, recurrent neural network based methods try to minimize a quantitative loss without reasoning about the temporal smoothness of human motion. Thus, these methods can suffer from unrealistic discontinuities. This is reinforced by the observation that these errors are larger in the High Variance tasks such as “Taking Photo” (both upright and kneeling poses) or “Directions” (high variance poses).

Our forecasting algorithm is inexpensive to compute while being faster and more accurate (for short horizons) than previous work. The method described in [12] (HMP) takes about 35 *ms* for one forward pass on a dedicated Nvidia

Titan X GPU. This translates to a maximum sampling rate of 28 Hz, assuming desktop-level hardware is available on-board. Note that this is the computational cost of just the HMP forecasting method. This is significant since it may be an additional bottleneck if used in conjunction with another 2D/3D pose recovery method for end-to-end pose recovery and forecasting over our forecasting method. Our method takes approximately 0.715 *ms* on an Intel i7-6700HQ CPU (laptop processor). This makes our method about 45× faster on cheaper and more accessible hardware.

It is important to note that HMP forecasts full body articulated 3D pose using about 34000× as many learned parameters, while we only model the torso plane. This makes the 45× running time speed-up less surprising. However, this difference in pose information demonstrates the power of our forecasting technique. Since HMP models the torso with more granularity and parameters than our method, it should yield much better performance in the medium to long term.

VI. CONCLUSION

We propose a novel end-to-end torso pose estimation and forecasting system which is relevant for rapid perception and re-planning loops of robot decision making in highly dynamic environments, such as the case of social navigation in an autonomous mobile, service robot.

We parameterized torso pose uniquely by the position and orientation of a torso plane (Equation 1). We evaluated the pose estimation quantitatively and compare against a state-of-the-art monocular approach, showing comparable results against a strong baseline. The evaluation was performed in a replicable manner using a publicly available dataset while also simulating the single viewpoint sensing of a mobile robot, thus allowing fair and easy benchmarking in the future.

In addition to torso pose estimation, our approach predictively models absolute torso position. We present a comparative quantitative evaluation and show that our simple filter and fit method outperforms complex recurrent neural network methods for the short-to-medium horizon case while being competitive over the long horizon case. For walking motions, it also accurately predicts the torso facing direction (plane azimuth) which is an important predictive cue of pedestrian trajectory intent. Further, our method is approximately $45\times$ faster on the torso plane forecasting task, implying suitability to navigation in human environments.

In future work, we would like to apply our method to tasks that require multi-person pose perception, like social navigation, to measure the intent prediction capability of torso pose. We imposed several constraints on this work, including a focus solely on the torso plane, which we would like to relax. We also only address the single human case, but this approach could be extended to multiple people in a top-down fashion using a suitable tracking approach.

REFERENCES

- [1] R. Kirby, R. Simmons, and J. Forlizzi, "Companion: A constraint optimizing method for person-acceptable navigation," in *IEEE RO-MAN*, pp. 607–612, September 2009.
- [2] I. Chatterjee and A. Steinfeld, "Performance of a low-cost, human-inspired perception approach for dense moving crowd navigation," in *IEEE RO-MAN*, pp. 578–585, Aug 2016.
- [3] E. Avrunin and R. Simmons, "Socially-appropriate approach paths using human data," in *The 23rd IEEE International Symposium on Robot and Human Interactive Communication*, pp. 1037–1042, Aug 2014.
- [4] M. Vázquez, E. J. Carter, B. McDorman, J. Forlizzi, A. Steinfeld, and S. E. Hudson, "Towards robot autonomy in group conversations: Understanding the effects of body orientation and gaze," in *HRI*, pp. 42–52, ACM, 2017.
- [5] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *CVPR*, 2017.
- [6] Hidalgo, Gines, "Openpose performance benchmark," 2017. <https://goo.gl/8Dk9HS>, Accessed: 05-16-2018.
- [7] J. Mainprice, R. Hayne, and D. Berenson, "Predicting human reaching motion in collaborative tasks using inverse optimal control and iterative re-planning," in *ICRA*, pp. 885–892, IEEE, 2015.
- [8] J. Lee and M. S. Ryoo, "Learning robot activities from first-person human videos using convolutional future regression," *Image*, vol. 500, p. 500, 2017.
- [9] K. Fragkiadaki, S. Levine, P. Felsen, and J. Malik, "Recurrent network models for human dynamics," in *ICCV*, pp. 4346–4354, IEEE, 2015.
- [10] H.-S. Fang, Y. Xu, W. Wang, X. Liu, and S.-C. Zhu, "Learning pose grammar to encode human body configuration for 3d pose estimation," 2018.
- [11] A. Jain, A. R. Zamir, S. Savarese, and A. Saxena, "Structural-rnn: Deep learning on spatio-temporal graphs," in *CVPR*, pp. 5308–5317, 2016.
- [12] J. Martinez, M. J. Black, and J. Romero, "On human motion prediction using recurrent neural networks," in *CVPR*, pp. 4674–4683, IEEE, 2017.
- [13] E. Jahangiri and A. L. Yuille, "Generating multiple diverse hypotheses for human 3d pose consistent with 2d joint detections," *arXiv preprint arXiv:1702.02258*, 2017.
- [14] H. Joo, T. Simon, X. Li, H. Liu, L. Tan, L. Gui, S. Banerjee, T. S. Godisart, B. Nabbe, I. Matthews, et al., "Panoptic studio: A massively multiview system for social interaction capture," *T-PAMI*, 2017.
- [15] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu, "Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments," *T-PAMI*, vol. 36, pp. 1325–1339, jul 2014.
- [16] Z. Liu, J. Zhu, J. Bu, and C. Chen, "A survey of human pose estimation: the body parts parsing based methods," *Journal of Visual Communication and Image Representation*, vol. 32, pp. 10–19, 2015.
- [17] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. V. Gehler, and B. Schiele, "Deepcut: Joint subset partition and labeling for multi person pose estimation," in *CVPR*, pp. 4929–4937, 2016.
- [18] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional pose machines," in *CVPR*, 2016.
- [19] K. Yamaguchi, A. C. Berg, L. E. Ortiz, and T. L. Berg, "Who are you with and where are you going?," in *CVPR*, pp. 1345–1352, IEEE, 2011.
- [20] B. D. Ziebart, N. Ratliff, G. Gallagher, C. Mertz, K. Peterson, J. A. Bagnell, M. Hebert, A. K. Dey, and S. Srinivasa, "Planning-based prediction for pedestrians," in *IROS*, pp. 3931–3936, IEEE, 2009.
- [21] A. Vemula, K. Muelling, and J. Oh, "Social attention: Modeling attention in human crowds," in *ICRA*, pp. 1–7, IEEE, 2018.
- [22] H. S. Koppula and A. Saxena, "Anticipating human activities using object affordances for reactive robotic response," *T-PAMI*, vol. 38, no. 1, pp. 14–29, 2016.
- [23] N. Hu, A. Bestick, G. Englebiene, R. Bajscy, and B. Kröse, "Human intent forecasting using intrinsic kinematic constraints," in *IROS*, pp. 787–793, IEEE, 2016.
- [24] H. Admoni and S. Srinivasa, "Predicting user intent through eye gaze for shared autonomy," in *AAAI Fall Symposium Series: Shared Autonomy in Research and Practice (AAAI Fall Symposium)*, pp. 298–303, 2016.
- [25] C.-M. Huang and B. Mutlu, "Anticipatory robot control for efficient human-robot collaboration," in *HRI*, pp. 83–90, IEEE Press, 2016.
- [26] K. M. Kitani, B. D. Ziebart, J. A. Bagnell, and M. Hebert, "Activity forecasting," in *ECCV*, pp. 201–214, Springer, 2012.
- [27] E. P. Fotiadis, M. Garzón, and A. Barrientos, "Human detection from a mobile robot using fusion of laser and vision information," *Sensors*, vol. 13, no. 9, pp. 11603–11635, 2013.
- [28] J. Miller, A. Hasfura, S.-Y. Liu, and J. P. How, "Dynamic arrival rate estimation for campus mobility on demand network graphs," in *IROS*, pp. 2285–2292, IEEE, 2016.
- [29] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *ECCV*, pp. 483–499, Springer, 2016.
- [30] V. Ramakrishna, D. Munoz, M. Hebert, J. A. Bagnell, and Y. Sheikh, "Pose machines: Articulated pose estimation via inference machines," in *ECCV*, pp. 33–47, Springer, 2014.
- [31] H.-S. Fang, S. Xie, Y.-W. Tai, and C. Lu, "RMPE: Regional multi-person pose estimation," in *ICCV*, 2017.
- [32] 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 *3DV*, 2017.
- [33] 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," *TOG*, vol. 36, no. 4, p. 44, 2017.
- [34] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis, "Coarse-to-fine volumetric prediction for single-image 3d human pose," in *CVPR*, pp. 1263–1272, IEEE, 2017.
- [35] X. Sun, J. Shang, S. Liang, and Y. Wei, "Compositional human pose regression," in *ICCV*, vol. 2, 2017.
- [36] D. Tome, C. Russell, and L. Agapito, "Lifting from the deep: Convolutional 3d pose estimation from a single image," *CVPR*, pp. 2500–2509, 2017.
- [37] X. Zhou, X. Sun, W. Zhang, S. Liang, and Y. Wei, "Deep kinematic pose regression," in *ECCV*, pp. 186–201, Springer, 2016.
- [38] J. Martinez, R. Hossain, J. Romero, and J. J. Little, "A simple yet effective baseline for 3d human pose estimation," in *IEEE ICCV*, vol. 206, p. 3, 2017.
- [39] C. Zimmermann, T. Welschehold, C. Dornhege, W. Burgard, and T. Brox, "3d human pose estimation in rgb-d images for robotic task learning," in *ICRA*, pp. 1986–1992, IEEE, 2018.
- [40] D. Holden, T. Komura, and J. Saito, "Phase-functioned neural networks for character control," *ACM TOG*, vol. 36, no. 4, p. 42, 2017.
- [41] Z. Fang and A. M. López, "Is the pedestrian going to cross? answering by 2d pose estimation," *arXiv preprint arXiv:1807.10580*, 2018.
- [42] Z. Fang, D. Vázquez, and A. M. López, "On-board detection of pedestrian intentions," *Sensors*, vol. 17, no. 10, p. 2193, 2017.
- [43] M. Vázquez, A. Steinfeld, and S. E. Hudson, "Parallel detection of conversational groups of free-standing people and tracking of their lower-body orientation," in *IROS*, pp. 3010–3017, IEEE, 2015.
- [44] T. Sekii, "Pose proposal networks," in *ECCV*, September 2018.