

Deep Learning-Based Human Pose Estimation: A Survey

Ce Zheng*, Wenhan Wu*, Taojiannan Yang, Sijie Zhu, Chen Chen, *Member, IEEE*, Ruixu Liu, Ju Shen, *Senior Member, IEEE*, Nasser Kehtarnavaz *Fellow, IEEE* and Mubarak Shah, *Fellow, IEEE*

Abstract—Human pose estimation aims to locate the human body parts and build human body representation (e.g., body skeleton) from input data such as images and videos. It has drawn increasing attention during the past decade and has been utilized in a wide range of applications including human-computer interaction, motion analysis, augmented reality, and virtual reality. Although the recently developed deep learning-based solutions have achieved high performance in human pose estimation, there still remain challenges due to insufficient training data, depth ambiguities, and occlusion. The goal of this survey paper is to provide a comprehensive review of recent deep learning-based solutions for both 2D and 3D pose estimation via a systematic analysis and comparison of these solutions based on their input data and inference procedures. More than 240 research papers since 2014 are covered in this survey. Furthermore, 2D and 3D human pose estimation datasets and evaluation metrics are included. Quantitative performance comparisons of the reviewed methods on popular datasets are summarized and discussed. Finally, the challenges involved, applications, and future research directions are concluded. We also provide a regularly updated project page: <https://github.com/zczcw/HL-HPE>

Index Terms—Survey of human pose estimation, 2D and 3D pose estimation, deep learning-based pose estimation, pose estimation datasets, pose estimation metrics

1 INTRODUCTION

HUMAN pose estimation (HPE), which has been extensively studied in computer vision literature, involves estimating the configuration of human body parts from input data captured by sensors, in particular images and videos. HPE provides geometric and motion information of the human body which has been applied to a wide range of applications (e.g., human-computer interaction, motion analysis, augmented reality (AR), virtual reality (VR), healthcare, etc.). With the rapid development of deep learning solutions in recent years, such solutions have been shown to outperform classical computer vision methods in various tasks including image classification [1], semantic segmentation [2], and object detection [3]. Significant progress and remarkable performance have already been made by employing deep learning techniques in HPE tasks. However, challenges such as occlusion, insufficient training data, and depth ambiguity still pose difficulties to be overcome. 2D HPE from images and videos with 2D pose annotations is easily achievable and high performance has been reached for the human pose estimation of a single person using deep learning techniques. More recently, attention has been paid

to highly occluded multi-person HPE in complex scenes. In contrast, for 3D HPE, obtaining accurate 3D pose annotations is much more difficult than its 2D counterpart. Motion capture systems can collect 3D pose annotation in controlled lab environments; however, they have limitations for in-the-wild environments. For 3D HPE from monocular RGB images and videos, the main challenge is depth ambiguities. In multi-view settings, viewpoints association is the key issue that needs to be addressed. Some works have utilized sensors such as depth sensor, inertial measurement units (IMUs), and radio frequency devices, but these approaches are usually not cost-effective and require special purpose hardware.

Given the rapid progress in HPE research, this article attempts to track recent advances and summarize their achievements in order to provide a clear picture of current research on deep learning-based 2D and 3D HPE.

1.1 Previous surveys and our contributions

Table 1 lists the related surveys and reviews previously reported on HPE. Among them, [4] [5] [6] [7] focus on the general field of visual-based human motion capture methods and their implementations including pose estimation, tracking, and action recognition. Therefore, pose estimation is only one of the topics covered in these surveys. The research works on 3D human pose estimation before 2012 are reviewed in [8]. The body parts parsing-based methods for single-view and multi-view HPE are reported in [9]. These surveys published during 2001-2015 mainly focus on conventional methods without deep learning. A survey on both traditional and deep learning-based methods related to HPE is presented in [10]. However, only a handful of deep learning-based approaches are included. The survey in [11] covers 3D HPE methods with RGB inputs. The survey in [13] only reviews 2D HPE

- * The first two authors are contributed equally.
- C. Zheng, W. Wu, T. Yang, S. Zhu and C. Chen are with the Department of Electrical and Computer Engineering, University of North Carolina at Charlotte, Charlotte, NC 28223.
E-mail: {czheng6, wwu25, tyang30, szhu3, chen.chen}@uncc.edu
- R. Liu and J. Shen are with the Department of Computer Science, University of Dayton, Dayton, OH 45469.
E-mail: {liur05, jshen1}@udayton.edu
- N. Kehtarnavaz is with the Department of Electrical and Computer Engineering, University of Texas at Dallas, Richardson, TX 75080.
E-mail: kehtar@utdallas.edu
- M. Shah is with the Center for Research in Computer Vision, University of Central Florida, Orlando, FL 32816.
E-mail: shah@crco.ucf.edu

TABLE 1: Previous reviews and surveys on HPE.

Title	Year	Venue	Brief Description
A survey of computer vision-based human motion capture [4]	2001	CVIU	A survey on human motion capture, including initialization, tracking, pose estimation, and recognition
A survey of advances in vision-based human motion capture and analysis [5]	2006	CVIU	A survey following [4] for human motion capture, summarizing the human motion capture methods from 2000 to 2006
Vision-based human motion analysis: An overview [6]	2007	CVIU	An overview of vision-based human motion methods and analysis based on markerless data
Advances in view-invariant human motion analysis: A review [7]	2010	TSMCS	A review of human motion methods based on human detection, view-invariant pose representation and estimation, and behavior understanding
Human pose estimation and activity recognition from multi-view videos: Comparative explorations of recent developments [8]	2012	JSTSP	A study of 3D human pose estimation and activity recognition based on multi-view approaches
A survey of human pose estimation: the body parts parsing based methods [9]	2015	JVCIR	A survey of body parts parsing-based methods for human pose estimation under single-view and multiple-view settings with different input sources (images, videos, and depth)
Human pose estimation from monocular images: A comprehensive survey [10]	2016	Sensors	A review of traditional and deep learning-based human pose estimation until 2016
3d human pose estimation: A review of the literature and analysis of covariates [11]	2016	CVIU	A review on 3D human pose estimation from RGB images and video sequences
Monocular human pose estimation: a survey of deep learning-based methods [12]	2020	CVIU	A survey of monocular human pose estimation using deep learning-based approaches
The progress of human pose estimation: a survey and taxonomy of models applied in 2D human pose estimation [13]	2020	IEEE Access	A summary of 2D human pose estimation methods and models until 2020

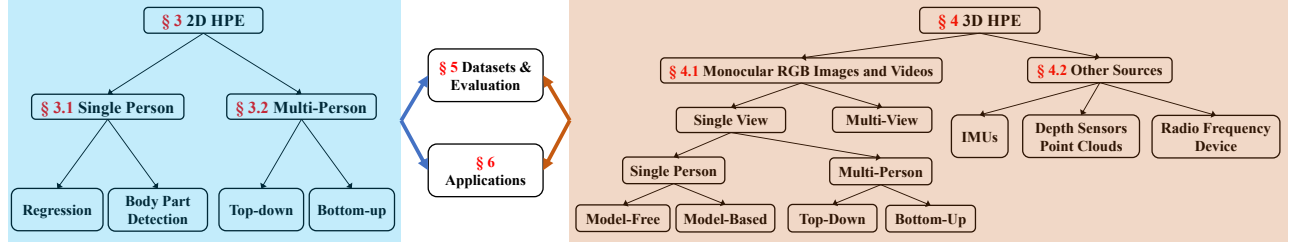


Fig. 1: Taxonomy of this survey.

methods and analyzes model interpretation. Monocular HPE from the classical to recent deep learning-based methods (till 2019) is summarized in [12]. However, it only covers 2D HPE and 3D single-view HPE from monocular images/videos. Also, no extensive performance comparison is given.

This survey aims to address the shortcomings of the previous surveys in terms of providing a systematic review of the recent deep learning-based solutions to 2D and 3D HPE but also covering other aspects of HPE including the performance evaluation of (2D and 3D) HPE methods on popular datasets, their applications, and comprehensive discussion. The key points that distinguish this survey from the previous ones are as follows:

- A comprehensive review of recent deep learning-based 2D and 3D HPE methods (up to 2020) is provided by categorizing them according to 2D or 3D scenario, single-view or multi-view, from monocular images/videos or other sources, and learning paradigm.
- Extensive performance evaluation of 2D and 3D HPE methods. We summarize and compare reported performances of promising methods on common datasets based on their categories. The comparison of results provides cues for the strength and weakness of different methods, revealing the research trends and future directions of HPE.
- An overview of a wide range of HPE applications, such as gaming, surveillance, AR/VR, and healthcare.
- An insightful discussion of 2D and 3D HPE is presented in terms of key challenges in HPE pointing to potential future research towards improving performance.

These contributions make our survey more comprehensive, up-to-date, and in-depth than previous survey papers.

1.2 Organization

In the following sections, we will cover various aspects of recent advances in HPE with deep learning.

We first overview the human body modeling techniques in § 2. Then, HPE is divided into two main categories: 2D HPE (§ 3) and 3D HPE (§ 4). Fig. 1 shows the taxonomy of deep learning methods for HPE. According to the number of people, 2D HPE methods are categorized into single-person and multi-person settings. For single-person methods (§ 3.1), there are two categories of deep learning-based methods: (1) regression methods, which directly build a mapping from input images to body joint coordinates by employing deep learning-based regressors; (2) body part detection methods, which consist of two steps: the first step involves generating heatmaps of keypoints (i.e., joints) for body part localization, and the second step involves assembling these detected keypoints into whole body pose or skeleton. For multi-person methods (§ 3.2), there are also two types of deep learning-based methods: (1) top-down methods, which construct human body poses by detecting the people first and then utilizing single-person HPE to predict the keypoints for each person; (2) bottom-up methods, which first detect body keypoints without knowing the number of people, then group the keypoints into individual poses.

3D HPE methods are classified according to the input source types: monocular RGB images and videos (§ 4.1), or other sensors (e.g., inertial measurement unit sensors, § 4.2). The majority of these methods use monocular RGB images and videos, and they are further divided into single-view and multi-view methods. Single-view methods are then separated by single-person versus multi-person. Multi-view settings are deployed mainly for multi-person pose estimation. Hence, single-person or multi-person is not specified in this category.

Next, depending on the 2D and 3D HPE pipelines, the datasets and evaluation metrics commonly used are summarized followed by a comparison of results of the promising methods (§ 5). In addition, various applications of HPE such as AR/VR are mentioned (§ 6). Finally, the paper ends by an insightful discussion of some promising directions for future research (§ 7).

2 HUMAN BODY MODELING

Human body modeling is an important aspect of HPE in order to represent keypoints and features extracted from input data. For example, most HPE methods use an N -joints rigid kinematic model. A human body is a sophisticated entity with joints and limbs, and contains body kinematic structure and body shape information. In typical methods, a model-based approach is employed to describe and infer human body pose, and render 2D and 3D poses. There are typically three types of models for human body modeling, i.e., kinematic model (used for 2D/3D HPE), planar model (used for 2D HPE) and volumetric model (used for 3D HPE), as shown in Fig. 2. In the following sections, a description of these models is provided covering different representations.

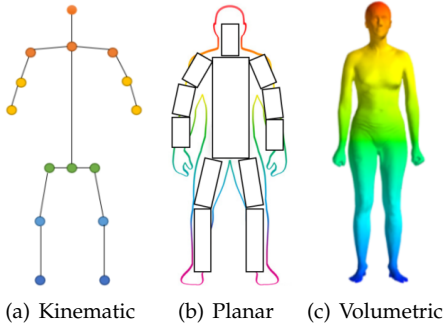


Fig. 2: Three types of models for human body modeling.

2.1 Kinematic model

The kinematic model, also called skeleton-based model [12] or kinematic chain model [14], as shown in Fig. 2 (a), includes a set of joint positions and the limb orientations to represent the human body structure. The pictorial structure model (PSM) [15] is a widely used graph model, which is also known as the tree-structured model. This flexible and intuitive human body model is successfully utilized in 2D HPE [16] [17] and 3D HPE [18] [19]. Although the kinematic model has the advantage of flexible graph-representation, it is limited in representing texture and shape information.

2.2 Planar model

Other than the kinematic model to capture the relations between different body parts, the planar model is used to represent the shape and appearance of a human body as shown in Fig. 2 (b). In the planar model, body parts are usually represented by rectangles approximating the human body contours. One example is the cardboard model [20], which is composed of body part rectangular shapes representing the limbs of a person. One of the early works [21] used the cardboard model in HPE. Another example is Active Shape Model (ASM) [22], which is widely used to capture the full human body graph and the silhouette deformations using principal component analysis [23] [24].

2.3 Volumetric models

With the increasing interest in 3D human reconstruction, many human body models have been proposed for a wide variety of human body shapes. We briefly discuss several

popular 3D human body models used in deep learning-based 3D HPE methods for recovering 3D human mesh. The volumetric model representation is depicted in Fig. 2 (c).

SMPL: Skinned Multi-Person Linear model [25] is a skinned vertex-based model which represents a broad range of human body shapes. SMPL can be modeled with natural pose-dependent deformations exhibiting soft-tissue dynamics. To learn how people deform with pose, there are 1786 high-resolution 3D scans of different subjects of poses with template mesh in SMPL to optimize the blend weights [26], pose-dependent blend shapes, the mean template shape, and the regressor from vertices to joint locations. SMPL is easy to deploy and compatible with existing rendering engines, therefore is widely adopted in 3D HPE methods.

DYNA: Dynamic Human Shape in Motion [27] model attempts to represent realistic soft-tissue motions for various body shapes. Motion related soft-tissue deformation is approximated by a low-dimensional linear subspace. In order to predict the low-dimensional linear coefficients of soft-tissue motion, the velocity and acceleration of the whole body, the angular velocities and accelerations of the body parts, and the soft-tissue shape coefficients are used. Moreover, DYNA leverages body mass index (BMI) to produce different deformations for people with different shapes.

Stitched Puppet Model [28] is a part-based graphical model integrated with a realistic body model. Different 3D body shapes and pose-dependent shape variations can be translated to the corresponding graph nodes representation. Each body part is represented by its own low-dimensional state space. The body parts are connected via pairwise potentials between nodes in the graph that “stitch” the parts together. In general, part connection via potential functions is performed by using message passing algorithms such as Belief Propagation (BP). To solve the problem that the state space of each part cannot be easily discretized to apply discrete BP, a max-product BP via a particle-based D-PMP model [29] is applied.

Frankenstein & Adam: The Frankenstein model [30] produces human motion parameters not only for body motion but also for facial expressions and hand gestures. This model is generated by blending models of the individual component meshes: SMPL [25] for the body, FaceWarehouse [31] for the face, and an artist rigged for the hand. All transform bones are merged into a single skeletal hierarchy while the native parameterization of each component is kept to express identity and motion variations. **The Adam model** [30] is optimized by the Frankenstein model using a large-scale capture of people’s clothes. With the ability to express human hair and clothing geometry, Adam is more suitable to represent human under real-world conditions.

GHUM & GHUML(ite): A fully trainable end-to-end deep learning pipeline is proposed in [32] to model statistical and articulated 3D human body shape and pose. GHUM is the moderate resolution version and GHUML is the low resolution version. GHUM and GHUML are trained by high resolution full-body scans (over 60,000 diverse human configurations in their dataset) in a deep variational auto-encoder framework. They are able to infer a host of components such as non-linear shape spaces, pose-space deformation correctives, skeleton joint center estimators, and blend skinning function [26].

3 2D HUMAN POSE ESTIMATION

2D HPE methods estimate the 2D position or spatial location of human body keypoints from images or videos. Traditional 2D HPE methods adopt different hand-crafted feature extraction techniques [33] [34] for body parts, and these early works describe human body as a stick figure to obtain global pose structures. Recently, deep learning-based approaches have achieved a major breakthrough in HPE by improving the performance significantly. In the following, we review deep learning-based 2D HPE methods with respect to single-person and multi-person scenarios.

3.1 2D single-person pose estimation

2D single-person pose estimation is used to localize human body joint positions when the input is a single-person image. If there are more than one person, the input image is cropped first so that there is only one person in each cropped patch (or sub-image). This process can be achieved automatically by an upper-body detector [35] or a full-body detector [3]. In general, there are two categories for single-person pipelines that employ deep learning techniques: regression methods and body part detection methods. Regression methods apply an end-to-end framework to learn a mapping from the input image to body joints or parameters of human body models [36]. The goal of body part detection methods is to predict approximate locations of body parts and joints [37] [38], which are normally supervised by heatmaps representation [39] [40]. Heatmap-based frameworks are now widely used in 2D HPE tasks. The general frameworks of 2D single-person HPE methods are depicted in Fig. 3.

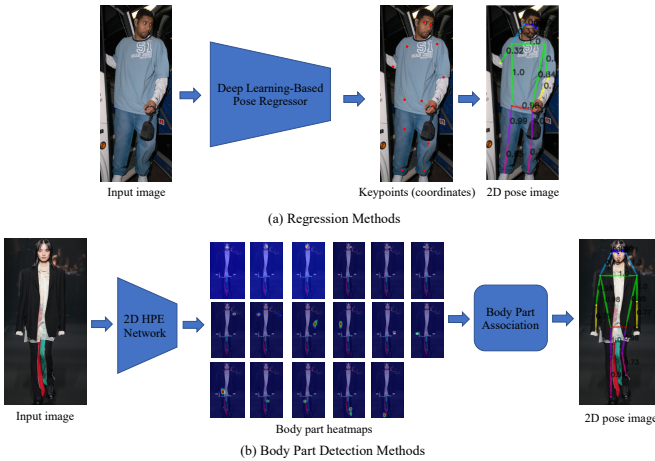


Fig. 3: Single-person 2D HPE frameworks. (a) Regression methods directly learn a mapping (via a deep neural network) from the original image to the kinematic body model and produce joint coordinates. (b) Body part detection methods predict body joint locations using the supervision of heatmaps.

3.1.1 Regression methods

There are many works based on the regression framework (e.g., [36] [41] [42] [43] [44] [45] [46] [47] [48] [49]) to predict joint coordinates from images as shown in Fig. 3 (a). Using AlexNet [1] as the backbone, Toshev and Szegedy

[36] proposed a cascaded deep neural network regressor named DeepPose to learn keypoints from images. Due to the impressive performance of DeepPose, the research paradigm of HPE began to shift from classic approaches to deep learning, in particular convolutional neural networks (CNNs). Based on GoogLeNet [50], Carreira et al. [42] proposed an Iterative Error Feedback (IEF) network, which is a self-correcting model to progressively change an initial solution by injecting the prediction error back to the input space. Sun et al. [43] introduced a structure-aware regression method called "compositional pose regression" based on ResNet-50 [51]. This method adopts a re-parameterized and bone-based representation that contains human body information and pose structure, instead of the traditional joint-based representation. Luvizon et al. [44] proposed an end-to-end regression approach for HPE using soft-argmax function to convert feature maps into joint coordinates in a fully differentiable framework.

A good feature that encodes rich pose information is critical for regression-based methods. One popular strategy to learn better feature representation is multi-task learning [52]. By sharing representations between related tasks (e.g., pose estimation and pose-based action recognition), the model can generalize better on the original task (pose estimation). Following this direction, Li et al. [46] proposed a heterogeneous multi-task framework that consists of two tasks: predicting joints coordinates from full images by building a regressor and detecting body parts from image patches using a sliding window. Fan et al. [47] proposed a Dual-Source (i.e., image patches and full images) Deep Convolutional Neural Network (DS-CNN) for two tasks: joint detection which determines whether a patch contains a body joint, and joint localization which finds the exact location of the joint in the patch. Each task corresponds to a loss function, and the combination of two tasks leads to improved results. Luvizon et al. [48] learned a multi-task network to jointly handle 2D/3D pose estimation and action recognition from video sequences.

3.1.2 Body part detection methods

Body part detection methods for HPE aim to train a body part detector to predict the positions of body joints. Recent detection methods tackle pose estimation as a heatmap prediction problem. Concretely, the goal is to estimate K heatmaps $\{H_1, H_2, \dots, H_K\}$ for a total of K keypoints. The pixel value $H_i(x, y)$ in each keypoint heatmap indicates the probability that the keypoint lies in the position (x, y) (see Fig. 3 (b)). The target (or ground-truth) heatmap is generated by a 2D Gaussian centered at the ground-truth joint location [39] [53]. Thus pose estimation networks are trained by minimizing the discrepancy (e.g., the Mean Squared-Error (MSE)) between the predicted heatmaps and target heatmaps.

Compared with joint coordinates, heatmaps provide richer supervision information by preserving the spatial location information to facilitate the training of convolutional networks. Therefore, there is a recent growing interest in leveraging heatmaps to represent the joint locations and developing effective CNN architectures for HPE, e.g., [53] [54] [39] [55] [56] [38] [40] [57] [58] [59] [60] [61] [62] [63] [64]. Tompson et al. [53] combined CNN-based body part detector with a part-based spatial-model into a unified

learning framework for 2D HPE. Lifshitz et al. [55] proposed a CNN-based method for predicting the locations of joints. It incorporates the keypoints votes and joint probabilities to determine the human pose representation. Wei et al. [40] introduced a convolutional networks-based sequential framework named Convolutional Pose Machines (CPM) to predict the locations of key joints with multi-stage processing (the convolutional networks in each stage utilize the 2D belief maps generated from previous stages and produce the increasingly refined predictions of body part locations). Newell et al. [38] proposed an encoder-decoder network named "stacked hourglass" (the encoder in this network squeezes features through bottleneck and then the decoder expands them) to repeat bottom-up and top-down processing with intermediate supervision. The stacked hourglass (SHG) network consists of consecutive steps of pooling and upsampling layers to capture information at every scale. Since then, complex variations of the SHG architecture were developed for HPE. Chu et al. [65] designed novel Hourglass Residual Units (HRUs), which extend the residual units with a side branch of filters with larger receptive fields, to capture features from various scales. Yang et al. [59] designed a multi-branch Pyramid Residual Module (PRM) to replace the residual unit in SHG, leading to enhanced invariance in scales of deep CNNs.

With the emergence of Generative Adversarial Networks (GANs) [66], they are explored in HPE to generate biologically plausible pose configurations and to discriminate the predictions with high confidence from those with low confidence, which could infer the potential poses for the occluded body parts. Chen et al. [67] constructed a structure-aware conditional adversarial network, named Adversarial PoseNet, which contains an hourglass network-based pose generator and two discriminators to discriminate against reasonable body poses from unreasonable ones. Chou et al. [68] built an adversarial learning-based network with two stacked hourglass networks sharing the same structure as discriminator and generator, respectively. The generator estimates the location of each joint, and the discriminator distinguishes the ground-truth heatmaps and predicted ones. Different from GANs-based methods that take HPE network as the generator and utilize the discriminator to provide supervision, Peng et al. [69] developed an adversarial data augmentation network to optimize data augmentation and network training by treating HPE network as a discriminator and using augmentation network as a generator to perform adversarial augmentations.

Besides these efforts in effective network design for HPE, body structure information is also investigated to provide more and better supervision information for building HPE networks. Yang et al. [70] designed an end-to-end CNN framework for HPE, which is able to find hard negatives by incorporating the spatial and appearance consistency among human body parts. A structured feature-level learning framework was proposed in [71] for reasoning the correlations among human body joints in HPE, which captures richer information of human body joints and improves the learning results. Ke et al. [72] designed a multi-scale structure-aware neural network, which combines multi-scale supervision, multi-scale feature combination, structure-aware loss information scheme, and a keypoint masking training

method to improve HPE models in complex scenarios. Tang et al. [73] built a hourglass-based supervision network, termed as Deeply Learned Compositional Model, to describe the complex and realistic relationships among body parts and learn the compositional pattern information (the orientation, scale and shape information of each body part) in human bodies. Tang and Wu [74] revealed that not all parts are related to each other, therefore introduced a Part-based Branches Network to learn representations specific to each part group rather than a shared representation for all parts.

Human poses in video sequences are (3D) spatio-temporal signals. Therefore, modeling the spatio-temporal information is important for HPE from videos. Jain et al. [75] designed a two-branch CNN framework to incorporate both color and motion features within frame pairs to build an expressive temporal-spatial model in HPE. Pfister et al. [76] proposed a convolutional network that is able to utilize temporal context information from multiple frames by using optical flow to align predicted heatmaps from neighbouring frames. Different from the previous video-based methods which are computationally intensive, Luo et al. [60] introduced a recurrent structure for HPE with Long Short-Term Memory (LSTM) [77] to capture temporal geometric consistency and dependency from different frames. This method results in a faster speed in training the HPE network for videos. Zhang et al. [78] introduced a key frame proposal network for capturing spatial and temporal information from frames and a human pose interpolation module for efficient video-based pose estimation.

3.2 2D multi-person pose estimation

Compared to single-person HPE, multi-person HPE is more difficult and challenging because it needs to figure out the number of people and their positions, and how to group keypoints for different people. In order to solve these problems, multi-person HPE methods can be classified into top-down and bottom-up methods. Top-down methods employ off-the-shelf person detectors to obtain a set of boxes (each corresponding to one person) from the input images, and then apply single-person pose estimators to each person box to generate multi-person poses. Different from top-down methods, bottom-up methods locate all the body joints in one image first and then group them to the corresponding subjects. In the top-down pipeline, the number of people in the input image will directly affect the computing time. The computing speed for bottom-up methods is usually faster than top-down methods since they do not need to detect the pose for each person separately. Fig. 4 shows the general frameworks for 2D multi-person HPE methods.

3.2.1 Top-down pipeline

In the top-down pipeline as shown in Fig. 4 (a), there are two important parts: a human body detector to obtain person bounding boxes and a single-person pose estimator to predict the locations of keypoints within these bounding boxes. A line of works focus on designing and improving the modules in HPE networks, e.g., [79] [80] [62] [81] [82] [83] [84] [85] [86] [87]. For example, in order to answer the question "how good could a simple method be" in building an HPE network, Xiao et al. [62] added a few deconvolutional layers in the ResNet

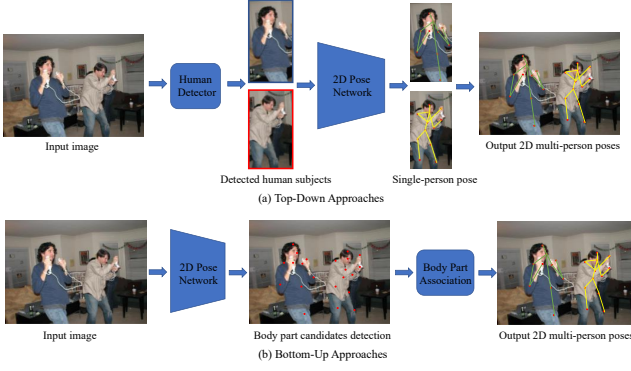


Fig. 4: Illustration of the multi-person 2D HPE frameworks. (a) Top-down approaches have two sub-tasks: (1) human detection and (2) pose estimation in the region of a single human; (b) Bottom-up approaches also have two sub-tasks: (1) detect all keypoints candidates of body parts and (2) associate body parts in different human bodies and assemble them into individual pose representations.

(backbone network) to build a simple yet effective structure to produce heatmaps for high-resolution representations. Sun et al. [81] presented a novel High-Resolution Net (HRNet) to learn reliable high resolution representations by connecting multi-resolution subnetworks in parallel and conducting repeated multi-scale fusions. To improve the accuracy of keypoint localization, Wang et al. [84] introduced a two-stage graph-based and model-agnostic framework, called Graph-PCNN. It consists of a localization subnet to obtain rough keypoint locations and a graph pose refinement module to get refined keypoints localization representations. In order to obtain more precise keypoints localization, Cai et al. [86] introduced a multi-stage network with a Residual Steps Network (RSN) module to learn delicate local representations by efficient intra-level feature fusion strategies, and a Pose Refine Machine (PRM) module to find a trade-off between local and global representations in the features.

Estimating poses under occlusion and truncation scenes often occurs in multi-person settings since the overlapping of limbs is inevitable. Human detectors may fail in the first step of top-down pipeline due to occlusion or truncation. Therefore, robustness to occlusion or truncation is an important aspect of the multi-person HPE approaches. Towards this goal, Iqbal and Gall [88] built a convolutional pose machine-based pose estimator to estimate the joint candidates. Then they used integer linear programming (ILP) to solve the joint-to-person association problem and obtain the human body poses even in presence of severe occlusions. Fang et al. [89] designed a novel regional multi-person pose estimation (RMPE) approach to improve the performance of HPE in complex scenes. Specifically, RMPE framework has three parts: Symmetric Spatial Transformer Network (to detect single person region within inaccurate bounding box), Parametric Pose Non-Maximum-Suppression (to solve the redundant detection problem), and Pose-Guided Proposals Generator (to augment training data). Papandreou et al. [79] proposed a two-stage architecture, consisting of a Faster R-CNN person detector to create bounding boxes for candidate human bodies and a keypoint estimator to predict the locations

of keypoints by using a form of heatmap-offset aggregation. The overall method works well in occluded and cluttered scenes. In order to alleviate the occlusion problem in HPE, Chen et al. [90] presented a Cascade Pyramid Network (CPN) which includes two parts: GlobalNet (a feature pyramid network to predict the invisible keypoints like eyes or hands) and RefineNet (a network to integrate all levels of features from the GlobalNet with a keypoint mining loss). Their results reveal that CPN has a good performance in predicting occluded keypoints. Su et al. [91] designed two modules, the Channel Shuffle Module and the Spatial & Channel-wise Attention Residual Bottleneck, to achieve channel-wise and spatial information enhancement for better multi-person pose estimation under occluded scenes. Qiu et al. [92] developed an Occluded Pose Estimation and Correction (OPEC-Net) module and an occluded pose dataset to solve the occlusion problem in crowd pose estimation. Umer et al. [93] proposed a keypoint correspondence framework to recover missed poses using temporal information of the previous frame in occluded scene. The network is trained using self-supervision in order to improve the pose estimation results in sparsely annotated video datasets.

3.2.2 Bottom-up pipeline

The bottom-up pipeline (e.g., [94] [95] [96] [17] [97] [98] [99] [100] [101] [102] [103]) has two main steps including body joint detection (i.e., extracting local features and predicting human body joint candidates) and joint candidates assembling for individual bodies (i.e., grouping joint candidates to build final pose representations with part association strategies) as illustrated in Fig. 4 (b).

Pishchulin et al. [94] proposed a Fast R-CNN based body part detector named DeepCut, which is one of the earliest two-step bottom-up approaches. It first detects all the body part candidates, then labels each part and assembles these parts using integer linear programming (ILP) to a final pose. However, DeepCut model is computationally expensive. To this end, Insafutdinov et al. [95] introduced DeeperCut to improve DeepCut by applying a stronger body part detector with a better incremental optimization strategy and image-conditioned pairwise terms to group body parts, leading to improved performance as well as a faster speed. Later, Cao et al. [17] built a detector named OpenPose, which uses Convolutional Pose Machines [40] (CPMs) to predict keypoints coordinates via heatmaps and Part Affinity Fields (PAFs, a set of 2D vector fields with vectormaps that encode the position and orientation of limbs) to associate the keypoints to each person. OpenPose largely accelerates the speed of the bottom-up multi-person HPE. Based on the OpenPose framework, Zhu et al. [104] improved the OpenPose structure by adding redundant edges to increase the connections between joints in PAFs and obtained better performance than the baseline approach. Although OpenPose-based methods have achieved impressive results on high resolution images, they have poor performance with low resolution images and occlusions. To address this problem, Kreiss et al. [100] proposed a bottom-up method called PifPaf that uses a Part Intensity Field (PIF) to predict the locations of body parts and a Part Association Field (PAF) to represent the joints association. This method outperformed previous OpenPose-based approaches on low resolution and occluded scenes. Motivated by OpenPose

[17] and stacked hourglass structure [38], Newell et al. [97] introduced a single-stage deep network to simultaneously obtain pose detections and group assignments. Following [97], Jin et al. [102] proposed a new differentiable Hierarchical Graph Grouping (HGG) method to learn the human part grouping. Based on [97] and [81], Cheng et al. [103] proposed a simple extension of HRNet, named Higher Resolution Network (HigherHRNet), which deconvolves the high-resolution heatmaps generated by HRNet to solve the scale variation challenge in bottom-up multi-person pose estimation.

Multi-task structures are also employed in bottom-up HPE methods. Papandreou et al. [105] introduced PersonLab to combine the pose estimation module and the person segmentation module for keypoints detection and association. PersonLab consists of short-range offsets (for refining heatmaps), mid-range offsets (for predicting the keypoints) and long-range offsets (for grouping keypoints into instances). Kocabas et al. [106] presented a multi-task learning model with a pose residual net, named MultiPoseNet, which can perform keypoints prediction, human detection and semantic segmentation tasks altogether.

3.3 2D HPE Summary

In summary, the performance of 2D HPE has been significantly improved with the blooming of deep learning techniques. In recent years, deeper and more powerful networks have promoted the performance in 2D single-person HPE such as DeepPose [36] and Stacked Hourglass Network [38], as well as in 2D multi-person HPE such as AlphaPose [89] and OpenPose [17].

Although these works have achieved sufficiently good performance in different 2D HPE scenarios, problems still remain. Regression and body part detection methods have their own advantages and limitations in 2D single-person HPE. Regression methods can learn a nonlinear mapping from input images to keypoint coordinates with an end-to-end framework, which offer a fast learning paradigm and a sub-pixel level prediction accuracy. However, they usually give sub-optimal solutions [44] due to the highly nonlinear problem. Body part detection methods, in particular heatmap-based frameworks, are more widely used in 2D HPE since (1) the probabilistic prediction of each pixel in heatmap can improve the accuracy of locating the keypoints; and (2) heatmaps provide richer supervision information by preserving the spatial location information. However, the precision of the predicted keypoints is dependent on the resolution of heatmaps. The computational cost and memory footprint are significantly increased when using high resolution heatmaps.

As for the top-down and bottom-up pipelines for 2D multi-person HPE, it is difficult to identify which method is better since both of them are widely used in recent works with their strengths and weaknesses. On one hand, top-down pipeline yields better results because it first detects each individual from the image using detection methods, then predicts the locations of keypoints using the single person-based approaches. In this case, the keypoint heatmap estimation within each detected person region is eased as the background is largely removed. On the other hand, bottom-up methods are generally faster than top-down methods, because they directly detect all the keypoints and group them

into individual poses using keypoint association strategies such as affinity linking [17], associative embedding [97], and pixel-wise keypoint regression [107].

There are several challenges in 2D HPE which need to be further addressed in future research. First is the reliable detection of individuals under significant occlusion, e.g., in crowd scenarios. The person detectors in top-down 2D HPE methods may fail to identify the boundaries of largely overlapped human bodies. Similarly, the difficulty of keypoint association is more pronounced for bottom-up approaches in occluded scenes.

The second challenge is computation efficiency. Although some methods like OpenPose [17] can achieve near real-time processing on special hardware with moderate computing power (e.g., 22 FPS on a machine with a Nvidia GTX 1080 Ti GPU), it is still difficult to implement the networks on resource-constrained devices. Real-world applications (e.g., online coaching, gaming, AR and VR) require more efficient HPE methods on commercial devices which can bring better interaction experience for users.

Another challenge lies in the limited data for rare poses. Although the size of current datasets for 2D HPE is large enough (e.g., COCO dataset [108]) for the normal pose estimation (e.g., standing, walking, running), these datasets have limited training data for unusual poses, e.g., falling. The data imbalance may cause model bias, resulting in poor performance on those poses. It would be useful to develop effective data generation or augmentation techniques to generate extra pose data for training more robust models.

4 3D HUMAN POSE ESTIMATION

3D HPE, which aims to predict locations of body joints in 3D space, has attracted much interest in recent years since it can provide extensive 3D structure information related to the human body. It can be applied to various applications (e.g., 3D movie and animation industries, virtual reality, and online 3D action prediction). Although significant improvements have recently been achieved in 2D HPE, 3D HPE still remains as a challenging task. Most existing research works tackle 3D HPE from monocular images or videos, which is an ill-posed and inverse problem due to projection of 3D to 2D where one dimension is lost. When multiple viewpoints are available or other sensors such as IMU and LiDAR are deployed, 3D HPE can be a well-posed problem employing information fusion techniques. Another limitation is that deep learning models are data-hungry and sensitive to the data collection environment. Unlike 2D human datasets where accurate 2D pose annotation can be easily obtained, collecting accurate 3D pose annotation is time-consuming and manual labeling is not practical. Also, datasets are usually collected from indoor environments with selected daily actions. Recent works [109] [110] [111] have validated the poor generalization of models trained with biased datasets by cross-dataset inference [112]. In this section, we first focus on 3D HPE from monocular RGB images and videos, and then cover 3D HPE based on other types of sensors.

4.1 3D HPE from monocular RGB images and videos

The monocular camera is the most widely used sensor for HPE in both 2D and 3D scenarios. Recent progress of

deep learning-based 2D HPE from monocular images and videos has enabled researchers to extend their works to 3D HPE. Specifically, deep learning-based 3D HPE methods are divided into two broad categories: single-view 3D HPE and multi-view 3D HPE.

4.1.1 Single-view 3D HPE

The reconstruction of 3D human poses from a single view of monocular images and videos is a nontrivial task that suffers from self-occlusions and other object occlusions, depth ambiguities, and insufficient training data. It is a severely ill-posed problem because different 3D human poses can be projected to a similar 2D pose projection. Moreover, for methods that build upon 2D joints, minor localization errors of the 2D body joints can lead to large pose distortion in the 3D space. Compared to the single-person scenario, the multi-person case is more complicated. Thus they are discussed separately in what follows.

A. Single-person 3D HPE

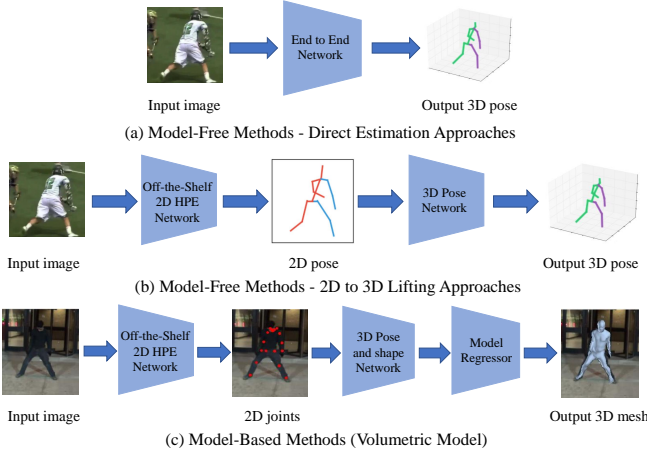


Fig. 5: Single-person 3D HPE frameworks. (a) Direct estimation approaches directly estimate the 3D human pose from 2D images. (b) 2D to 3D lifting approaches leverage the predicted 2D human pose (intermediate representation) for 3D pose estimation. (c) Model-based methods incorporate parametric body models to recover high-quality 3D human mesh. The 3D pose and shape parameters inferred by the 3D pose and shape network are fed into the model regressor to reconstruct 3D human mesh. Part of the figure is from [113].

Single-person 3D HPE approaches can be classified into model-free and model-based categories based on whether they employ a human body model (as listed in Section 2) to estimate 3D human pose or not.

Model-free methods. The model-free methods do not employ human body models to reconstruct 3D human representation. These methods can be further divided into two classes: (1) Direct estimation approaches, and (2) 2D to 3D lifting approaches.

Direct estimation: As shown in Fig. 5(a), direct estimation methods infer the 3D human pose from 2D images without intermediately estimating 2D pose representation, e.g., [114] [115] [116] [43] [117] [118] [119]. One of the early deep learning approaches was proposed by Li and Chan [114]. They

employed a shallow network to train the body part detection with sliding windows and the pose coordinate regression synchronously. A follow-up approach was proposed by Li et al. [115] where the image-3D pose pairs were used as the network input. The score network can assign high scores to the correct image-3D pose pairs and low scores to other pairs. However, these approaches are highly inefficient because they require multiple forward network inferences. Sun et al. [43] proposed a structure-aware regression approach. Instead of using joint-based representation, they adopted a bone-based representation with more stability. A compositional loss was defined by exploiting 3D bone structure with bone-based representation that encodes long range interactions between the bones. Tekin et al. [116] encoded the structural dependencies between joints by learning a mapping of 3D pose to a high-dimensional latent space. The learned high-dimensional pose representation can enforce structural constraints of the 3D pose. Pavlakos et al. [117] [118] introduced a volumetric representation to convert the highly non-linear 3D coordinate regression problem to a manageable form in a discretized space. The voxel likelihoods for each joint in the volume were predicted by a convolutional network. Ordinal depth relations of human joints were used to alleviate the need for accurate 3D ground truth pose.

2D to 3D lifting: Motivated by the recent success of 2D HPE, 2D to 3D lifting approaches that infer 3D human pose from the intermediately estimated 2D human pose have become a popular 3D HPE solution as illustrated in Fig. 5(b). Benefiting from the excellent performance of state-of-the-art 2D pose detectors, 2D to 3D lifting approaches generally outperform direct estimation approaches. In the first stage, off-the-shelf 2D HPE models are employed to estimate 2D pose, and then in the second stage 2D to 3D lifting is used to obtain 3D pose. Chen and Ramanan [120] deployed a nearest neighbor matching of the predicted 2D pose and 3D pose from a library. However, 3D HPE could fail when the 3D pose is not conditionally independent of the image given the 2D pose. Martinez et al. [121] proposed a simple but effective fully connected residual network to regress 3D joint locations based on the 2D joint locations. Despite achieving the state-of-the-art results at that time, the method could fail due to reconstruction ambiguity of over-reliance on the 2D pose detector [118]. Tekin et al. [122] and Zhou et al. [123] utilized 2D heatmaps instead of 2D pose as intermediate representations for estimating 3D pose. Moreno-Noguer [124] inferred the 3D human pose via distance matrix regression where the distances of 2D and 3D body joints were encoded into two Euclidean Distance Matrices (EDMs). EDMs are invariant to in-plane image rotations and translations, as well as scaling invariance when applying normalization operations. Wang et al. [125] developed a Pairwise Ranking Convolutional Neural Network (PRCNN) to predict the depth ranking of pairwise human joints. Then, a coarse-to-fine pose estimator was used to regress the 3D pose from 2D joints and the depth ranking matrix. Jahangiri and Yuille [126], Sharma et al. [127], and Li and Lee [128] first generated multiple diverse 3D pose hypotheses then applied ranking networks to select the best 3D pose.

Given that a human pose can be represented as a graph where the joints are the nodes and the bones are the edges, Graph Convolutional Networks (GCNs) have been applied

to the 2D-to-3D pose lifting problem by showing promising performance [129] [130] [131] [132] [133]. Choi et al. [131] proposed Pose2Mesh, which is a GCN-based method to refine the intermediate 3D pose from its PoseNet. With GCN, the MeshNet regresses the 3D coordinates of mesh vertices with graphs constructed from the mesh topology. Ci et al. [129] proposed a generic framework, named Locally Connected Network (LCN), which leverages both fully connected network and GCN to encode the relationship between local joint neighborhoods. LCN can overcome the limitations of GCN that weight sharing scheme harms pose estimation model’s representation ability, and the structure matrix lacks flexibility to support customized node dependence. Zhao et al. [130] also tackled the limitation of the shared weight matrix of convolution filters for all the nodes in GCN. A Semantic-GCN was proposed to investigate the semantic information and relationship. The semantic graph convolution (SemGConv) operation is used to learn channel-wise weights for edges. Both local and global relationships among nodes are captured since SemGConv and non-local layers are interleaved.

3D HPE datasets are usually collected from controlled environments with selected daily motions. It is difficult to obtain the 3D pose annotations for in-the-wild data. Thus 3D HPE for in-the-wild data with unusual poses and occlusions is still a challenge. To this end, a group of 2D to 3D lifting methods pay attention to estimate the 3D human pose from in-the-wild images without 3D pose annotations such as [109] [134] [135] [110] [111]. Zhou et al. [109] proposed a weakly supervised transfer learning method that uses 2D annotations of in-the-wild images as weak labels. 3D pose estimation module was connected with intermediate layers of the 2D pose estimation module. For in-the-wild images, 2D pose estimation module performed a supervised 2D heatmap regression and a 3D bone length constraint induced loss was applied in the weakly supervised 3D pose estimation module. Habibie et al. [134] tailored a projection loss to refine the 3D human pose without 3D annotation. A 3D-2D projection module was designed to estimate the 2D body joint locations with the predicted 3D pose from earlier network layer. The projection loss was used to update the 3D human pose without requiring 3D annotations. Inspired by [136], Chen et al. [135] proposed an unsupervised lifting network based on the closure and invariance lifting properties with a geometric self-consistency loss for the lift-reproject-lift process. Closure means for a lifted 3D skeleton, after random rotation and re-projection, the resulting 2D skeleton will lie within the distribution of valid 2D pose. Invariance means when changing the viewpoint of 2D projection from a 3D skeleton, the re-lifted 3D skeleton should be the same.

Instead of estimating 3D human pose from monocular images, videos can provide temporal information to improve accuracy and robustness of 3D HPE, e.g., [137] [138] [139] [140] [141] [142] [143] [144]. Hossain and Little [145] proposed a recurrent neural network using a Long Short-Term Memory (LSTM) unit with shortcut connections to exploit temporal information from sequences of human pose. Their method exploits the past events in a sequence-to-sequence network to predict temporally consistent 3D pose. Noticing that the complementary property between spatial constraints and temporal correlations is usually ignored

by prior work, Dabral et al. [139], Cai et al. [142], and Li et al. [146] exploited the spatial-temporal relationships and constraints (e.g., bone-length constraint and left-right symmetry constraint) to improve 3D HPE performance from sequential frames. Pavlo et al. [140] proposed a temporal convolution network to estimate 3D pose over 2D keypoints from consecutive 2D sequences. However, their method is based on the assumption that prediction errors are temporally non-continuous and independent, which may not hold in presence of occlusions [141]. Based on [140], Chen et al. [147] added bone direction module and bone length module to ensure human anatomy temporal consistency across video frames, while Liu et al. [148] utilized the attention mechanism to recognize significant frames and model long-range dependencies in large temporal receptive fields. Zeng et al. [133] employed the split-and-recombine strategy to address the rare and unseen pose problem. The human body is first split into local regions for processing through separate temporal convolutional network branches, then the low-dimensional global context obtained from each branch is combined to maintain global coherence.

Model-based methods. Model-based methods incorporate parametric body models as noted in Section 2 (such as kinematic model and volumetric model) to estimate human pose and shape as shown in Fig. 5(c).

The kinematic model is an articulated body representation by connected bones and joints with kinematic constraints, which has gained increasing attention in 3D HPE in recent years. Many methods leverage prior knowledge based on the kinematic model such as skeletal joints connectivity information, joints rotation properties, and fixed bone-length ratios for plausible pose estimation, e.g., [149] [19] [150] [151] [152] [153] [154] [155]. Zhou et al. [149] embedded a kinematic model into a network as kinematic layers to enforce the orientation and rotation constraints. Nie et al. [150] and Lee et al. [156] employed a skeleton-LSTM network to leverage joint relations and connectivity. Observing that human body parts have a distinct degree of freedom (DOF) based on the kinematic structure, Wang et al. [151] and Nie et al. [154] proposed bidirectional networks to model the kinematic and geometric dependencies of the human skeleton. Kundu et al. [152] [157] designed a kinematic structure preservation approach by inferring local-kinematic parameters with energy-based loss and explored 2D part segments based on the parent-relative local limb kinematic model. Xu et al. [153] demonstrated that noisy 2D joint is one of the key obstacles for accurate 3D pose estimation. Hence a 2D pose correction module was employed to refine unreliable 2D joints based on the kinematic structure. Zanfir et al. [158] introduced a kinematic latent normalizing flow representation (a sequence of invertible transformations applied to the original distribution) with differentiable semantic body part alignment loss functions.

Compared with the kinematic model, which produces human poses or skeletons, volumetric models can recover high-quality human mesh, providing extra shape information of human body. As one of the most popular volumetric models, the SMPL model [25] has been widely used in 3D HPE, e.g., [159] [160] [161] [162] [163] [164] [165] [166] [167] [168], because it is compatible with existing rendering engines. Tan et al. [161], Tung et al. [162], Pavlakos et al.

[169], and Omran et al. [170] regressed SMPL parameters to reconstruct 3D human mesh. Instead of predicting SMPL parameters, Kolotouros et al. [171] regressed the locations of the SMPL mesh vertices using a Graph-CNN architecture. Zhu et al. [172] combined the SMPL model with a hierarchical mesh deformation framework to enhance the flexibility of free-form 3D deformation. Kundu et al. [173] included a color-recovery module in the SMPL model to obtain vertex color via reflectional symmetry. Arnab et al. [113] pointed out that methods using the SMPL model usually fail on the in-the-wild data. They employed the bundle adjustment method to cope with occlusion, unusual poses and object blur. Doersch and Zisserman [165] proposed a transfer learning method to regress SMPL parameters by training on the synthetic human video dataset SURREAL [174]. Kocabas et al. [175] included the large-scale motion capture dataset AMASS [176] for adversarial training of their SMPL-based method named VIBE (Video Inference for Body Pose and Shape Estimation). VIBE leveraged AMASS to discriminate between real human motions and predicted pose by pose regression module. Since low-resolution visual content is more common in real-world scenarios than the high-resolution visual content, existing well-trained models may fail when resolution is degraded. Xu et al. [177] introduced the contrastive learning scheme into self-supervised resolution-aware SMPL-based network. The self-supervised contrastive learning scheme uses a self-supervision loss and a contrastive feature loss to enforce the feature and scale consistency.

There are several extended SMPL-based models to address the limitations of the SMPL model such as high computational complexity, and lack of hands and facial landmarks. Bogo et al. [159] proposed SMPLify to estimate 3D human mesh, which fits the SMPL model to the detected 2D joints and minimizes the re-projection error. An extended version of SMPLify was presented by Lassner et al. [160]. The running time is reduced by employing a random forest regression to regress SMPL parameters, but it still cannot achieve real-time throughput. Kanazawa et al. [178] further proposed an adversarial learning approach to directly infer SMPL parameters in real-time. Pavlakos et al. [179] introduced a new model, named SMPL-X, that can also predict fully articulated hands and facial landmarks. Following the SMPLify method, they also proposed SMPLify-X, which is an improved version learned from AMASS dataset [176]. Hassan et al. [163] further extended SMPLify-X to PROX – a method enforcing Proximal Relationships with Object eXclusion by adding 3D environmental constraints. Kolotouros et al. [164] integrated the regression-based and optimization-based SMPL parameter estimation methods to a new one named SPIN (SMPL oPtimization IN the loop) while employing SMPLify in the training loop. Osman et al. [180] upgraded SMPL to STAR by training with additional 10,000 scans for better model generalization. The number of model parameters is reduced to 20% of that of SMPL.

Instead of using the SMPL-based model, other volumetric models have also been used for recovering 3D human mesh, e.g., [181] [182] [183] [184]. Chen et al. [182] introduced a Cylinder Man Model to generate occlusion labels for 3D data and performed data augmentation. A pose regularization term was introduced to penalize wrong estimated occlusion labels. Xiang et al. [183] utilized the Adam model [30] to

reconstruct the 3D motions. A 3D human representation, named 3D Part Orientation Fields (POFs), was introduced to encode the 3D orientation of human body parts in the 2D space. Wang et al. [185] presented a new Bone-level Skinned Model of human mesh, which decouples bone modelling and identity-specific variations by setting bone lengths and joint angles. Fisch and Clark [186] introduced an orientation keypoints model which can compute full 3-axis joint rotations including yaw, pitch, and roll for 6D HPE.

B. Multi-person 3D HPE

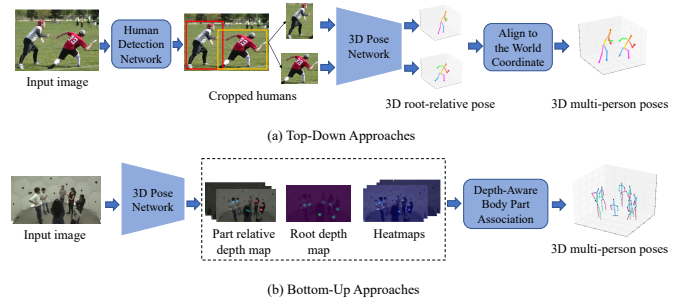


Fig. 6: Illustration of the multi-person 3D HPE frameworks. (a) Top-Down approaches first detect single-person regions by human detection network. For each single-person region, individual 3D pose can be estimated by 3D pose network. Then all 3D poses are aligned to the world coordinate. (b) Bottom-Up approaches first estimate all body joints and depth maps, then associate body parts to each person according to the root depth and part relative depth. Part of the figure is from [187].

For 3D multi-person HPE from monocular RGB images or videos, similar categories as 2D multi-person HPE are noted here: top-down approaches and bottom-up approaches as shown in Fig. 6 (a) and Fig. 6 (b), respectively. The comparison between 2D top-down and bottom-up approaches in Section 3.2 is also applicable for the 3D case.

Top-down approaches. Top-down approaches of 3D multi-person HPE first perform human detection to detect each individual person. Then for each detected person, absolute root (center joint of the human) coordinate and 3D root-relative pose are estimated by 3D pose networks. Based on the absolute root coordinate of each person and their root-relative pose, all poses are aligned to the world coordinate. Rogez et al. [188] localized candidate regions of each person to generate potential poses, and used a regressor to jointly refine the pose proposals. This localization-classification-regression method, named LCR-Net, performed well on the controlled environment datasets but could not generalize well to in-the-wild images. To address this issue, Rogez et al. [189] proposed LCR-Net++ by using synthetic data augmentation for the training data to improve performance. Zanfir et al. [190] added semantic segmentation to the 3D multi-person HPE module with scene constraints. Additionally, the 3D temporal assignment problem was tackled by the Hungarian matching method for video-based multi-person 3D HPE. Moon et al. [191] introduced a camera distance-aware approach that the cropped human images were fed into their developed RootNet to estimate the camera-centered coordinates of human body’s roots. Then the root-relative 3D

pose of each cropped human was estimated by the proposed PoseNet. Benzine et al. [192] proposed a single-shot approach named PandaNet (Pose estimAtioN and Detection Anchor-based Network). A low-resolution anchor-based representation was introduced to avoid the occlusion problem. A pose-aware anchor selection module was developed to address the overlapping problem by removing the ambiguous anchors. An automatic weighting of losses associated with different scales was used to handle the imbalance issue of different sizes of people. Li et al. [193] tackled the lack of global information in the top-down approaches. They adopted a Hierarchical Multi-person Ordinal Relations method to leverage body level semantic and global consistency for encoding the interaction information hierarchically.

Bottom-up approaches. In contrast to the top-down approaches, bottom-up approaches first produce all body joint locations and depth maps, then associate body parts to each person according to the root depth and part relative depth. A key challenge of bottom-up approaches is how to group human body joints belonging to each person. Zanfir et al. [194] formulated the person grouping problem as a binary integer programming (BIP) problem. A limb scoring module was used to estimate candidate kinematic connections of detected joints and a skeleton grouping module assembled limbs into skeletons by solving the BIP problem. Nie et al. [101] proposed a Single-stage multi-person Pose Machine (SPM) to define the unique identity root joint for each person. The body joints were aligned to each root joint by using the dense displacement maps. However, this method is limited in that only paired 2D images and 3D pose annotations could be used for supervised learning. Without paired 2D images and 3D pose annotations, Kundu et al. [195] proposed a frozen network to exploit the shared latent space between two diverse modalities under a practical deployment paradigm such that the learning could be cast as a cross-model alignment problem. Fabbri et al. [196] developed a distance-based heuristic for linking joints in a multi-person setting. Specifically, starting from detected heads (i.e., the joint with the highest confidence), the remaining joints are connected by selecting the closest ones in terms of 3D Euclidean distance.

Another challenge of bottom-up approaches is occlusion. To cope with this challenge, Metha et al. [197] developed an Occlusion-Robust Pose-Maps (ORPM) approach to incorporate redundancy into the location-maps formulation, which facilitates person association in the heatmaps especially for occluded scenes. Zhen et al. [187] leveraged a depth-aware part association algorithm to assign joints to individuals by reasoning about inter-person occlusion and bone-length constraints. Mehta et al. [198] quickly inferred intermediate 3D pose of visible body joints regardless of the accuracy. Then the completed 3D pose is reconstructed by inferring occluded joints using learned pose priors and global context. The final 3D pose was refined by applying temporal coherence and fitting the kinematic skeletal model.

Comparison of top-down and bottom-up approaches. Top-down approaches usually achieve promising results by relying on the state-of-the-art person detection methods and single-person pose estimation methods. But the computational complexity and the inference time may become excessive with the increase in the number of humans, especially in crowded scenes. Moreover, since top-down

approaches first detect the bounding box for each person, global information in the scene may get neglected. The estimated depth of cropped region may be inconsistent with the actual depth ordering and the predicted human bodies may be placed in overlapping positions. On the contrary, the bottom-up approaches enjoy linear computation and time complexity. However, if the goal is to recover 3D body mesh, it is not straightforward for the bottom-up approaches to reconstruct human body meshes. For top-down approaches, after detecting each individual person, human body mesh of each person can be easily recovered by incorporating the model-based 3D single-person HPE estimator. While for the bottom-up approaches, additional model regressor module is needed to reconstruct human body meshes based on the final 3D poses.

4.1.2 Multi-view 3D HPE

The partial occlusion is a challenging problem for 3D HPE in the single-view setting. The natural solution to overcome this problem is to estimate 3D human pose from multiple views, since the occluded part in one view may become visible in other views. In order to reconstruct the 3D pose from multiple views, the association of corresponding location between different cameras needs to be resolved.

A group of methods [199] [200] [201] [202] [203] used body models to tackle the association problem by optimizing model parameters to match the model projection with the 2D pose. The widely used 3D pictorial structure model [204] is such a model. However, these methods usually need large memory and expensive computational cost, especially for multi-person 3D HPE under multi-view settings. Rhodin et al. [205] employed a multi-view consistency constraint in the network, however it requires a large amount of 3D ground-truth training data. To overcome this limitation, Rhodin et al. [206] further proposed an encoder-decoder framework to learn the geometry-aware 3D latent representation from multi-view images and background segmentation without 3D annotations. Chen et al. [207], Dong et al. [202], Chen et al. [208], Mitra et al. [209], Iqbal et al. [210], Zhang et al. [211], and Huang et al. [212] proposed multi-view matching frameworks to reconstruct 3D human pose across all viewpoints with consistency constraints. Pavlakos et al. [199] and Zhang et al. [213] aggregated the 2D keypoint heatmaps of multi-view images into a 3D pictorial structure model based on all the calibrated camera parameters. However, when multi-view camera environments change, the model needs to be retrained. Liang et al. [201] and Habermann et al. [214] inferred the non-rigid 3D deformation parameters to reconstruct a 3D human body mesh from multi-view images. Kadkhodamohammadi and Padoy [215], Qiu et al. [200], and Kocabas et al. [216] employed epipolar geometry to match paired multi-view poses for 3D pose reconstruction and generalized their methods to new multi-view camera environments. It should be noted that matching each pair of views separately without the cycle consistency constraint may lead to incorrect 3D pose reconstructions [202]. Tu et al. [203] aggregated all the features in each camera view in the 3D voxel space to avoid incorrect estimation in each camera view. A cuboid proposal network and a pose regression network were designed to localize all people and to estimate the 3D pose, respectively. When given sufficient viewpoints

(more than ten), it is not practical to use all viewpoints for 3D pose estimation. Pirinen et al. [217] proposed a self-supervised reinforcement learning approach to select a small set of viewpoints to reconstruct the 3D pose via triangulation.

Besides accuracy, the lightweight architecture, fast inference time, and efficient adaptation to new camera settings also need to be taken into consideration in multi-view HPE. In contrast to [202] which matched all view inputs together, Chen et al. [218] applied an iterative processing strategy to match 2D poses of each view with the 3D pose while the 3D pose was updated iteratively. Compared to the previous methods whose running time may explode with the increase in the number of cameras, their method has linear time complexity. Remelli et al. [219] encoded images of each view into a unified latent representation so that the feature maps were disentangled from camera viewpoints. As a lightweight canonical fusion, these 2D representations are lifted to the 3D pose using a GPU-based Direct Linear Transform to accelerate the processing. In order to improve the generalization ability of multi-view fusion scheme, Xie et al. [220] proposed a pre-trained multi-view fusion model (MetaFuse), which can be efficiently adapted to new camera settings with few labeled data. They deployed the model-agnostic meta-learning framework to learn the optimal initialization of the generic fusion model for adaptation.

4.2 3D HPE from other sources

Although monocular RGB camera is the most common device used for 3D HPE, other sensors (e.g., depth sensor, IMUs, and radio frequency device) are also used for this purpose.

Depth and point cloud sensors: Depth sensors have gained more attention recently for conducting 3D computer vision tasks due to their low-cost and increased utilization. As one of the key challenges in 3D HPE, the depth ambiguity problem can be alleviated by utilizing depth sensors. Yu et al. [221] presented a single-view and real-time method named DoubleFusion to estimate 3D human pose from a single depth sensor without using images. The inner body layer was used to reconstruct 3D shape by volumetric representation, and the outer layer updated the body shape and pose by fusing more geometric details. Xiong et al. [222] proposed an Anchor-to-Joint regression network (A2J) using depth images. 3D joints positions were estimated by integrating estimated multiple anchor points with global-local spatial context information. Kadkhodamohammadi et al. [223] used multi-view RGB-D cameras to capture color images with depth information in the real operating room environments. A random forest-based prior was deployed to incorporate priori environment information. The final 3D pose was estimated by multi-view fusion and RGB-D optimization. Zhi et al. [224] reconstructed detailed meshes with high-resolution albedo texture from RGB-D video.

Compared with depth images, point clouds can provide more information. The state-of-the-art point cloud feature extraction techniques, PointNet [225] and PointNet++ [226], have demonstrated excellent performance for classification and segmentation tasks. Jiang et al. [227] combined PointNet++ with the SMPL body model to regress 3D human pose. A modified PointNet++ with a graph aggregation module can extract more useful unordered features. After mapping

into ordered skeleton joint features by an attention module, a skeleton graph module extracts ordered features to regress SMPL parameters for accurate 3D pose estimation. Wang et al. [228] presented PointNet++ with a spatial-temporal mesh attention convolution method to predict 3D human meshes with refinement.

IMUs with monocular images: Wearable Inertial Measurement Units (IMUs) can track the orientation and acceleration of specific human body parts by recording motions without object occlusions and clothes obstructions. However, the drifting problem may occur overtime when using IMUs. Marcard et al. [229] proposed the Sparse Inertial Poser (SIP) to reconstruct human pose from 6 IMUs attached to the human body. The collected information was fitted into the SMPL body model with coherence constraints to obtain accurate results. Marcard et al. [230] further associated 6-17 IMU sensors with a hand-held moving camera for in-the-wild 3D HPE. A graph-based optimization method was introduced to assign each 2D person detection to a 3D pose candidate from long-range frames. Huang et al. [231] addressed the limitation of the Sparse Inertial Poser (SIP) method [229]. Multiple pose parameters can generate the same IMU orientation, also collecting IMUs data is time-consuming. Thus, a large synthetic dataset was created by placing virtual sensors on the SMPL mesh to obtain orientations and accelerations from motion capture sequences of AMASS dataset [176]. A bi-directional RNN framework was proposed to map IMU orientations and accelerations to SMPL parameters with past and future information. Zhang et al. [232] introduced an orientation regularized pictorial structure model to estimate 3D pose from multi-view heatmaps associated with IMUs orientation. Huang et al. [233] proposed DeepFuse, a two-stage approach by fusing IMUs data with multi-view images. The first stage only processes multi-view images to predict a volumetric representation and the second stage uses IMUs to refine the 3D pose by an IMU-bone refinement layer.

Radio frequency device: Radio frequency (RF) based sensing technology has also been used to localize people. The ability to traverse walls and to bounce off human bodies in the WiFi range without carrying wireless transmitters is the major advantage for deploying a RF-based sensing system. Also, privacy can be preserved due to non-visual data. However, RF signals have a relatively low spatial resolution compared to visual camera images and the RF systems have shown to generate coarse 3D pose estimation. Zhao et al. [234] proposed a RF-based deep learning method, named RF-Pose, to estimate 2D pose for multi-person scenarios. Later the extended version, named RF-Pose3D [235], can estimate 3D skeletons for multi-person. Based on these, Zhao et al. [236] presented a temporal adversarial training method with multi-headed attention module, named RF-Avatar, to recover a full 3D body mesh using the SMPL body model.

Other sensors/sources: Besides using the aforementioned sensors, Isogawa et al. [237] estimated 3D human pose from the 3D spatio-temporal histogram of photons captured by a non-line-of-sight (NLOS) imaging system. Tome et al. [238] tackled the egocentric 3D pose estimation via a fish-eye camera. Saini et al. [239] estimated human motion using images captured by multiple Autonomous Micro Aerial Vehicles (MAVs). Clever et al. [240] focused on the HPE of the rest position in bed from pressure images which were

collected by a pressure sensing mat.

4.3 3D HPE Summary

3D HPE has made significant advancements in recent years. Since a large number of 3D HPE methods apply the 2D to 3D lifting strategy, the performance of 3D HPE has been improved considerably due to the progress made in 2D HPE. Some 2D HPE methods such as OpenPose [17], CPN [90], AlphaPose [89], and HRNet [81] have been extensively used as 2D pose detector in 3D HPE methods. Besides the 3D pose, some methods also recover 3D human mesh from images or videos, e.g., [164] [175] [241] [242]. However, despite the progress made so far, there are still several challenges.

One challenge is the model generalization. High-quality 3D ground truth pose annotations depend on motion capture systems which cannot be easily deployed in random environment. Therefore, the existing datasets are mainly captured in constrained scenes. The state-of-the-art methods can achieve promising results on these datasets, but their performance degrades when applied to in-the-wild data. It is possible to leverage gaming engines to generate synthetic datasets with diverse poses and complex scenes, e.g., SURREAL dataset [174] and GTA-IM dataset [243]. However, learning from synthetic data may not achieve the desired performance due to a gap between synthetic and real data distributions.

Same as 2D HPE, robustness to occlusion and computation efficiency are two key challenges for 3D HPE as well. The performance of current 3D HPE methods drops considerably in crowded scenarios due to severe mutual occlusion and possibly low resolution content of each person. 3D HPE is more computation demanding than 2D HPE. For example, 2D to 3D lifting approaches rely on 2D poses as intermediate representations for inferring 3D poses. Therefore, it is critical to develop computationally efficient 2D HPE pipelines while maintaining high accuracy for pose estimation.

5 DATASETS AND EVALUATION METRICS

Datasets are very much needed in conducting HPE. They are also necessary to provide a fair comparison among different algorithms. Collecting a comprehensive and universal dataset poses challenges due to the complexity and variations of application scenes. A number of datasets have been collected to evaluate and compare results based on different metrics. In this section, we present the traditional datasets utilized in HPE, as well as more recent ones used for 2D and 3D deep learning-based HPE methods. In addition to these datasets with different features and task requirements, this section also covers several commonly used evaluation metrics for both 2D and 3D HPE. The results achieved by existing methods on the popular datasets are summarized as well.

5.1 Datasets for 2D HPE

There are many 2D human pose datasets before deep learning found its way into human pose estimation. These datasets are of two types: (1) upper body pose datasets including Buffy Stickmen [244], ETHZ PASCAL Stickmen [245], We Are Family [246], Video Pose 2 [247] and Sync. Activities [248]; and (2) full-body pose datasets including PASCAL Person Layout [249], Sport [250] and UIUC people [251]. However,

TABLE 2: Datasets for 2D HPE.

Image-based datasets						
Name	Year	Single-Person / Multi-Person	Joints	Size		
				Train	Val	Test
FLIC [252]	2013	Single	10	5k	-	1k
FLIC-full [252]	2013	Single	10	20k	-	-
FLIC-plus [53]	2014	Single	10	17k	-	-
MPII [253]	2014	Single	16	29k	-	12k
		Multiple	16	3.8k	-	1.7k
LSP [16]	2010	Single	14	1k	-	1k
LSP-extended [254]	2011	Single	14	10k	-	-
COCO2016 [108]	2016	Multiple	17	45k	22k	80k
COCO2017 [108]	2017	Multiple	17	64k	2.7k	40k
AIC-HKD [255]	2017	Multiple	14	210k	30k	60k
CrowdPose [256]	2019	Multiple	14	10k	2k	8k
Video-based datasets						
Name	Year	Single-Person / Multi-Person	Joints	Size		
				Train	Val	Test
Penn Action [257]	2013	Single	13	1k	-	1k
J-HMDB [258]	2013	Single	15	0.6k	-	0.3k
PoseTrack [259]	2017	Multiple	15	292	50	208

only a few recent works use these 2D HPE datasets because they have many limitations such as lack of diverse object movements and small number of images. Since deep learning-based approaches are fueled by large amounts of training data, only the large-scale 2D HPE datasets are reviewed in this section. They are summarized under two different categories (image-based and video-based) in Table 2.

5.1.1 Image-based datasets

Frames Labeled In Cinema (FLIC) Dataset [252] is one of the early image-based 2D HPE datasets, which contains 5,003 images collected automatically from Hollywood movies. Around 4,000 images are used as the training set and the rest are used as the testing set. The FLIC dataset uses a body part detector named Poselets [260] to obtain about 20K person candidates from every tenth frame of 30 popular Hollywood movies. The subjects in these images have different kinds of poses. The full set of frames harvested from the movies is called the FLIC-full dataset. It is a superset of the original FLIC dataset and contains 20,928 occluded, non-frontal samples. A new FLIC-based dataset named FLIC-plus was introduced in [53] by removing all the images that contain the same scene with the test set in the FLIC dataset. Dataset Link: <https://bensapp.github.io/flic-dataset.html>

Leeds Sports Pose (LSP) Dataset [16] has 2,000 annotated images from Flickr and 8 sports tags covering different sports including athletics, badminton, baseball, gymnastics, parkour, soccer, tennis, and volleyball. In the LSP dataset, every person's full body is labeled with a total of 14 joints. In addition, the Leeds Sports Pose Extended dataset (LSP-extended) [254] extends the LSP dataset and is only used for training. LSP-extended dataset has over 10,000 images from Flickr. In most recent research, LSP and LSP-extended datasets have been used for single-person HPE. Dataset Link: <https://sam.johnson.io/research/lsp.html>

Max Planck Institute for Informatics (MPII) Human Pose Dataset [253] is a popular dataset for evaluation of articulated HPE. The dataset includes around 25,000 images containing over 40,000 individuals with annotated body joints. Based on [261], the images were systematically collected by a two-level hierarchical method to capture everyday human activities. The entire dataset covers 410 human activities and all the images are labeled. Each image was extracted from a YouTube video and provided with preceding and following un-annotated frames. Moreover, rich annotations including

body part occlusions, 3D torso and head orientations are labeled by workers on Amazon Mechanical Turk. Images in MPII are suitable for 2D single-person or multi-person HPE. Dataset Link: <http://human-pose.mpi-inf.mpg.de/#>

Microsoft Common Objects in Context (COCO) Dataset [108] is the most widely used large-scale dataset. It has more than 330,000 images and 200,000 labeled subjects with keypoints, and each individual person is labeled with 17 joints. The COCO dataset is not only proposed for pose estimation and analysis, but also used for object detection and image segmentation in natural environments, recognition in context, etc. There are two versions of the COCO datasets for HPE: COCO keypoints 2016 and COCO keypoints 2017, which are hosted by COCO 2016 Keypoints Detection Challenge and COCO 2017 Keypoints Detection Challenge, respectively. The difference between COCO 2016 and COCO 2017 lies in the training, validation and test split. The COCO dataset has been extensively used in multi-person HPE works. In addition, Jin et al. [262] proposed COCO-WholeBody Dataset with whole-body annotations for HPE. Dataset Link: <https://cocodataset.org/#home>

AI Challenger Human Keypoint Detection (AIC-HKD) Dataset [255] is currently the largest training dataset for 2D HPE. It has 300,000 annotated images for keypoint detection. There are 210,000 images for training, 30,000 images for validation and over 600,000 images for testing. The images were collected from internet search engines and mainly focused on daily activities of people. Dataset Link: <https://challenger.ai/>

CrowdPose Dataset [256] is one of the latest dataset for 2D HPE in crowded and occlusive scenarios. This dataset contains 20,000 images selected from 30,000 images with *Crowd Index* (a measurement satisfies uniform distribution to judge the crowding level in images). The training, validation and testing datasets have 10,000 images, 2,000 images and 8,000 images separately. Dataset Link: <https://github.com/Jeff-sjtu/CrowdPose>

5.1.2 Video-based datasets

Penn Action Dataset [257] consists of 2,326 video sequences with 15 different actions and human joint annotations. The videos contain frames with annotations from sports actions: baseball pitch, baseball swing, tennis forehand, tennis serve, bench press, bowling, clean and jerk, golf swing, jump rope, jumping jacks, pull up, push up, sit up, squat, and strum guitar. The annotations for images were labeled using Amazon Mechanical Turk. Dataset Link: <http://dreamdragon.github.io/PennAction/>

Joint-annotated Human Motion Database (J-HMDB) [258] is a fully annotated video dataset for action recognition, human detection and HPE. There are 21 action categories including brush hair, catch, clap, climb stairs, golf, jump, kick ball, pick, pour, pull-up, push, run, shoot ball, shoot bow, shoot gun, sit, stand, swing baseball, throw, walk, and wave. There are 928 video clips comprising 31,838 annotated frames. A 2D articulated human puppet model is applied to generate all the annotations based on Amazon Mechanical Turk. The 70 percent images in the J-HMDB dataset are used for training and the rest images are used for testing. Dataset Link: <http://jhmdb.is.tue.mpg.de/>

PoseTrack Dataset [259] is a large-scale dataset for multi-person pose estimation and articulated tracking in video analysis. Each person in a video has a unique track ID with annotations. PoseTrack contains 1,356 video sequences, around 46,000 annotated video frames and 276,000 body pose annotations for training, validation and testing. Dataset Link: <https://posetrack.net/>

5.2 Evaluation Metrics for 2D HPE

It is difficult to precisely evaluate the performance of HPE because there are many features and requirements that need to be considered (e.g., upper/full human body, single/multiple pose estimation, the size of human body). As a result, many evaluation metrics have been used for 2D HPE. Here we summarize the commonly used ones.

Percentage of Correct Parts (PCP) [263] is a measure commonly used in early works on 2D HPE, which evaluates stick predictions to report the localization accuracy for limbs. The localization of limbs is determined when the distance between the predicted joint and ground truth joint is less than a fraction of the limb length (between 0.1 to 0.5). In some works, the PCP measure is also referred to as PCP@0.5, where the threshold is 0.5. This measure is used on the LSP dataset for single-person HPE evaluation. However, PCP has not been widely implemented in latest works because it penalizes the limbs with short length which are hard to detect. The performance of a model is considered better when it has a higher PCP measure. In order to address the drawbacks of PCP, Percentage of Detected Joints (PDJ) is introduced, where a prediction joint is considered as detected if the distance between predicted joints and true joints is within a certain fraction of the torso diameter [36].

Percentage of Correct Keypoints (PCK) [264] is also used to measure the accuracy of localization of different keypoints within a given threshold. The threshold is set to 50 percent of the head segment length of each test image and it is denoted as PCKh@0.5. PCK is referred to as PCK@0.2 when the distance between detected joints and true joints is less than 0.2 times the torso diameter. The higher the PCK value, the better model performance is regarded.

Average Precision (AP) and Average Recall (AR). AP measure is an index to measure the accuracy of keypoints detection according to precision (the ratio of true positive results to the total positive results) and recall (the ratio of true positive results to the total number of ground truth positives). AP computes the average precision value for recall over 0 to 1. AP has several similar variants. For example, Average Precision of Keypoints (APK) is introduced in [264]. Mean Average Precision (mAP), which is the mean of average precision over all classes, is a widely used metric on the MPII and PoseTrack datasets. Average Recall (AR) is another metric used in the COCO keypoint evaluation [265]. Object Keypoint Similarity (OKS) plays the similar role as the Intersection over Union (IoU) in object detection and is used for AP or AR. This measure is computed from the scale of the subject and the distance between predicted points and ground truth points. The COCO evaluation usually uses mAP across 10 OKS thresholds as the evaluation metric.

TABLE 3: Comparison of different 2D single-person HPE methods on the LSP dataset using PCP measure (for the individual limbs like torso, head, legs, arms and the whole body) with "Person-Centric" annotations. The best two scores are marked in red and blue, respectively.

Leeds Sports Pose (LSP)									
	Year	Method	Torso	Upper Leg	Lower Leg	Upper arm	Forearm	Head	Total
Body detection	2016	[55]	97.3	88.9	84.5	80.4	71.4	94.7	84.2
	2016	[40]	98.0	92.2	89.1	85.8	77.9	95.0	88.3
	2016	[56]	97.7	92.4	89.3	86.7	79.7	95.2	88.9
	2017	[65]	98.4	95.0	92.8	88.5	81.2	95.7	90.9
Regression	2016	[42]	95.3	81.8	73.3	66.7	51.0	84.4	72.5
	2019	[44]	98.2	93.6	91.0	86.6	78.2	96.8	89.4

TABLE 4: Comparison of different methods on the MPII dataset for 2D single-person HPE using PCKh@0.5 measure (i.e., the threshold is equal to 50 percent of the head segment length of each test image). The best two scores are marked in red and blue, respectively.

Max Planck Institute for Informatics (MPII)										
	Year	Method	Head	Shoulder	Elbow	Wrist	Hip	Knee	Ankle	Total
Body detection	2016	[58]	97.7	95.0	88.2	83.0	87.9	82.6	78.4	88.1
	2016	[95]	96.8	95.2	89.3	84.4	88.4	83.4	78.0	88.5
	2016	[40]	97.8	95.0	88.7	84.0	88.4	82.8	79.4	88.5
	2016	[56]	97.9	95.1	89.9	85.3	89.4	85.7	81.7	89.7
	2016	[38]	98.2	96.3	91.2	87.1	90.1	87.4	83.6	90.9
	2017	[65]	98.5	96.3	91.9	88.1	90.6	88.0	85.0	91.5
	2017	[68]	98.2	96.8	92.2	88.0	91.3	89.1	84.9	91.8
	2017	[67]	98.1	96.5	92.5	88.5	90.2	89.6	86.0	91.9
	2017	[59]	98.5	96.7	92.5	88.7	91.1	88.6	86.0	92.0
	2018	[72]	98.5	96.8	92.7	88.4	90.6	89.3	86.3	92.1
	2018	[73]	98.4	96.9	92.6	88.7	91.8	89.4	86.2	92.3
	2018	[62]	98.5	96.6	91.9	87.6	91.1	88.1	84.1	91.5
	2019	[81]	98.6	96.9	92.8	89.0	91.5	89.0	85.7	92.3
	2019	[63]	98.6	97.0	92.8	88.8	91.7	89.8	86.6	92.5
	2019	[82]	98.4	97.1	93.2	89.2	92.0	90.1	85.5	92.6
	2020	[64]	-	-	-	-	-	-	-	92.7
	2020	[86]	98.5	97.3	93.9	89.9	92.0	90.6	86.8	93.0
Regression	2017	[188]	-	-	-	-	-	-	-	74.2
	2016	[42]	95.7	91.7	81.7	72.4	82.8	73.2	66.4	81.3
	2017	[43]	97.5	94.3	87.0	81.2	86.5	78.5	75.4	86.4
	2019	[49]	98.3	96.4	91.5	87.4	90.0	87.1	83.7	91.1
	2019	[44]	98.1	96.6	92.0	87.5	90.6	88	82.7	91.2

Note: [95], [81], [82], [86], [188] are 2D multi-person HPE methods, which are also applied to the single-person case here.

TABLE 5: Comparison of different 2D multi-person HPE methods on the full testing set of the MPII dataset using mAP measure. The best two scores are marked in red and blue, respectively.

Max Planck Institute for Informatics (MPII)										
	Year	Method	Head	Shoulder	Elbow	Wrist	Hip	Knee	Ankle	Total
Top-down	2016	[88]	58.4	53.9	44.5	35.0	42.2	36.7	31.1	43.1
	2017	[89]	88.4	86.5	78.6	70.4	74.4	73.0	65.8	76.7
Bottom-up	2016	[95]	78.4	72.5	60.2	51.0	57.2	52.0	45.4	59.5
	2017	[96]	88.8	87.0	75.9	64.9	74.2	68.8	60.5	74.3
	2017	[97]	92.1	89.3	78.9	69.8	76.2	71.6	64.7	77.5
	2017	[17]	91.2	87.6	77.7	66.8	75.4	68.9	61.7	75.6
	2018	[98]	91.8	89.5	80.4	69.6	77.3	71.7	65.5	78.0
	2019	[101]	89.7	87.4	80.4	72.4	76.7	74.9	68.3	78.5

5.3 Performance Comparison of 2D HPE Methods

In Tables 3 ~ 6, we have summarized the performance of different 2D HPE methods on the popular datasets together with the relevant and commonly used evaluation metrics. For comparison on the LSP dataset, the PCP measure is employed to evaluate the performance of body detection-based and regression-based methods in Table 3. Table 4 shows the comparison results for different 2D HPE methods on the MPII dataset using PCKh@0.5 measure. It is worth noting that body detection methods generally have better performance than regression methods, thus gaining more popularity in recent 2D HPE research. In Table 5, the mAP comparison on the full testing set of the MPII dataset is reported. Table 6 presents

the experimental results of different 2D HPE methods on the test-dev set of the COCO dataset, together with a summary of the experiment settings (extra data, backbones in models, input images size) and AP scores for each approach.

5.4 Datasets for 3D HPE

In contrast to numerous 2D human pose datasets with high-quality annotation, acquiring accurate 3D annotation for 3D HPE datasets is a challenging task that requires motion capture systems such as MoCap and wearable IMUs. Due to this requirement, many 3D pose datasets are created in constrained environments. Here, the widely used 3D pose datasets under different settings are summarized in Table 7.

TABLE 6: Comparison of different 2D multi-person HPE methods on the test-dev set of the COCO dataset using AP measure (AP.5 means AP at OKS = 0.50, AP.75 means AP at OKS = 0.75, AP(M) is used for medium objects, AP(L) is used for large objects). The best two scores are marked in red and blue, respectively.

Microsoft Common Objects in Context (COCO)										
	Year	Method	Extra data	Backbone	Input size	AP	AP.5	AP.75	AP(M)	AP(L)
Top-down	2017	[79]	no	ResNet-101	353×257	64.9	85.5	71.3	62.3	70.0
	2017	[79]	yes	ResNet-101	353×257	68.5	87.1	75.5	65.8	73.3
	2018	[90]	no	ResNet-Inception	384×288	72.1	91.4	80.0	68.7	77.2
	2018	[62]	no	ResNet-152	384×288	73.7	91.9	81.1	70.3	80.0
	2019	[81]	no	HRNet-W32	384×288	74.9	92.5	82.8	71.3	80.9
	2019	[81]	no	HRNet-W48	384×288	75.5	92.5	83.3	71.9	81.5
	2019	[81]	yes	HRNet-W48	384×288	77.0	92.7	84.5	73.4	83.1
	2019	[82]	yes	4xResNet-50	384×288	77.1	93.8	84.6	73.4	82.3
	2020	[87]	no	HRNet-W48	384×288	76.2	92.5	83.6	72.5	82.4
	2020	[87]	yes	HRNet-W48	384×288	77.4	92.6	84.6	73.6	83.7
	2020	[86]	no	4xRSN-50	384×288	78.6	94.3	86.6	75.5	83.3
Bottom-up	2017	[17]	no	-	-	61.8	84.9	67.5	57.1	68.2
	2017	[97]	no	Hourglass	512×512	65.5	86.8	72.3	60.6	72.6
	2018	[105]	no	ResNet-152	1401×1401	68.7	89.0	75.4	64.1	75.5
	2019	[99]	no	ResNet-101	800×800	64.8	87.8	71.1	60.4	71.5
	2019	[101]	no	Hourglass	384×384	66.9	88.5	72.9	62.6	73.1
	2020	[102]	no	Hourglass	512×512	67.6	85.1	73.7	62.7	74.6
	2020	[103]	no	HRNet-W48	640×640	70.5	89.3	77.2	66.6	75.8

TABLE 7: Datasets for 3D HPE.

Dataset	Year	Capture system	Environment	Size	Single person	Multi-person	Single view	Multi-view
HumanEva [266]	2010	Marker-based MoCap	Indoor	6 subject, 7 actions, 40k frames	Yes	No	Yes	Yes
Human3.6M [267]	2014	Marker-based MoCap	Indoor	11 subjects, 17 actions, 3.6M frames	Yes	No	Yes	Yes
CMU Panoptic [268]	2016	Marker-less MoCap	Indoor	8 subjects, 1.5M frames	Yes	Yes	Yes	Yes
MPI-INF-3DHP [269]	2017	Marker-less MoCap	Indoor and outdoor	8 subjects, 8 actions, 1.3M frames	Yes	No	Yes	Yes
TotalCapture [270]	2017	Marker-based MoCap with IMUs	Indoor	5 subjects, 5 actions, 1.9M frames	Yes	No	Yes	Yes
3DPW [230]	2018	Hand-held cameras with IMUs	Indoor and outdoor	7 subjects, 51k frames	Yes	Yes	Yes	No
MuPoTS-3D [197]	2018	Marker-less MoCap	Indoor and outdoor	8 subjects, 8k frames	Yes	Yes	Yes	Yes
AMASS [176]	2019	Marker-based MoCap	Indoor and outdoor	300 subjects, 9M frames	Yes	No	Yes	Yes
NBA2K [271]	2020	NBA2K19 game engine	Indoor	27 subjects, 27k poses	Yes	No	Yes	No
GTA-IM [243]	2020	GTA game engine	Indoor	1M frames	Yes	No	Yes	No
Occlusion-Person [213]	2020	Unreal Engine 4 game engine	Indoor	73k frames	Yes	No	Yes	Yes

HumanEva Dataset [266] contains 7 calibrated video sequences (4 gray-scale and 3 color) with ground truth 3D annotation captured by a commercial MoCap system from ViconPeak. The database consists of 4 subjects performing 6 common actions (walking, jogging, gesturing, throwing and catching a ball, boxing, and combo) in a $3m \times 2m$ area. Dataset Link: <http://humaneva.is.tue.mpg.de/>

Human3.6M [267] is the most widely used indoor dataset for 3D HPE from monocular images and videos. There are 11 professional actors (6 males and 5 females) performing 17 activities (e.g., smoking, taking photo, talking on the phone) from 4 different views in an indoor laboratory environment. This dataset contains 3.6 million 3D human poses with 3D ground truth annotation captured by accurate marker-based MoCap system. There are 3 protocols with different training and testing data splits. Protocol #1 uses images of subjects S1, S5, S6, and S7 for training, and images of subjects S9 and S11 for testing. Protocol #2 uses the same training-testing split as Protocol #1, but the predictions are further post-processed by a rigid transformation before comparing to the ground-truth. Protocol #3 uses images of subjects S1, S5, S6, S7, and S9 for training, and images of subjects S11 for testing. Dataset Link: <http://vision.imar.ro/human3.6m/>

MPI-INF-3DHP [269] is a dataset captured by a commercial marker-less MoCap system in a multi-camera studio. There are 8 actors (4 males and 4 females) performing 8 human activities including walking, sitting, complex exercise posed, and dynamic actions. More than 1.3 million frames

from 14 cameras were recorded in a green screen studio which allows automatic segmentation and augmentation. Dataset Link: <http://gvv.mpi-inf.mpg.de/3dhp-dataset/>

TotalCapture Dataset [270] contains fully synchronised videos with IMU and Vicon labeling for over 1.9 million frames. There are 13 sensors placed on key body parts such as head, upper and lower back, upper and lower limbs, and feet. The data was collected indoors with 8 calibrated full HD video cameras at 60 Hz measuring roughly $4 \times 6 m^2$. There are 5 actors (4 males and 1 female) performing actions, repeated 3 times, including walking, running, and freestyle. Dataset Link: <https://cvssp.org/data/totalcapture/>

CMU Panoptic Dataset [268] contains 65 sequences (5.5 hours) with 1.5 million of 3D skeletons of multiple people scenes. This dataset was captured by a marker-less motion capture system with 480 VGA camera views, more than 30 HD views, 10 RGB-D sensors, and a calibrated hardware-based synchronization system. The test set contains 9,600 frames from HD cameras for 4 activities (Ultimatum, Mafia, Haggling, and Pizza). Dataset Link: domedb.perception.cs.cmu.edu/

3DPW Dataset [230] was collected by hand-held cameras with IMUs in natural scenes capturing daily activities (e.g., shopping in the city, going up-stairs, doing sports, drinking coffee, and taking the bus). There are 60 video sequences (more than 51,000 frames) in this dataset and the corresponding 3D poses were computed by wearable IMUs. The test set contains 9,600 frames from HD cameras for 4

activities (Ultimatum, Mafia, Haggling, and Pizza). Dataset Link: <https://virtualhumans.mpi-inf.mpg.de/3DPW/>

MuCo-3DHP Dataset [197] is a multi-person 3D training set composed by the MPI-INF-3DHP single-person dataset with ground truth 3D pose from multi-view marker-less motion capture system. Background augmentation and shading-aware foreground augmentation of person appearance were applied to enable data diversity. Dataset Link: <http://gvv.mpi-inf.mpg.de/projects/SingleShotMultiPerson/>

MuPoTS-3D Dataset [197] is a multi-person 3D test set and its ground-truth 3D poses were captured by a multi-view marker-less MoCap system containing 20 real-world scenes (5 indoor and 15 outdoor). There are challenging samples with occlusions, drastic illumination changes, and lens flares in some of the outdoor footage. More than 8,000 frames were collected in the 20 sequences by 8 subjects. Dataset Link: <http://gvv.mpi-inf.mpg.de/projects/SingleShotMultiPerson/>

AMASS Dataset [176] was created by unifying 15 different optical marker-based MoCap datasets and using the SMPL model to represent human motion sequences. This large dataset contains more than 40 hours of motion data in 8,593 sequences of 9 million frames sampled at 60 Hz. More than 11,000 motions were recorded over 300 subjects. Dataset Link: <https://amass.is.tue.mpg.de/>

NBA2K Dataset [271] was extracted from the NBA2K19 video games by intercepting calls between the game engine and the graphic card using RenderDoc. The synthetic dataset contains 27,144 basketball poses spanning 27 subjects. The 3D poses of 35 keypoints and the corresponding RGB images are provided in this dataset with high quality. Dataset Link: <https://github.com/luyangzhu/NBA2K-dataset>

GTA-IM Dataset [243] is a GTA Indoor Motion dataset collected from Grand Theft Auto (GTA) video game by the GTA game engine. There are one million RGB-D frames of 1920×1080 resolution with ground-truth 3D human pose of 98 joints, covering various actions including sitting, walking, climbing, and opening the door. Each scene contains several settings such as living rooms, bedrooms and kitchens that emphasize human-scene interactions. Dataset Link: <https://people.eecs.berkeley.edu/~zhcao/hmp/>

Occlusion-Person Dataset [213] is a multi-view synthetic dataset with occlusion labels for the joints in images. UnrealCV [272] was used to render multi-view images and depth maps from 3D models. A total of 8 cameras were used every 45 degrees on a circle of two meters radius. This dataset contains 73K frames with 20.3% of the joints occluded. The ground truth 3D annotation and occlusion labels are also provided. Dataset Link: https://github.com/zhezh/occlusion_person

5.5 Evaluation Metrics for 3D HPE

MPJPE (Mean Per Joint Position Error) is the most widely used evaluation metric to assess the performance of 3D HPE. MPJPE is computed by using the Euclidean distance between the estimated 3D joints and the ground truth positions as follows:

$$MPJPE = \frac{1}{N} \sum_{i=1}^N \|J_i - J_i^*\|_2, \quad (1)$$

where N is the number of joints, J_i and J_i^* are the ground truth position and the estimated position of the i_{th} joint, respectively.

MPJPE, also called Reconstruction Error, is the MPJPE after rigid alignment by a post-processing between the estimated pose and the ground truth pose.

NMPJPE is defined as the MPJPE after normalizing the predicted positions in scale to the reference [205].

MPVE (Mean Per Vertex Error) [169] measures the Euclidean distances between the ground truth vertices and the predicted vertices as follows:

$$MPVE = \frac{1}{N} \sum_{i=1}^N \|V_i - V_i^*\|_2, \quad (2)$$

where N is the number of vertices, V is the ground truth vertices, and V^* is the estimated vertices.

3DPCK is a 3D extended version of the Percentage of Correct Keypoints (PCK) metric used in 2D HPE evaluation. An estimated joint is considered as correct if the distance between the estimation and the ground-truth is within a certain threshold. Generally the threshold is set to 150mm.

Summary. As pointed out by Ji et al. [273], low MPJPE does not always indicate an accurate pose estimation as it depends on the predicted scale of human shape and skeleton. Although 3DPCK is more robust to incorrect joints, it cannot evaluate the precision of correct joints. Also, existing metrics are designed to evaluate the precision of an estimated pose in a single frame. However, the temporal consistency and smoothness of reconstructed human pose cannot be examined over continuous frames by existing evaluation metrics. Designing frame-level evaluation metrics that can evaluate 3D HPE performance with temporal consistency and smoothness remains an open problem.

5.6 Performance Comparison of 3D HPE Methods

Tables 8 ~ 11 provide performance comparison of different 3D HPE methods on the widely used datasets corresponding to the single-view single-person, single-view multi-person, and multi-view scenarios.

In Table 8, most 3D single-view single-person HPE models successfully estimate 3D human pose on the Human3.6M dataset with remarkable precision. Although the Human3.6M dataset has a large size of training and testing data, it only contains 11 actors (6 male and 5 female) performing 17 activities such as eating, discussion, smoking, and taking photo. When estimating 3D pose on the in-the-wild data with more complex scenarios, the performance of these methods are degraded. It is also observed that model-based methods perform on par with model-free methods. Utilizing temporal information when video data is available can improve the performance. 3D single-view multi-person HPE is a harder task than 3D single-view single-person HPE due to more severe occlusion. As shown in Table 9 and Table 10, good progress has been made in single-view multi-person HPE methods in recent years. By comparing the results from Table 8 and Table 11, it is evident that the performance (e.g., MPJPE under Protocol 1) of multi-view 3D HPE methods has improved compared to single-view 3D HPE methods using the same dataset and evaluation metric.

TABLE 8: Comparison of different 3D single-view single-person HPE approaches on the Human3.6M dataset. The best two scores are marked in red and blue, respectively. Here in model-free approaches, “Direct” indicates the method directly estimating 3D pose without 2D pose representation. “Lifting” indicates the method lifting the 2D pose representation to the 3D space (i.e., 3D pose). “Temporal” means the method using temporal information.

Model-free methods			Protocol 1 MPJPE ↓																	Protocol 2 PMPJPE ↓	
Model	Year	Method	Dir.	Disc.	Eat.	Greet	Phone	Photo	Pose	Purch.	Sit	SitD.	Somme	Wait	WalkD.	Walk	WalkT.	Avg.	Avg.		
Direct	2017	[117]	67.4	71.9	66.7	69.1	72.0	77.0	65.0	68.3	83.7	96.5	71.7	65.8	74.9	59.1	63.2	71.9	51.9		
Direct	2018	[118]	48.5	54.4	54.4	52.0	59.4	65.3	49.9	52.9	65.8	71.1	56.6	52.9	60.9	44.7	47.8	56.2	41.8		
Lifting	2017	[121]	51.8	56.2	58.1	59.0	69.5	78.4	55.2	58.1	74.0	94.6	62.3	59.1	65.1	49.5	52.4	62.9	47.7		
Lifting	2017	[122]	54.2	61.4	60.2	61.2	79.4	78.3	63.1	81.6	70.1	107.3	69.3	70.3	74.3	51.8	63.2	69.7	-		
Lifting	2017	[109]	54.8	60.7	58.2	71.4	62.0	65.5	53.8	55.6	75.2	111.6	64.2	66.1	51.4	63.2	55.3	64.9	-		
Lifting	2018	[110]	51.5	58.9	50.4	57.0	62.1	65.4	49.8	52.7	69.2	85.2	57.4	58.4	43.6	60.1	47.7	58.6	37.7		
Lifting	2018	[125]	49.2	55.5	53.6	53.4	63.8	67.7	50.2	51.9	70.3	81.5	57.7	51.5	58.6	44.6	47.2	57.8	42.9		
Lifting	2019	[123]	34.4	42.4	36.6	42.1	38.2	39.8	34.7	40.2	45.6	60.8	39.0	42.6	42.0	29.8	31.7	39.9	27.9		
Lifting	2019	[134]	54.0	65.1	58.5	62.9	67.9	75.0	54.0	60.6	82.7	98.2	63.3	61.2	66.9	50.0	56.5	65.7	49.2		
Lifting	2019	[128]	43.8	48.6	49.1	49.8	57.6	61.5	45.9	48.3	62.0	73.4	54.8	50.6	56.0	43.4	45.5	52.7	42.6		
Lifting	2019	[127]	48.6	54.5	54.2	55.7	62.6	72.0	50.5	54.3	70.0	78.3	58.1	55.4	61.4	45.2	49.7	58.0	-		
Lifting	2019	[129]	46.8	52.3	44.7	50.4	52.9	68.9	49.6	46.4	60.2	78.9	51.2	50.0	54.8	40.4	43.3	52.7	42.2		
Lifting	2019	[130]	47.3	60.7	51.4	60.5	61.1	49.9	47.3	68.1	86.2	55.0	67.8	61.0	42.1	60.6	45.3	57.6	-		
Temporal	2018	[139]	44.8	50.4	44.7	49.0	52.9	61.4	43.5	45.5	63.1	87.3	51.7	48.5	52.2	37.6	41.9	52.1	36.3		
Temporal	2019	[140]	45.2	46.7	43.3	45.6	48.1	55.1	44.6	44.3	57.3	65.8	47.1	44.0	49.0	32.8	33.9	46.8	36.5		
Temporal	2019	[142]	44.6	47.4	45.6	48.8	50.8	59.0	47.2	43.9	57.9	61.9	49.7	46.6	51.3	37.1	39.4	48.8	39.0		
Temporal	2020	[147]	41.4	43.5	40.1	42.9	46.6	51.9	41.7	42.3	53.9	60.2	45.4	41.7	46.0	31.5	32.7	44.1	35.0		
Temporal	2020	[148]	41.8	44.8	41.1	44.9	47.4	54.1	43.4	42.2	56.2	63.6	45.3	43.5	45.3	31.3	32.2	45.1	35.6		
Model-based methods			Protocol 1 MPJPE ↓																	Protocol 2 PMPJPE ↓	
Model	Year	Method	Dir.	Disc.	Eat.	Greet	Phone	Photo	Pose	Purch.	Sit	SitD.	Somme	Wait	WalkD.	Walk	WalkT.	Avg.	Avg.		
SMPL	2016	[159]	62.0	60.2	67.8	76.5	92.1	77.0	73.0	75.3	100.3	137.3	83.4	77.3	79.7	86.8	81.7	69.3	-		
SMPL	2018	[170]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	59.9		
SMPL	2019	[171]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	50.1		
SMPL	2019	[164]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	41.1		
SMPL	2019	[113]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	77.8		
SMPL	2019	[175]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	65.6		
CylinderMan	2019	[182]	38.3	41.3	46.1	40.1	41.6	51.9	41.8	40.9	51.5	58.4	42.2	44.6	41.7	33.7	30.1	42.9	32.8		
Adam	2019	[183]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	58.3		
Kinematic	2018	[156]	43.8	51.7	48.8	53.1	52.2	74.9	52.7	44.6	56.9	74.3	56.7	66.4	47.5	68.4	45.6	55.8	46.2		
Kinematic	2019	[151]	44.7	48.9	47.0	49.0	56.4	67.7	48.7	47.0	63.0	78.1	51.1	50.1	54.5	40.1	43.0	52.6	40.7		
Kinematic	2020	[157]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	50.8		
Kinematic	2020	[153]	37.4	43.5	42.7	42.7	46.6	59.7	41.3	45.1	52.7	60.2	45.8	43.1	47.7	33.7	37.1	45.6	36.2		

TABLE 9: Comparison of different 3D single-view multi-person HPE approaches on the MuPoTS-3D dataset. The best two scores are marked in red and blue, respectively.

MuPoTS-3D				
	Year	Method	3DPCK ↑	
			All people	Matched people
Top down	2019	[189]	70.6	74.0
	2019	[191]	81.8	82.5
	2020	[166]	69.1	72.2
Bottom up	2018	[197]	65.0	69.8
	2019	[198]	70.4	-
	2020	[192]	72.0	-
	2020	[187]	73.5	80.5

TABLE 10: Comparison of different 3D single-view multi-person HPE approaches on the CMU Panoptic dataset. Ultimatum, Mafia, Haggling, and pizza denote four activities. The best two scores are marked in red and blue, respectively.

CMU Panoptic						
	Year	Method	MPJPE ↓			
			Haggle	Mafia	Ultim.	Pizza
Top down	2018	[190]	140.0	165.9	150.7	156.0
	2020	[166]	129.6	133.5	153.0	156.7
	2020	[193]	50.9	50.5	50.7	68.2
Bottom up	2018	[194]	72.4	78.8	66.8	94.3
	2020	[196]	45.0	95.0	58.0	79.0
	2020	[192]	40.6	37.6	31.3	55.8
	2020	[192]	40.6	37.6	31.3	55.8

5.7 Conference Workshops and Challenges for HPE

Due to the increasing interest in HPE, workshops and challenges on HPE are held in conjunction with computer vision conference venues like CVPR, ICCV and ECCV. These workshops aim at gathering researchers and practitioners work on HPE to discuss the current state-of-the-art as well as future directions to concentrate on. In Table 12, we summarize the relevant 2D and 3D HPE workshops and challenges in this research field from 2017 to 2020.

TABLE 11: Comparison of different 3D multi-view HPE approaches on the Human3.6M dataset. The best two scores are marked in red and blue, respectively.

Human3.6M					
Year	Method	Use extra 3D data	Protocol 1		Protocol 2
			MPJPE ↓	Normalized MPJPE ↓	PMPJPE ↓
2017	[199]	No	56.9	-	-
2018	[215]	Yes	57.9	-	-
2019	[201]	Yes	79.9	-	45.1
2019	[216]	No	60.6	60.0	47.5
2019	[200]	No	31.2	-	-
2019	[200]	Yes	26.2	-	-
2020	[219]	No	30.2	-	-
2020	[220]	No	29.3	-	-
2020	[213]	Yes	19.5	-	-

6 APPLICATIONS

In this section, we review related works of exploring HPE for a few popular applications (Fig. 7).

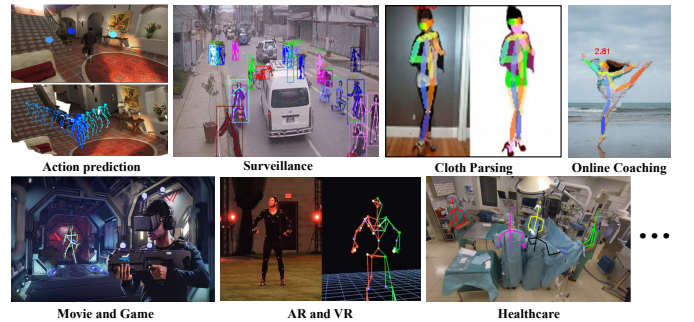


Fig. 7: Various applications of human pose estimation.

Action recognition, prediction, detection, and tracking: Pose information has been utilized as cues for various applications such as action recognition, prediction, detection,

TABLE 12: Conference workshops and challenges for 2D and 3D HPE.

ICCV 2017	PoseTrack Challenge: Human Pose Estimation and Tracking in the Wild	https://posetrack.net/workshops/iccv2017/
ECCV 2018	PoseTrack Challenge: Articulated People Tracking in the Wild	https://posetrack.net/workshops/eccv2018/
CVPR 2018	3D humans 2018: 1st International workshop on Human pose, motion, activities and shape	https://project.inria.fr/humans2018/#
CVPR 2019	3D humans 2019: 2nd International workshop on Human pose, motion, activities and shape	https://sites.google.com/view/humans3d/
CVPR 2019	Workshop On Augmented Human: Human-centric Understanding	https://vuhcs.github.io/vuhcs-2019/index.html
CVPR 2020	Towards Human-Centric Image/Video Synthesis	https://vuhcs.github.io/
ECCV 2020	3D poses in the wild challenge	https://virtualhumans.mpi-inf.mpg.de/3DPW_Challenge/
ACM Multimedia 2020	Large-scale Human-centric Video Analysis in Complex Events	http://humaninevents.org/

and tracking. Angelini et al. [274] proposed a real-time action recognition method using a pose-based algorithm. Yan et al. [275] leveraged the dynamic skeleton modality of pose for action recognition. Markovitz et al. [276] studied human pose graphs for anomaly detection of human actions in videos. Cao et al. [243] used the predicted 3D pose for long-term human motion prediction. Sun et al. [277] proposed a view-invariant probabilistic pose embedding for video alignment.

Pose-based video surveillance enjoys the advantage of preserving privacy by monitoring through pose and human mesh representation instead of human sensitive identities. Das et al. [278] embedded video with pose to identify activities of daily living for monitoring human behavior.

Action correction and online coaching: Some activities such as dancing, sporting, and professional training require precise human body control to strictly react as the standard pose. Normally personal trainers are responsible for the pose correction and action guidance in a face-to-face manner. With the help of 3D HPE and action detection, AI personal trainers can make coaching more convenient by simply setting up cameras without personal trainer presenting. Wang et al. [279] designed an AI coaching system with a pose estimation module for personalized athletic training assistance.

Clothes parsing: The e-commerce trends have brought about a noticeable impact on various aspects including clothes purchases. Clothing product in pictures can no longer satisfy customers’ demands, and customers hope to see the reliable appearance as they wear their selected clothes. Clothes parsing and pose transfer [280] make it possible by inferring the 3D appearance of a person wearing a specific clothes. HPE can provide plausible human body regions for cloth parsing. Moreover, the recommendation system can be upgraded by evaluating appropriateness based on the inferred reliable 3D appearance of customers with selected items. Patel et al. [281] achieved clothing prediction from 3D pose, shape and garment style.

Animation, movie, and gaming: Motion capture is the key component to present characters with complex movements and realistic physical interactions in industries of animation, movie, and gaming. The motion capture devices are usually expensive and complicated to set up. HPE can provide realistic pose information while alleviating the demand for professional high-cost equipment [282] [283].

AR and VR: Augmented Reality (AR) technology aims to enhance the interactive experience of digital objects into the real-world environment. The objective of Virtual Reality (VR) technology is to provide an immersive experience for the users. AR and VR devices use human pose information as input to achieve their goals of different applications. A cartoon character can be generated in real-world scenes to replace the real person. Weng et al. [284] created 3D character animation from single photo with the help of 3D

pose estimation and human mesh recovery. Zhang et al. [285] presented a pose-based system that converts broadcast tennis match videos into interactive and controllable video sprites. The players in the video sprites preserve the techniques and styles as real professional players.

Healthcare: HPE provides quantitative human motion information that physicians can diagnose some complex diseases, create rehabilitation training, and operate physical therapy. Lu et al. [286] designed a pose-based estimation system for assessing Parkinson’s disease motor severity. Gu et al. [287] developed a pose-based physical therapy system that patients can be evaluated and advised at home. Furthermore, such a system can be established to detect abnormal action and to predict the following actions ahead of time. Alerts are sent immediately if the system determines that danger may occur. Chen et al. [288] used the HPE algorithms for fall detection monitoring in order to provide immediate assistant. Also, HPE methods can provide reliable posture labels of patients in hospital environments to augment research on neural correlates to natural behaviors [289].

7 CONCLUSION AND FUTURE DIRECTIONS

In this survey, we have presented a systematic overview of recent deep learning-based 2D and 3D HPE methods. A comprehensive taxonomy and performance comparison of these methods have been covered. Despite great success, there are still many challenges as discussed in Sections 3.3 and 4.3. Here, we further point out a few promising future directions to promote advances in HPE research.

- **Domain adaptation for HPE.** For some applications such as estimating human pose from infant images [290] or artwork collections [291], there are not enough training data with ground truth annotations. Moreover, data for these applications exhibit different distributions from that of the standard pose datasets. HPE methods trained on existing standard datasets may not generalize well across different domains. The recent trend to alleviate the domain gap is utilizing GAN-based learning approaches. Nonetheless, how to effectively transfer the human pose knowledge to bridge domain gaps remains unaddressed.
- **Human body models** such as SMPL, SMPLify, SMPL-X, GHUM & GHUML, and Adam are used to model human mesh representation. However, these models have a huge number of parameters. How to reduce the number of parameters while preserving the reconstructed mesh quality is an intriguing problem. Also, different people have various deformations of body shape. A more effective human body model may utilize other information such as BMI [180] and silhouette [292] for better generalization.
- **Most existing methods ignore human interaction with 3D scenes.** There are strong human-scene relationship

constraints that can be explored such as a human subject cannot be simultaneously present in the locations of other objects in the scene. The physical constraints with semantic cues can provide reliable and realistic 3D HPE.

- 3D HPE is employed in visual tracking and analysis. Existing 3D human pose and shape reconstruction from videos are not smooth and continuous. One reason is that the evaluation metrics such as MPJPE cannot evaluate the smoothness and the degree of realisticness. Appropriate frame-level evaluation metrics focusing on temporal consistency and motion smoothness should be developed.
- Existing well-trained networks pay less attention to resolution mismatch. The training data of HPE networks are usually high resolution images or videos, which may lead to inaccurate estimation when predicting human pose from low resolution input. The contrastive learning scheme [293] (e.g., the original image and its low resolution version as a positive pair) might be helpful for building resolution-aware HPE networks.
- Deep neural networks in vision tasks are vulnerable to adversarial attacks. The imperceptible noise can significantly affect the performance of HPE. There are few works [294] [295] that consider adversarial attack for HPE. The study of defense against adversarial attacks can improve the robustness of HPE networks and facilitate real-world pose-based applications.
- Human body parts may have different movement patterns and shapes due to the heterogeneity of the human body. A single shared network architecture may not be optimal for estimating all body parts with various degrees of freedom. Neural Architecture Search (NAS) [296] can search the optimal architecture for estimating each body part [297], [298]. Also, NAS can be used for discovering efficient HPE network architectures to reduce the computational cost [299]. It is also worth exploring multi-objective NAS in HPE when multiple objectives (e.g, latency, accuracy and energy consumption) have to be met.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NeurIPS*, 2012.
- [2] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *CVPR*, 2015.
- [3] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE TPAMI*, 2016.
- [4] T. B. Moeslund and E. Granum, "A survey of computer vision-based human motion capture," *CVIU*, 2001.
- [5] T. B. Moeslund, A. Hilton, and V. Krüger, "A survey of advances in vision-based human motion capture and analysis," *CVIU*, 2006.
- [6] R. Poppe, "Vision-based human motion analysis: An overview," *CVIU*, 2007.
- [7] X. Ji and H. Liu, "Advances in view-invariant human motion analysis: a review," *IEEE TSMC*, 2009.
- [8] M. B. Holte, C. Tran, M. M. Trivedi, and T. B. Moeslund, "Human pose estimation and activity recognition from multi-view videos: Comparative explorations of recent developments," *IEEE Journal of Selected Topics in Signal Processing*, 2012.
- [9] Z. Liu, J. Zhu, J. Bu, and C. Chen, "A survey of human pose estimation: the body parts parsing based methods," *JVCIR*, 2015.
- [10] W. Gong, X. Zhang, J. González, A. Sobral, T. Bouwmans, C. Tu, and E.-h. Zahzah, "Human pose estimation from monocular images: A comprehensive survey," *Sensors*, 2016.
- [11] N. Sarafianos, B. Boteanu, B. Ionescu, and I. A. Kakadiaris, "3d human pose estimation: A review of the literature and analysis of covariates," *CVIU*, 2016.
- [12] Y. Chen, Y. Tian, and M. He, "Monocular human pose estimation: A survey of deep learning-based methods," *CVIU*, 2020.
- [13] T. L. Munea, Y. Z. Jembre, H. T. Weldegebrgel, L. Chen, C. Huang, and C. Yang, "The progress of human pose estimation: A survey and taxonomy of models applied in 2d human pose estimation," *IEEE Access*, 2020.
- [14] E. Mariniou, D. Papava, and C. Sminchisescu, "Pictorial human spaces: How well do humans perceive a 3d articulated pose?" in *ICCV*, 2013.
- [15] S. Zuffi, O. Freifeld, and M. J. Black, "From pictorial structures to deformable structures," in *CVPR*, 2012.
- [16] S. Johnson and M. Everingham, "Clustered pose and nonlinear appearance models for human pose estimation," in *BMVC*, 2010.
- [17] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *CVPR*, 2017.
- [18] X. Chen and A. L. Yuille, "Parsing occluded people by flexible compositions," in *CVPR*, 2015.
- [19] 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," *ACM TOG*, 2017.
- [20] S. X. Ju, M. J. Black, and Y. Yacoob, "Cardboard people: A parameterized model of articulated image motion," in *FG*, 1996.
- [21] H. Jiang, "Finding human poses in videos using concurrent matching and segmentation," in *ACCV*, 2010.
- [22] T. F. Cootes, C. J. Taylor, D. H. Cooper, and J. Graham, "Active shape models-their training and application," *CVIU*, 1995.
- [23] O. Freifeld, A. Weiss, S. Zuffi, and M. J. Black, "Contour people: A parameterized model of 2d articulated human shape," in *CVPR*, 2010.
- [24] R. Urtasun and P. Fua, "3d human body tracking using deterministic temporal motion models," in *ECCV*, 2004.
- [25] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, "SMPL: A skinned multi-person linear model," *ACM TOG*, 2015.
- [26] L. Kavan, "Part i: direct skinning methods and deformation primitives," in *ACM SIGGRAPH*, 2014.
- [27] G. Pons-Moll, J. Romero, N. Mahmood, and M. J. Black, "Dyna: A model of dynamic human shape in motion," *ACM TOG*, 2015.
- [28] S. Zuffi and M. J. Black, "The stitched puppet: A graphical model of 3D human shape and pose," in *CVPR*, 2015.
- [29] J. Pacheco, S. Zuffi, M. Black, and E. Sudderth, "Preserving modes and messages via diverse particle selection," in *ICML*, 2014.
- [30] H. Joo, T. Simon, and Y. Sheikh, "Total capture: A 3d deformation model for tracking faces, hands, and bodies," in *CVPR*, 2018.
- [31] C. Cao, Y. Weng, S. Zhou, Y. Tong, and K. Zhou, "Facewarehouse: A 3d facial expression database for visual computing," *IEEE TVCG*, 2013.
- [32] H. Xu, E. G. Bazavan, A. Zafir, W. T. Freeman, R. Sukthankar, and C. Sminchisescu, "Ghum & ghumi: Generative 3d human shape and articulated pose models," in *CVPR*, 2020.
- [33] M. Dantone, J. Gall, C. Leistner, and L. Van Gool, "Human pose estimation using body parts dependent joint regressors," in *CVPR*, 2013.
- [34] G. Gkioxari, B. Hariharan, R. Girshick, and J. Malik, "Using k-poselets for detecting people and localizing their keypoints," in *CVPR*, 2014.
- [35] A. S. Micilotta, E.-J. Ong, and R. Bowden, "Real-time upper body detection and 3d pose estimation in monoscopic images," in *ECCV*, 2006.
- [36] A. Toshev and C. Szegedy, "DeepPose: Human pose estimation via deep neural networks," in *CVPR*, 2014.
- [37] X. Chen and A. L. Yuille, "Articulated pose estimation by a graphical model with image dependent pairwise relations," in *NeurIPS*, 2014.
- [38] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *ECCV*, 2016.
- [39] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler, "Efficient object localization using convolutional networks," in *CVPR*, 2015.
- [40] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional pose machines," in *CVPR*, 2016.
- [41] T. Pfister, K. Simonyan, J. Charles, and A. Zisserman, "Deep convolutional neural networks for efficient pose estimation in gesture videos," in *ACCV*, 2014.

- [42] J. Carreira, P. Agrawal, K. Fragkiadaki, and J. Malik, "Human pose estimation with iterative error feedback," in *CVPR*, 2016.
- [43] X. Sun, J. Shang, S. Liang, and Y. Wei, "Compositional human pose regression," in *ICCV*, 2017.
- [44] D. C. Luvizon, H. Tabia, and D. Picard, "Human pose regression by combining indirect part detection and contextual information," *Computers & Graphics*, 2019.
- [45] A. Nibali, Z. He, S. Morgan, and L. Prendergast, "Numerical coordinate regression with convolutional neural networks," *arXiv preprint arXiv:1801.07372*, 2018.
- [46] S. Li, Z.-Q. Liu, and A. B. Chan, "Heterogeneous multi-task learning for human pose estimation with deep convolutional neural network," in *CVPR Workshops*, 2014.
- [47] X. Fan, K. Zheng, Y. Lin, and S. Wang, "Combining local appearance and holistic view: Dual-source deep neural networks for human pose estimation," in *CVPR*, 2015.
- [48] D. C. Luvizon, D. Picard, and H. Tabia, "2d/3d pose estimation and action recognition using multitask deep learning," in *CVPR*, 2018.
- [49] F. Zhang, X. Zhu, and M. Ye, "Fast human pose estimation," in *CVPR*, 2019.
- [50] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015.
- [51] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [52] S. Ruder, "An overview of multi-task learning in deep neural networks," *arXiv preprint arXiv:1706.05098*, 2017.
- [53] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler, "Joint training of a convolutional network and a graphical model for human pose estimation," in *NeurIPS*, 2014.
- [54] V. Ramakrishna, D. Munoz, M. Hebert, J. A. Bagnell, and Y. Sheikh, "Pose machines: Articulated pose estimation via inference machines," in *ECCV*, 2014.
- [55] I. Lifshitz, E. Fetaya, and S. Ullman, "Human pose estimation using deep consensus voting," in *ECCV*, 2016.
- [56] A. Bulat and G. Tzimiropoulos, "Human pose estimation via convolutional part heatmap regression," in *ECCV*, 2016.
- [57] G. Gkioxari, A. Toshev, and N. Jaitly, "Chained predictions using convolutional neural networks," in *ECCV*, 2016.
- [58] V. Belagiannis and A. Zisserman, "Recurrent human pose estimation," in *FG*, 2017.
- [59] W. Yang, S. Li, W. Ouyang, H. Li, and X. Wang, "Learning feature pyramids for human pose estimation," in *ICCV*, 2017.
- [60] Y. Luo, J. Ren, Z. Wang, W. Sun, J. Pan, J. Liu, J. Pang, and L. Lin, "Lstm pose machines," in *CVPR*, 2018.
- [61] B. Debnath, M. O'Brien, M. Yamaguchi, and A. Behera, "Adapting mobilenets for mobile based upper body pose estimation," in *AVSS*, 2018.
- [62] B. Xiao, H. Wu, and Y. Wei, "Simple baselines for human pose estimation and tracking," in *ECCV*, 2018.
- [63] H. Zhang, H. Ouyang, S. Liu, X. Qi, X. Shen, R. Yang, and J. Jia, "Human pose estimation with spatial contextual information," *arXiv preprint arXiv:1901.01760*, 2019.
- [64] B. Artacho and A. Savakis, "Unipose: Unified human pose estimation in single images and videos," in *CVPR*, 2020.
- [65] X. Chu, W. Yang, W. Ouyang, C. Ma, A. L. Yuille, and X. Wang, "Multi-context attention for human pose estimation," in *CVPR*, 2017.
- [66] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *NeurIPS*, 2014.
- [67] Y. Chen, C. Shen, X.-S. Wei, L. Liu, and J. Yang, "Adversarial posenet: A structure-aware convolutional network for human pose estimation," in *ICCV*, 2017.
- [68] C.-J. Chou, J.-T. Chien, and H.-T. Chen, "Self adversarial training for human pose estimation," in *APSIPA ASC*, 2018.
- [69] X. Peng, Z. Tang, F. Yang, R. S. Feris, and D. Metaxas, "Jointly optimize data augmentation and network training: Adversarial data augmentation in human pose estimation," in *CVPR*, 2018.
- [70] W. Yang, W. Ouyang, H. Li, and X. Wang, "End-to-end learning of deformable mixture of parts and deep convolutional neural networks for human pose estimation," in *CVPR*, 2016.
- [71] X. Chu, W. Ouyang, H. Li, and X. Wang, "Structured feature learning for pose estimation," in *CVPR*, 2016.
- [72] L. Ke, M.-C. Chang, H. Qi, and S. Lyu, "Multi-scale structure-aware network for human pose estimation," in *ECCV*, 2018.
- [73] W. Tang, P. Yu, and Y. Wu, "Deeply learned compositional models for human pose estimation," in *ECCV*, 2018.
- [74] W. Tang and Y. Wu, "Does learning specific features for related parts help human pose estimation?" in *CVPR*, 2019.
- [75] A. Jain, J. Tompson, Y. LeCun, and C. Bregler, "Modeep: A deep learning framework using motion features for human pose estimation," in *ACCV*, 2014.
- [76] T. Pfister, J. Charles, and A. Zisserman, "Flowing convnets for human pose estimation in videos," in *ICCV*, 2015.
- [77] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, 1997.
- [78] Y. Zhang, Y. Wang, O. Camps, and M. Sznai, "Key frame proposal network for efficient pose estimation in videos," *arXiv preprint arXiv:2007.15217*, 2020.
- [79] G. Papandreou, T. Zhu, N. Kanazawa, A. Toshev, J. Tompson, C. Bregler, and K. Murphy, "Towards accurate multi-person pose estimation in the wild," in *CVPR*, 2017.
- [80] S. Huang, M. Gong, and D. Tao, "A coarse-fine network for keypoint localization," in *ICCV*, 2017.
- [81] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *CVPR*, 2019.
- [82] W. Li, Z. Wang, B. Yin, Q. Peng, Y. Du, T. Xiao, G. Yu, H. Lu, Y. Wei, and J. Sun, "Rethinking on multi-stage networks for human pose estimation," *arXiv preprint arXiv:1901.00148*, 2019.
- [83] G. Moon, J. Y. Chang, and K. M. Lee, "Posefix: Model-agnostic general human pose refinement network," in *CVPR*, 2019.
- [84] J. Wang, X. Long, Y. Gao, E. Ding, and S. Wen, "Graph-pcnn: Two stage human pose estimation with graph pose refinement," *arXiv preprint arXiv:2007.10599*, 2020.
- [85] J. Huang, Z. Zhu, F. Guo, and G. Huang, "The devil is in the details: Delving into unbiased data processing for human pose estimation," in *CVPR*, 2020.
- [86] Y. Cai, Z. Wang, Z. Luo, B. Yin, A. Du, H. Wang, X. Zhou, E. Zhou, X. Zhang, and J. Sun, "Learning delicate local representations for multi-person pose estimation," *arXiv preprint arXiv:2003.04030*, 2020.
- [87] F. Zhang, X. Zhu, H. Dai, M. Ye, and C. Zhu, "Distribution-aware coordinate representation for human pose estimation," in *CVPR*, 2020.
- [88] U. Iqbal and J. Gall, "Multi-person pose estimation with local joint-to-person associations," in *ECCV*, 2016.
- [89] H.-S. Fang, S. Xie, Y.-W. Tai, and C. Lu, "Rmpe: Regional multi-person pose estimation," in *ICCV*, 2017.
- [90] Y. Chen, Z. Wang, Y. Peng, Z. Zhang, G. Yu, and J. Sun, "Cascaded pyramid network for multi-person pose estimation," in *CVPR*, 2018.
- [91] K. Su, D. Yu, Z. Xu, X. Geng, and C. Wang, "Multi-person pose estimation with enhanced channel-wise and spatial information," in *CVPR*, 2019.
- [92] L. Qiu, X. Zhang, Y. Li, G. Li, X. Wu, Z. Xiong, X. Han, and S. Cui, "Peeking into occluded joints: A novel framework for crowd pose estimation," *arXiv preprint arXiv:2003.10506*, 2020.
- [93] R. Umer, A. Doering, B. Leibe, and J. Gall, "Self-supervised keypoint correspondences for multi-person pose estimation and tracking in videos," *arXiv preprint arXiv:2004.12652*, 2020.
- [94] 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*, 2016.
- [95] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele, "Deepercut: A deeper, stronger, and faster multi-person pose estimation model," in *ECCV*, 2016.
- [96] E. Insafutdinov, M. Andriluka, L. Pishchulin, S. Tang, E. Levinkov, B. Andres, and B. Schiele, "Arttrack: Articulated multi-person tracking in the wild," in *CVPR*, 2017.
- [97] A. Newell, Z. Huang, and J. Deng, "Associative embedding: End-to-end learning for joint detection and grouping," in *NeurIPS*, 2017.
- [98] M. Fieraru, A. Khoreva, L. Pishchulin, and B. Schiele, "Learning to refine human pose estimation," in *CVPR Workshops*, 2018.
- [99] Z. Tian, H. Chen, and C. Shen, "Directpose: Direct end-to-end multi-person pose estimation," *arXiv preprint arXiv:1911.07451*, 2019.
- [100] S. Kreiss, L. Bertoni, and A. Alahi, "Pifpaf: Composite fields for human pose estimation," in *CVPR*, 2019.
- [101] X. Nie, J. Feng, J. Zhang, and S. Yan, "Single-stage multi-person pose machines," in *ICCV*, 2019.

- [102] S. Jin, W. Liu, E. Xie, W. Wang, C. Qian, W. Ouyang, and P. Luo, "Differentiable hierarchical graph grouping for multi-person pose estimation," *arXiv preprint arXiv:2007.11864*, 2020.
- [103] B. Cheng, B. Xiao, J. Wang, H. Shi, T. S. Huang, and L. Zhang, "Higherhrnet: Scale-aware representation learning for bottom-up human pose estimation," *arXiv preprint arXiv:1908.10357*, 2019.
- [104] X. Zhu, Y. Jiang, and Z. Luo, "Multi-person pose estimation for posetrack with enhanced part affinity fields," in *ICCV PoseTrack Workshop*, 2017.
- [105] G. Papandreou, T. Zhu, L.-C. Chen, S. Gidaris, J. Tompson, and K. Murphy, "Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model," in *ECCV*, 2018.
- [106] M. Kocabas, S. Karagoz, and E. Akbas, "Multiposenet: Fast multi-person pose estimation using pose residual network," in *ECCV*, 2018.
- [107] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv preprint arXiv:1904.07850*, 2019.
- [108] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*, 2014.
- [109] X. Zhou, Q. Huang, X. Sun, X. Xue, and Y. Wei, "Towards 3d human pose estimation in the wild: A weakly-supervised approach," in *ICCV*, 2017.
- [110] W. Yang, W. Ouyang, X. Wang, J. Ren, H. Li, and X. Wang, "3d human pose estimation in the wild by adversarial learning," in *CVPR*, 2018.
- [111] B. Wandt and B. Rosenhahn, "Repnet: Weakly supervised training of an adversarial reprojection network for 3d human pose estimation," in *CVPR*, 2019.
- [112] S. Li, L. Ke, K. Pratama, Y.-W. Tai, C.-K. Tang, and K.-T. Cheng, "Cascaded deep monocular 3d human pose estimation with evolutionary training data," in *CVPR*, 2020.
- [113] A. Arnab, C. Doersch, and A. Zisserman, "Exploiting temporal context for 3d human pose estimation in the wild," in *CVPR*, 2019.
- [114] S. Li and A. B. Chan, "3d human pose estimation from monocular images with deep convolutional neural network," in *ACCV*, 2014.
- [115] S. Li, W. Zhang, and A. B. Chan, "Maximum-margin structured learning with deep networks for 3d human pose estimation," in *ICCV*, 2015.
- [116] B. Tekin, I. Katircioglu, M. Salzmann, V. Lepetit, and P. Fua, "Structured prediction of 3d human pose with deep neural networks," in *BMVC*, 2016.
- [117] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis, "Coarse-to-fine volumetric prediction for single-image 3D human pose," in *CVPR*, 2017.
- [118] G. Pavlakos, X. Zhou, and K. Daniilidis, "Ordinal depth supervision for 3D human pose estimation," in *CVPR*, 2018.
- [119] G. Moon and K. M. Lee, "I2l-meshnet: Image-to-lixel prediction network for accurate 3d human pose and mesh estimation from a single rgb image," in *ECCV*, 2020.
- [120] C.-H. Chen and D. Ramanan, "3d human pose estimation = 2d pose estimation + matching," *CVPR*, 2017.
- [121] J. Martinez, R. Hossain, J. Romero, and J. J. Little, "A simple yet effective baseline for 3d human pose estimation," in *ICCV*, 2017.
- [122] B. Tekin, P. Márquez-Neila, M. Salzmann, and P. Fua, "Learning to fuse 2d and 3d image cues for monocular body pose estimation," in *ICCV*, 2017.
- [123] K. Zhou, X. Han, N. Jiang, K. Jia, and J. Lu, "Hemlets pose: Learning part-centric heatmap triplets for accurate 3d human pose estimation," in *ICCV*, 2019.
- [124] F. Moreno-Noguer, "3d human pose estimation from a single image via distance matrix regression," in *CVPR*, 2017.
- [125] M. Wang, X. Chen, W. Liu, C. Qian, L. Lin, and L. Ma, "Drpose3d: Depth ranking in 3d human pose estimation," in *IJCAI*, 2018.
- [126] E. Jahangiri and A. L. Yuille, "Generating multiple diverse hypotheses for human 3d pose consistent with 2d joint detections," in *ICCV Workshops*, 2017.
- [127] S. Sharma, P. T. Varigonda, P. Bindal, A. Sharma, and A. Jain, "Monocular 3d human pose estimation by generation and ordinal ranking," in *ICCV*, 2019.
- [128] C. Li and G. H. Lee, "Generating multiple hypotheses for 3d human pose estimation with mixture density network," in *CVPR*, 2019.
- [129] H. Ci, C. Wang, X. Ma, and Y. Wang, "Optimizing network structure for 3d human pose estimation," in *ICCV*, 2019.
- [130] L. Zhao, X. Peng, Y. Tian, M. Kapadia, and D. N. Metaxas, "Semantic graph convolutional networks for 3d human pose regression," in *CVPR*, 2019.
- [131] H. Choi, G. Moon, and K. M. Lee, "Pose2mesh: Graph convolutional network for 3d human pose and mesh recovery from a 2d human pose," in *ECCV*, 2020.
- [132] K. Liu, R. Ding, Z. Zou, L. Wang, and W. Tang, "A comprehensive study of weight sharing in graph networks for 3d human pose estimation," in *ECCV*, 2020.
- [133] A. Zeng, X. Sun, F. Huang, M. Liu, Q. Xu, and S. Lin, "Srnet: Improving generalization in 3d human pose estimation with a split-and-recombine approach," in *ECCV*, 2020.
- [134] I. Habibie, W. Xu, D. Mehta, G. Pons-Moll, and C. Theobalt, "In the wild human pose estimation using explicit 2d features and intermediate 3d representations," in *CVPR*, 2019.
- [135] C.-H. Chen, A. Tyagi, A. Agrawal, D. Drover, R. MV, S. Stojanov, and J. M. Rehg, "Unsupervised 3d pose estimation with geometric self-supervision," in *CVPR*, 2019.
- [136] D. Drover, C.-H. Chen, A. Agrawal, A. Tyagi, and C. Phuoc Huynh, "Can 3d pose be learned from 2d projections alone?" in *ECCV*, 2018.
- [137] X. Zhou, M. Zhu, K. Derpanis, and K. Daniilidis, "Sparseness meets deepness: 3d human pose estimation from monocular video," in *CVPR*, 2016.
- [138] X. Zhou, M. Zhu, G. Pavlakos, S. Leonardos, K. G. Derpanis, and K. Daniilidis, "Monocap: Monocular human motion capture using a cnn coupled with a geometric prior," *IEEE TPAMI*, 2018.
- [139] R. Dabral, A. Mundhada, U. Kusupati, S. Afaq, A. Sharma, and A. Jain, "Learning 3d human pose from structure and motion," in *ECCV*, 2018.
- [140] D. Pavlo, C. Feichtenhofer, D. Grangier, and M. Auli, "3d human pose estimation in video with temporal convolutions and semi-supervised training," in *CVPR*, 2019.
- [141] Y. Cheng, B. Yang, B. Wang, Y. Wending, and R. Tan, "Occlusion-aware networks for 3d human pose estimation in video," in *ICCV*, 2019.
- [142] Y. Cai, L. Ge, J. Liu, J. Cai, T. Cham, J. Yuan, and N. M. Thalmann, "Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks," in *ICCV*, 2019.
- [143] J. Wang, S. Yan, Y. Xiong, and D. Lin, "Motion guided 3d pose estimation from videos," in *ECCV*, 2020.
- [144] B. Tekin, A. Rozantsev, V. Lepetit, and P. Fua, "Direct prediction of 3d body poses from motion compensated sequences," in *CVPR*, 2016.
- [145] M. Rayat Imtiaz Hossain and J. J. Little, "Exploiting temporal information for 3d human pose estimation," in *ECCV*, 2018.
- [146] Z. Li, X. Wang, F. Wang, and P. Jiang, "On boosting single-frame 3d human pose estimation via monocular videos," in *ICCV*, 2019.
- [147] T. Chen, C. Fang, X. Shen, Y. Zhu, Z. Chen, and J. Luo, "Anatomy-aware 3d human pose estimation in videos," *arXiv preprint arXiv:2002.10322*, 2020.
- [148] R. Liu, J. Shen, H. Wang, C. Chen, S.-c. Cheung, and V. Asari, "Attention mechanism exploits temporal contexts: Real-time 3d human pose reconstruction," in *CVPR*, 2020.
- [149] X. Zhou, X. Sun, W. Zhang, S. Liang, and Y. Wei, "Deep kinematic pose regression," in *ECCV*, 2016.
- [150] B. X. Nie, P. Wei, and S. Zhu, "Monocular 3d human pose estimation by predicting depth on joints," in *ICCV*, 2017.
- [151] J. Wang, S. Huang, X. Wang, and D. Tao, "Not all parts are created equal: 3d pose estimation by modeling bi-directional dependencies of body parts," in *ICCV*, 2019.
- [152] J. N. Kundu, S. Seth, M. Rahul, M. Rakesh, V. B. Radhakrishnan, and A. Chakraborty, "Kinematic-structure-preserved representation for unsupervised 3d human pose estimation," in *AAAI*, 2020.
- [153] J. Xu, Z. Yu, B. Ni, J. Yang, X. Yang, and W. Zhang, "Deep kinematics analysis for monocular 3d human pose estimation," in *CVPR*, 2020.
- [154] Q. Nie, Z. Liu, and Y. Liu, "Unsupervised human 3d pose representation with viewpoint and pose disentanglement," in *ECCV*, 2020.
- [155] G. Georgakis, R. Li, S. Karanam, T. Chen, J. Kosecka, and Z. Wu, "Hierarchical kinematic human mesh recovery," in *ECCV*, 2020.
- [156] K. Lee, I. Lee, and S. Lee, "Propagating lstm: 3d pose estimation based on joint interdependency," in *ECCV*, 2018.

- [157] J. N. Kundu, S. Seth, V. Jampani, M. Rakesh, R. V. Babu, and A. Chakraborty, "Self-supervised 3d human pose estimation via part guided novel image synthesis," in *CVPR*, 2020.
- [158] A. Zanfir, E. G. Bazavan, H. Xu, B. Freeman, R. Sukthankar, and C. Sminchisescu, "Weakly supervised 3d human pose and shape reconstruction with normalizing flows," in *ECCV*, 2020.
- [159] 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 *ECCV*, 2016.
- [160] 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," in *CVPR*, 2017.
- [161] I. B. Vince Tan and R. Cipolla, "Indirect deep structured learning for 3d human body shape and pose prediction," in *BMVC*, 2017.
- [162] H.-Y. F. Tung, H.-W. Tung, E. Yumer, and K. Fragkiadaki, "Self-supervised learning of motion capture," in *NeurIPS*, 2017.
- [163] M. Hassan, V. Choutas, D. Tzionas, and M. J. Black, "Resolving 3D human pose ambiguities with 3D scene constraints," in *ICCV*, 2019.
- [164] N. Kolotouros, G. Pavlakos, M. Black, and K. Daniilidis, "Learning to reconstruct 3d human pose and shape via model-fitting in the loop," in *ICCV*, 2019.
- [165] C. Doersch and A. Zisserman, "Sim2real transfer learning for 3d human pose estimation: motion to the rescue," in *NeurIPS*, 2019.
- [166] W. Jiang, N. Kolotouros, G. Pavlakos, X. Zhou, and K. Daniilidis, "Coherent reconstruction of multiple humans from a single image," in *CVPR*, 2020.
- [167] X. Xu, H. Chen, F. Moreno-Noguer, L. A. Jeni, and F. De la Torre, "3d human shape and pose from a single low-resolution image with self-supervised learning," in *ECCV*, 2020.
- [168] T. Zhang, B. Huang, and Y. Wang, "Object-occluded human shape and pose estimation from a single color image," in *CVPR*, 2020.
- [169] G. Pavlakos, L. Zhu, X. Zhou, and K. Daniilidis, "Learning to estimate 3D human pose and shape from a single color image," in *CVPR*, 2018.
- [170] M. Omran, C. Lassner, G. Pons-Moll, P. V. Gehler, and B. Schiele, "Neural body fitting: Unifying deep learning and model-based human pose and shape estimation," in *3DV*, 2018.
- [171] N. Kolotouros, G. Pavlakos, and K. Daniilidis, "Convolutional mesh regression for single-image human shape reconstruction," in *CVPR*, 2019.
- [172] H. Zhu, X. Zuo, S. Wang, X. Cao, and R. Yang, "Detailed human shape estimation from a single image by hierarchical mesh deformation," in *CVPR*, 2019.
- [173] J. N. Kundu, M. Rakesh, V. Jampani, R. M. Venkatesh, and R. V. Babu, "Appearance consensus driven self-supervised human mesh recovery," in *ECCV*, 2020.
- [174] G. Varol, J. Romero, X. Martin, N. Mahmood, M. J. Black, I. Laptev, and C. Schmid, "Learning from synthetic humans," in *CVPR*, 2017.
- [175] M. Kocabas, N. Athanasiou, and M. J. Black, "Vibe: Video inference for human body pose and shape estimation," in *CVPR*, 2020.
- [176] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black, "Amass: Archive of motion capture as surface shapes," in *ICCV*, 2019.
- [177] X. Xu, H. Chen, F. Moreno-Noguer, L. A. Jeni, and F. D. la Torre, "3d human shape and pose from a single low-resolution image with self-supervised learning," in *ECCV*, 2020.
- [178] A. Kanazawa, M. J. Black, D. W. Jacobs, and J. Malik, "End-to-end recovery of human shape and pose," in *CVPR*, 2018.
- [179] G. Pavlakos, V. Choutas, N. Ghorbani, T. Bolkart, A. A. A. Osman, D. Tzionas, and M. J. Black, "Expressive body capture: 3d hands, face, and body from a single image," in *CVPR*, 2019.
- [180] A. A. A. Osman, T. Bolkart, and M. J. Black, "STAR: A spare trained articulated human body regressor," in *ECCV*, 2020.
- [181] A. Qammar and A. A. Argyros, "Mocapnet: Ensemble of snn encoders for 3d human pose estimation in rgb images," in *BMVC*, 2019.
- [182] Y. Cheng, B. Yang, B. Wang, W. Yan, and R. T. Tan, "Occlusion-aware networks for 3d human pose estimation in video," in *ICCV*, 2019.
- [183] D. Xiang, H. Joo, and Y. Sheikh, "Monocular total capture: Posing face, body, and hands in the wild," in *CVPR*, 2019.
- [184] S. Saito, T. Simon, J. Saragih, and H. Joo, "Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization," in *CVPR*, 2020.
- [185] H. Wang, R. A. Guler, I. Kokkinos, G. Papandreou, and S. Zafeiriou, "Blsm: A bone-level skinned model of the human mesh," in *ECCV*, 2020.
- [186] M. Fisch and R. Clark, "Orientation keypoints for 6d human pose estimation," *arXiv preprint arXiv:2009.04930*, 2020.
- [187] J. Zhen, Q. Fang, J. Sun, W. Liu, W. Jiang, H. Bao, and X. Zhou, "Smap: Single-shot multi-person absolute 3d pose estimation," in *ECCV*, 2020.
- [188] G. Rogez, P. Weinzaepfel, and C. Schmid, "Lcr-net: Localization-classification-regression for human pose," in *CVPR*, 2017.
- [189] G. Rogez, P. Weinzaepfel, and C. Schmid, "Lcr-net++: Multi-person 2d and 3d pose detection in natural images," *IEEE TPAMI*, 2019.
- [190] A. Zanfir, E. Marinoiu, and C. Sminchisescu, "Monocular 3d pose and shape estimation of multiple people in natural scenes: The importance of multiple scene constraints," in *CVPR*, 2018.
- [191] G. Moon, J. Chang, and K. M. Lee, "Camera distance-aware top-down approach for 3d multi-person pose estimation from a single rgb image," in *ICCV*, 2019.
- [192] A. Benzine, F. Chabot, B. Luvison, Q. C. Pham, and C. Achard, "Pandanet: Anchor-based single-shot multi-person 3d pose estimation," in *CVPR*, 2020.
- [193] J. Li, C. Wang, W. Liu, C. Qian, and C. Lu, "Hmor: Hierarchical multi-person ordinal relations for monocular multi-person 3d pose estimation," in *ECCV*, 2020.
- [194] A. Zanfir, E. Marinoiu, M. Zanfir, A.-I. Popa, and C. Sminchisescu, "Deep network for the integrated 3d sensing of multiple people in natural images," in *NeurIPS*, 2018.
- [195] J. N. Kundu, A. Revanur, G. V. Waghmare, R. M. Venkatesh, and R. V. Babu, "Unsupervised cross-modal alignment for multi-person 3d pose estimation," in *ECCV*, 2020.
- [196] M. Fabbri, F. Lanzi, S. Calderara, S. Alletto, and R. Cucchiara, "Compressed volumetric heatmaps for multi-person 3d pose estimation," in *CVPR*, 2020.
- [197] D. Mehta, O. Sotnychenko, F. Mueller, W. Xu, S. Sridhar, G. Pons-Moll, and C. Theobalt, "Single-shot multi-person 3d pose estimation from monocular rgb," in *3DV*, 2018.
- [198] D. Mehta, O. Sotnychenko, F. Mueller, W. Xu, M. Elgharib, P. Fua, H.-P. Seidel, H. Rhodin, G. Pons-Moll, and C. Theobalt, "Xnect: Real-time multi-person 3d motion capture with a single rgb camera," *ACM TOG*, 2020.
- [199] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis, "Harvesting multiple views for marker-less 3d human pose annotations," in *CVPR*, 2017.
- [200] H. Qiu, C. Wang, J. Wang, N. Wang, and W. Zeng, "Cross view fusion for 3d human pose estimation," in *ICCV*, 2019.
- [201] J. Liang and M. C. Lin, "Shape-aware human pose and shape reconstruction using multi-view images," in *ICCV*, 2019.
- [202] J. Dong, W. Jiang, Q. Huang, H. Bao, and X. Zhou, "Fast and robust multi-person 3d pose estimation from multiple views," in *CVPR*, 2019.
- [203] H. Tu, C. Wang, and W. Zeng, "Voxelpose: Towards multi-camera 3d human pose estimation in wild environment," in *ECCV*, 2020.
- [204] M. Burenus, J. Sullivan, and S. Carlsson, "3d pictorial structures for multiple view articulated pose estimation," in *CVPR*, 2013.
- [205] H. Rhodin, J. Spörri, I. Katircioglu, V. Constantin, F. Meyer, E. Müller, M. Salzmann, and P. Fua, "Learning monocular 3d human pose estimation from multi-view images," in *CVPR*, 2018.
- [206] H. Rhodin, M. Salzmann, and P. Fua, "Unsupervised geometry-aware representation for 3d human pose estimation," in *ECCV*, 2018.
- [207] X. Chen, K. Lin, W. Liu, C. Qian, and L. Lin, "Weakly-supervised discovery of geometry-aware representation for 3d human pose estimation," in *CVPR*, 2019.
- [208] H. Chen, P. Guo, P. Li, G. H. Lee, and G. Chirikjian, "Multi-person 3d pose estimation in crowded scenes based on multi-view geometry," in *ECCV*, 2020.
- [209] R. Mitra, N. B. Gundavarapu, A. Sharma, and A. Jain, "Multiview-consistent semi-supervised learning for 3d human pose estimation," in *CVPR*, 2020.
- [210] U. Iqbal, P. Molchanov, and J. Kautz, "Weakly-supervised 3d human pose learning via multi-view images in the wild," in *CVPR*, 2020.
- [211] Y. Zhang, L. An, T. Yu, X. Li, K. Li, and Y. Liu, "4d association graph for realtime multi-person motion capture using multiple video cameras," in *CVPR*, 2020.

- [212] C. Huang, S. Jiang, Y. Li, Z. Zhang, J. Traish, C. Deng, S. Ferguson, and R. Y. D. Xu, "End-to-end dynamic matching network for multi-view multi-person 3d pose estimation," in *ECCV*, 2020.
- [213] Z. Zhang, C. Wang, W. Qiu, W. Qin, and W. Zeng, "Adafuse: Adaptive multiview fusion for accurate human pose estimation in the wild," *IJCV*, 2020.
- [214] M. Habermann, W. Xu, M. Zollhoefer, G. Pons-Moll, and C. Theobalt, "Deepcap: Monocular human performance capture using weak supervision," in *CVPR*, 2020.
- [215] A. Kadkhodamohammadi and N. Padoy, "A generalizable approach for multi-view 3d human pose regression," *arXiv preprint arXiv:1804.10462*, 2018.
- [216] M. Kocabas, S. Karagoz, and E. Akbas, "Self-supervised learning of 3d human pose using multi-view geometry," in *CVPR*, 2019.
- [217] A. Pirinen, E. Gärtner, and C. Sminchisescu, "Domes to drones: Self-supervised active triangulation for 3d human pose reconstruction," in *NeurIPS*, 2019.
- [218] L. Chen, H. Ai, R. Chen, Z. Zhuang, and S. Liu, "Cross-view tracking for multi-human 3d pose estimation at over 100 fps," in *CVPR*, 2020.
- [219] E. Remelli, S. Han, S. Honari, P. Fua, and R. Wang, "Lightweight multi-view 3d pose estimation through camera-disentangled representation," in *CVPR*, 2020.
- [220] R. Xie, C. Wang, and Y. Wang, "Metafuse: A pre-trained fusion model for human pose estimation," in *CVPR*, 2020.
- [221] T. Yu, J. Zhao, Z. Zheng, K. Guo, Q. Dai, H. Li, G. Pons-Moll, and Y. Liu, "Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor," *IEEE TPAMI*, 2019.
- [222] F. Xiong, B. Zhang, Y. Xiao, Z. Cao, T. Yu, J. Zhou Tianyi, and J. Yuan, "A2j: Anchor-to-joint regression network for 3d articulated pose estimation from a single depth image," in *ICCV*, 2019.
- [223] A. Kadkhodamohammadi, A. Gangi, M. de Mathelin, and N. Padoy, "A multi-view rgb-d approach for human pose estimation in operating rooms," in *WACV*, 2017.
- [224] T. Zhi, C. Lassner, T. Tung, C. Stoll, S. G. Narasimhan, and M. Vo, "Texmesh: Reconstructing detailed human texture and geometry from rgb-d video," in *ECCV*, 2020.
- [225] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *CVPR*, 2017.
- [226] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in *NeurIPS*, 2017.
- [227] H. Jiang, J. Cai, and J. Zheng, "Skeleton-aware 3d human shape reconstruction from point clouds," in *ICCV*, 2019.
- [228] K. Wang, J. Xie, G. Zhang, L. Liu, and J. Yang, "Sequential 3d human pose and shape estimation from point clouds," in *CVPR*, 2020.
- [229] T. von Marcard, B. Rosenhahn, M. J. Black, and G. Pons-Moll, "Sparse inertial poser: Automatic 3d human pose estimation from sparse imus," in *Computer Graphics Forum*, 2017.
- [230] T. von Marcard, R. Henschel, M. J. Black, B. Rosenhahn, and G. Pons-Moll, "Recovering accurate 3d human pose in the wild using imus and a moving camera," in *ECCV*, 2018.
- [231] Y. Huang, M. Kaufmann, E. Aksan, M. J. Black, O. Hilliges, and G. Pons-Moll, "Deep inertial poser: Learning to reconstruct human pose from sparse inertial measurements in real time," *ACM TOG*, 2018.
- [232] Z. Zhang, C. Wang, W. Qin, and W. Zeng, "Fusing wearable imus with multi-view images for human pose estimation: A geometric approach," in *CVPR*, 2020.
- [233] F. Huang, A. Zeng, M. Liu, Q. Lai, and Q. Xu, "Deepfuse: An imu-aware network for real-time 3d human pose estimation from multi-view image," in *WACV*, 2020.
- [234] M. Zhao, T. Li, M. A. Alsheikh, Y. Tian, H. Zhao, A. Torralba, and D. Katabi, "Through-wall human pose estimation using radio signals," in *CVPR*, 2018.
- [235] M. Zhao, Y. Tian, H. Zhao, M. A. Alsheikh, T. Li, R. Hristov, Z. Kabelac, D. Katabi, and A. Torralba, "Rf-based 3d skeletons," in *2018 Conference of the ACM Special Interest Group on Data Communication*, 2018.
- [236] M. Zhao, Y. Liu, A. Raghu, T. Li, H. Zhao, A. Torralba, and D. Katabi, "Through-wall human mesh recovery using radio signals," in *ICCV*, 2019.
- [237] M. Isogawa, Y. Yuan, M. O'Toole, and K. M. Kitani, "Optical non-line-of-sight physics-based 3d human pose estimation," in *CVPR*, 2020.
- [238] D. Tome, P. Peluse, L. Agapito, and H. Badino, "xr-egopose: Egocentric 3d human pose from an hmd camera," in *ICCV*, 2019.
- [239] N. Saini, E. Price, R. Tallamraju, R. Enciclaud, R. Ludwig, I. Martinović, A. Ahmad, and M. Black, "Markerless outdoor human motion capture using multiple autonomous micro aerial vehicles," in *ICCV*, 2019.
- [240] H. M. Clever, Z. Erickson, A. Kapusta, G. Turk, K. Liu, and C. C. Kemp, "Bodies at rest: 3d human pose and shape estimation from a pressure image using synthetic data," in *CVPR*, 2020.
- [241] W. Zeng, W. Ouyang, P. Luo, W. Liu, and X. Wang, "3d human mesh regression with dense correspondence," in *CVPR*, 2020.
- [242] K. Zhou, B. L. Bhatnagar, and G. Pons-Moll, "Unsupervised shape and pose disentanglement for 3d meshes," *arXiv preprint arXiv:2007.11341*, 2020.
- [243] Z. Cao, H. Gao, K. Mangalam, Q. Cai, M. Vo, and J. Malik, "Long-term human motion prediction with scene context," in *ECCV*, 2020.
- [244] V. Ferrari, M. Marin-Jimenez, and A. Zisserman, "Progressive search space reduction for human pose estimation," in *CVPR*, 2008.
- [245] M. Eichner, V. Ferrari, and S. Zurich, "Better appearance models for pictorial structures," in *BMVC*, 2009.
- [246] M. Eichner and V. Ferrari, "We are family: Joint pose estimation of multiple persons," in *ECCV*, 2010.
- [247] B. Sapp, D. Weiss, and B. Taskar, "Parsing human motion with stretchable models," in *CVPR*, 2011.
- [248] M. Eichner and V. Ferrari, "Human pose co-estimation and applications," *IEEE TPAMI*, 2012.
- [249] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *IJCV*, 2010.
- [250] Y. Wang, D. Tran, and Z. Liao, "Learning hierarchical poselets for human parsing," in *CVPR*, 2011.
- [251] L.-J. Li and L. Fei-Fei, "What, where and who? classifying events by scene and object recognition," in *ICCV*, 2007.
- [252] B. Sapp and B. Taskar, "Modex: Multimodal decomposable models for human pose estimation," in *CVPR*, 2013.
- [253] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele, "2d human pose estimation: New benchmark and state of the art analysis," in *CVPR*, 2014.
- [254] S. Johnson and M. Everingham, "Learning effective human pose estimation from inaccurate annotation," in *CVPR*, 2011.
- [255] J. Wu, H. Zheng, B. Zhao, Y. Li, B. Yan, R. Liang, W. Wang, S. Zhou, G. Lin, Y. Fu *et al.*, "Ai challenger: A large-scale dataset for going deeper in image understanding," *arXiv preprint arXiv:1711.06475*, 2017.
- [256] J. Li, C. Wang, H. Zhu, Y. Mao, H.-S. Fang, and C. Lu, "Crowdpose: Efficient crowded scenes pose estimation and a new benchmark," in *ICCV*, 2019.
- [257] W. Zhang, M. Zhu, and K. G. Derpanis, "From actemes to action: A strongly-supervised representation for detailed action understanding," in *ICCV*, 2013.
- [258] H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black, "Towards understanding action recognition," in *ICCV*, 2013.
- [259] M. Andriluka, U. Iqbal, E. Ensafutdinov, L. Pishchulin, A. Milan, J. Gall, and S. B., "PoseTrack: A benchmark for human pose estimation and tracking," in *CVPR*, 2018.
- [260] L. Bourdev and J. Malik, "Poselets: Body part detectors trained using 3d human pose annotations," in *ICCV*, 2009.
- [261] B. E. Ainsworth, W. L. Haskell, S. D. Herrmann, N. Meckes, D. R. Bassett Jr, C. Tudor-Locke, J. L. Greer, J. Vezina, M. C. Whitt-Glover, and A. S. Leon, "2011 compendium of physical activities: a second update of codes and met values," *Medicine & Science in Sports & Exercise*, 2011.
- [262] S. Jin, L. Xu, J. Xu, C. Wang, W. Liu, C. Qian, W. Ouyang, and P. Luo, "Whole-body human pose estimation in the wild," *arXiv preprint arXiv:2007.11858*, 2020.
- [263] M. Eichner, M. Marin-Jimenez, A. Zisserman, and V. Ferrari, "2d articulated human pose estimation and retrieval in (almost) unconstrained still images," *IJCV*, 2012.
- [264] Y. Yang and D. Ramanan, "Articulated human detection with flexible mixtures of parts," *IEEE TPAMI*, 2012.
- [265] MSCOCO keypoint detection evaluation metric. (2016). [Online]. Available: <https://cocodataset.org/#keypoints-eval>

- [266] L. Sigal, A. Balan, and M. J. Black, "HumanEva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion," *IJCV*, 2010.
- [267] 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 TPAMI*, 2014.
- [268] H. Joo, T. Simon, X. Li, H. Liu, L. Tan, L. Gui, S. Banerjee, T. S. Godisart, B. Nabbe, I. Matthews, T. Kanade, S. Nobuhara, and Y. Sheikh, "Panoptic studio: A massively multiview system for social interaction capture," *IEEE TPAMI*, 2017.
- [269] 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.
- [270] M. Trumble, A. Gilbert, C. Malleson, A. Hilton, and J. Collomosse, "Total capture: 3d human pose estimation fusing video and inertial sensors," in *BMVC*, 2017.
- [271] L. Zhu, K. Rematas, B. Curless, S. Seitz, and I. Kemelmacher-Shlizerman, "Reconstructing nba players," in *ECCV*, 2020.
- [272] W. Qiu, F. Zhong, Y. Zhang, S. Qiao, Z. Xiao, T. S. Kim, and Y. Wang, "Unrealcv: Virtual worlds for computer vision," *ACM MM Open Source Software Competition*, 2017.
- [273] X. Ji, Q. FANG, J. DONG, Q. SHUAI, W. JIANG, and X. ZHOU, "A survey on monocular 3d human pose estimation," *Virtual Reality & Intelligent Hardware*, 2020.
- [274] F. Angelini, Z. Fu, Y. Long, L. Shao, and S. M. Naqvi, "Actionxpose: A novel 2d multi-view pose-based algorithm for real-time human action recognition," *arXiv preprint arXiv:1810.12126*, 2018.
- [275] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *AAAI*, 2018.
- [276] A. Markovitz, G. Sharir, I. Friedman, L. Zelnik-Manor, and S. Avidan, "Graph embedded pose clustering for anomaly detection," in *CVPR*, 2020.
- [277] J. J. Sun, J. Zhao, L.-C. Chen, F. Schroff, H. Adam, and T. Liu, "View-invariant probabilistic embedding for human pose," in *ECCV*, 2020.
- [278] S. Das, S. Sharma, R. Dai, F. Br mond, and M. Thonnat, "Vpn: Learning video-pose embedding for activities of daily living," in *ECCV*, 2020.
- [279] J. Wang, K. Qiu, H. Peng, J. Fu, and J. Zhu, "Ai coach: Deep human pose estimation and analysis for personalized athletic training assistance," in *ACM MM*, 2019.
- [280] Y. Li, C. Huang, and C. C. Loy, "Dense intrinsic appearance flow for human pose transfer," in *CVPR*, 2019.
- [281] C. Patel, Z. Liao, and G. Pons-Moll, "Tailornet: Predicting clothing in 3d as a function of human pose, shape and garment style," in *CVPR*, 2020.
- [282] N. S. Willett, H. V. Shin, Z. Jin, W. Li, and A. Finkelstein, "Pose2pose: pose selection and transfer for 2d character animation," in *IUI*, 2020.
- [283] J. Liu, H. Fu, and C.-L. Tai, "Posetween: Pose-driven tween animation," in *ACM Symposium on User Interface Software and Technology*, 2020.
- [284] C. Weng, B. Curless, and I. Kemelmacher-Shlizerman, "Photo wake-up: 3d character animation from a single photo," in *CVPR*, 2019.
- [285] H. Zhang, C. Scutto, M. Agrawala, and K. Fatahalian, "Vid2player: Controllable video sprites that behave and appear like professional tennis players," *arXiv preprint arXiv:2008.04524*, 2020.
- [286] M. Lu, K. Poston, A. Pfefferbaum, E. V. Sullivan, L. Fei-Fei, K. M. Pohl, J. C. Niebles, and E. Adeli, "Vision-based estimation of mds-updrs gait scores for assessing parkinson's disease motor severity," *arXiv preprint arXiv:2007.08920*, 2020.
- [287] Y. Gu, S. Pandit, E. Saraei, T. Nordahl, T. Ellis, and M. Betke, "Home-based physical therapy with an interactive computer vision system," in *ICCV Workshops*, 2019.
- [288] W. Chen, Z. Jiang, H. Guo, and X. Ni, "Fall detection based on key points of human-skeleton using openpose," *Symmetry*, 2020.
- [289] K. Chen, P. Gabriel, A. Alasfour, C. Gong, W. K. Doyle, O. Devinsky, D. Friedman, P. Dugan, L. Melloni, T. Thesen *et al.*, "Patient-specific pose estimation in clinical environments," *JTEHM*, 2018.
- [290] X. Huang, N. Fu, S. Liu, K. Vyas, A. Farnoosh, and S. Ostadabbas, "Invariant representation learning for infant pose estimation with small data," *arXiv preprint arXiv:2010.06100*, 2020.
- [291] P. Madhu, A. Villar-Corrales, R. Kosti, T. Bendschus, C. Reinhardt, P. Bell, A. Maier, and V. Christlein, "Enhancing human pose estimation in ancient vase paintings via perceptually-grounded style transfer learning," *arXiv preprint arXiv:2012.05616*, 2020.
- [292] Z. Li, A. Heyden, and M. Oskarsson, "A novel joint points and silhouette-based method to estimate 3d human pose and shape," *arXiv preprint arXiv:2012.06109*, 2020.
- [293] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," *arXiv preprint arXiv:2002.05709*, 2020.
- [294] J. Liu, N. Akhtar, and A. Mian, "Adversarial attack on skeleton-based human action recognition," *arXiv preprint arXiv:1909.06500*, 2019.
- [295] N. Jain, S. Shah, A. Kumar, and A. Jain, "On the robustness of human pose estimation," in *CVPR Workshops*, 2019.
- [296] T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," *JMLR*, 2019.
- [297] Z. Chen, Y. Huang, H. Yu, B. Xue, K. Han, Y. Guo, and L. Wang, "Towards part-aware monocular 3d human pose estimation: An architecture search approach," in *ECCV*, 2020.
- [298] X. Gong, W. Chen, Y. Jiang, Y. Yuan, X. Liu, Q. Zhang, Y. Li, and Z. Wang, "Autopose: Searching multi-scale branch aggregation for pose estimation," *arXiv preprint arXiv:2008.07018*, 2020.
- [299] W. Zhang, J. Fang, X. Wang, and W. Liu, "Efficientpose: Efficient human pose estimation with neural architecture search," *arXiv preprint arXiv:2012.07086*, 2020.