# Dynamic Graph Reasoning for Multi-person 3D Pose Estimation

Zhongwei Qiu*
University of Science and Technology Beijing
qiuzhongwei@xs.ustb.edu.cn

Qiansheng Yang*
Baidu
yangqiansheng@baidu.com

Jian Wang
Baidu
wangjian33@baidu.com

Dongmei Fu
University of Science and Technology Beijing
fdm_ustb@ustb.edu.cn

## ABSTRACT

Multi-person 3D pose estimation is a challenging task because of occlusion and depth ambiguity, especially in the cases of crowd scenes. To solve these problems, most existing methods explore modeling body context cues by enhancing feature representation with graph neural networks or adding structural constraints. However, these methods are not robust for their single-root formulation that decoding 3D poses from a root node with a pre-defined graph. In this paper, we propose GR-M3D, which models the **M**ulti-person **3D** pose estimation with dynamic **G**raph **R**easoning. The decoding graph in GR-M3D is predicted instead of pre-defined. In particular, It firstly generates several data maps and enhances them with a scale and depth aware refinement module (SDAR). Then multiple root keypoints and dense decoding paths for each person are estimated from these data maps. Based on them, dynamic decoding graphs are built by assigning path weights to the decoding paths, while the path weights are inferred from those enhanced data maps. And this process is named dynamic graph reasoning (DGR). Finally, the 3D poses are decoded according to dynamic decoding graphs for each detected person. GR-M3D can adjust the structure of the decoding graph implicitly by adopting soft path weights according to input data, which makes the decoding graphs be adaptive to different input persons to the best extent and more capable of handling occlusion and depth ambiguity than previous methods. We empirically show that the proposed bottom-up approach even outperforms top-down methods and achieves state-of-the-art results on three 3D pose datasets.

## CCS CONCEPTS

• **Computing methodologies** → **Object recognition**.

## KEYWORDS

Human Pose Estimation, Multi-person, Graph Reasoning

*Equal Contribution

## 1 INTRODUCTION

The goal of multi-person 3D human pose estimation is to estimate the 3D coordinates of human joints from multiple human bodies in a monocular RGB image, which is a challenging and fundamental task. Recently, 3D human pose estimation has drawn a lot of attention because of its broad applications, such as human behavior understanding [22, 48], human-object interaction detection [20], and athletic training assistance [44]. Although remarkable progresses [24, 28, 37, 41, 43] have been achieved in 3D pose estimation, the challenges of depth ambiguity and occlusions remain.

Existing multi-person 3D pose estimation methods mainly include top-down, bottom-up, and single-stage strategies. Top-down strategy [28, 37, 38, 43] firstly predicts the human bounding box and the absolute depth of root points in each person and then conducts single-person 3D pose estimation in the bounding box. Bottom-up strategy [26, 27, 51, 54] firstly estimates the 3D coordinates for each human joint in an image and then assigns them to different human instances. Single-stage strategy [31, 55] predicts keypoints offsets for each place in the image to generate 3D poses. The top-down methods are more accurate but more costly since the human detection and the repeated stages of extracting features for each person. The bottom-up and single-stage approaches are more efficient but uncompetitive in accuracy. The occlusions, non-uniform scales, and variable depths of each person are more difficult to handle when the input is a whole image, which are our interests.

For 3D human pose estimation, the common way to decode 3D pose is based on heatmaps [23, 28, 41, 43, 54], which decodes 3D coordinates from 3D heatmaps by an isolated structure as Figure 1 (a). To mitigate occlusion problem, [23, 54] model context information in coordinate level and feature level, respectively. Recent works [31, 55] bridge these two levels to some extent by estimating the 3D keypoints offsets, directly decoding keypoints coordinates from root point with 3D offsets. For example, based on single root point, CenterNet [55] and SPM [31] estimate other 3D keypoints by star structure propagation as Figure 1 (b), which suffers the long-distance transmission problem. SPM further improves the information propagation path by tree structure decoding, as shown in Figure 1 (c). The tree structure brings accumulation error since the

(a) Isolated     (b) Star Structure     (c) Tree Structure

● Root Node   ◉ Target Node
→ Message Propagation Path       (d) Dynamic Graph Reasoning

**Figure 1: (a) Heatmap-based approaches [28, 41] decode human poses by locating the point with maximum confidence on 3D heatmaps. Each node is isolated when decoding. (b) Star graph approaches [31, 55] start from single-root node, broadcasting information to other nodes by star graph, which suffers from long-distance information transmission problem. (c) Tree graph approach [31] starts from single-root node, hierarchically transferring the message from the parent node to the child node, which suffers from the occlusion problem and cumulative error. (d) Our dynamic graph reasoning (DGR) extracts dense decoding paths from multiple root keypoints, further reasons the best dynamic decoding graphs for each person in the image from the predicted dense decoding paths. The final 3D poses are decoded with the guidance of reasoned dynamic decoding graphs.**

hierarchical decoding scheme, which is serious when the middle nodes are occluded.

The single-root decoding strategy, such as star and tree structure decoding, is unreliable since the depth ambiguity. Firstly, it is hard to calculate the localization of a new joint when its root joint is occluded. An incorrect starting point will result in a wrong whole 3D pose prediction. Secondly, the predicted depth from single root points is not robust since lacking the global context of the human body. Thirdly, the decoding graph is fixed for each person, which may fail to build the message propagation path when the middle nodes are occluded.

Based on the above analysis, we propose a dynamic graph reasoning method (DGR). As shown in Figure 1 (d), each joint in a person serves as a root keypoint to generate dense decoding paths. DGR reasons the dynamic decoding graph for each person from these dense decoding paths. Each message propagation path in the dynamic decoding graph is calculated by combining starting point, target point, and message path score, which is adaptively adjusted according to the occlusions in the input image. This mechanism enables the best decoding paths for each target joint. Although some works [23, 52, 53] use graph convolutional network or local human structure to learn better structural features, the learned body graph is also fixed. In the inference phase, this learned structural graph

can't be adjusted adaptively with the occlusion of the input image. But our DGR can self-adaptively adjust the decoding graph due to the multi-root decoding and dynamic graph reasoning mechanism.

To better predict the depth and build a reliable message propagation graph, we further propose a scale and depth aware refinement (SDAR) module. Inspired by the intuition of near big far small, which is a basic perspective principle, SADR concatenates multiple initial scale features and depth features and then generates refined scale and depth features. The refined scale-aware feature and depth-aware feature are beneficial to generating good root keypoints and building robust pose graphs.

Our main contributions can be summarized as follows:

- We argue that decoding 3D poses from a single root point is unreliable and propose a novel dynamic graph reasoning (DGR) method to decode multi-person 3D poses.
- We propose SADR to generate better root keypoints and reliable decoding graph representations by aggregating global depth and scale information.
- The proposed bottom-up approach outperforms even top-down methods and achieves new state-of-the-art results on three widely-used benchmarks: Human3.6M, MuPoTS-3D, and CMU Panoptic datasets.

## 2 RELATED WORK

### 2.1 2D Human Pose Estimation

Mainstream multi-person 2D pose estimation approaches include top-down and bottom-up methods. Top-down approaches [4, 39, 47] firstly conduct human detection. Then, they crop images and perform single-person human pose estimation for each human instance. Xiao *et al.* [47] propose a simple baseline for 2D pose estimation, which uses ResNet as the backbone and follows several up-sample layers to generate heatmaps. Sun *et al.* [39] propose HRNet to generate high-resolution representation. Bottom-up approaches [3, 17, 18] estimate the keypoints for all human instances in an image and then group them into multiple instances. Cao *et al.* [3] propose part affinity fields to group keypoints. Kocabas *et al.* [17] propose Multiposenet, a framework to finish detection, pose estimation, and grouping at the same time. The accuracy of bottom-up approaches is lower than top-down approaches since the different scales and low resolution of persons.

Some approaches [31, 42, 46, 55] discard the method of locating from heatmaps, but learn 2D offset to decode pose. Nie *et al.* [31] propose SPM to predict root points of the human body and offsets of each keypoints. The coordinates of keypoints can be obtained from the root points and offsets. DirectPose [42] and Point-set anchors [46] use deformable convolution to align the features of pose. These offset-based methods provide insight for 3D pose estimation.

### 2.2 3D Human Pose Estimation

For single person cases, the main-stream methods follow the architecture of top-down 2D pose estimation methods. They change the 2D pose regression heads to 3D heads, including *inferring* and *lifting* methods. The *inferring* methods [16, 25, 33, 40, 41] directly regress 3D pose from learned 3D heatmaps. For example, Sun *et al.* [41] conduct integral operation on 3D heatmaps to obtain 3D pose coordinates. The *lifting* methods [24, 29, 49, 52] firstly estimate

2D pose by 2D pose estimator, and then lift 2D pose to 3D pose by a sample neural network, such as SRNet [52].

Recently, some works [28, 37, 38, 43, 51] study multi-person 3D pose estimation. Rogez *et al.* [37, 38] locate the human bounding box and generate a set of anchor poses for each human. The anchor poses are refined to the final pose by a regression module. Moon *et al.* [28] propose a top-down pipeline as multi-person 2D pose estimation. They locate the human bounding box and the depth of root keypoints in each bounding box and then conduct single-person 3D pose estimation in the bounding box. Zanfir *et al.* [51] propose MubyNet, which estimates keypoints and limb core at the same time and then integrates limb score to group keypoints into different persons. Wang *et al.* [43] propose hierarchical multi-person ordinal relations as an additional loss for depth learning. Cheng *et al.* [5] integrate top-down and bottom-up networks to relieve the problems of occlusion and close interactions.

### 2.3 Graph Reasoning in Pose Estimation

Since human joints can be naturally seemed as a graph structure, many works [2, 6, 23, 31, 35, 36, 53] try to utilize this information. Zhao *et al.* [53] take the human joints as the nodes of graph and then build a semantic graph convolutional network to learn joints relation. More works try to learn a better graph representation for human pose, such as spatial-temporal graph convolutional network [2], dynamic graph convolutional network [36], context pose [23] and so on. Other works [6, 52] study the local or global structure of the human body. However, these methods just learn graph representation for the human pose, the learned human graph is fixed in the inference phase. In this paper, different from previous works, we propose a graph reasoning approach, which can self-adaptively adjust the graph in the inference phase.

## 3 APPROACH

Let $p$ denotes the coordinate of a target joint, then the offset-based methods [31, 55] can be formulated as:

$$p = p^r + \Delta p \tag{1}$$

where $p^r$ is the coordinate of a root joint, such as body center or the parent joint of $p$ based on the human skeleton. $\Delta p$ is a learned offset from the root joint to the target joint.

However, predicting 3D pose from single root is not robust since depth ambiguities and occlusions. Thus, we propose a new graph-based method, which takes full advantage of the context, scale, and depth information of multiple roots, and makes the 3D pose prediction more robust for handing depth ambiguity and occlusion.

### 3.1 Framework

The overview of our dynamic graph reasoning is illustrated in Figure 2. First of all, four maps and deep feature map $\mathbf{F}$ are extracted from the backbone network. These maps are denoted as $\mathbf{M}_h^I$, $\mathbf{M}_d^I$, $\mathbf{M}_s^I$ and $\mathbf{M}_o^I$ of shapes $(K+1, H, W)$, $(K, H, W)$, $(2, H, W)$ and $(3K, H, W)$ respectively. $K$ is the number of joint categories, $H$ and $W$ are the height and width of these maps. $\mathbf{M}_h^I$ is the heat map for $K$ joints and one body center (the midpoint of left hip and right hip). $\mathbf{M}_d^I$ is the depth map for $K$ joints. $\mathbf{M}_s^I$ is a map where each pixel preserves the 2D offset from itself to the corresponding body center.

And it can be regarded as scale map since the 2D offset indicates the scale of a person. $\mathbf{M}_o^I$ is a 3D offset map, where each pixel records the 3D offsets from its 3D location to the corresponding $K$ joints.

Based on those maps, our method outputs multi-person 3D pose by carrying out the following three parts successively, which are Scale and Depth Aware Refinement (SDAR), Multi-person Root Keypoints Decoding (MRKD), and Dynamic Graph Reasoning (DGR). Firstly, SDAR refines the four maps mentioned above by integrating scale and depth information, and the refined maps are denoted as $\mathbf{M}_h$, $\mathbf{M}_d$, $\mathbf{M}_s$, $\mathbf{M}_o$. Secondly, MRKD decodes multi-person root keypoints from those refined maps. $\mathbf{M}_h$, $\mathbf{M}_d$, $\mathbf{M}_s$ are used to decode 3D root keypoints and $\mathbf{M}_o$ is used to decode dense decoding paths. Finally, based on the root keypoints and predicted dense decoding paths, DGR reasons the dynamic decoding graphs for each person in the input image, to further decode robust multi-person 3D poses as equation 1 with the guidance of dynamic decoding graphs. The details of the three parts and the training process are introduced in the following.

### 3.2 Scale and Depth Aware Refinement

As demonstrated in Figure 3, SDAR firstly computes refined scale map and depth map by the following operations:

$$\mathbf{M}_s^R = Conv_s(Concat(\mathbf{M}_h^I, \mathbf{M}_s^I)) \tag{2}$$

$$\mathbf{M}_d^R = Conv_d(Concat(\mathbf{M}_o^I, \mathbf{M}_d^I)) \tag{3}$$

where $Conv_s$ and $Conv_d$ are sequential modules including convolution and normalization layers, and $Concat$ is an operator for concatenating tensors in channel dimension. Equation 2 merges heat map into scale map, which makes $\mathbf{M}_s$ can provide both scale information and attention cue for downstream operations. In Equation 3, depth map is refined to capture more global depth context from $M_o^I$.

Once got the refined scale map and depth map, we multiply them with the feature map $\mathbf{F}$ respectively, and then sum up these products as $\mathbf{F}_{sd}$. Based on this enhanced feature, the refined heat map and 3D offset map can be predicted by the following calculation:

$$\mathbf{M}_h^R, \mathbf{M}_o^R = Split(Conv_f(\mathbf{F}_{sd}))$$
$$\mathbf{F}_{sd} = \mathbf{M}_s \odot \mathbf{F} + \mathbf{M}_d \odot \mathbf{F} \tag{4}$$

where $Conv_f$ is also a sequential convolution module, $Split$ is an operator for splitting tensor in channel dimension, and $\odot$ means element-wise multiplication.

The final outputs are:

$$\mathbf{M}_* = \mathbf{M}_*^I + \mathbf{M}_*^R \tag{5}$$

where $*$ is wildcard character for $\{h, s, d, o\}$.

In the following, $\mathbf{M}_h$, $\mathbf{M}_s$ and $\mathbf{M}_d$ will serve for multi-person root keypoints decoding, while $\mathbf{M}_h$ and $\mathbf{M}_o$ will contribute for dynamic graph reasoning. After the refining by SDAR as Equation 4, the final multi-person pose predicted by DGR could be more precise since the heatmap and 3D offset map had been enhanced by scale and depth references.

**Figure 2: The framework of our approach includes: 1) Given an image, backbone network processes deep feature ($F$) and four maps ($M_h^I, M_s^I, M_d^I, M_o^I$). 2) Scale and Depth Aware Refinement (SDAR) module aggregates scale and depth information to enhance feature $F$ and refine the four maps with the scale and depth context. 3) Multiple root keypoints are decoded from $M_h, M_s, M_d$ by Multi-person Root Keypoints Decoding (MRKD). 4) Based on the multi-root keypoints, we conduct dynamic graph reasoning on predicted dense paths to generate dynamic decoding graphs for each person, to further decode final 3D poses.**

## 3.3 Multi-person Root Keypoints Decoding

Before carrying out graph reasoning, we decode multi-person root 3D keypoints out of the maps outputted from SDAR. This process is called multi-person root keypoints decoding (MRKD). At first, $N \times K$ independent keypoints and $N$ body centers can be detected from $\mathbf{M}_h$ by extracting local maximums, where $N$ represents the number of persons and is obtained by finding $N$ high-confidence points according to a confidence threshold. Meanwhile, the depth of each keypoint is the value of $\mathbf{M}_d$ at the corresponding 2D location. And then, we assign these keypoints to $N$ persons. Let $[p_1, ..., p_m, ..., p_M]$ as the $M = N \times K$ predicted keypoints and $[c_1, ..., c_n, ..., c_N]$ as the $N$ predicted body centers from $\mathbf{M}_h$, where $p_m$ and $c_n$ are both 2D coordinates. Then each $p_m$ is corresponding to a unique $c_n$ by solving a distance matrix $E^{M \times N}$ with hungarian algorithm, while the definition of each element in $E$ is:

$$E[m, n] = \mathcal{E}(\widetilde{c}_m, c_n)$$
$$\widetilde{c}_m = \mathbf{M}_s|_{p_m} + p_m \tag{6}$$

where $\mathcal{E}$ represents calculating euclidean distance, and $\widetilde{c}_m$ is a regressed 2D body center coordinate from $p_m$ depending on the semantic information at $p_m$, while $\mathbf{M}_s|_{p_m}$ means the 2D offset from $p_m$ to this center, which can be assigned as the 2D vector at point $p_m$ on $\mathbf{M}_s$.

## 3.4 Dynamic Graph Reasoning

After obtaining root keypoints, the decoding graphs are inferred by conducting dynamic graph reasoning (DGR) for the final 3D pose estimation. Supposing that a person can be regarded as a undirected acyclic graph, denoted as $G(\mathbb{P}, \mathbb{E})$, where $\mathbb{P} = \{p^i, i \in [1, K]\}$ is the set of joints for this person, $p^i$ is the 2D coordinate of the $i^{th}$ joint. Here $\mathbb{P}$ is initialized with the root joints detected from $M_h$ and $M_s$. While, $\mathbb{E} = \{e^{ij}, i \in [1, K], j \in [1, K]\}$ is the set of dense decoding paths and $e^{ij}$ means the decoding path from the $i^{th}$ joint to the $j^{th}$ joint. And it is valued by the 3D offset which is expressed as:

$$e^{ij} = \mathbf{M}_o|_{p^i}^j \tag{7}$$

where $\mathbf{M}_o|_{p^i}^j$ means the 3D offset to the $j^{th}$ target joint from $\mathbf{M}_o$ at the position of root joint $p^i$.

The goal of DGR is to reason the best decoding paths $\hat{\mathbb{E}}$ from dense decoding paths set $\mathbb{E}$, to further construct decoding graph $G(\mathbb{P}, \hat{\mathbb{E}})$ for each person in the input image. Intuitively, we can directly pick or drop candidate paths for generating $\hat{\mathbb{E}}$. However, this hard selecting manner may not be optimal. Here, we adopt a soft manner that assigning a weight on each candidate path for reasoning the decoding graph. And the weights for all paths are inferred from heat map $\mathbf{M}_h$ and offset map $\mathbf{M}_o$. While the weighted

decoding paths can be expressed as

$$\hat{\mathbb{E}} = \{(e^{ij}, \mathcal{W}(p^i, p^j)), i \in [1, K], j \in [1, K]\} \tag{8}$$

where $\mathcal{W}(p^i, p^j)$ means the corresponding path weight for $e^{ij}$, which can be calculated as:

$$\mathcal{W}(p^i, p^j) = \mathbf{M}_h|_{p^i}^i \mathcal{R}(p^i, p^j)\mathbf{M}_h|_{p^j}^j \tag{9}$$

where $\mathbf{M}_h|_{p^i}^i$ and $\mathbf{M}_h|_{p^j}^j$ serve as confidence scores which are the heat values of the $i^{th}$ joint at point $p^i$ on $\mathbf{M}_h$ and the $j^{th}$ joint at point $p^j$. While $\mathcal{R}(p^i, p^j)$ is a bone confidence formulated as:

$$R(p^i, p^j) = \exp(-(\frac{||\mathbf{M}_o|_{p^i}^j||_2}{||\mathbf{M}_o|_{p^h}^c||_2} + \gamma(i, j))) \tag{10}$$

where $h$ and $c$ mean the joint index of head-top and mid-hip, respectively. $\mathbf{M}_o|_{p^h}^c$ is the 3D offset from the head-top to mid-hip joint. $\gamma(i, j) = ||\frac{\sigma(i,j)}{\sigma(h,c)}||_2$ is a priori ratio, and $\sigma(i, j)$ is the average bone length between the $i^{th}$ and the $j^{th}$ joint, counted from the training dataset, while $\sigma(h, c)$ is the average bone length between the $h^{th}$ joint and the $c^{th}$ joint. Thus, $R(p^i, p^j)$ indicates the confidence of predicted edge in the human structural graph, which appears as a production of instance-level propagation confidence and statistical priori propagation confidence.

Given a pair of predicted root joint set $\mathbb{P}$ and decoding graph $G(\mathbb{P}, \hat{\mathbb{E}})$, the final 3D pose of a certain person can be decoded. Concretely, we firstly extend $p^i$ in $\mathbb{P}$ as $p_{3d}^i$ by concatenating $p^i$ with its corresponding depth $d^i$ predicted from $M_d$. And the extended root joint set is denoted as $\mathbb{P}_{3d}$. Then we decode $G(\mathbb{P}_{3d}, \hat{\mathbb{E}})$ and update the 3D coordinates of each joint as:

$$\hat{p}_{3d}^j = \frac{\sum_i^K \mathcal{W}(p^i, p^j)(p_{3d}^i + e^{ij})}{\sum_i^K \mathcal{W}(p^i, p^j)}, \quad j = [1, K] \tag{11}$$
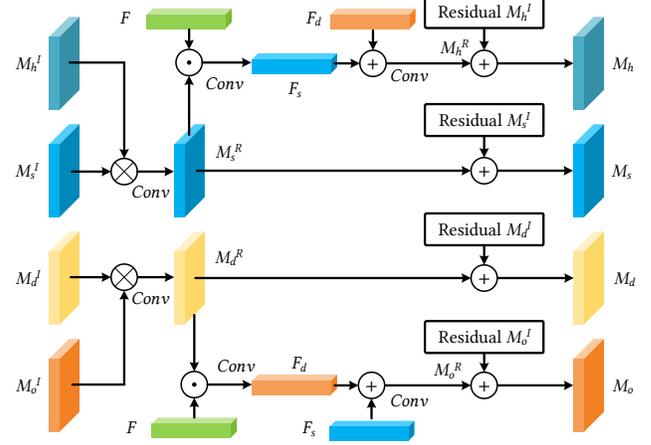
where $\hat{p}_{3d}^j$ represents the estimated 3D coordinates of the $j^{th}$ target joint.

There are two reasons that the DGR can predict better 3D poses: 1) Compared with the single-root decoding mechanism, the depth value predicted by the multi-root decoding mechanism are more robust. 2) For each person, DGR builds a dynamic decoding graph, which can be adjusted by the path weights according to the deep feature of input image adaptively, even in the inference phase. As stated in Equation 9, the graph path confidence is determined by start point confidence, edge confidence, and endpoint confidence, which perceives the occlusion of the inputs image. This makes the 3D pose decoding more robust to occlusion conditions.

## 3.5 Loss Function

For training the network, we impose supervision both on the input and output maps of SDAR. The total loss is :

$$\mathcal{L} = \mathcal{L}_h + \alpha \mathcal{L}_s + \beta \mathcal{L}_d + \mathcal{L}_o, \tag{12}$$



Figure 3: The architecture of scale and depth refinement (SDAR). SDAR aggregates scale and depth information to learn scale-aware and depth-aware maps, further to refine keypoints heatmaps and 3D offsets. ⊗ denotes concatenate operation, ⊕ denotes element-wise addition, ⊙ denotes element-wise multiplication.

where $\alpha = 0.1$ and $\beta = 0.1$ are loss weights. and

$$\begin{aligned}
\mathcal{L}_h &= \mathcal{L}_{MSE}(\mathbf{M}_h^I, \overline{\mathbf{M}}_h) + \mathcal{L}_{MSE}(\mathbf{M}_h, \overline{\mathbf{M}}_h) \\
\mathcal{L}_s &= \mathcal{L}_{L1}(\mathbf{M}_s^I, \overline{\mathbf{M}}_s) + \mathcal{L}_{L1}(\mathbf{M}_s, \overline{\mathbf{M}}_s) \\
\mathcal{L}_d &= \mathcal{L}_{L1}(\delta(\mathbf{M}_d^I), \overline{\mathbf{M}}_d) + \mathcal{L}_{L1}(\delta(\mathbf{M}_d), \overline{\mathbf{M}}_d) \\
\mathcal{L}_o &= \mathcal{L}_{L1}(\mathbf{M}_o^I, \overline{\mathbf{M}}_o) + \mathcal{L}_{L1}(\mathbf{M}_o, \overline{\mathbf{M}}_o)
\end{aligned} \tag{13}$$

where $\overline{\mathbf{M}}_*$ represents ground truth map. $\mathcal{L}_{MSE}$ is standard MSE loss. $\mathcal{L}_{L1}$ is L1 loss, and only pixels around body joints on $\overline{\mathbf{M}}_s$, $\overline{\mathbf{M}}_d$ and $\overline{\mathbf{M}}_o$ are active for training. For training depth map, output transformation $\delta(x) = 1/sigmoid(x) - 1$ is applied on $\mathbf{M}_d^I$ and $\mathbf{M}_d$ before computing loss, following [8].

## 4 EXPERIMENTS

### 4.1 Datasets and evaluation metrics

**MuCo-3DHP and MuPoTS-3D** MuCo-3DHP is a large-scale indoor multi-person training dataset [25], including 400K frames for training. For testing, MuPoTS-3D, which consists of real images with various camera poses, are collected from indoor and outdoor scenes. There are 20 real-world scenes in the MuPoTS-3D dataset, which are labeled with ground-truth 3D poses for multiple person subjects, making it a convincing benchmark to test the generalization ability of the 3D pose model.

The widely-used evaluation metric for multi-person 3D pose estimation is $3DPCK$. If the Euclidean distance between predicted and ground-truth is smaller than the threshold (150mm), the prediction is marked as a correct prediction. $PCK_{rel}$ measures relative pose accuracy after root alignment, and $PCK_{abs}$ measures absolute pose accuracy without root alignment. The area under the curve of $3DPCK$ over various thresholds is defined as $AUC$. To evaluate in the crowded scenes, we use *Crowd Index* to split hard cases and

**Table 1: The ablation study of SDAR and DGR on MuPoTS-3D dataset. The star and tree decoding approaches are introduced in [31]. The backbone of all results in this table is ResNet-34.**

| Model | Enhance Feature | Decoding Graph | 3DPCK | Δ |
|-------|-----------------|----------------|-------|---|
| Baseline | × | Star | 75.4 | - |
| Baseline | × | Tree | 76.7 | ↑1.7% |
| Ours(GR-M3D) | × | DGR | 77.9 | ↑3.3% |
| Ours(GR-M3D) | SDAR | DGR | **78.6** | ↑4.2% |

**Table 2: The ablation study of our approach (GR-M3D) based on different backbone on MuPoTS-3D. HG2 denotes Hourglass with 2 blocks. RN denotes ResNet. HRN32 denotes HRNet-32. Baseline is with star structure.**

| Model | Backbone | Params(MB) | Time(s/img) | 3DPCK | Δ |
|-------|----------|------------|-------------|-------|---|
| Baseline | RN34 | 32.1 | 0.022 | 75.4 | ↑ 4.2% |
| GR-M3D | | 34.4 | 0.026 | 78.6 | |
| Baseline | RN50 | 40.6 | 0.029 | 77.1 | ↑ 3.1% |
| GR-M3D | | 42.9 | 0.034 | 79.5 | |
| Baseline | RN101 | 59.6 | 0.041 | 79.1 | ↑ 2.8% |
| GR-M3D | | 61.9 | 0.045 | 81.3 | |
| Baseline | HRN32 | 40.1 | 0.061 | 80.6 | ↑ 2.9% |
| GR-M3D | | 42.4 | 0.065 | 82.9 | |
| Baseline | HG2 | 192.1 | 0.151 | 82.2 | ↑ 2.9% |
| GR-M3D | | 196.6 | 0.158 | **84.6** | |

**Table 3: The ablation study of GR-M3D on MuPoTS-3D dataset with different crowd indexes. Backbone is ResNet-50 here. Baseline is with star structure.**

| Crowd Index | > 0.0 | > 0.3 | > 0.5 | > 0.7 |
|-------------|-------|-------|-------|-------|
| Baseline | 77.1 | 76.4 | 74.8 | 73.2 |
| GR-M3D | **79.5**↑2.7% | **78.8**↑3.1% | **77.9**↑4.1% | **77.5**↑5.9% |

easy cases. $0 \leq Crowd\ Index < 1$ is introduced in [19]. The bigger *Crowd Index* means the serious occlusion by human bodies.

**Human3.6M** Human3.6M [13] is the largest indoor benchmark for single-person 3D pose estimation, which consists of 3.6M video frames. Collectors use the motion capture system to obtain the ground-truth 3D poses. MPJPE and PA-MPJPE are widely used to measure the accuracy of the 3D root-relative pose. They calculate the euclidean distance between predicted and ground-truth 3D joint coordinates after root joint alignment and further rigid alignment (i.e., Procrustes analysis [11]).

**CMU Panoptic** CMU Panoptic [15] is a large-scale multi-person 3D human pose estimation dataset, captured by multiple cameras. It's challenging to recover the multi-person 3D pose since heavy mutual occlusion. Following the settings of [50, 54], 9600 images of two cameras (16, 30) from four activities (Haggling, Mafia, Ultimatum, Pizza) as the testing set, and 160k images from different

videos as training set. For fair comparison with [50, 54], MPJPE is the evaluation metric.

## 4.2 Implement details

GR-M3D is trained on 8 V100 GPUs with a batch size of 16/GPU and input resolution is $512 \times 512$. Adam optimizer is adopted and the initial learning rate is 5e-4, which decreases 10× at 60 and 90 epochs. The total epoch number is 110. Random flip, random occlusion, rotation, and color jittering are used, and the range of rotation is $[-\pi, \pi]$. The confidence threshold for root ketpoints is 0.5. Unless otherwise specified, the backbone of GR-M3D is Hourglass network. Following previous works [9, 28, 54], additional 2D images in MPII [1] and COCO [21] are mixed into 3D datasets for training.

## 4.3 Ablation study

Our baseline is the recently-developed offset-based method. They learn keypoint 3D offset and then decode human pose from single root keypoint by the corresponding star or tree pose graphs [31, 55].

**Effectiveness of SDAR and DGR** For the ablation study, we use ResNet-34 as the backbone. The ablation study of SDAR and DGR is shown in table 1. We can observe that decoding the 3D pose with the tree graph is better than the star graph, which gains a relative improvement of 1.7%. Our DGR achieves 77.9 3DPCK, which gains a relative improvement of 3.3%, compared with the star graph decoding method. We further use the scale and depth aware refinement module to improve the learned keypoints heatmaps, scale maps, depth maps, and 3D offset maps. SDAR achieves 78.6 3DPCK based on our DGR decoding approach. The whole gain of GR-M3D with SDAR and DGR is 3.2 3DPCK, a relative improvement of 4.2%, compared with the star structure decoding approach [31].

We also make a comparison of star, tree, and DGR decoding in Figure 4. As shown in Figure 4, the decoding methods based on star and tree structure fail in occlusion cases. Our DGR performs well on these challenging cases since the multi-root-based dynamic graph reasoning mechanism.

**Different backbones** To evaluate the effectiveness of SDAR and DGR on different backbones, we use several widely used networks ( Hourglass [30] , ResNet [12], HRNet-32 [39]) as backbone. As shown in table 2, our GR-M3D outperforms baseline over a relative of 3% in different backbones. GR-M3D achieves 84.6 3DPCK on the MuPoTS-3D dataset, which is based on the Hourglass backbone. GR-M3D is lightweight and fast. GR-M3D only increases 2MB parameters, compared with baseline methods. GR-M3D achieves real-time based on ResNet-34 and ResNet-50. Even with the heavy backbones (ResNet-101, HRNet-32), GR-M3D achieves 15-23 fps, which indicates that GR-M3D has broad application potential.

**Evaluation in crowded scenes** To evaluate the ability of GR-M3D to handle the occlusion cases, which are often seen in the crowded scenes, we evaluate GR-M3D on MuPoTS-3D with different *Crowd Index*. The bigger value of *Index* means more crowded cases in the image. That represents the pose variances and occlusions are more serious. As shown in table 3, with the increase of *Crowd Index*, the 3DPCK of baseline decreases. At the same time, our GR-M3D outperforms baseline in all *Crowd Index*. Even in the most challenging cases, with a crowd index of over 0.7, GR-M3D outperforms baseline with a relative improvement of 5.9%. These

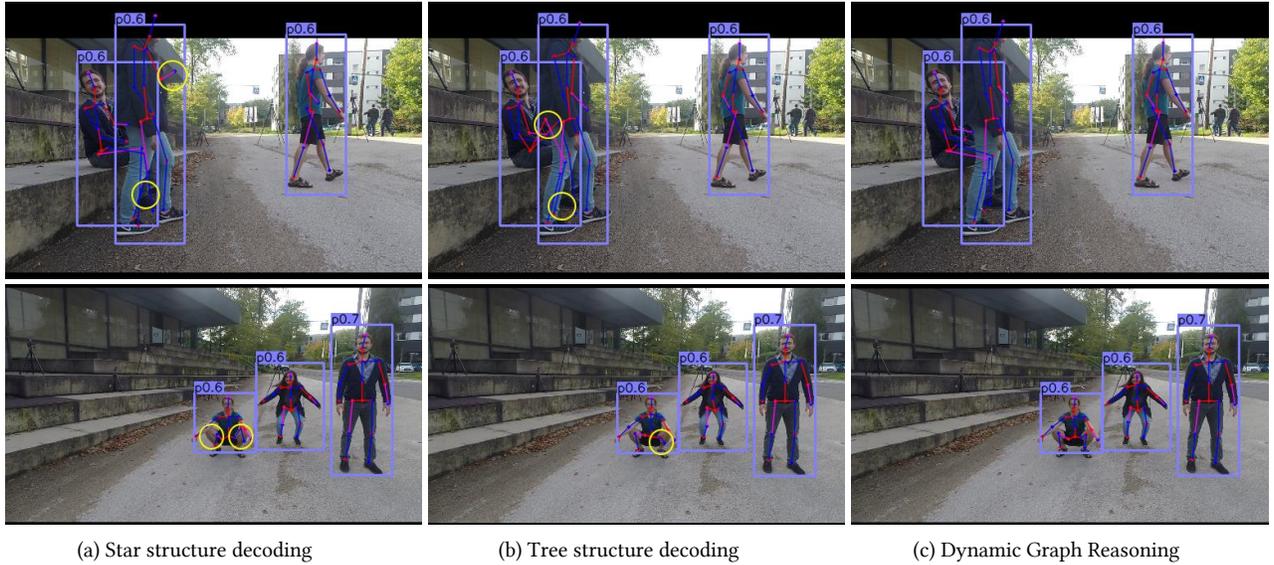| (a) Star structure decoding | (b) Tree structure decoding | (c) Dynamic Graph Reasoning |

Figure 4: The comparison of (a) star structure decoding, (b) tree structure decoding, and (c) our dynamic graph reasoning (DGR). DGR generates better results on these cases with occlusions or strange poses than star and tree structure decoding.

Table 4: Comparison with state-of-the-art approaches on Human3.6M dataset.

| MPJPE(mm) | Direct. | Discuss | Eating | Greet | Phone | Photo | Pose | Purch. | Sitting | SitingD | Smoke | Wait | WalkD | Walk | WalkT | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Jahangiri *et al.* [14] | 74.4 | 66.7 | 67.9 | 75.2 | 77.3 | 70.6 | 64.5 | 95.6 | 127.3 | 79.6 | 79.1 | 73.4 | 67.4 | 71.8 | 72.8 | 77.6 |
| Mehta *et al.* [25] | 57.5 | 68.6 | 59.6 | 67.3 | 78.1 | 56.9 | 69.1 | 98.0 | 117.5 | 69.5 | 82.4 | 68.0 | 55.3 | 76.5 | 61.4 | 72.9 |
| Martinez *et al.* [24] | 51.8 | 56.2 | 58.1 | 59.0 | 69.5 | 55.2 | 58.1 | 74.0 | 94.6 | 62.3 | 78.4 | 59.1 | 49.5 | 65.1 | 52.4 | 62.9 |
| Sun *et al.* [40] | 52.8 | 54.8 | 54.2 | 54.3 | 61.8 | 53.1 | 53.6 | 71.7 | 86.7 | 61.5 | 67.2 | 53.4 | 47.1 | 61.6 | 63.4 | 59.1 |
| Pavlakos *et al.* [32] | 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 |
| Sun *et al.* [41] | 47.5 | 47.7 | 49.5 | 50.2 | 51.4 | 43.8 | 46.4 | 58.9 | 65.7 | **49.4** | 55.8 | 47.8 | 38.9 | 49.0 | 43.8 | 49.6 |
| Moon *et al.* [28] | 50.5 | 55.7 | 50.1 | 51.7 | 53.9 | 46.8 | 50.0 | 61.9 | 68.0 | 52.5 | 55.9 | 49.9 | 41.8 | 56.1 | 46.9 | 53.3 |
| Cai *et al.* [2] | 46.5 | 48.8 | 47.6 | 50.9 | 52.9 | 61.3 | 48.3 | 45.8 | 59.2 | 64.4 | 51.2 | 48.4 | 53.5 | 39.2 | 41.2 | 50.6 |
| Zeng *et al.* [52] | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 49.9 |
| Wehrbein *et al.* [45] | 38.5 | 42.6 | 39.9 | 41.7 | 46.5 | 51.6 | 39.9 | 40.8 | 49.5 | 56.8 | 45.3 | 46.4 | 46.8 | 37.8 | 40.4 | 44.3 |
| Ma *et al.* [23] | **36.3** | 42.8 | 39.5 | **40.0** | **43.9** | 48.8 | 36.7 | **44.0** | 51.0 | 63.1 | 44.3 | 40.6 | 44.4 | 34.9 | **36.7** | 43.4 |
| **Ours(GR-M3D)** | 37.1 | **40.4** | **39.3** | 41.2 | **43.1** | **43.2** | **31.8** | 44.7 | **47.2** | 59.9 | **41.1** | **37.2** | **42.1** | **33.7** | 37.6 | **41.3** |
| PA MPJPE(mm) | Direct. | Discuss | Eating | Greet | Phone | Photo | Pose | Purch. | Sitting | SitingD | Smoke | Wait | WalkD | Walk | WalkT | Avg |
| Martinez *et al.* [24] | 39.5 | 43.2 | 46.4 | 47.0 | 51.0 | 41.4 | 40.6 | 56.5 | 69.4 | 49.2 | 56.0 | 45.0 | 38.0 | 49.5 | 43.1 | 47.7 |
| Fang *et al.* [10] | 38.2 | 41.7 | 43.7 | 44.9 | 48.5 | 40.2 | 38.2 | 54.5 | 64.4 | 47.2 | 55.3 | 44.3 | 36.7 | 47.3 | 41.7 | 45.7 |
| Sun *et al.* [41] | 36.9 | 36.2 | 40.6 | 40.4 | 41.9 | 34.9 | 35.7 | 50.1 | 59.4 | 40.4 | 44.9 | 39.0 | 30.8 | 39.8 | 36.7 | 40.6 |
| Cai *et al.* [2] | 36.8 | 38.7 | 38.2 | 41.7 | 40.7 | 46.8 | 37.9 | 35.6 | 47.6 | 51.7 | 41.3 | 36.8 | 42.7 | 31.0 | 34.7 | 40.2 |
| Ma *et al.* [23] | 30.5 | 34.9 | 32.0 | 32.2 | 35.0 | 37.8 | 28.6 | 32.6 | 40.8 | 52.0 | 35.0 | 31.9 | 35.6 | 26.6 | 28.5 | 34.6 |
| Moon *et al.* [28] | 31.0 | 30.6 | 39.9 | 35.5 | 34.8 | **30.2** | 32.1 | 35.0 | 43.8 | **35.7** | 37.6 | 30.1 | 24.6 | 35.7 | 29.3 | 34.0 |
| Wehrbein *et al.* [45] | 27.9 | 31.4 | 29.7 | 30.2 | 34.9 | 37.1 | 27.3 | 28.2 | 39.0 | 46.1 | 34.2 | 32.3 | 33.6 | 26.1 | 27.5 | 32.4 |
| **Ours(GR-M3D)** | **24.4** | **26.3** | **25.4** | **26.5** | **30.6** | 31.4 | **24.3** | **29.7** | **30.2** | 36.4 | **27.5** | **23.4** | **22.9** | **24.8** | **25.2** | **27.3** |

results show that our GR-M3D can handle the problems of pose variances and heavy occlusions.

## 4.4 Comparison with state-of-the-art methods

Following the settings of SOTA methods, we report the results of GR-M3D on the Human3.6M dataset. As shown in table 4, GR-M3D outperforms SOTA methods and achieves gains of 5% and 16% in MPJPE and PA MPJPE based on Hourglass backbone, respectively.

For multi-person cases, the commonly used dataset is the MuCo-3DHP and MuPoTS-3D dataset. We train GR-M3D on the MuCo-3DHP and test on the MuPoTS-3D. The results are shown in table 5. Based on Hourglass backbone, GR-M3D achieves the state-of-the-art results and obtains relative improvements of 2.5%, 6.5%, and 3.3% in $PCK_{rel}$, $PCK_{abs}$, and $AUC_{rel}$, respectively.
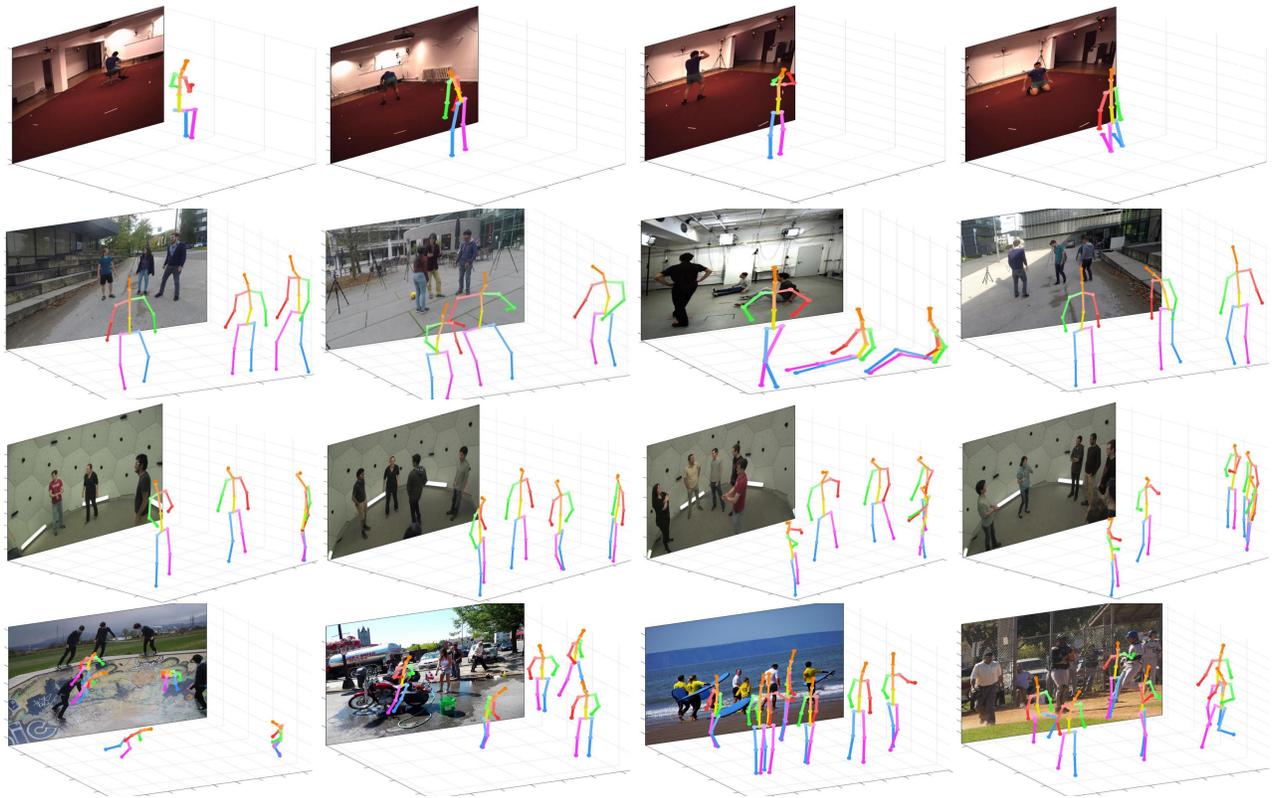
**Figure 5: Visualization of the predicted 3D poses by GR-M3D.**

**Table 5: Comparison on MuPoTS-3D, a 3D multi-person pose estimation dataset. "-" shows that the results are not available. GR-M3D outperforms the SOTA methods. Bigger is better.**

| Methods | Category | $PCK_{rel}$ | $PCK_{abs}$ | $AUC_{rel}$ |
|---|---|---|---|---|
| Lcr-net [37] | | 62.4 | - | - |
| Lcr-net++ [38] | | 74.0 | - | - |
| HG-RCNN [7] | Top-down | 74.2 | - | - |
| HMOR [43] | | 82.0 | - | - |
| PoseNet [28] | | **82.5** | **31.8** | **40.9** |
| ORPM [27] | | 69.8 | - | - |
| Xnect [26] | Bottom-up | 75.8 | - | - |
| SMAP [54] | | **80.5** | **38.7** | **42.7** |
| **Ours(GR-M3D)** | Bottom-up | 84.6↑2.5% | 41.2↑6.5% | 44.1↑3.3% |

**Table 6: Comparison with state-of-the-art methods on CMU Panoptic dataset, 3D multi-person pose estimation dataset. ∗ denotes refining results by an extra network. Lower is better.**

| MPJPE(mm) | Hagg. | Mafia | Ultim. | Pizza | Avg. |
|---|---|---|---|---|---|
| DMHS [34] | 217.9 | 187.3 | 193.6 | 221.3 | 203.4 |
| SemanticFB [50] | 140.0 | 165.9 | 150.7 | 156.0 | 153.4 |
| PoseNet [28] | 89.6 | 91.3 | 79.6 | 90.1 | 87.6 |
| MubyNet [51] | 72.4 | 78.8 | 66.8 | 94.3 | 78.1 |
| SMAP [54] | 71.8 | 72.5 | 65.9 | 82.1 | 73.1 |
| LoCO [9] | 45.0 | 95.0 | 58.0 | 79.0 | 69.0 |
| SMAP [54]* | 63.1 | 60.3 | 56.6 | 67.1 | 61.8 |
| **Ours(GR-M3D)** | **57.1** | **58.3** | **53.4** | **62.7** | **57.9↓ 6.3%** |

## 4.5 Generalization in the wild

The images in MuPoTS-3D are collected in the wild scenes. Table 5 has shown that GR-M3D outperforms the state-of-the-art methods, which demonstrate the generalization ability of GR-M3D in the wild scenes. We also conduct experiments on COCO [21], a larger scale 2D pose estimation dataset. All of the images are collected in challenging, uncontrolled conditions. We directly predict the 3D pose on the COCO images since no 3D pose annotations. The model is based on Hourglass and trained on MuCo-3DHP. The results are shown in Figure 5. The visualization results in Figure 5 are from COCO, Human3.6M, MuPoTS-3D, and CMU Panoptic datasets, respectively. GR-M3D performs well on these challenging

CMU Panoptic dataset is another widely-used benchmark for multi-person 3D pose estimation. Following [51, 54], we take experiments on this dataset and show results in table 6. Compared with state-of-the-art methods, GR-M3D obtains a relative improvement of 6.3% in MPJPE based on the Hourglass backbone.

cases, even on the images collected in the outdoor scenes in the COCO dataset, which shows the generalization ability of GR-M3D.

## 5 CONCLUSION

We propose a novel bottom-up approach (GR-M3D) for 3D multi-person pose estimation, which mitigates occlusions and depth ambiguity by capturing more global context information. We firstly design a scale and depth refinement (SDAR) module to enhance the learned feature maps, to further generate better root keypoints and build robust message propagation paths. DGR reasons the dynamic decoding graphs from the predicted message propagation paths to decoding 3D poses. GR-M3D outperforms previous works and achieves state-of-the-art results on three widely-used benchmarks.

## 6 ACKNOWLEDGEMENT

# REFERENCES

[1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2014. 2d human pose estimation: New benchmark and state of the art analysis. In *CVPR*. 3686–3693.

[2] Yujun Cai, Liuhao Ge, Jun Liu, Jianfei Cai, Tat-Jen Cham, Junsong Yuan, and Nadia Magnenat Thalmann. 2019. Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks. In *ICCV*. 2272–2281.

[3] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. 2017. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*. 7291–7299.

[4] Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gang Yu, and Jian Sun. 2018. Cascaded pyramid network for multi-person pose estimation. In *CVPR*. 7103–7112.

[5] Yu Cheng, Bo Wang, Bo Yang, and Robby T Tan. 2021. Monocular 3D Multi-Person Pose Estimation by Integrating Top-Down and Bottom-Up Networks. In *CVPR*. 7649–7659.

[6] Hai Ci, Chunyu Wang, Xiaoxuan Ma, and Yizhou Wang. 2019. Optimizing network structure for 3d human pose estimation. In *ICCV*. 2262–2271.

[7] Rishabh Dabral, Nitesh B Gundavarapu, Rahul Mitra, Abhishek Sharma, Ganesh Ramakrishnan, and Arjun Jain. 2019. Multi-person 3d human pose estimation from monocular images. In *3DV*. IEEE, 405–414.

[8] David Eigen, Christian Puhrsch, and Rob Fergus. 2014. Depth Map Prediction from a Single Image using a Multi-Scale Deep Network. In *NeurIPS*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger (Eds.), Vol. 27. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2014/file/7bccfde7714a1ebadf06c5f4cea752c1-Paper.pdf

[9] Matteo Fabbri, Fabio Lanzi, Simone Calderara, Stefano Alletto, and Rita Cucchiara. 2020. Compressed volumetric heatmaps for multi-person 3d pose estimation. In *CVPR*. 7204–7213.

[10] Hao-Shu Fang, Yuanlu Xu, Wenguan Wang, Xiaobai Liu, and Song-Chun Zhu. 2018. Learning pose grammar to encode human body configuration for 3d pose estimation. In *AAAI*, Vol. 32.

[11] John C Gower. 1975. Generalized procrustes analysis. *Psychometrika* 40, 1 (1975), 33–51.

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*. 770–778.

[13] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. 2013. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *TPAMI* 36, 7 (2013), 1325–1339.

[14] Ehsan Jahangiri and Alan L Yuille. 2017. Generating multiple diverse hypotheses for human 3d pose consistent with 2d joint detections. In *CVPRW*. 805–814.

[15] Hanbyul Joo, Tomas Simon, Xulong Li, Hao Liu, Lei Tan, Lin Gui, Sean Banerjee, Timothy Godisart, Bart Nabbe, Iain Matthews, et al. 2017. Panoptic studio: A massively multiview system for social interaction capture. *TPAMI* 41, 1 (2017), 190–204.

[16] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. 2018. End-to-end recovery of human shape and pose. In *CVPR*. 7122–7131.

[17] Muhammed Kocabas, Salih Karagoz, and Emre Akbas. 2018. Multiposenet: Fast multi-person pose estimation using pose residual network. In *ECCV*. 417–433.

[18] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. 2019. Pifpaf: Composite fields for human pose estimation. In *CVPR*. 11977–11986.

[19] Jiefeng Li, Can Wang, Hao Zhu, Yihuan Mao, Hao-Shu Fang, and Cewu Lu. 2019. Crowdpose: Efficient crowded scenes pose estimation and a new benchmark. In *CVPR*. 10863–10872.

[20] Yong-Lu Li, Xinpeng Liu, Han Lu, Shiyi Wang, Junqi Liu, Jiefeng Li, and Cewu Lu. 2020. Detailed 2d-3d joint representation for human-object interaction. In *CVPR*. 10166–10175.

[21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *ECCV*. Springer, 740–755.

[22] Ziyu Liu, Hongwen Zhang, Zhenghao Chen, Zhiyong Wang, and Wanli Ouyang. 2020. Disentangling and unifying graph convolutions for skeleton-based action recognition. In *CVPR*. 143–152.

[23] Xiaoxuan Ma, Jiajun Su, Chunyu Wang, Hai Ci, and Yizhou Wang. 2021. Context Modeling in 3D Human Pose Estimation: A Unified Perspective. In *CVPR*. 6238–6247.

[24] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little. 2017. A simple yet effective baseline for 3d human pose estimation. In *ICCV*. 2640–2649.

[25] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. 2017. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *3DV*. IEEE, 506–516.

[26] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Mohamed Elgharib, Pascal Fua, Hans-Peter Seidel, Helge Rhodin, Gerard Pons-Moll, and Christian Theobalt. 2019. Xnect: Real-time multi-person 3d human pose estimation with a single rgb camera. *arXiv preprint arXiv:1907.00837* (2019).

[27] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Srinath Sridhar, Gerard Pons-Moll, and Christian Theobalt. 2018. Single-shot multi-person 3d pose estimation from monocular rgb. In *3DV*. IEEE, 120–130.

[28] Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. 2019. Camera distance-aware top-down approach for 3d multi-person pose estimation from a single rgb image. In *ICCV*. 10133–10142.

[29] Francesc Moreno-Noguer. 2017. 3d human pose estimation from a single image via distance matrix regression. In *CVPR*. 2823–2832.

[30] Alejandro Newell, Kaiyu Yang, and Jia Deng. 2016. Stacked hourglass networks for human pose estimation. In *ECCV*. Springer, 483–499.

[31] Xuecheng Nie, Jiashi Feng, Jianfeng Zhang, and Shuicheng Yan. 2019. Single-stage multi-person pose machines. In *ICCV*. 6951–6960.

[32] Georgios Pavlakos, Xiaowei Zhou, and Kostas Daniilidis. 2018. Ordinal depth supervision for 3d human pose estimation. In *CVPR*. 7307–7316.

[33] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. 2017. Coarse-to-fine volumetric prediction for single-image 3D human pose. In *CVPR*. 7025–7034.

[34] Alin-Ionut Popa, Mihai Zanfir, and Cristian Sminchisescu. 2017. Deep multitask architecture for integrated 2d and 3d human sensing. In *CVPR*. 6289–6298.

[35] Zhongwei Qiu, Kai Qiu, Jianlong Fu, and Dongmei Fu. 2019. Learning recurrent structure-guided attention network for multi-person pose estimation. In *ICME*. IEEE, 418–423.

[36] Zhongwei Qiu, Kai Qiu, Jianlong Fu, and Dongmei Fu. 2020. Dgcn: Dynamic graph convolutional network for efficient multi-person pose estimation. In *AAAI*, Vol. 34. 11924–11931.

[37] Gregory Rogez, Philippe Weinzaepfel, and Cordelia Schmid. 2017. Lcr-net: Localization-classification-regression for human pose. (2017), 3433–3441.

[38] Gregory Rogez, Philippe Weinzaepfel, and Cordelia Schmid. 2019. Lcr-net++: Multi-person 2d and 3d pose detection in natural images. *TPAMI* 42, 5 (2019), 1146–1161.

[39] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. 2019. Deep high-resolution representation learning for human pose estimation. In *CVPR*. 5693–5703.

[40] Xiao Sun, Jiaxiang Shang, Shuang Liang, and Yichen Wei. 2017. Compositional human pose regression. In *ICCV*. 2602–2611.

[41] Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. 2018. Integral human pose regression. In *ECCV*. 529–545.

[42] Zhi Tian, Hao Chen, and Chunhua Shen. 2019. Directpose: Direct end-to-end multi-person pose estimation. *arXiv preprint arXiv:1911.07451* (2019).

[43] Can Wang, Jiefeng Li, Wentao Liu, Chen Qian, and Cewu Lu. 2020. Hmor: Hierarchical multi-person ordinal relations for monocular multi-person 3d pose estimation. In *ECCV*. Springer, 242–259.

[44] Jianbo Wang, Kai Qiu, Houwen Peng, Jianlong Fu, and Jianke Zhu. 2019. AI coach: Deep human pose estimation and analysis for personalized athletic training assistance. In *ACM MM*. 374–382.

[45] Tom Wehrbein, Marco Rudolph, Bodo Rosenhahn, and Bastian Wandt. 2021. Probabilistic monocular 3d human pose estimation with normalizing flows. In *ICCV*. 11199–11208.

[46] Fangyun Wei, Xiao Sun, Hongyang Li, Jingdong Wang, and Stephen Lin. 2020. Point-set anchors for object detection, instance segmentation and pose estimation. In *ECCV*. Springer, 527–544.

[47] Bin Xiao, Haiping Wu, and Yichen Wei. 2018. Simple baselines for human pose estimation and tracking. In *ECCV*. 466–481.

[48] Sijie Yan, Yuanjun Xiong, and Dahua Lin. 2018. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *AAAI*.

[49] Wei Yang, Wanli Ouyang, Xiaolong Wang, Jimmy Ren, Hongsheng Li, and Xiaogang Wang. 2018. 3d human pose estimation in the wild by adversarial learning. In *CVPR*. 5255–5264.

[50] Andrei Zanfir, Elisabeta Marinoiu, and Cristian Sminchisescu. 2018. Monocular 3d pose and shape estimation of multiple people in natural scenes-the importance of multiple scene constraints. In *CVPR*. 2148–2157.

[51] Andrei Zanfir, Elisabeta Marinoiu, Mihai Zanfir, Alin-Ionut Popa, and Cristian Sminchisescu. 2018. Deep network for the integrated 3d sensing of multiple people in natural images. *NeurIPS* 31 (2018), 8410–8419.

[52] Ailing Zeng, Xiao Sun, Fuyang Huang, Minhao Liu, Qiang Xu, and Stephen Lin. 2020. Srnet: Improving generalization in 3d human pose estimation with a split-and-recombine approach. In *ECCV*. Springer, 507–523.

[53] Long Zhao, Xi Peng, Yu Tian, Mubbasir Kapadia, and Dimitris N Metaxas. 2019. Semantic graph convolutional networks for 3d human pose regression. In *CVPR*. 3425–3435.

[54] Jianan Zhen, Qi Fang, Jiaming Sun, Wentao Liu, Wei Jiang, Hujun Bao, and Xiaowei Zhou. 2020. Smap: Single-shot multi-person absolute 3d pose estimation. In *ECCV*. Springer, 550–566.

[55] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. 2019. Objects as points. *arXiv preprint arXiv:1904.07850* (2019).

## SUPPLEMENTARY MATERIAL

In this supplementary material, we introduce the algorithm details in Section A. To better understand the dynamic graph reasoning (DGR) in our approach, we show the visualization of the dynamic decoding graph in Section B. The limitations and failure cases of our approaches are shown in Section C. More visualization cases are shown in Section D.

## A  ALGORITHM DETAILS

The algorithm details of our GR-M3D are shown in Algorithm 1. In Algorithm 1, the backbone network can be ResNet [12], HRNet [39], and Hourglass [30], etc. Unless otherwise specified, the backbone of GR-M3D is the Hourglass network. *conv* means the convolutional layers with a kernel size of $1 \times 1$ to generate the four data maps.

---

**Algorithm 1** GR-M3D with dynamic graph reasoning (DGR)

---

**Input:** $I$: Input image; $K$: Joint number; $N$: Person number; $\phi(\cdot)$: Backbone network; $F$: Deep features; *conv*: Convolution layers with kernel size of $1 \times 1$; $SDAR(\cdot)$: The Scale and Depth Aware Refinement module.

**Output:** $\mathbb{P}_{3d}$: $\{\hat{p}_{3d}^j, j \in [1, K]\}$;

1: $F = \phi(I)$;
2: $M_h^I, M_s^I, M_d^I, M_o^I = conv(F)$;
3: $M_h, M_s, M_d, M_o = SDAR(M_h^I, M_s^I, M_d^I, M_o^I, F)$;
4: Obtain center points set $C = \{c_1, c_2, ..., c_N\}$ from $M_h$;
5: Obtain 2D keypoints set $P_{2D} = \{p_1, p_2, ..., p_M\}$ from $M_h$;
6: Assign $P_{2D}$ to $C$ as Eq.(6);
7: Obtain dense decoding paths $\mathbb{E} = \{e^{ij}, i \in [1, K], j \in [1, K]\}$;
8: Calculate dynamic decoding graphs $\hat{\mathbb{E}} = \{(e^{ij}, \mathcal{W}(p^i, p^j)), i \in [1, K], j \in [1, K]\}$ according to Eq.(8), Eq.(9), and Eq.(10);
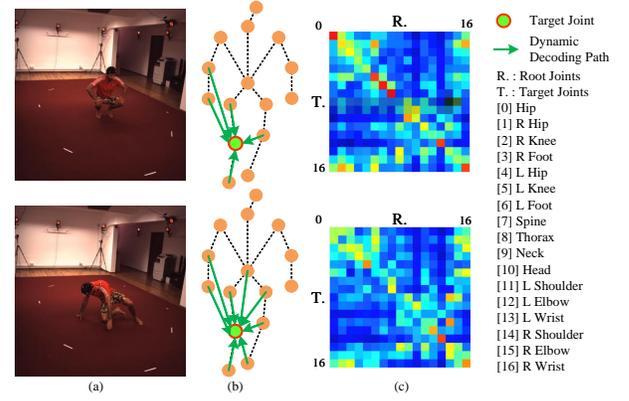9: Calculate 3D coordinates $\hat{p}_{3d}^j$ as Eq.(11).

---

## B  DYNAMIC DECODING GRAPH

To better understand the proposed dynamic graph reasoning (DGR) in GR-M3D for 3D human pose estimation, we visualize the path weights in dynamic decoding graphs of different people in Figure 6.
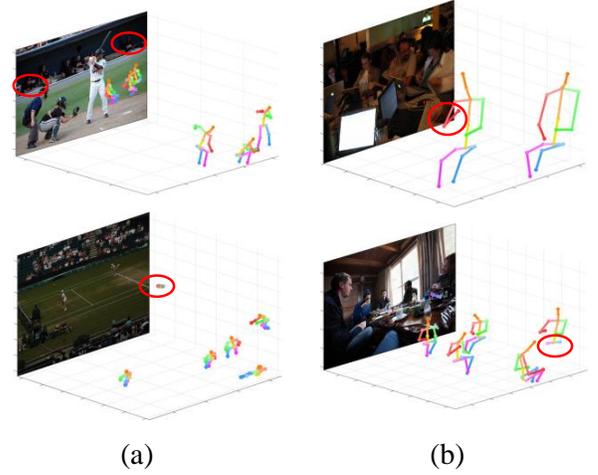
As shown in Figure 6, for the two input cases performed by the same person, our GR-M3D predicted different decoding graphs as Figure 6 (c). As shown in Figure 6 (b), for the cases in the second row, the decoding graph contains more decoding paths with high weights ($> 0.1$) from the hip, left hip, and left foot since occlusion. Compared with these two cases, it shows that the predicted decoding graph is self-adapting due to the dynamic graph reasoning (DGR) mechanism.

## C  LIMITATIONS AND FAILURE CASES

We discuss the limitations of the proposed GR-M3D in this section and show some failure cases in Figure 7. As shown in Figure 7 (a), for some small instances, our approach can not handle and misses them. Due to the input is whole image for GR-M3D, the small instances with limited resolution lack sufficient information for GR-M3D to handle. As shown in Figure 7 (b), for some cases with heavy occlusions, it is hard for GM-M3D to tackle.
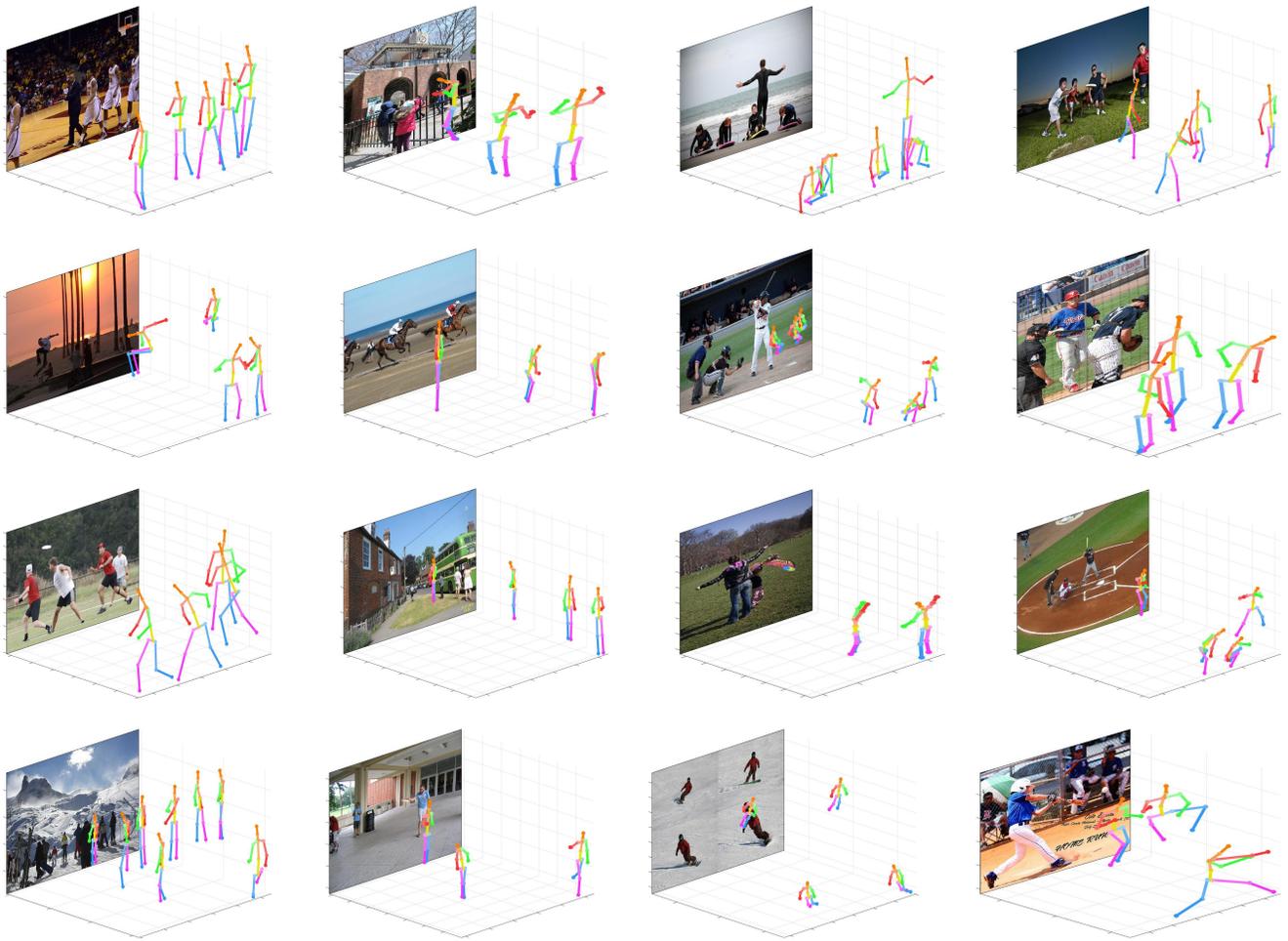


Figure 6: The visualization of dynamic decoding graphs. (a) Different actions from the same person, (b) The illustration of dynamic decoding paths (Green arrows show the decoding paths that weight is greater than 0.1), (c) The weight matrix of the predicted dynamic decoding graph(DDG). The DDGs for the two cases from the same person are different, which demonstrates that our method can self-adaptively estimate the best decoding graph for each person according to the different inputs.



Figure 7: The visualization of failure cases from COCO [21] dataset. (a) Some missed instances since they are too small, (b) Some wrong predictions since heavy occlusions.

## D  MORE VISUALIZATION RESULTS

To demonstrate the generalization ability of proposed GR-M3D, the more visualization results on COCO [21] and MuPoTS-3D [25] datasets are shown in Figure 8 and Figure 9. The images in COCO [21] dataset are collected in unconstraint and in-the-wild conditions. There is no 3D pose annotations in COCO dataset. The images in MuPoTS-3D [25] dataset are collected in constraint and in-the-wild conditions. All the results in Figure 8 and Figure 9 are predicted by GR-M3D, which is trained on MuCo-3DHP dataset [25].

**Figure 8: Visualization of the predicted 3D poses by GR-M3D on COCO [21] dataset.**

As shown in Figure 8, our GR-M3D performs well on these in-the-wild images with strange poses, occlusions, and variational backgrounds. As shown in Figure 9, our GR-M3D still can handle these images with different human actions, changing camera viewpoint, and occlusions. These results show that our GR-M3D has a strong generalization ability on handling these in-the-wild cases.

**Figure 9: Visualization of the predicted 3D poses by GR-M3D on MuPoTS-3D [25] dataset.**