PRML: PROGRESSIVE MULTI-TASK LEARNING FOR MONOCULAR 3D HUMAN POSE ESTIMATION

Anonymous authors

Paper under double-blind review

ABSTRACT

The lifting-based framework has dominated the field of monocular 3D human pose estimation by leveraging the well-detected 2D pose as an intermediate representation. However, it neglects different initial states between 2D pose and per-joint depth. The initial state of the 2D pose is well-detected, but the per-joint depth is unknown and needs to be learned from scratch. The lifting-based framework encodes the well-detected 2D pose and unknown per-joint depth in an entangled feature space, explicitly introducing depth uncertainty to the well-detected 2D pose. To address this limitation, we present a progressive multi-task learning pose estimation framework named **PrML**. First, PrML introduces two task branches to refine the well-detected 2D pose features and to learn the per-joint depth features. This dual-branch design reduces the explicit influence of uncertain depth features on 2D pose features. Second, PrML employs a task-aware decoder to indirectly supplement the complementary information between the refined 2D pose features and learned per-joint depth features. This step establishes the connection between 2D pose and per-joint depth, compensating for the lack of interaction caused by the dual-branch design. We conduct theoretical analysis from the perspective of mutual information and arrive at a loss to supervise this feature complementary process. Finally, we use two regression heads to regress the 2D pose and per-joint depth, respectively, and concatenate them to obtain the final 3D pose. Extensive experiments show that PrML outperforms the conventional lifting-based framework with fewer parameters on two widely used datasets: Human3.6M and MPI-INF-3DHP. Code is available at https://anonymous.4open.science/r/PrML.

031 032 033

034

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

1 INTORDUCTION

Monocular 3D human pose estimation has been a crucial problem in computer vision, which aims to locate the 3D joint positions of a human body (Moon & Lee, 2020; Pavlakos et al., 2018; Chen et al., 2021). Nowadays, monocular 3D human pose estimation finds widespread applications in 037 various scenarios, including motion prediction (Liu et al., 2021b; 2022b), action recognition (Zhang et al., 2022a), and human-robot interaction (Gong et al., 2022; Ye et al., 2021). Existing monocular 3D human pose estimation methods can be categorized as the end-to-end manner and lifting-based 040 manner. The end-to-end approaches (Kanazawa et al., 2018; Pavlakos et al., 2017; Sun et al., 2018) 041 directly estimate the 3D pose from the input image without the intermediate 2D pose representation. 042 Different from the end-to-end manner, lifting-based methods (Martinez et al., 2017; Liu et al., 2020) 043 first obtain 2D pose using 2D pose detector (Newell et al., 2016; Chen et al., 2018) and then lift 044 the 2D pose in image coordinate to the 3D pose in camera coordinate. These lifting-based methods usually outperform the end-to-end manner and dominate the monocular 3D human pose estimation.

Recent lifting-based methods (Zheng et al., 2021; Li et al., 2022b; Zhang et al., 2022b; Yu et al., 2023;
Zhu et al., 2023; Peng et al., 2024) for monocular 3D human pose estimation focus on designing
various spatio-temporal encoders. As shown in Figure 1 left, they project the 2D pose into an
entangled feature space and regress the 3D pose from it. This lifting process neglects the different
initial states of 2D pose and per-joint depth. It encodes the well-detected 2D pose and unknown
per-joint depth in an entangled feature space, which introduces a main limitation: the high
uncertainty of the per-joint depth may erode the 2D pose. It is well-known that the monocular 3D
human pose estimation task is an ill-posed problem and inherently suffers from depth ambiguity (Li
et al., 2022b; Ma et al., 2021b; Wehrbein et al., 2021). One 2D pose possibly corresponds to multiple

065

066

067

068

069

083

084

085

086 087



Figure 1: Given a 2D pose in the image coordinate, we aim to estimate the 3D pose in the camera coordinate. Left: Conventional lifting-based framework directly projects the 2D pose in an entangled feature space and regression the 3D pose from it. Right: Our proposed progressive multi-task learning framework. The 2D pose and per-joint depth features are learned separately in the first step. In the second step, we perform feature interaction to supplement the complementary information. Finally, we regress the 2D pose and per-joint depth and concatenate them to obtain the final 3D pose.



Figure 2: Qualitative Comparison of 2D Pose (Ground Truth, Input, MotionBERT (Zhu et al., 2023) and Ours). We project the 2D pose in the camera coordinate (part of the output 3D pose) back to the image coordinate for comparison. The powerful lifting-based method MotionBERT gets a 2D pose worse than the input, which contradicts our intuition. In contrast, our framework obtains a 2D pose better than the input. Please refer to Appendix **F** for more qualitative and quantitative comparisons.

088 3D poses, where the lifting process is inherently ambiguous (Yu et al., 2021). To validate the impact of depth uncertainty on the 2D pose, we project the 2D pose in the camera coordinate (part of output 089 3D pose) back to the image coordinate and compare it with the ground truth 2D pose and input 2D 090 pose. As shown in Figure 2, despite learning through multiple spatio-temporal encoders, the 2D pose 091 of the powerful lifting method MotionBERT (Zhu et al., 2023) is even worse than the original input 092 2D pose. This observation provides empirical evidence that directly encoding the well-detected 2D 093 pose features and the unknown per-joint depth features in an entangled feature space will inevitably 094 introduce explicit uncertainty to the 2D pose and cause erosion. To provide more empirical support for the high uncertainty of per-joint depth, we conduct quantitative comparisons of Mean Per Joint 096 Position Error (MPJPE) across different axes for different hard actions (Zeng et al., 2021) with MotionBERT (Zhu et al., 2023), GLA-GCN (Yu et al., 2023) and KTPFormer (Peng et al., 2024). As 098 shown in Figure 3, the MPJPE of per-joint depth is significantly higher than the MPJPE of 2D pose and accounts for the majority of the overall MPJPE. These quantitative findings highlight the high 099 uncertainty of per-joint depth compared to the well-detected 2D pose. 100

Motivated by these qualitative and quantitative observations, we propose a progressive multi-task learning pose estimation framework named **PrML** to address this limitation. As shown in Figure 1 right, the first step of PrML introduces two task branches: refining the well-detected 2D pose features and learning the per-joint depth features. The dual-branch design brings two benefits. First, learning the features of the 2D pose and per-joint depth separately avoids the explicit impact of uncertain depth features on the 2D pose. Second, the model parameters are not shared across two task branches, which makes the training more targeted. After dual-branch learning, we obtain the refined 2D pose features and learned depth features, mitigating the uncertainty of depth features. In light of this, the second



Figure 3: Quantitative Comparison of Mean Per Joint Position Error (MPJPE) of different axes for all actions and three hard actions (Zeng et al., 2021) with lifting-based methods (Zhu et al., 2023; Yu et al., 2023; Peng et al., 2024). The MPJPE of the Z-axis (per-joint depth) is significantly higher than the X-Y axes (2D pose) and accounts for the majority of the overall (3D pose) MPJPE. Our proposed framework achieves better results across different axes than the lifting-based framework.

125 step of PrML employs a task-aware decoder to indirectly supplement the complementary information 126 between the refined 2D pose features and the learned per-joint depth features. This step compensates for the information loss caused by the dual-branch structure and establishes the connection between 127 2D pose and per-joint depth. We also conduct theoretical analysis from the perspective of mutual 128 information (Becker, 1996) and arrive at a loss to supervise this feature complementary process. 129 Finally, we regress the 2D pose and per-joint depth, respectively, and concatenate them to obtain the 130 final 3D pose. As shown in Figure 2, our framework could reduce the erosion of the 2D pose caused 131 by depth uncertainty. The quantitative results in Figure 3 also demonstrate our framework performs 132 favorably across different parts (2D pose and per-joint depth) of the 3D pose. Extensive experiments 133 on two widely used monocular 3D human pose estimation benchmarks (i.e., Human3.6M (Ionescu 134 et al., 2013) and MPI-INF-3DHP (Mehta et al., 2017)) demonstrate that the proposed progressive 135 multi-task learning framework outperforms conventional lifting-based framework in terms of accuracy 136 and robustness with fewer parameters. The key contributions of this paper are as follows:

- We tackle an overlooked different initial states between the well-detected 2D pose and the unknown per-joint depth of the lifting-based framework and present a novel progressive multi-task learning pose estimation framework named PrML to address it.
- We propose a task-aware decoder to indirectly supplement the complementary information between 2D pose and per-joint depth after task learning. We also conduct theoretical analysis from the perspective of mutual information to explicitly supervise this feature complementary process.
- Our framework achieves state-of-the-art results on Human3.6M and MPI-INF-3DHP datasets with fewer parameters. These results demonstrate the potential of the progressive multi-task learning framework for future monocular 3D human pose estimation research.
- 145 146 147

148 149

137

138

139

140

141

142

143

144

119

120

121

122

123 124

2 RELATED WORK

Monocular 3D Human Pose Estimation. Existing methods for monocular 3D human pose estimation 150 can be categorized as end-to-end and lifting-based. End-to-end approaches (Kanazawa et al., 2018; 151 Pavlakos et al., 2017; Sun et al., 2018) directly estimate the 3D pose from the input image without the 152 intermediate 2D pose representation. With the reliable achievement of 2D human pose detectors (Chen 153 et al., 2018; Newell et al., 2016; Sun et al., 2019), lifting-based methods (Fang et al., 2018; Martinez 154 et al., 2017; Zhao et al., 2019; Liu et al., 2020) first obtain 2D pose representations in the image and then lift the 2D joint coordinates to 3D space. Recently, Transformers (Vaswani et al., 2017) 156 have been applied to various visual tasks (Dosovitskiy et al., 2021; Carion et al., 2020). For the 157 monocular 3D human pose estimation task, PoseFormer (Zheng et al., 2021) introduces transformer 158 architecture to leverage spatial and temporal dependency. MHFormer (Li et al., 2022b) addresses the depth ambiguity by learning multiple pose hypotheses and MixSTE (Zhang et al., 2022b) 159 constructs a mixed spatiotemporal transformer to capture the temporal motion of different body joints. STCFormer (Tang et al., 2023) decomposed spatio-temporal attention and integrated the 161 structure-enhanced positional embedding. On the other hand, MotionBERT (Zhu et al., 2023) trains

a unified model for multiple downstream tasks. In (Peng et al., 2024), KTPFormer uses two prior attention modules to facilitate pose estimation. Moreover, MotionAGFormer (Soroush Mehraban, 2024) using two parallel transformer and GCNFormer streams to better learn the underlying 3D structure. However, these methods are developed within the conventional lifting-based framework. In contrast, we propose a progressive multi-task learning framework to estimate 3D human pose.

167 Multi-Task Learning. Multi-Task Learning (MTL) (Caruana, 1997) is a learning paradigm in 168 machine learning, and it aims to leverage useful information contained in multiple related tasks to 169 help improve the generalization performance of all the tasks (Zhang & Yang, 2021). Numerous models 170 have been explored (Vandenhende et al., 2020; Brüggemann et al., 2021) within the MTL framework. 171 Moreover, existing approaches analyze the optimization of multi-task learning by designing multi-task 172 loss (Liu et al., 2021a; Li et al., 2022a) or gradient manipulations (Yu et al., 2020; Wang et al., 2020). MTL has been widely used in computer vision, such as image classification (Rebuffi et al., 2017), 173 semantic segmentation (Hoyer et al., 2021; Li et al., 2023), and dense prediction (Proesmans et al., 174 2022; Hoyer et al., 2021). In (Iqbal et al., 2018), they introduce a novel scale and translation invariant 175 2.5D pose representation contain 2D pose and depth. Our approach is motivated by these former 176 attempts but from the perspective of decomposing the single 3D human pose estimation task into two 177 sub-tasks and learning them in a progressive manner, which is a novel and unexplored question. 178

179 **Mutual Information.** Mutual Information plays an important role in the representation learning. As the pioneering work among mutual information methods, Linsker (Linsker, 1988) proposes to 180 maximize mutual information between the input and output. Designing optimization objectives based 181 on mutual information maximization has been extensively studied (Becker, 1992; Wiskott & Se-182 jnowski, 2002). For human pose estimation, CV-MIM (Zhao et al., 2021) introduces a representation 183 learning method to disentangle pose-dependent and view-dependent factors from 2D human poses. 184 FAMI-Pose (Liu et al., 2022a) designs a mutual information loss to maximize the complementary 185 information between temporal frames. TDMI (Feng et al., 2023) proposes to minimize the mutual 186 information between useful and noisy constituents of the raw features. To the best of our knowledge, 187 we are the first to introduce mutual information loss to the monocular 3D human pose estimation task.

188 189

3 RETHINKING LIFTING-BASED MONOCULAR 3D HUMAN POSE ESTIMATION

190 191

Since SimpleBaseline (Martinez et al., 2017) proposes the 2D-to-3D lifting framework, numerous methods (Pavllo et al., 2019; Zhang et al., 2022b; Li et al., 2022b; Shan et al., 2022; Zhao et al., 2023b; Shan et al., 2023; Zhu et al., 2023; Tang et al., 2023; Peng et al., 2024; Mehraban et al., 2024) have been developed within this framework. These lifting-based methods usually outperform the end-to-end manner (Kanazawa et al., 2018; Pavlakos et al., 2017; Sun et al., 2018) and have been the dominant paradigm in monocular 3D human pose estimation for a long time.

The ensuing question is why lifting-based methods perform better than end-to-end approaches. We 198 argue that this is mainly attributed to leveraging the 2D pose as an intermediate representation. First, 199 there exists a high relevance between 2D pose and 3D pose. Regressing 3D pose directly from 200 raw images is a highly nonlinear and challenging problem (Pavlakos et al., 2017). This difficulty 201 also exists in 2D human pose estimation (Pfister et al., 2015; Tompson et al., 2014). In contrast, 202 with the widespread usage of 2D human pose detectors (Chen et al., 2018; He et al., 2017; Newell 203 et al., 2016; Sun et al., 2019), lifting-based methods could leverage the well-detected 2D pose, 204 which contributes to its 3D counterpart and make network training easy. Second, the 2D pose is 205 exceptionally lightweight regarding memory cost compared to raw image. This property enables 206 lifting-based methods to leverage long-term temporal clues to address the occlusion and achieve advanced accuracy. (e.g., 243 frames for MixSTE (Zhang et al., 2022b), MotionBERT (Zhu et al., 207 2023), and KTPFormer (Peng et al., 2024); large as 351 frames for MHFormer (Li et al., 2022b)) 208

Once we have a well-detected 2D pose, lifting it directly to 3D space is natural and simple. However, these lifting-based methods neglect different initial states between 2D pose and per-joint depth and encode the well-detected 2D pose features and unknown per-joint depth features in an entangled feature space. This leads to the fact that despite these methods (Zhu et al., 2023; Peng et al., 2024; Li et al., 2022b) striving to design various encoders to leverage the well-detected 2D pose, the 2D pose itself is inevitably eroded by the uncertainty of depth features (see Figure 2). This paper presents a progressive multi-task learning framework that addresses the different initial states between 2D pose and per-joint depth and provides a new choice for future monocular 3D human pose estimation.



Figure 4: Overview of the proposed progressive multi-task learning framework PrML, which comprises a shared bottom, a 2D pose branch, a depth branch, and a task-aware decoder.

4 METHODOLOGY

238

239

240 241

242 243

244 245

246

247

248 249

250

251

262 263 264

266

4.1 **PROBLEM FORMULATION**

Given a 2D pose sequence $X \in \mathbb{R}^{T \times J \times C_{in}}$, the goal of monocular 3D human pose estimation is to estimate the 3D pose sequence $\overline{Y} \in \mathbb{R}^{T \times J \times C_{out}}$. Here, T refers to the number of input frames, and J refers to the number of joints. C_{in} and C_{out} denote the dimension of the input and output.

4.2 MULTI-TASK LEARNING BRANCH

One of the widely used multi-task learning models is proposed by Caruana (Caruana, 1997; 1993), 252 which has a shared-bottom model structure that substantially reduces the risk of overfitting (Ma et al., 253 2018). We first use a linear embedding layer to project the 2D pose sequence into high-dimensional 254 features. Then, we employ the DSTFormer block proposed by MotionBERT (Zhu et al., 2023) as our shared bottom to extract general features $F \in \mathbb{R}^{T \times J \times C}$. The DSTFormer block is composed of 255 256 spatial-temporal and temporal-spatial branches. The outputs of two branches are adaptively fused by 257 an attention regressor. Next, we add the learnable 2D pose position embedding and per-joint depth position embedding to F to obtain the 2D pose features $F_{2D} \in \mathbb{R}^{T \times J \times C}$ and the per-joint depth features $F_D \in \mathbb{R}^{T \times J \times C}$ respectively. C denotes the feature dimension. Subsequently, the 2D pose 258 259 branch and depth branch repeat the temporal transformer encoder (TF_T) and spatial transformer 260 encoder (TF_S) for N times to refine the 2D pose features and learn depth features separately as: 261

$$F_{2D}^{n} = TF_{S}(TF_{T}(F_{2D}^{n-1})) \quad F_{D}^{n} = TF_{S}(TF_{T}(F_{D}^{n-1})) \qquad n = 1...N$$
(1)

265 4.3 TASK-AWARE DECODER

Coarse Alignment. We first use a fully connected layer $\mathcal{F}_{FC}(\cdot)$ and softmax function $Softmax(\cdot)$ to compute the coarse alignment parameters $\Theta = (\theta_1, \theta_2) \in \mathbb{R}^{T \times J \times 2}$ to project the 2D pose features and per-joint depth features into a shared feature space and obtain the shared features $F_S \in \mathbb{R}^{T \times J \times C}$. This coarse alignment operation mitigates the uncertainty erosion by avoiding direct interaction between the 2D pose feature and the per-joint depth feature and can be expressed as follows:

$$F_{S} = \underbrace{Softmax(\mathcal{F}_{FC}(F_{2D}^{N} \oplus F_{D}^{N}))}_{\Theta} \begin{pmatrix} F_{2D}^{N} \\ F_{D}^{N} \end{pmatrix} = \theta_{1}F_{2D}^{N} + \theta_{2}F_{D}^{N}$$
(2)

Feature Complement. The shared features F_S obtained after coarse alignment have the 2D pose and per-joint depth features but lack precise and targeted support information. Thus, we introduce the learnable 2D pose bias $B_{2D} \in \mathbb{R}^{T \times J \times C}$ and per-joint depth bias $B_D \in \mathbb{R}^{T \times J \times C}$ to address this issue. We concatenate the task bias B_{2D} and B_D with the shared features F_S to obtain the 2D pose support features $S_{2D} \in \mathbb{R}^{T \times J \times 2C}$ and depth support features $S_D \in \mathbb{R}^{T \times J \times 2C}$. Then, we utilize Multi-Head Cross-Attention (Vaswani et al., 2017) (TF_C) with the original features acting as the query and the support features serving as the key and value for feature supplementation. Technically, we get the enhanced 2D pose features \tilde{F}_{2D} and enhanced depth features \tilde{F}_D as follows:

$$\widetilde{F}_{2D} = TF_C(F_{2D}^N, S_{2D}) \quad \widetilde{F}_D = TF_C(F_D^N, S_D)$$
(3)

4.4 MUTUAL INFORMATION OBJECTIVE

 Mutual Information. Mutual information (MI) is an important measurement to quantify the statistical dependency of two random variables. Given two random variables x and y, p(x, y) represents the joint probability distribution between x and y, while p(x) and p(y) represent their marginal distributions. The mutual information between two random variables x and y is defined as:

$$\mathcal{I}(X;Y) = \int_{Y} \int_{X} p(x,y) \log\left(\frac{p(x,y)}{p(x)p(y)}\right) dxdy \tag{4}$$

Mutual Information Loss. Within the task-aware decoder, our goal is to explicitly supervise the feature complementary process. This mutual information objective can be formulated as follows:

$$\max\left[\mathcal{I}(Y_{2D}; S_{2D} \mid B_{2D}) + \mathcal{I}(Y_D; S_D \mid B_D)\right]$$
(5)

where Y_{2D} and Y_D denote the 2D pose and per-joint depth of the 3D pose label. Intuitively, optimizing this objective will maximize the mutual information between the support feature and the label to ehance the feature complementary. Due to the notorious difficulty of the conditional MI computations especially in neural networks (Hjelm et al., 2018; Tian et al., 2021), we factorize Equation 5 as:

$$\mathcal{I}(Y_D; S_D \mid B_D) = \mathcal{I}(Y_D; S_D) - \mathcal{I}(S_D; B_D) + \int_{Y_D} \underbrace{D_{KL}(P_{(S_D, B_D)|Y_D} \parallel P_{S_D|Y_D} P_{B_D|Y_D})}_{\text{KL Divergence} \ge 0} dP_{Y_D}$$

$$\geq \mathcal{I}(Y_D; S_D) - \mathcal{I}(S_D; B_D) \tag{6}$$

Since both $\mathcal{I}(Y_D; S_D)$ and $\mathcal{I}(S_D; B_D)$ are non-negative, the $\mathcal{I}(Y_D; S_D) - \mathcal{I}(S_D; B_D)$ will result in negative values during training. This will yield negative values during training, leading to vanishing gradients problem and preventing training from converging. Therefore, we simplified the implementation of mutual information by calculating only the first term $\mathcal{I}(Y_D; S_D)$. Finally, we obtain two simplified mutual information optimization objectives as follows:

$$\mathcal{L}_{MI} = \lambda_{2D} \mathcal{I}(Y_{2D}; S_{2D}) + \lambda_D \mathcal{I}(Y_D; S_D)$$
(7)

312 The λ_{2D} and λ_D serve as hyper-parameters in our framework to balance different objects.

4.5 REGRESSION HEAD AND LOSS FUNCTION

We use two regression heads (MLP) to regress the 2D pose $\overline{Y}_{2D} \in \mathbb{R}^{T \times J \times 2}$ and per-joint depth $\overline{Y}_D \in \mathbb{R}^{T \times J \times 1}$ respectively and concatenate them to generate the 3D pose sequence \overline{Y} . Losses are independently calculated for 2D pose and per-joint depth as Equation 8. For the 2D pose, we use L2 loss to minimize the errors between predictions and ground truth. For the per-joint depth, we use mean absolute error loss to minimize the errors between the estimated per-joint depth and label.

$$\mathcal{L}_{3D} = \underbrace{\frac{1}{JT} \sum_{j=1}^{J} \sum_{t=1}^{T} \left\| Y_{2D}^{j,t} - \overline{Y}_{2D}^{j,t} \right\|_{2}}_{JT} + \underbrace{\frac{1}{JT} \sum_{j=1}^{J} \sum_{t=1}^{T} \left| Y_{D}^{j,t} - \overline{Y}_{D}^{j,t} \right|}_{JT}}_{JT}$$
(8)

2D Pose Optimization Objective

Depth Optimization Objective

Where $Y_{2D}^{j,t}$ and $Y_D^{j,t}$ are the 2D pose and per-joint depth of 3D pose label. $\overline{Y}_{2D}^{j,t}$ and $\overline{Y}_D^{j,t}$ are the predicted results of the *j*-th joint in *t*-th frame. In addition, the temporal consistency loss \mathcal{L}_T from (Hossain & Little, 2018) is introduced to produce smooth poses. The total loss \mathcal{L} is defined as follows:

$$\mathcal{L} = \mathcal{L}_{3D} + \lambda_T \mathcal{L}_T + \lambda_{MI} \mathcal{L}_{MI} \tag{9}$$

328

330 331

332 333

334

where λ_T and λ_{MI} are hyper-parameters to balance the ratio of different loss terms.

5 EXPERIMENTS

5.1 EXPERIMENT SETTING

335 We evaluate our model on two large-scale monocular 3D human pose estimation datasets: Hu-336 man3.6M (Ionescu et al., 2013) and MPI-INF-3DHP (Mehta et al., 2017). For the Human3.6M 337 dataset, we report the MPJPE (Mean Per Joint Position Error) and P-MPJPE (Procrustes-MPJPE) as evaluation metrics as prior methods (Li et al., 2022b; Zhu et al., 2023; Zhang et al., 2022b; Zhao 338 et al., 2023b). For the MPI-INF-3DHP dataset, similar to existing approaches (Shan et al., 2022; 339 Tang et al., 2023; Chen et al., 2023; Zhu et al., 2023), we use ground truth 2D pose as input and 340 report MPJPE, Percentage of Correct Keypoint (PCK) with the threshold of 150mm, and Area Under 341 Curve (AUC) as the evaluation metrics. Please refer to Appendix B for implementation details. 342

Table 1: Results on Human3.6M in millimeters (mm) under MPJPE using 2D pose detected by
 SH (Newell et al., 2016) following MotionBERT (Zhu et al., 2023). T is the length of the input 2D pose sequence. The best result is shown in bold, and the second-best result is underlined.

0.40	_																	
346	MPJPE	Т	Dir.	Disc.	Eat	Greet	Phone	Photo	Pose	Pur.	Sit	SitD.	Smoke	Wait	WalkD.	Walk	WalkT.	Avg
347	MHFormer (Li et al., 2022b)	81	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	44.5
	MixSTE (Zhang et al., 2022b)	81	39.8	43.0	38.6	40.1	43.4	50.6	40.6	41.4	52.2	56.7	43.8	40.8	43.9	29.4	30.3	42.4
348	P-STMO (Shan et al., 2022)	81	41.7	44.5	41.0	42.9	46.0	51.3	42.8	41.3	54.9	61.8	45.1	42.8	43.8	30.8	30.7	44.1
	PoseFormerV2 (Zhao et al., 2023b)	81	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	46.0
349	STCFormer (Tang et al., 2023)	81	40.6	43.0	38.3	40.2	43.5	52.6	40.3	40.1	51.8	57.7	42.8	39.8	42.3	28.0	29.5	42.0
050	GLA-GCN (Yu et al., 2023)	81	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
350	MotionBERT (Zhu et al., 2023)	81	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
054	KTPFormer (Peng et al., 2024)	81	39.1	<u>41.9</u>	37.3	40.1	44.0	51.3	39.8	41.0	<u>51.4</u>	56.0	<u>43.0</u>	41.0	42.6	28.8	<u>29.5</u>	41.8
331	MotionAGFormer (Soroush Mehraban, 2024)	81	41.9	42.7	40.4	<u>37.6</u>	45.6	51.3	41.0	<u>38.0</u>	54.1	58.8	45.5	<u>40.4</u>	<u>39.8</u>	29.4	31.0	42.5
352	PrML (Ours)	81	<u>39.7</u>	41.4	39.4	35.5	43.1	<u>50.7</u>	<u>40.0</u>	37.2	51.1	56.0	43.7	40.5	39.1	28.7	28.8	41.0
002	MHFormer (Li et al., 2022b)	351	39.2	43.1	40.1	40.9	44.9	51.2	40.6	41.3	53.5	60.3	43.7	41.1	43.8	29.8	30.6	43.0
353	MixSTE (Zhang et al., 2022b)	243	37.6	40.9	37.3	39.7	42.3	49.9	40.1	39.8	51.7	55.0	42.1	39.8	41.0	27.9	27.9	40.9
	P-STMO (Shan et al., 2022)	243	38.9	42.7	40.4	41.1	45.6	49.7	40.9	39.9	55.5	59.4	44.9	42.2	42.7	29.4	29.4	42.8
354	PoseFormerV2 (Zhao et al., 2023b)	243	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	45.2
	STCFormer (Tang et al., 2023)	243	39.6	41.6	37.4	38.8	43.1	51.1	39.1	39.7	51.4	57.4	41.8	38.5	40.7	27.1	28.6	41.0
355	GLA-GCN (Yu et al., 2023)	243	41.3	44.3	40.8	41.8	45.9	54.1	42.1	41.5	57.8	62.9	45.0	42.8	45.9	29.4	29.9	44.4
	MotionBERT (Zhu et al., 2023)	243	<u>36.3</u>	38.7	38.6	33.6	42.1	50.1	36.2	35.7	<u>50.1</u>	56.6	41.3	<u>37.4</u>	<u>37.7</u>	25.6	26.5	39.2
356	KTPFormer (Peng et al., 2024)	243	37.3	39.2	35.9	37.6	42.5	48.2	38.6	39.0	51.4	55.9	41.6	39.0	40.0	27.0	27.4	40.1
	MotionAGFormer (Soroush Mehraban, 2024)	243	36.8	38.5	35.9	33.0	41.1	48.6	38.0	34.8	49.0	51.4	40.3	37.4	36.3	27.2	27.2	38.4
357	PrML (Ours)	243	36.0	38.2	<u>37.3</u>	<u>33.5</u>	40.4	46.9	<u>37.5</u>	34.6	48.9	<u>52.9</u>	<u>40.7</u>	36.6	<u>36.7</u>	<u>26.1</u>	26.1	38.2

Table 2: Results on Human3.6M in millimeters (mm) under MPJPE using ground truth 2D pose. T is
the number of input frames. Seq2seq refers to estimating 3D pose sequences rather than only the
center frame. MACs/frames represents multiply-accumulate operations for each output frame. The
best result is shown in bold, and the second-best result is underlined.

302									
363	Method	Venue	Framework	Seq2Seq	Т	Parameter	MACs	MACs/frame	MPJPE
	MHFormer (Li et al., 2022b)	CVPR'22	Lifting-Based	×	351	30.9M	7.1G	7096M	30.5
364	MixSTE (Zhang et al., 2022b)	CVPR'22	Lifting-Based	\checkmark	243	33.6M	139.0G	572M	21.6
265	P-STMO (Shan et al., 2022)	ECCV'22	Lifting-Based	×	243	6.2M	0.7G	740M	29.3
305	PoseFormerV2 (Zhao et al., 2023b)	CVPR'23	Lifting-Based	×	243	14.3M	0.5G	528M	-
366	STCFormer (Tang et al., 2023)	CVPR'23	Lifting-Based	\checkmark	243	4.7M	19.6G	80M	21.3
	GLA-GCN (Yu et al., 2023)	ICCV'23	Lifting-Based	×	243	1.3M	1.5G	1556M	21.0
367	MotionBERT (Zhu et al., 2023)	ICCV'23	Lifting-Based	\checkmark	243	42.5M	174.7G	719M	17.8
000	KTPFormer (Peng et al., 2024)	CVPR'24	Lifting-Based	\checkmark	243	33.7M	69.5G	286M	19.0
300	MotionAGFormer (Soroush Mehraban, 2024)	WACV'24	Lifting-Based	\checkmark	243	19.0M	78.3G	322M	17.3
369	PrML (Ours)	-	Multi-Task Learning	\checkmark	243	13.0M	49.3G	203M	17.2

³⁷⁰ 371

372

358

5.2 COMPARISON WITH STATE-OF-THE-ART METHODS

Human3.6M. We compare our method with several state-of-the-art techniques on the Human3.6M
dataset. For fair comparisons, only the results of models without extra pre-training on additional
data are included. Table 1 summarizes the performance comparisons in terms of MPJPE of all 15
actions, and the number of the input frames T is also given for each method. Our method achieved
state-of-the-art performance with an MPJPE of 38.2mm with T = 243. It is worth noting that our
method in the case of T = 81 input frames still achieves state-of-the-art performance with an MPJPE

378 error of 41.0mm and even surpasses the performance of several methods with a higher number of input 379 frames. For example, this result outperforms P-STMO (Shan et al., 2022) (41.0mm v.s. 42.8mm), 380 PoseformerV2 (Zhao et al., 2023b) (41.0mm v.s. 45.2mm) with 243 frames, and MHFormer (Li et al., 381 2022b) even with 351 frames (41.0mm v.s. 43.0mm). These results demonstrate the effectiveness of 382 PrML. To further validate the effectiveness of the multi-task learning framework, we also report the model parameters, MACs (Multiply-Accumulate Operations), and MPJPE using 2D ground truth 383 as input. As shown in Table 2, our method with T = 243 achieves the best performance with an 384 MPJPE of 17.2mm, which outperforms the lifting-based framework with faster inference speed. For 385 example, this result outperforms MotionBERT (Zhu et al., 2023) (17.2mm v.s. 17.8mm). Due to 386 space limitations, we present the results of P-MPJPE (Procrustes-MPJPE) in Appendix C. 387

388 MPI-INF-3DHP. To demonstrate the generalization capability of our 389 model, we also evaluate our model 390 on the challenging MPI-INF-3DHP 391 dataset, which includes more com-392 plex scenes and motions. Following 393 previous works (Zheng et al., 2021; 394 Zhang et al., 2022b; Shan et al., 2022; 395 Tang et al., 2023; Li et al., 2022b; 396 2024), we use ground truth 2D pose 397 as input and set the number of input 398 frames as 9, 27, or 81. As observed 399 in Table 3, our method with T = 81achieves the best performance with 400

Table 3: Results on MPI-INF-3DHP dataset under PCK, AUC, and MPJPE using ground truth 2D pose as input. T is the number of input frames. Seq2seq refers to estimating 3D pose sequence. (*) indicate our re-implementation.

Method	Т	Seq2Seq	PCK↑	AUC↑	MPJPE↓
MHFormer (Li et al., 2022b)	9	×	93.8	63.3	58.0
MixSTE (Zhang et al., 2022b)	27	\checkmark	94.4	66.5	54.9
P-STMO (Shan et al., 2022)	81	×	97.9	75.8	32.2
PoseFormerV2 (Zhao et al., 2023b)	81	×	97.9	78.8	27.8
GLA-GCN (Yu et al., 2023)	81	×	98.5	79.1	27.8
STCFormer (Tang et al., 2023)	81	\checkmark	98.7	83.9	23.1
MotionBERT* (Zhu et al., 2023)	81	\checkmark	98.7	85.6	16.5
KTPFormer (Peng et al., 2024)	81	\checkmark	98.9	<u>85.9</u>	16.7
MotionAGFormer (Soroush Mehraban, 2024)	81	\checkmark	98.2	85.3	<u>16.2</u>
PrML (Ours)	9	\checkmark	98.1	82.2	23.3
PrML (Ours)	27	\checkmark	98.6	85.8	18.1
PrML (Ours)	81	\checkmark	98.9	86.9	15.7

the PCK of 98.9%, AUC of 86.9%, and MPJPE of 15.7mm. Similar to the previous findings, our method with T = 9, 27 input frames still outperforms the previous state-of-the-art methods and achieves the MPJPE of 23.3mm and 18.1mm, respectively. More remarkably, our method with T = 9input frames outperforms the GLA-GCN (Yu et al., 2023) with T = 81 input frames, despite having only one-ninth of the input frames (23.3mm v.s. 27.8mm, 9 frames vs. 81 frames).

Robustness to Noisy 2D Pose. Benefiting from the 2D pose branch, our method not only preserves 406 well-detected 2D pose features but also allows us to handle noisy 2D pose input. To demonstrate that 407 the inclusion of the 2D pose branch helps improve the robustness of the proposed method, we make 408 the pose estimation task more challenging by adding zero-mean Gaussian noise to the ground-truth 409 2D pose on the Human3.6M (Ionescu et al., 2013) and MPI-INF-3DHP (Mehta et al., 2017). As 410 shown in Figure 5, the experimental evidence reveals that our proposed PrML suffers from less 411 performance drop as the standard deviation of Gaussian noise (sigma) increases compared with the 412 powerful lifting-based method MotionBERT (Zhu et al., 2023) while being more efficient. 413



Figure 5: Comparisons of PrML and MotionBERT (Zhu et al., 2023) in terms of robustness to noise on Human3.6M and MPI-INF-3DHP datasets. Zero-mean Gaussian noise of standard deviation sigma is added to ground truth 2D pose, and we show their performance drop (Δ MPJPE, in millimeters) as sigma increases. The size of markers indicates the computational cost of models.

428 5.3 ABLATION STUDY

429

423

424

425

426

427

430 We perform extensive ablation studies focused on analyzing the contribution of each component in 431 our proposed PrML. Experiments are conducted on the Human3.6M (Ionescu et al., 2013) dataset with T = 243 as the number of input frames and MPJPE is used as the evaluation metric. Analysis on Effectiveness of Components. The results in Section 5.2 have demonstrated that our framework achieves better results compared to the conventional lifting-based framework in terms of accuracy and robustness with fewer parameters. In this section, we show how to construct our proposed progressive multi-task learning pose estimation framework step by step.

- **Baseline:** Most multi-task learning models have a shared-bottom model structure following (Caruana, 1997) which substantially reduces the risk of overfitting (Ma et al., 2018). We follow this design and use the DSTFormer block from MotionBERT (Zhu et al., 2023) as the shared bottom. As shown in Table 4, this structure achieves 50.6mm MPJPE and serves as our baseline.
- **Multi-Task Branch:** To reduce the explicit influence of uncertain depth features on the welldetected 2D pose features, we introduce the multi-task branch design to refine the 2D pose features and learn the depth features separately. Based on the shared bottom, we incorporate the 2D pose branch and depth branch respectively bringing 6.2mm and 9.0mm error reduction (see also Table 4). With both 2D pose and depth branches, we achieve the MPJPE of 39.1mm.
- **Task-Aware Decoder:** As shown in Table 4, we achieve a reduction in MPJPE from 50.6mm to 46.2mm by adding task-aware decoder to the shared bottom. The performance is further improved to 38.6mm when task-aware decoder is introduced in conjunction with the multi-task branch.
- **PrML:** By introducing the mutual information loss to explicitly supervise the feature complement and the learning of bias, our TAD module further resulted in a 2.2mm error reduction (from 46.2mm to 44.0mm). After incorporating all components, we obtain the complete version of our PrML, which achieves the best performance with an MPJPE of 38.2mm.

Model Setting	Shared Bottom	2D Pose Branch	Depth Branch	Task-Aware Decoder	Mutual Information	MPJPE \downarrow
Baseline (Shared Bottom)	\checkmark	-	-	-	-	50.6
+ 2D Pose Branch Only	\checkmark	\checkmark	-	-	-	44.4 (-6.2)
+ Depth Branch Only	\checkmark	-	\checkmark	-	-	41.6 (-9.0)
+ TAD Only	\checkmark	-	-	\checkmark	-	46.2 (-4.4)
+ TAD + MI Loss	\checkmark	-	-	\checkmark	\checkmark	44.0 (-6.6)
+ Multi-Task Branch	\checkmark	\checkmark	\checkmark	-	-	39.1 (-11.5)
+ Multi-Task Branch + TAD	\checkmark	\checkmark	\checkmark	\checkmark	-	38.6 (-12.0)
PrML	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	38.2 (-12.4)

Table 4: Analysis of the effectiveness of each component within PrML.

Analysis on Task-Aware Decoder (TAD). We first analyze the effectiveness of each operation in TAD and report performance in Table 5. By incorporating coarse alignment and task bias separately, we achieve the results of 39.5mm and 39.1mm. When both of them are incorporated, we obtain the best performance with an MPJPE of 38.2mm. We also examine the impact of different task biases within TAD. As shown in Table 6, the best results are achieved when both the 2D pose bias and per-joint bias are introduced and made learnable during the training process.

Table	e 5: Analysis of	various desig	Table 6: Analysis of task biases within TAD.						
Step	Feature Complement	Coarse Alignment	Task Bias	MPJPE	Step	2D Pose Bias	Per-Joint Depth Bias	Learnable	MPJPE
1	√	-	-	39.5	1	√	-	~	38.7
2	\checkmark	\checkmark	-	39.1	2	-	\checkmark	\checkmark	38.5
3	\checkmark	-	\checkmark	38.7	3	\checkmark	\checkmark	-	38.5
TAD	\checkmark	\checkmark	\checkmark	38.2	TAD	\checkmark	\checkmark	\checkmark	38.2

Multi-Task Learning for Lifting-Based Methods. An alternative to our framework is to directly copy a branch from the lifting-based meth-ods to construct the multi-task learn-ing framework. However, this would increase the number of parameters substantially and lead to an unfair comparison. In light of this, we per-

Table 7: Analysis on the generalization of multi-task learningframework. (*) denotes our re-implementation.

Method	Framework	\mid MPJPE \downarrow
MixSTE* (Zhang et al., 2022b)	Lifting-Based	40.9
+ Two Regression Heads	Multi-Task Learning	40.0 (-0.9)
MotionBERT* (Zhu et al., 2023)	Lifting-Based	39.8
+ Two Regression Heads	Multi-Task Learning	38.9 (- <mark>0.9</mark>)
CA-PF* (Zhao et al., 2023a)	Lifting-Based	41.4
+ Two Regression Heads	Multi-Task Learning	40.2 (-1.2)

form an embarrassingly simple transformation: replacing the original single regression head of
the lifting-based framework with two regression heads. This replacement transforms the liftingbased framework into a hard parameter sharing multi-task learning framework (Ruder, 2017). As
shown in Table 7, such simple modification leads to performance improvement in both multi-frames
(MixSTE (Zhang et al., 2022b), MotionBERT (Zhu et al., 2023)) and single-frame (CA-PF (Zhao et al., 2023a)) lifting-based methods. Due to space limitations, we present more details in Appendix C.1.

486 5.4VISUALIZATION 487

488 Feature Distribution Visualization. To further analyze the difference between the 2D pose features 489 and per-joint features, we utilize t-SNE (Van der Maaten & Hinton, 2008) to visualize their feature distributions. We select samples from Human3.6M (Ionescu et al., 2013) and visualize the features 490 before regression heads (i.e., F_{2D} and F_D). As shown in Figure 6, the distributions of the 2D features 491 and depth features are different across various situations. These qualitative results provide strong 492 evidence that we should not simply encode the 2D pose features and per-joint depth features in an 493 entangled feature space. Due to space limitations, we present more visualization in our Appendix F. 494



504

505

506 507

509

510

521 522

523

524 525 526

527



Figure 6: Feature distributions visualization of 2D features (green) and depth features (purple) using t-SNE (Van der Maaten & Hinton, 2008) method on Human3.6M (Ionescu et al., 2013) dataset. The distributions of the 2D features and depth features are different across various situations.

3D Human Pose Estimation Visualization. We present 3D human pose estimation results by our proposed PrML and MotionBERT (Zhu et al., 2023). As shown in Figure 7, our method generalizes well to in-the-wild videos including self-occlusion and fast motion.

511									
512		900 manada 900 *	000 mmmmm 000 . 1	- 000	- 00 00	PR	PR	PR	. PR
513	Input 2D Pose								-2
514		A	A	N I	N				
515		TTTE DA				+TI	+	HTTTTTT	HTTTED.
516	MotionBERT					5	5	5	5
517	mouoliblitti								
518									
519									
520	PrML								A B
521		A BEEN	A BEEN	A BERT	A BEEN	H BY BY	A CS SA	A BER	HAS 3

Figure 7: Qualitative comparisons of PrML with MotionBERT (Zhu et al., 2023). The green cycle indicates locations where our method achieves better results. See Appendix F for more comparison.

CONCLUSION 6

528 This work presents a novel progressive multi-task learning framework named PrML for monocular 529 3D human pose estimation. Our framework addresses the limitation of the lifting-based framework 530 that neglects different initial states between a well-detected 2D pose and an unknown per-joint depth. 531 PrML first learns 2D pose features and per-joint depth features separately by multi-task branch design and then employs a task-aware decoder to indirectly supplement information between the refined 532 2D pose features and the learned per-joint depth features. We also propose a mutual information 533 loss to supervise the feature complementary process. Extensive quantitative experimental results on 534 the Human3.6M and MPI-INF-3DHP datasets show that our PrML outperforms the conventional lifting-based framework in terms of accuracy and robustness with fewer parameters. 536

Future Work. The core contribution of our work is providing a new framework for monocular 3D human pose estimation. To this end, we use the widely used spatial and temporal transformer as our 538 encoder to ensure a fair comparison with the lifting-based framework. It will be novel and interesting to design specific encoders for different tasks to extend our framework in future research.

540	REFERENCES
541	

550

551

554

586

542	Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation:
543	New benchmark and state of the art analysis. In CVPR, June 2014.

- Suzanna Becker. An information-theoretic unsupervised learning algorithm for neural networks. University of Toronto, 1992.
- Suzanna Becker. Mutual information maximization: models of cortical self-organization. *Network: Computation in neural systems*, 7(1):7, 1996.
 - David Brüggemann, Menelaos Kanakis, Anton Obukhov, Stamatios Georgoulis, and Luc Van Gool. Exploring relational context for multi-task dense prediction. In *ICCV*, pp. 15869–15878, 2021.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey
 Zagoruyko. End-to-end object detection with transformers. In *ECCV*, pp. 213–229, 2020.
- R Caruana. Multitask learning: A knowledge-based source of inductive bias1. In *ICML*, pp. 41–48.
 Citeseer, 1993.
 - Rich Caruana. Multitask learning. *Machine learning*, 28:41–75, 1997.
- Hanyuan Chen, Jun-Yan He, Wangmeng Xiang, Zhi-Qi Cheng, Wei Liu, Hanbing Liu, Bin Luo,
 Yifeng Geng, and Xuansong Xie. HDFormer: High-order directed transformer for 3D human pose
 estimation. In *IJCAI*, pp. 581–589, 2023.
- Tianlang Chen, Chen Fang, Xiaohui Shen, Yiheng Zhu, Zhili Chen, and Jiebo Luo. Anatomy-aware
 3D human pose estimation with bone-based pose decomposition. *IEEE TCSVT*, 32(1):198–209, 2021.
- Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gang Yu, and Jian Sun. Cascaded
 pyramid network for multi-person pose estimation. In *CVPR*, pp. 7103–7112, 2018.
- Hai Ci, Chunyu Wang, Xiaoxuan Ma, and Yizhou Wang. Optimizing network structure for 3D human pose estimation. In *ICCV*, pp. 2262–2271, 2019.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.
 In *ICLR*, 2021.
- Hao-Shu Fang, Yuanlu Xu, Wenguan Wang, Xiaobai Liu, and Song-Chun Zhu. Learning pose grammar to encode human body configuration for 3D pose estimation. In *AAAI*, volume 32, 2018.
- Runyang Feng, Yixing Gao, Xueqing Ma, Tze Ho Elden Tse, and Hyung Jin Chang. Mutual information-based temporal difference learning for human pose estimation in video. In *CVPR*, pp. 17131–17141, 2023.
- Lin Geng Foo, Tianjiao Li, Hossein Rahmani, Qiuhong Ke, and Jun Liu. Unified pose sequence modeling. In *CVPR*, pp. 13019–13030, 2023.
- Jia Gong, Zhipeng Fan, Qiuhong Ke, Hossein Rahmani, and Jun Liu. Meta agent teaming active
 learning for pose estimation. In *CVPR*, pp. 11079–11089, 2022.
- 587 Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, pp. 2961–2969, 2017.
- ⁵⁸⁹ R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *ICLR*, 2018.
- 593 Mir Rayat Imtiaz Hossain and James J Little. Exploiting temporal information for 3D human pose estimation. In ECCV, pp. 68–84, 2018.

594 595	Lukas Hoyer, Dengxin Dai, Yuhua Chen, Adrian Koring, Suman Saha, and Luc Van Gool. Three ways to improve semantic segmentation with self-supervised depth estimation. In <i>CVPR</i> , pp.
596	11130–11140 2021
597	11150 11110, 2021.
598	Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale
599	datasets and predictive methods for 3D human sensing in natural environments. IEEE TPAMI, 36
600	(7):1325–1339, 2013.
601	Uner John Daule Malakanan Thomas David Lucasa Call and Ian Kauta Hand access the
602	via latent 2.5d heatman regression. In ECCV Sentember 2018
603	via latent 2.50 heatinap regression. In ECCV, September 2018.
604	Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of
605	human shape and pose. In CVPR, pp. 7122–7131, 2018.
606	
607	Diederik P Kingma. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
608	Wei-Hong Li Xialei Liu and Hakan Bilen. Learning multiple dense prediction tasks from partially
609	annotated data. In <i>CVPR</i> , pp. 18879–18889, 2022a.
610	
611	Wenhao Li, Hong Liu, Hao Tang, Pichao Wang, and Luc Van Gool. MHFormer: Multi-hypothesis
612	transformer for 3D human pose estimation. In CVPR, pp. 13147–13156, 2022b.
613	Wanhaa Li Mangunan Lin Hang Lin Diahaa Wang Lialun Cai and Nien Saka Haunglass takanigan
614	for efficient transformer based 3D human pose estimation. In CVPR 2024
615	for encient transformer based 3D numan pose estimation. In CVT R, 2024.
616	Xiangtai Li, Henghui Ding, Haobo Yuan, Wenwei Zhang, Jiangmiao Pang, Guangliang Cheng, Kai
617	Chen, Ziwei Liu, and Chen Change Loy. Transformer-based visual segmentation: A survey. arXiv
618	preprint arXiv:2304.09854, 2023.
619	Delah Linghan Calf angenization in a generated actived. Commuter 21(2):105, 117, 1088
620	Raiph Linsker. Self-organization in a perceptual network. Computer, 21(3):105–117, 1988.
621	Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for
622	multi-task learning. NeurIPS, 34:18878–18890, 2021a.
623	
624	Kenkun Liu, Rongqi Ding, Zhiming Zou, Le Wang, and Wei Tang. A comprehensive study of weight
625	sharing in graph networks for 3D human pose estimation. In ECCV, pp. 318–334, 2020.
626	Zhenguang Liu, Pengxiang Su, Shuang Wu, Xuaniing Shen, Haipeng Chen, Yanbin Hao, and Meng
627	Wang. Motion prediction using trajectory cues. In <i>ICCV</i> , pp. 13299–13308, 2021b.
628	
620	Zhenguang Liu, Runyang Feng, Haoming Chen, Shuang Wu, Yixing Gao, Yunjun Gao, and Xiang
630	Wang. Temporal feature alignment and mutual information maximization for video-based human
622	pose estimation. In <i>CVPR</i> , pp. 11006–11016, 2022a.
632	Zhenguang Liu, Shuang Wu, Shuyuan Jin, Shouling Ji, Oi Liu, Shijian Lu, and Li Cheng. Investigating
624	pose representations and motion contexts modeling for 3D motion prediction. <i>IEEE TPAMI</i> . 45(1):
635	681–697, 2022b.
636	
637	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint
638	arXiv:1/11.05101, 2017.
639	Jiagi Ma Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H Chi, Modeling task relationships
640	in multi-task learning with multi-gate mixture-of-experts. In ACM KDD, pp. 1930–1939 2018
641	6
642	Xiaoxuan Ma, Jiajun Su, Chunyu Wang, Hai Ci, and Yizhou Wang. Context modeling in 3d human
643	pose estimation: A unified perspective. In CVPR, pp. 6238–6247, June 2021a.
644	Visovuon Ma Jisiun Su Chunya Wang Hai Ci and Vizhay Wang Contact modeling in 2D human
645	nose estimation: A unified perspective In CVPR pp. 6238-6247 2021b
646	pose estimation. It unified perspective. In CVI A, pp. 0250-0247, 20210.
	Inlight Martinez, Payet Hossein, Javier Pomero, and James I. Little. A simple yet offective baseline

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

648 649 650	Soroush Mehraban, Yiqian Qin, and Babak Taati. Evaluating recent 2d human pose estimators for 2d-3d pose lifting. In 2024 IEEE 18th International Conference on Automatic Face and Gesture Recognition (FG), pp. 1–5. IEEE, 2024.
651 652 653 654	Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. Monocular 3D human pose estimation in the wild using improved CNN supervision. In <i>3DV</i> , pp. 506–516, 2017.
655 656	Gyeongsik Moon and Kyoung Mu Lee. I2I-meshnet: Image-to-lixel prediction network for accurate 3D human pose and mesh estimation from a single rgb image. In <i>ECCV</i> , pp. 752–768, 2020.
657 658 659	Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In <i>ECCV</i> , pp. 483–499, 2016.
660 661	Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. Coarse-to-fine volumetric prediction for single-image 3D human pose. In <i>CVPR</i> , pp. 7025–7034, 2017.
662 663 664	Georgios Pavlakos, Xiaowei Zhou, and Kostas Daniilidis. Ordinal depth supervision for 3D human pose estimation. In <i>CVPR</i> , pp. 7307–7316, 2018.
665 666	Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3D human pose estimation in video with temporal convolutions and semi-supervised training. In <i>CVPR</i> , pp. 7753–7762, 2019.
667 668 669	Jihua Peng, Yanghong Zhou, and PY Mok. Ktpformer: Kinematics and trajectory prior knowledge- enhanced transformer for 3d human pose estimation. <i>arXiv preprint arXiv:2404.00658</i> , 2024.
670 671	Tomas Pfister, James Charles, and Andrew Zisserman. Flowing convnets for human pose estimation in videos. In <i>ICCV</i> , pp. 1913–1921, 2015.
673 674 675	Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. <i>arXiv preprint arXiv:2307.01952</i> , 2023.
676 677	Marc Proesmans, Dengxin Dai, and Luc Van Gool. Multi-task learning for dense prediction tasks: A survey. <i>IEEE TPAMI</i> , 44(7), 2022.
678 679 680	Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Learning multiple visual domains with residual adapters. <i>NeurIPS</i> , 30, 2017.
681 682	Sebastian Ruder. An overview of multi-task learning in deep neural networks. <i>arXiv preprint arXiv:1706.05098</i> , 2017.
683 684 685	Wenkang Shan, Zhenhua Liu, Xinfeng Zhang, Shanshe Wang, Siwei Ma, and Wen Gao. P-STMO: Pre-trained spatial temporal many-to-one model for 3D human pose estimation. In <i>ECCV</i> , 2022.
686 687 688 689	Wenkang Shan, Zhenhua Liu, Xinfeng Zhang, Zhao Wang, Kai Han, Shanshe Wang, Siwei Ma, and Wen Gao. Diffusion-based 3D human pose estimation with multi-hypothesis aggregation. <i>arXiv</i> preprint arXiv:2303.11579, 2023.
690 691	Babak Taati Soroush Mehraban, Vida Adeli. Motionagformer: Enhancing 3d human pose estimation with a transformer-genformer network. In <i>WACV</i> , 2024.
692 693	Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In <i>CVPR</i> , pp. 5693–5703, 2019.
695 696	Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. Integral human pose regression. In <i>ECCV</i> , pp. 529–545, 2018.
697 698 699	Zhenhua Tang, Zhaofan Qiu, Yanbin Hao, Richang Hong, and Ting Yao. 3D human pose estimation with spatio-temporal criss-cross attention. In <i>CVPR</i> , pp. 4790–4799, 2023.
700 701	Xudong Tian, Zhizhong Zhang, Shaohui Lin, Yanyun Qu, Yuan Xie, and Lizhuang Ma. Farewell to mutual information: Variational distillation for cross-modal person re-identification. In <i>CVPR</i> , pp. 1522–1531, 2021.

702 703 704	Jonathan J Tompson, Arjun Jain, Yann LeCun, and Christoph Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. <i>NeurIPS</i> , 27, 2014.
705	Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. JMLR, 9(11), 2008.
706 707 708	Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Mti-net: Multi-scale task interaction networks for multi-task learning. In <i>ECCV</i> , pp. 527–543. Springer, 2020.
709 710	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>NeurIPS</i> , pp. 5998–6008, 2017.
711 712 713	Zirui Wang, Yulia Tsvetkov, Orhan Firat, and Yuan Cao. Gradient vaccine: Investigating and improv- ing multi-task optimization in massively multilingual models. <i>arXiv preprint arXiv:2010.05874</i> , 2020.
714 715 716	Tom Wehrbein, Marco Rudolph, Bodo Rosenhahn, and Bastian Wandt. Probabilistic monocular 3D human pose estimation with normalizing flows. In <i>ICCV</i> , pp. 11199–11208, 2021.
717 718	Laurenz Wiskott and Terrence J Sejnowski. Slow feature analysis: Unsupervised learning of invariances. <i>Neural computation</i> , 14(4):715–770, 2002.
719 720 721 722	Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 10371–10381, 2024.
723 724	Mang Ye, He Li, Bo Du, Jianbing Shen, Ling Shao, and Steven CH Hoi. Collaborative refining for person re-identification with label noise. <i>IEEE TIP</i> , 31:379–391, 2021.
725 726 727 728	Bruce XB Yu, Zhi Zhang, Yongxu Liu, Sheng-hua Zhong, Yan Liu, and Chang Wen Chen. Gla-gcn: Global-local adaptive graph convolutional network for 3D human pose estimation from monocular video. In <i>ICCV</i> , pp. 8818–8829, 2023.
729 730	Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. <i>NeurIPS</i> , 33:5824–5836, 2020.
731 732 733 734	Zhenbo Yu, Bingbing Ni, Jingwei Xu, Junjie Wang, Chenglong Zhao, and Wenjun Zhang. Towards alleviating the modeling ambiguity of unsupervised monocular 3d human pose estimation. In <i>ICCV</i> , pp. 8651–8660, 2021.
735 736	Ailing Zeng, Xiao Sun, Lei Yang, Nanxuan Zhao, Minhao Liu, and Qiang Xu. Learning skeletal graph neural networks for hard 3d pose estimation. In <i>ICCV</i> , pp. 11436–11445, 2021.
737 738 720	Can Zhang, Tianyu Yang, Junwu Weng, Meng Cao, Jue Wang, and Yuexian Zou. Unsupervised pre-training for temporal action localization tasks. In <i>CVPR</i> , pp. 14031–14041, 2022a.
740 741 742	Jinlu Zhang, Zhigang Tu, Jianyu Yang, Yujin Chen, and Junsong Yuan. MixSTE: Seq2seq mixed spatio-temporal encoder for 3D human pose estimation in video. In <i>CVPR</i> , pp. 13232–13242, 2022b.
743	Yu Zhang and Qiang Yang. A survey on multi-task learning. IEEE TKDE, 34(12):5586–5609, 2021.
744 745 746	Long Zhao, Xi Peng, Yu Tian, Mubbasir Kapadia, and Dimitris N Metaxas. Semantic graph convolutional networks for 3D human pose regression. In <i>CVPR</i> , pp. 3425–3435, 2019.
747 748 749 750	Long Zhao, Yuxiao Wang, Jiaping Zhao, Liangzhe Yuan, Jennifer J Sun, Florian Schroff, Hartwig Adam, Xi Peng, Dimitris Metaxas, and Ting Liu. Learning view-disentangled human pose representation by contrastive cross-view mutual information maximization. In <i>CVPR</i> , pp. 12793–12802, 2021.
751 752 753	Qitao Zhao, Ce Zheng, Mengyuan Liu, and Chen Chen. A single 2D pose with context is worth hundreds for 3D human pose estimation. In <i>NeurIPS</i> , 2023a.
754 755	Qitao Zhao, Ce Zheng, Mengyuan Liu, Pichao Wang, and Chen Chen. PoseFormerV2: Exploring frequency domain for efficient and robust 3D human pose estimation. In <i>CVPR</i> , pp. 8877–8886, 2023b.

756 757 758	Ce Zheng, Sijie Zhu, Matias Mendieta, Taojiannan Yang, Chen Chen, and Zhengming Ding. 3D human pose estimation with spatial and temporal transformers. In <i>ICCV</i> , pp. 11656–11665, 2021.
759	Wentao Zhu, Xiaoxuan Ma, Zhaoyang Liu, Libin Liu, Wayne Wu, and Yizhou Wang. MotionBERT:
760	A unified perspective on learning human motion representations. In <i>ICCV</i> , pp. 15085–15099, 2023.
761	
762	
763	
764	
765	
766	
767	
768	
769	
770	
771	
772	
773	
774	
775	
776	
777	
778	
779	
780	
781	
782	
783	
704	
700	
787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
807	
808	
009	

810 A APPENDIX 811

The Appendix is organized as follows:

- Section **B**: Experiment Setting
- Section C: Additional Experiment Analysis
- Section D: Limitation
- Section E: Additional Analysis of Mutual Information
- Section F: Additional Visualization
- 822 823 824

825 826

827

812

813 814

815

816 817

818 819

820 821

B EXPERIMENT SETTING

B.1 DATASETS AND EVALUATION METRICS

Human3.6M (Ionescu et al., 2013) is the most popular benchmark for indoor 3D human pose estimation, which contains approximately 3.6 million frames captured by 4 cameras at different views. This dataset contains 11 subjects performing 15 typical actions (e.g., walking and sitting). To ensure a fair comparison, we follow previous methods (Zheng et al., 2021; Zhang et al., 2022b; Tang et al., 2023; Zhu et al., 2023) by using subjects 1, 5, 6, 7, and 8 for model training and subjects 9 and 11 for evaluation.

MPI-INF-3DHP (Mehta et al., 2017) is a recently proposed large-scale challenging dataset with both indoor and outdoor scenes. The training set comprises 8 subjects, covering 8 activities, ranging from walking and sitting to complex exercise poses and dynamic actions. The test set covers 7 activities, containing three scenes: green screen, non-green screen, and outdoor environments. It complements existing test sets with more diverse motions (standing/walking, sitting/reclining, exercise, sports (dynamic poses), on the floor, dancing/miscellaneous).

840 Evaluation Metrics. For the Human3.6M dataset, we use two common evaluation metrics: MPJPE 841 and P-MPJPE. MPJPE (Mean Per Joint Position Error) is computed as the mean Euclidean distance 842 between the estimated joints and the ground truth in millimeters after aligning their root joints (hip). 843 P-MPJPE (Procrustes-MPJPE) is the MPJPE after the estimated joints align to the ground truth via a 844 rigid transformation. For the MPI-INF-3DHP dataset, following previous works (Shan et al., 2022; 845 Tang et al., 2023; Chen et al., 2023; Zhu et al., 2023), we use ground truth 2D pose as input and 846 report MPJPE, Percentage of Correct Keypoint (PCK) with the threshold of 150mm, and Area Under 847 Curve (AUC) as the evaluation metrics.

848 849

850 851

B.2 IMPLEMENTATION DETAILS

Our model is implemented using PyTorch and executed on a server equipped with 2 NVIDIA 3090 852 GPUs. We apply horizontal flipping augmentation for both training and testing following (Tang 853 et al., 2023; Zhu et al., 2023; Foo et al., 2023; Zhao et al., 2023a). For model training, we set each 854 mini-batch as 16 sequences. The network parameters are optimized using AdamW (Loshchilov & 855 Hutter, 2017) optimizer over 90 epochs with a weight decay of 0.01. The initial learning rate is set to 856 5e-4 with an exponential learning rate decay schedule and the decay factor is 0.99. In the experiments 857 on Human3.6M, two kinds of input are utilized, including the 2D ground truth and the Stacked 858 Hourglass (Newell et al., 2016) 2D pose detection, following (Zhu et al., 2023; Ci et al., 2019). For 859 MPI-INF-3DHP, 2D ground truth is used following previous works (Zheng et al., 2021; Zhang et al., 860 2022b; Shan et al., 2022; Zhu et al., 2023). While our proposed framework is capable of adapting to 861 input sequences of any length, to be fair, we choose specific input sequence lengths (denoted as T) for two datasets to compare our method with other approaches that have a certain 2D input length (Zheng 862 et al., 2021; Zhang et al., 2022b; Shan et al., 2022; Tang et al., 2023): Human3.6M (T = 81, 243), 863 MPI-INF-3DHP (T = 9, 27, 81).



Figure 8: Previous pipeline (MixSTE (Zhang et al., 2022b) or MotionBERT (Zhu et al., 2023)): Using a single regression head to directly estimate the 3D pose. Improved pipeline: Replacing the original single regression head with two regression heads to transform the existing lifting framework into a hard parameter sharing multi-task learning framework (Ruder, 2017).

C ADDITIONAL EXPERIMENT ANALYSIS

C.1 MULTI-TASK LEARNING FOR LIFTING-BASED METHODS.

As shown in Figure 8, the 2D-to-

889 3D lifting process employs a sin-890 gle regression head to directly es-891 timate the 3D pose after extracting spatio-temporal information. We 892 transform the existing 2D-to-3D lift-893 ing framework into a hard parame-894 ter sharing multi-task learning frame-895 work (Ruder, 2017) by replacing the 896 original single regression head with 897 two regression heads. We regress the 898

Table 8: Analysis on the generalization of multi-task learning framework. (*) denotes our re-implementation.

Method	Framework	MPJPE
MixSTE* (Zhang et al., 2022b)	Lifting-Based	40.9
+ Two Regression Heads	Multi-Task Learning	40.0↓0.9
MotionBERT* (Zhu et al., 2023)	Lifting-Based	39.8
+ Two Regression Heads	Multi-Task Learning	38.9↓0.9
CA-PF* (Zhao et al., 2023a)	Lifting-Based	41.4
+ Two Regression Heads	Multi-Task Learning	$40.2 \downarrow 1.2$

2D pose and per-joint depth separately. Then, we concatenate them together to obtain the 3D pose. As
shown in Table 8, such simple modification leads to improvement in multi-frames (MixSTE (Zhang
et al., 2022b), MotionBERT (Zhu et al., 2023)) and single-frame (CA-PF (Zhao et al., 2023a))
lifting-based methods.

Table 9: Analysis on the initial distribution of Task Bias within Task-Aware Decoder. G and L represent Gaussian distribution and Laplacian distribution respectively.

Step	2D Pose Bias	Per-Joint Depth Bias	MPJPE
1	L	G	38.9
2	G	L	38.6
3	L	L	38.4
Ours	G	G	38.2

Table 10: Analysis on various micro designs within transformer block. S and T denote Spatial Transformer and Temporal Transformer.

Step	S	Т	T–S	$S{\rightarrow}T$	$T {\rightarrow} S$	MPJPE
1	✓	.(40.6
3		v	\checkmark			39.0
4				\checkmark		38.3
Ours					\checkmark	38.2

912

913

902 903

904

905

906

879

880

881

882 883 884

885

887

914

915

C.2 ANALYSIS ON TASK BIAS WITHIN TASK-AWARE DECODER.

In this section, we investigate the impact of different initial distributions of task bias on the model.
 As shown in Table 9, we use Gaussian or Laplace distribution as our initial distributions. We achieve a result of 38.9mm when initializing the 2D pose bias with Laplace distribution and the

per-joint depth with Gaussian distribution. The MPJPE decreased from 38.9mm to 38.6mm when
 we initialized the 2D pose bias with Gaussian distribution and the per-joint depth with Laplace
 distribution. Experimental results revealed that employing Gaussian distribution for both biases
 yielded the best performance.

C.3 ANALYSIS ON STRUCTURE DESIGN OF TRANSFORMER BLOCK.

We conducted experiments on Human3.6M for five different types, including $S, T, S - T, S \rightarrow T$, $T \rightarrow S$ in the Transformer layer, where S and T denote Spatial Transformer Encoder and Temporal Transformer Encoder, respectively. S - T is a two-stream layer that uses Spatial Transformer and Temporal Transformer Encoder simultaneously. $S \rightarrow T$ means Spatial Transformer first and then Temporal Transformer; the same goes for $T \rightarrow S$. The results in Table 10 show that the $T \rightarrow S$ variant has improved by 2.4mm (from 40.6mm to 38.2mm), 1.4mm (from 39.6mm to 38.2mm), 0.8mm (from 39.0mm to 38.2mm), and 0.1mm (from 38.3mm to 38.2mm), respectively, compared to the other four variants.

932 933

934

923

924

C.4 ANALYSIS ON HYPERPARAMETER SETTING.

935 We consider the dimension C of hidden feature, and the weight λ_{MI} of mutual information loss as free parameters. 936 Table 11 shows the results of different hyperparameter 937 settings. We divide the configurations into 2 groups by 938 row, and allocate different values for one hyperparameter 939 in each group, while keeping the other hyperparameter 940 fixed, to evaluate the impact of each configuration. The 941 best results were obtained when the C was 128 and λ_{MI} 942 was 0.01. Based on the results shown in the table, we 943 chose the combination of C = 128, and λ_{MI} = 0.01 as our 944 configuration. 945

Table 11: Analysis on the	e hyperparame-
ter setting.	

Dimension (C)	λ_{MI}	MPJPE
64	0.01	40.4
96	0.01	39.3
128	0.01	38.2
128	1	38.6
128	0.1	38.5
128	0.01	38.2

946 C.5 Additional Quantitative Results.947

In this section, we provide the results of P-MPJPE on the Human3.6M (Ionescu et al., 2013) dataset.
As shown in Table 12, our proposed PrML achieved the best performance with the MPJPE of 38.2mm and the third-best results with the P-MPJPE of 32.3mm.

Table 12: Results on Human3.6M in millimeters (mm) under MPJPE. T is the number of input frames.
Seq2seq refers to estimating 3D pose sequences rather than only the center frame. MACs/frames
represents multiply-accumulate operations for each output frame. The best result is shown in bold, and the second-best result is underlined.

Method	Venue	Framework	Seq2Seq	Т	Parameter	MACs	MACs/frame	MPJPE	P-MPJPE
MHFormer (Li et al., 2022b)	CVPR'22	Lifting-Based	×	351	30.9M	7.1G	7096M	43.0	34.4
MixSTE (Zhang et al., 2022b)	CVPR'22	Lifting-Based	\checkmark	243	33.6M	139.0G	572M	40.9	32.6
P-STMO (Shan et al., 2022)	ECCV'22	Lifting-Based	×	243	6.2M	0.7G	740M	42.8	34.4
PoseFormerV2 (Zhao et al., 2023b)	CVPR'23	Lifting-Based	×	243	14.3M	0.5G	528M	45.2	35.6
STCFormer (Tang et al., 2023)	CVPR'23	Lifting-Based	\checkmark	243	4.7M	19.6G	80M	41.0	32.0
GLA-GCN (Yu et al., 2023)	ICCV'23	Lifting-Based	×	243	1.3M	1.5G	1556M	44.4	34.8
MotionBERT (Zhu et al., 2023)	ICCV'23	Lifting-Based	\checkmark	243	42.5M	174.7G	719M	39.2	-
KTPFormer (Peng et al., 2024)	CVPR'24	Lifting-Based	\checkmark	243	33.7M	69.5G	286M	40.1	31.9
MotionAGFormer (Soroush Mehraban, 2024)	WACV'24	Lifting-Based	\checkmark	243	19.0M	78.3G	322M	38.4	32.5
PrML (Ours)	-	Multi-Task Learning	~	243	13.0M	49.3G	203M	38.2	32.3

961 962

951

963

964

965

C.6 ATTEMPT TO USE IMAGE FEATURE.

As the inherent information entropy of the input 2D pose sequence, the performance of methods only using 2D pose sequences as input gradually approaches the theoretical limit. Incorporating image features is a potential solution to this general limitation. There are several single-frame lifting-based methods that utilize image features, such as ContextPose (Ma et al., 2021a) and CA-PF (Zhao et al., 2023a). Unfortunately, integrating these methods into existing multi-frame lifting methods or our PrML is challenging due to their complex model architectures. Therefore, we employ one cross-attention layer to fuse image features, making a preliminary attempt to fuse image features into our

972 framework. Due to the lightweight nature of 2D poses, traditional lifting-based frameworks typically 973 input a 2D pose sequence. However, inputting a long-term image sequence is clearly impractical. 974 How to effectively extend image input from a single frame to multiple frames remains an open 975 question. To tackle this problem, we employ VAE to compress the original images into latent features, 976 allowing for multi-frame input. We use the pre-trained VAE (Kingma, 2013) from SDXL (Podell et al., 2023). Although we have compressed the images using VAE, processing 243 frames remains 977 computationally expensive. Therefore, we leverage VAE to incorporate 15 frames image feature, 978 making a preliminary attempt to fuse multi-frame image features into our framework. 979

980 As shown in the Table 13, despite using simple cross-attention to fuse image features, we still observe 981 performance improvement. When we incorporated image features into the 2D branch, we observed 982 a significant improvement. This is likely due to the fact that image features can easily capture 2D spatial information. However, when we added image features to the depth branch, the improvement 983 was limited. This is because we used a simple cross-attention mechanism to process the images, and 984 extracting depth information from RGB images is a challenging task. Similarly, when we incorporated 985 image features into all branches, the improvement was modest. These experiments demonstrate the 986 potential for further extending our framework. 987

988 Moreover, the insights from multi-task learning can be further extended. Using 2D pose detectors 989 to estimate 2D poses has been widely used, so why not similarly utilize powerful depth estimation networks (like DepthAnything (Yang et al., 2024)) to estimate a relative depth from the image to 990 facilitate 3D human pose estimation as well? We can integrate depth features into existing single-991 frame methods that utilize image features to validate this idea. As shown in Table 14, by incorporating 992 depth features from DepthAnything (Yang et al., 2024), we can reduce the error of CA-PF by 1.9mm, 993 which demonstrates the effectiveness of depth features. 994

995 In summary, how to leverage image features to facilitate multi-frame 3D human pose estimation is a non-trivial question and our future research direction. 996

997

998 999 1000 1001 1002

Table 13: Results on the Hu- Table 14: Results on the Hu- Table 15: Results on the 3DPW in man3.6M in millimeters (mm) man3.6M in millimeters (mm) un- millimeters (mm) under MPJPE

under MPJPE. I is the number			der MPJPE.		and P-MPJPE.		
out frame	s.		Method	MPJPE	Method	MPJPE	P-MPJPE
d	Т	MPJPE	CA-PF (Zhao et al., 2023a)	41.4	MotionBERT (Zhu et al., 2023) MotionBERT (Zhu et al., 2023)	85.5 76.9	50.2 47.2
e Feature for	15 2D 15	49.7 47.2	+ Depui Feature (Tang et al., 2024)	39.5	PrML	73.9	49.1
d ge Feature for	T 15 2D 15	MPJPE 49.7 47.2	+ Depth Feature (Yang et al., 2024)	41.4 39.5	MotionBERT (Zhu et al., 2023) PrML	-	76.9 73.9

+ Image Feature for Depth 15 49.6 49.5 + Image Feature for All 15

1004 1005

ADDITIONAL EXPERIMENT ON 3DPW. C.7

1008 We also conduct an experiment on the 3DPW dataset following the experiment setting of Motion-BERT (Zhu et al., 2023). As other methods we compared did not conduct experiments on 3DPW, we 1010 are unable to provide their performance in the table. As shown in the Table 15, the MPJPE of our 1011 PrML is not only lower than the MotionBERT trained from scratch but also lower than its finetune variant. PrML's P-MPJPE is still lower than the MotionBERT trained from scratch. 1012

- 1013 D LIMITATION. 1014
- 1015

1020

Although we propose a novel framework for monocular 3D human pose estimation, our input remains 1016 the same as conventional lifting-based methods: a 2D pose sequence. Such 2D joint coordinates 1017 undoubtedly cause visual representation reduction compared to raw images. This inherent information 1018 entropy of the 2d pose sequence limits the theoretical upper bound on performance. 1019

ADDITIONAL ANALYSIS OF MUTUAL INFORMATION 1021 E

Calculating the conditional mutual information loss is notoriously difficult, especially in neural 1023 networks (Hjelm et al., 2018; Tian et al., 2021). The approximate mutual information loss is much 1024 easier to compute. Given $A \ge B$, maximizing B can maximize the lower bound of A. By maximizing 1025 the approximate mutual information, we can maximize the lower bound of the conditional mutual

information. This is in line with our original goal: provide as much relevant information as possible. In summary, although it is a mathematically suboptimal result, it is the best result we can achieve.

We use Y to represent the label, S to represent task support features and B to represent task bias. Optimizing this objective will maximize the mutual information between task support features and the label to support the task. By the definition of condition mutual information, we have:

$$\begin{aligned} \mathcal{I}(Y;S \mid B) &= \int_{B} \left(\int_{S} \int_{Y} p_{Y,S\mid B}(y,s|b) \log \left(\frac{p_{Y,S\mid B}(y,s|b)}{p_{Y\mid B}(y|b)p_{S\mid B}(s|b)} \right) dyds \right) p_{B}(b)db \\ &= \int_{S} \int_{Y} p(y,s) \log \left(\frac{p(y,s)}{p(y)p(s)} \right) dyds - \int_{B} \int_{S} p(s,b) \log \left(\frac{p(s,b)}{p(s)p(b)} \right) dsdb \\ &+ \int_{Y} \left(\int_{B} \int_{S} p_{S,B\mid Y}(s,b|y) \log \left(\frac{p_{S,B\mid Y}(s,b|y)}{p_{S\mid Y}(s|y)p_{B\mid Y}(b|y)} \right) dsdb \right) p_{Y}(y)dy \\ &= \mathcal{I}(Y;S) - \mathcal{I}(S;B) + \int_{Y} \underbrace{D_{KL}(P_{(S,B)\mid Y} \|P_{S\mid Y}P_{B\mid Y})}_{KL \text{ Divergence } \geq 0} dP_{Y} \\ &= \mathcal{I}(Y;S) - \mathcal{I}(S;B) + \mathbb{E}_{Y} [D_{KL}(P_{(S,B)\mid Y} \|P_{S\mid Y}P_{B\mid Y})] \\ &\geq \mathcal{I}(Y;S) - \mathcal{I}(S;B) \end{aligned}$$

The optimization objective thus becomes:

$$\max \mathcal{I}(Y; S \mid B) \longrightarrow \max \mathcal{I}(Y; S) - \mathcal{I}(S; B)$$
(11)

However, since both $\mathcal{I}(Y; S)$ and $\mathcal{I}(S; B)$ are non-negarive (KL Divergence ≥ 0), the $\mathcal{I}(Y; S)$ – $\mathcal{I}(S;B)$ will result in negative values during training. This will cause the training process to fail to converge. Therefore, we simplified the implementation of mutual information by calculating only the first term. Our final optimization objective becomes:

$$\max \mathcal{I}(Y; S \mid B) \longrightarrow \max \mathcal{I}(Y; S) \tag{12}$$

> F ADDITIONAL VISUALIZATION

Table 16: Results on Human3.6M. (PCKh@0.5)

To further comparing the accuracy of the projected 2D pose, we report the PCKh@0.5 following MPII (Andriluka et al., 2014). The table shows that our PrML achieves better accuracy in the projected 2D poses than the MotionBERT.

(1 CIVIE 0.5)							
Method	PCKh@0.5						
MotionBERT	95.9						
PrML	96.4						



Figure 9: Qualitative Comparison of 2D Pose (Ground Truth, Input, MotionBERT (Zhu et al., 2023) and Ours). We project the 2D pose in the camera coordinate system back to the image coordinate system for comparison. The powerful lifting-based method MotionBERT gets a 2D pose worse than the input, which contradicts our intuition. In contrast, our proposed framework obtains a 2D pose better than the input.



1123Figure 10: Feature distributions visualization of 2D feature (green) and depth feature (purple) using1124t-SNE (Van der Maaten & Hinton, 2008) method on Human3.6M (Ionescu et al., 2013) dataset.1125(Random 243 frames poses, all 17 joints are visualized)



Figure 11: Feature distributions visualization of 2D feature (green) and depth feature (purple) using
t-SNE (Van der Maaten & Hinton, 2008) method on Human3.6M (Ionescu et al., 2013) dataset.
(Continuous 243 frames poses, all 17 joints are visualized)



Figure 12: Qualitative comparisons of PrML with MotionBERT (Zhu et al., 2023). The green cycle indicates locations where our method achieves better results.



Figure 13: Feature distributions visualization of 2D features (green) and depth features (purple) and entangled features (orange) using t-SNE (Van der Maaten & Hinton, 2008) method on Human3.6M (Ionescu et al., 2013) dataset. The 2D features and depth features are extracted from our PrML (Multi-Task Learning). The entangled features are extracted from the lifting-based method MotionBERT (Zhu et al., 2023). Visualization results show that the distribution of entangled features is more compact, while our 2D pose features and per-joint depth features exhibit a more diverse distribution, indicating the potential to learn various representations.

1214

- 1240
- 1241



Figure 14: Feature distributions visualization of 2D feature, depth feature and entangled feature
using t-SNE (Van der Maaten & Hinton, 2008) method on Human3.6M (Ionescu et al., 2013) dataset.
(Random 243 frames poses, all 17 joints are visualized)



Figure 15: Feature distributions visualization of 2D feature, depth feature and entangled feature
using t-SNE (Van der Maaten & Hinton, 2008) method on Human3.6M (Ionescu et al., 2013) dataset.
(Continuous 243 frames poses, all 17 joints are visualized)



Figure 16: Feature distributions visualization of 2D feature, depth feature, and entangled feature in their own feature space using t-SNE (Van der Maaten & Hinton, 2008) method on Human3.6M (Ionescu et al., 2013) dataset. (Random 243 frames poses, all 17 joints are visualized)