

# POSE PRIORS FROM LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

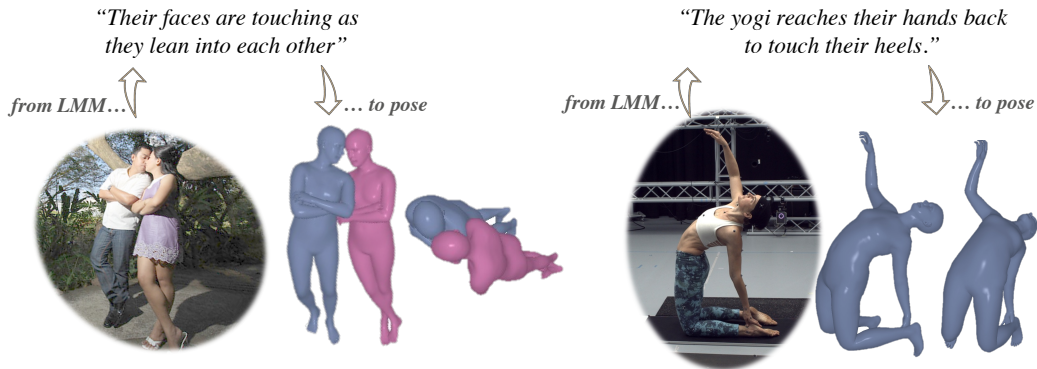


Figure 1: **Optimizing contacts in 3D human pose.** Our approach leverages the semantic priors of a Large Multimodal Model (LMM) by converting natural language descriptions of individuals in an image into mathematical constraints. We can then optimize the 3D pose estimates using these constraints. These examples show image descriptions generated by an LMM and corresponding refined pose estimates.

## ABSTRACT

We present a pose optimization method that enforces accurate physical contact constraints when estimating the 3D pose of humans. Our central insight is that since language is often used to describe physical interaction, large pretrained text-based models can act as priors on pose estimation. We can thus leverage this insight to improve pose estimation by converting natural language descriptors, generated by a large multimodal model (LMM), into tractable losses to constrain the 3D pose optimization. Despite its simplicity, our method produces surprisingly compelling pose reconstructions of people in close contact, correctly capturing the semantics of the social and physical interactions. We demonstrate that our method rivals more complex state-of-the-art approaches that require expensive human annotation of contact points and training specialized models. Moreover, unlike previous approaches, our method provides a unified framework for resolving self-contact and person-to-person contact.<sup>1</sup>

## 1 INTRODUCTION

In 3D pose estimation, the dominant forms of labeled training data are 3D pose (obtained via motion capture) and 2D keypoints (Goel et al., 2023). Scaling the amount of this data is expensive, as it requires either special technology or fine-grained human effort. On the other hand, in 2D computer vision tasks such as recognition and segmentation, prior work shows that natural language supervision provides a path to strong performance (Radford et al., 2021; Xu et al., 2022a). Can we use language-based models to improve 3D pose estimation?

Our focus in this work is enabling computers to correctly perceive physical contact when estimating pose (See Figure 1). Specifically, we aim to build a system that takes as input a single view of people during close physical interaction or one person in a pose that involves self-contact and produces accurate 3D mesh reconstructions of each person as output. This setting is challenging for state-of-the-art pose regression models, as some body parts are frequently occluded by other ones, and also

<sup>1</sup>Our code will be publicly available at the time of publication.

054 challenging for pose optimization methods relying on 2D keypoints, which do not convey contact  
055 points. Previously proposed approaches address these issues by curating task-specific datasets via  
056 motion capture or human-annotated points of contact between body parts (Muller et al., 2021; Fieraru  
057 et al., 2021; Müller et al., 2023).

058 As physical contact is a universal human social signal, humans developed extensive terminology for  
059 its particularities. Detailed descriptions of touch in different contexts are widely discussed in texts that  
060 range from love-song lyrics such as Paul Anka’s “Put your head on my shoulder” to Shakespeare’s  
061 “See how she leans her cheek upon her hand.” (Romeo and Juliet). It touches on subjects from love to  
062 meditative poses.

063 Our main insight is that since written language discusses our physical interactions (hugs, kisses, fist  
064 fights, yoga poses, etc.) at great length, we should be able to extract a semantic prior on humans’  
065 poses from a pretrained large multimodal model (LMM) (Achiam et al., 2023; Liu et al., 2023;  
066 Dai et al., 2023). Just like a prior trained on motion capture data, this language-based prior can tell  
067 us which contacts are most likely in poses and interactions. Through this approach, we avoid the  
068 time-consuming and expensive collection of training data involving motion capture or annotated self  
069 and cross-person contacts that previous refinement methods require.

070 This insight leads us to a simple framework, which we call ProsePose. We prompt a pre-trained  
071 LMM, with the image and request as output a formatted list of contact constraints between body parts.  
072 We then convert this list of constraints into a loss function that can be optimized jointly with other  
073 common losses, such as 2D keypoint loss, to refine the initial estimates of a pose regression model.  
074 The prompt provides an intuitive way for the system designer to adapt the generated constraints to  
075 their setting (e.g. if they want to focus on yoga or dance).

076 We show in experiments on three 2-person interaction datasets and one dataset of complex yoga  
077 poses that ProsePose produces more accurate reconstructions than previous approaches that do not  
078 use a large amount of task-specific data for training. These results indicate that LMMs, without any  
079 additional finetuning, offer a useful prior for pose reconstruction.

080 In summary, (1) we show that LMMs have implicit semantic knowledge of poses that is useful for  
081 pose estimation, and (2) we formulate a novel framework that converts free-form natural language  
082 responses from a pre-trained LMM into tractable loss functions that can be used for pose optimization.

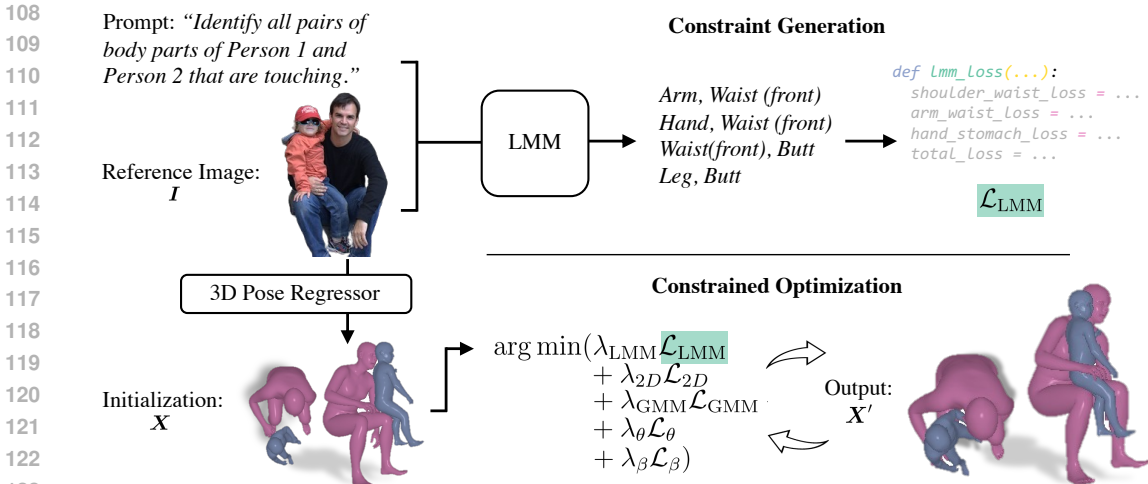
## 084 2 RELATED WORK

### 086 2.1 3D HUMAN POSE RECONSTRUCTION

087 Reconstructing 3D human poses from single images is an active area of research. Prior works have  
088 explored using optimization-based approaches (Pavlakos et al., 2019a; Guan et al., 2009; Lassner et al.,  
089 2017; Pavlakos et al., 2019b; Rempe et al., 2021) or pure regression (Kanazawa et al., 2018; Arnab  
090 et al., 2019; Guler & Kokkinos, 2019; Joo et al., 2021; Kolotouros et al., 2019) to estimate the 3D  
091 body pose given a single image. HMR2 (Goel et al., 2023) is a recent state-of-the-art regression model  
092 in this line of work. Building on these monocular reconstruction approaches, some methods have  
093 looked into reconstructing multiple individuals jointly from a single image. These methods (Zanfir  
094 et al., 2018; Jiang et al., 2020; Sun et al., 2021) use deep networks to reason about multiple people in  
095 a scene to directly output multi-person 3D pose predictions. BEV (Sun et al., 2022) accounts for the  
096 relative proximity of people explicitly using relative depth annotations to reason about proxemics  
097 when predicting and placing each individuals in the scene (e.g. depth of people with respect to one  
098 another). However, approaches in both categories generally do not accurately capture physical contact  
099 between parts of a single person or between people (Müller et al., 2023; Muller et al., 2021).

### 100 2.2 CONTACT INFERENCE IN 3D POSE RECONSTRUCTION

101 3D pose reconstruction is especially challenging when there is self-contact or inter-person contact.  
102 This has motivated a line of work on pose reconstruction approaches tailored for this setting. Muller  
103 et al. (2021) focuses on predicting self contact regions for 3D pose estimation by leveraging a dataset  
104 with collected contact annotations to model complex poses such as arm on hip or crossed arms.  
105 Fieraru et al. (2020) introduces the first dataset with hand-annotated ground truth contact labels  
106 between two people. REMIPS (Fieraru et al., 2021) and BUDDI (Müller et al., 2023) train models on  
107 the person-to-person contact maps in this data in order to improve 3D pose estimation of multiple  
people from a single image. CloseInt (Huang et al., 2024) trains a physics-guided diffusion model on



124 Figure 2: **LMM-guided Pose Estimation** Our method takes as input an image of one or two people in contact.  
125 We first obtain initial pose estimates for each person from a pose regressor. Then we use an LMM to generate  
126 contact constraints, each of which is a pair of body parts that should be touching. This list of contacts is converted  
127 into a loss function  $\mathcal{L}_{LMM}$ . We optimize the pose estimates using  $\mathcal{L}_{LMM}$  and other losses to produce a refined  
128 estimate of each person’s pose that respects the predicted contacts.

129 two-person motion capture data for this task. However, contact annotations, which are crucial for  
130 these approaches, are difficult and expensive to acquire. Our method does not require any training on  
131 such annotations. Instead, we leverage an LMM’s implicit knowledge about pose to constrain pose  
132 optimization to capture both self- and person-to-person contact.

### 132 2.3 LANGUAGE PRIORS ON HUMAN POSE

133 There exists a plethora of text to 3D human pose and motion datasets (Punnakkal et al., 2021; Guo  
134 et al., 2022; Plappert et al., 2016), which have enabled work focused on generating 3D motion  
135 sequences of a single person performing a general action (Tevet et al., 2023; Jiang et al., 2023; Zhang  
136 et al., 2023). This line of work has been extended to generating the motion of two people conditioned  
137 on text (Shafir et al., 2023; Liang et al., 2023).

138 PoseScript (Delmas, Ginger and Weinzaepfel, Philippe and Lucas, Thomas and Moreno-Noguer,  
139 Francesc and Rogez, Grégory, 2022) is a method for generating a single person’s pose from fine-  
140 grained descriptions. They leverage a library of predefined pose descriptors, from which they form  
141 detailed textual annotations for their motion capture dataset. By training a model on this data, they can  
142 generate various plausible poses. PoseFix (Delmas, Ginger and Weinzaepfel, Philippe and Moreno-  
143 Noguer, Francesc and Rogez, Grégory, 2023) considers the problem of modifying a pose given a  
144 fine-grained description of the desired change, and introduces a labeled dataset for this task. The  
145 PoseFix method then trains a model on this data to predict the modified pose given the initial pose  
146 and description. PoseGPT (Feng et al., 2023), like our work, focuses on the problem of monocular 3D  
147 reconstruction of people. PoseGPT is a pose regressor that uses language as part of its training data.  
148 However, PoseGPT does not produce better pose estimates than previous state-of-the-art regressors  
149 (i.e. regressors that do not use language) and applies only to the one-person setting.

150 Our work differs from previous work on language and pose in several ways. First, whereas all prior  
151 work trains a model on data with pairs of language and pose, which is expensive to collect, our  
152 method leverages the existing knowledge in an LMM to reason about pose. Second, prior work in  
153 this area focuses on either the one-person or the two-person setting. In contrast, our work presents a  
154 single framework to reason about physical contacts within or between poses. Finally, in scenes with  
155 physical contact, we show that our method improves the pose estimates of state-of-the-art regressors.

## 156 3 GUIDING POSE OPTIMIZATION WITH AN LMM

157  
158 Given an image, our goal is to estimate the 3D body pose of individuals in the image while capturing  
159 the self and cross-person contact points. While we cannot trivially use natural language responses  
160 (hug, kiss) to directly optimize 3D body poses, we leverage the key insight that LMMs understand  
161 *how* to articulate a given pose (arms around waist, lips touching). We propose a method to structure  
these articulations into constraints and convert them into loss functions.

More concretely, our framework, illustrated by Figure 2, takes as input the image  $I$  and the bounding boxes  $B$  of the subjects of interest. In the first stage, the image is passed to a pose regressor to obtain a rough estimate of the 3D pose  $X^p$  for each individual  $p$  in the image. In the second stage, we prompt a LMM with the image and a set of instructions in order to generate a list of self- or inter-person contact constraints, which we then convert into a loss function (Sec. 3.4). Finally, in the third stage, we jointly optimize the generated loss function with several other pre-defined loss terms (Sec. 3.4). We refer to our framework as **ProsePose**.

### 3.1 PRELIMINARIES

While our approach scales in principle to an arbitrary number of individuals, we focus our description on the two-person case to keep the exposition simple. We also demonstrate results on the one-person case, which is simply an extension of the two-person case. In particular, we apply our method to the one-person case by setting  $X^0 = X^1$ . Please see § 6 for details on the differences between the two-person and one-person cases.

**Large Multimodal Models** An LMM is a model that takes as input an image and a text prompt and produces text output that answers the prompt based on the image. Our framework is agnostic to the architecture of the LMM. LMMs are typically trained to respond to wide variety of instructions (Liu et al., 2023; Dai et al., 2023), but at the same time, LMMs are prone to hallucination (Leng et al., 2023; Li et al., 2023). Handling cases of hallucination is a key challenge when using LMMs, and we mitigate this issue by aggregating information across several samples from the LMM.

**Pose representation.** We use a human body model (Pavlakos et al., 2019a) to represent each person  $p \in \{0, 1\}$ . The body model is composed of a pose parameter that defines the joint rotations  $\theta \in \mathbb{R}^{d_\theta \times 3}$ , where  $d_\theta$  is the number of joints, and a shape parameter  $\beta \in \mathbb{R}^{d_\beta}$ , where  $d_\beta$  is the dimensions of the shape parameter. We can apply a global rotation  $\Phi \in \mathbb{R}^3$  and translation  $t \in \mathbb{R}^3$  to place each person in the world coordinate space. The full set of parameters for each person is denoted by  $X^p = [\theta^p, \beta^p, \Phi^p, t^p]$ . For simplicity, we refer to the parameter set  $(X^0, X^1)$  as  $X$ .

These parameters can be plugged into a differentiable function that maps to a mesh consisting of  $d_v$  vertices  $V \in \mathbb{R}^{d_v \times 3}$ . From the mesh, we can obtain a subset of the vertices representing the 3D locations of the body’s joints  $J \in \mathbb{R}^{d_j \times 3}$ . From these joints, we can calculate the 2D keypoints  $K_{proj}$  by projecting the 3D joints to 2D using the camera intrinsics  $\Pi$  predicted from (Pavlakos et al., 2019a).

$$K_{proj} = \Pi(J) \in \mathbb{R}^{d_j \times 2}. \quad (1)$$

**Vertex regions.** In order to define contact constraints between body parts, we define a set of *regions* of vertices. Prior work on contact has partitioned the body in to fine-grained regions (Fieraru et al., 2020). However, since our constraints are specified by a LMM trained on natural language, the referenced body parts are often coarser in granularity. We therefore update the set of regions to reflect this language bias by combining these fine-grain regions into larger, more commonly referenced body parts such as arm, shoulder (front&back), back, and waist (front&back). Please see § 6.2 for a visualization of the coarse regions. Formally, we write  $R \in \mathbb{R}^{d_r \times 3}$  to denote a region with  $d_r$  vertices, which is part of the full mesh ( $R \subset V$ ).

**Constraint definition.** A contact constraint specifies which body parts from two meshes should be touching. Using the set of coarse regions, we define contact constraints as pairs of coarse regions  $c = (R_a, R_b)$  between a region  $R_a$  of one mesh and  $R_b$  of the other mesh, as shown in Figure 3. For instance, (“hand”, “arm”) indicates a hand should touch an arm.

### 3.2 POSE INITIALIZATION

We obtain a rough initial estimate of the 3D pose from a regression-based method. The regressor takes as input the image  $I$  and outputs estimates for the body model parameters  $\theta, \beta, r$ , and  $t$  for each subject.

### 3.3 CONSTRAINT GENERATION WITH A LMM

Our method strives to enforce contact constraints for the estimated 3D poses. Our key insight is to leverage a LMM to identify regions of contact between different body parts on the human body surface. As shown in Figure 2 (top), we prompt the LMM with an image and ask it to output a list of

all plausible regions that are in contact. However, we cannot simply use natural language descriptions to directly optimize a 3D mesh. As such, we propose a framework to convert these constraints into a loss function.

**LMM-based constraint generation.** Given the image  $I$ , we first use the bounding boxes  $B$  to crop the part containing the subjects. We then use an image segmentation model to mask any extraneous individuals. While cropping and masking the image may remove information, we find the LMMs are relatively robust to missing context, and more importantly, this allows us to indicate which individuals to focus on. Given the segmented image, we ask the LMM to generate a set  $C = \{c_1, \dots, c_m\}$  of all pairs of body parts that are touching, where  $m$  is the total number of constraints the LMM generates for the image.

In the prompt, we specify the full set of coarse regions to pick from. We find that LMMs fail to reliably reference the left and right limbs correctly or consistently, so we designed this set of coarse regions such that they do not disambiguate the chirality of the hands, arms, legs, feet, and shoulders. Instead, the two hands are grouped together, the two arms are grouped together, etc. Nevertheless, if the LMM uses “left” or “right” to reference a region, despite the instruction to not do so, we directly use the part of the region with the specified chirality rather than considering both possibilities.

Motivated by the chain-of-thought technique, which has been shown to improve language model performance on reasoning tasks (Wei et al., 2022), we ask the LMM to write its reasoning or describe the pose before listing the constraints. For the full prompt used in each setting, please refer to § 6.

We sample  $N$  responses from the LMM, yielding  $N$  sets of constraints  $\{C_1, C_2, \dots, C_N\}$ . The next step is to convert each constraint set  $C_j$ , where  $j \in \{1, 2, \dots, N\}$ , into a loss term.

**Loss function generation.** We first filter out contact pairs that occur fewer than  $f$  times across all constraint sets, where  $f$  is a hyperparameter. Then for each contact pair  $c = (R_a, R_b)$  in  $C_j$ , we define  $dist(c)$  as the minimum distance between the two regions:

$$dist(c) = \min \|v_a - v_b\|_2 \quad \forall v_a \in R_a, \forall v_b \in R_b \tag{2}$$

where  $\{v_a, v_b\} \in \mathbb{R}^3$ . In practice, the number of vertices in each region can be very large. To make this computation tractable, we first take a random sample of vertices from  $R_a$  and from  $R_b$  before computing distances between pairs of vertices in these samples. Furthermore, since the ordering of the people in the LMM constraints is unknown (i.e. does  $R_a$  come from the mesh defined by parameter  $X^0$  or  $X^1$ ), we compute the overall loss for both possibilities and take the minimum. We use  $c^\top = (R_b, R_a)$  to denote the flipped ordering. We then sum over all constraints in the list  $C_j$ :

$$dist_{sum}(C_j) = \min \left( \sum_{c \in C_j} dist(c), \sum_{c \in C_j} dist(c^\top) \right) \tag{3}$$

Each constraint set sampled from the LMM is likely to contain noise or hallucination. To mitigate the effect of this, we average over all  $N$  losses corresponding to each constraint set to obtain the overall LMM loss. This technique is similar to self-consistency (Wang et al., 2022), which is commonly used for code generation tasks. Concretely, the overall LMM loss is defined as

$$\mathcal{L}_{LMM} = \frac{1}{N} \sum_{j=1}^N dist_{sum}(C_j) \tag{4}$$

If a constraint set  $C_j$  is empty (i.e. the LMM does not suggest any contact pairs), then we set  $dist_{sum}(C_j) = 0$ . If there are several such constraint sets, we infer that the LMM has low confidence about the contact points (if any) in the image. To handle these cases, we set a threshold  $t$  and if the number of empty constraint sets is at least as large as  $t$ , we gracefully backoff to the appropriate baseline optimization procedure (described in Sections 4.1 and 4.2 for each setting).



Figure 3: **Notation.** Given an image  $I$ , we can lift each individual into corresponding 3D meshes  $V$ . We define contact constraints  $c$  as pairs of regions  $(R_a, R_b)$  in contact. The loss is defined in terms of the distance between the vertices  $(v_a, v_b)$  on the mesh.

### 3.4 CONSTRAINED POSE OPTIMIZATION

Drawing from previous optimization-based approaches (Müller et al., 2023; Bogo et al., 2016; Pavlakos et al., 2019b), we employ several additional losses in the optimization. We then minimize the joint loss to obtain a refined subset of the body model parameters  $\mathbf{X}' = [\boldsymbol{\theta}', \beta', \mathbf{t}']$ :

$$[\boldsymbol{\theta}', \beta', \mathbf{t}'] = \arg \min(\lambda_{\text{LMM}}\mathcal{L}_{\text{LMM}} + \lambda_{\text{GMM}}\mathcal{L}_{\text{GMM}} + \lambda_{\beta}\mathcal{L}_{\beta} + \lambda_{\theta}\mathcal{L}_{\theta} + \lambda_{2D}\mathcal{L}_{2D} + \lambda_P\mathcal{L}_P)$$

Following Müller et al. (2023), we divide the optimization into two stages. In the first stage, we optimize all three parameters. In the second stage, we optimize only  $\boldsymbol{\theta}$  and  $\mathbf{t}$ , keeping the shape  $\beta$  fixed. Here, we detail all of the remaining losses used in the optimization.

**Pose and shape priors.** We compute a loss  $\mathcal{L}_{\text{GMM}}$  based on the Gaussian Mixture pose prior of Bogo et al. (2016) and a shape loss  $\mathcal{L}_{\beta} = \|\beta\|_2^2$ , which penalizes extreme deviations from the body model’s mean shape.

**Initial pose loss.** To ensure we do not stray too far from the initialization, we penalize large deviations from the initial pose  $\mathcal{L}_{\theta} = \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^2$ .

**2D keypoint loss.** Similar to BUDDI (Müller et al., 2023), for each person in the image, we obtain pseudo ground truth 2D keypoints and their confidences from OpenPose (Cao et al., 2019) and ViTPose (Xu et al., 2022b). Given this pseudo ground truth, we merge all the keypoints into  $\mathbf{K} \in \mathbb{R}^{d_j \times 2}$ , and their corresponding confidences into  $\gamma \in \mathbb{R}^{d_j}$ . From the predicted  $\mathbf{X}'$ , we can compute the 2D projection of each 3D joint location using Equation 3.1. Then, the 2D keypoint loss is defined as:

$$\mathcal{L}_{2D} = \sum_{j=1}^{d_j} \gamma(\mathbf{K}_{proj} - \mathbf{K})^2 \quad (5)$$

**Interpenetration loss.** To prevent parts of one mesh from being in the interior of the other, we add an interpenetration loss. Generically, given two sets of vertices  $\mathbf{V}_0$  and  $\mathbf{V}_1$ , we use winding numbers to compute the subset of  $\mathbf{V}_0$  that intersects  $\mathbf{V}_1$ , which we denote as  $\mathbf{V}_{0,1}$ . Similarly,  $\mathbf{V}_{1,0}$  is the subset of  $\mathbf{V}_1$  that intersects  $\mathbf{V}_0$ . The interpenetration loss is then defined as

$$\mathcal{L}_P = \sum_{x \in \mathbf{V}_{0,1}} \min_{v_1 \in \mathbf{V}_1} \|x - v_1\|_2^2 + \sum_{y \in \mathbf{V}_{1,0}} \min_{v_0 \in \mathbf{V}_0} \|y - v_0\|_2^2 \quad (6)$$

Due to computational cost, this loss is computed on low-resolution versions of the two meshes (roughly 1000 vertices per mesh).

## 4 EXPERIMENTS

We conduct experiments on several datasets in the two-person and one-person settings. In this section, we first provide important implementation details and a description of the metrics that we use to evaluate our method and previous approaches. We then present quantitative and qualitative results showing that ProsePose refines pose estimates to capture semantically relevant contact in each setting.

**Implementation details.** Following prior work on two-person pose estimation (Müller et al., 2023), we use BEV (Sun et al., 2022) to initialize the poses since it was trained to predict both the body pose parameters and the placement of each person in the scene. However, on the single person yoga poses, we find that the pose parameter estimates of HMR2 (Goel et al., 2023) are much higher quality, so we initialize the body pose using HMR2.

We use the SMPL-X (Pavlakos et al., 2019a) body model and GPT4-V (Achiam et al., 2023) as the LMM with temperature = 0.7 when sampling from it.<sup>2</sup> We also include results when using LLaVA as the LMM in § 7.4. We use Segment Anything (Kirillov et al., 2023) as the segmentation model, used to remove extraneous people in the image (we only apply this step for FlickrCI3D, since other datasets are from motion capture). Unless otherwise specified, we set  $N = 20$  samples. For all

<sup>2</sup>We access GPT4-V, specifically the gpt-4-vision-preview model, via the OpenAI API: [platform.openai.com](https://platform.openai.com). We use the “high” detail setting for image input.

Table 1: **Two-person Results.** Joint PA-MPJPE (lower is better) and Avg. PCC (higher is better). For FlickrCI3D, PA-MPJPE is computed using the pseudo-ground-truth fits. **Bold** indicates best method without contact supervision in each column.

	Hi4D	FlickrCI3D		CHI3D	
	PA-MPJPE $\downarrow$	PA-MPJPE $\downarrow$	PCC $\uparrow$	PA-MPJPE $\downarrow$	PCC $\uparrow$
<i>Without contact supervision</i>					
BEV (Sun et al., 2022)	144	106	64.8	<b>96</b>	71.4
Heuristic	116	67	77.8	105	74.1
ProsePose	<b>93</b>	<b>58</b>	<b>79.9</b>	100	<b>75.8</b>
<i>With contact supervision</i>					
BUDDI (Müller et al., 2023)	89	65.9	81.9	68	78.6

of our 2-person experiments,  $f = 1$ , while  $f = 10$  in the 1-person setting. We set  $t = 2$  for the experiment on the CHI3D dataset and  $t = N$  for all other experiments. We set  $\lambda_{\text{LMM}} = 1000$  in the 2-person experiments, and  $\lambda_{\text{LMM}} = 10000$  in the 1-person setting. In the two-person case, all other loss coefficients are taken directly from Müller et al. (2023). In the one-person case, we find that removing the GMM pose prior and doubling the weight on the initial pose loss improves optimization dramatically, likely because the complex yoga poses are out of distribution for the GMM prior. These hyperparameters and our prompts were chosen based on experiments on the validation sets. Furthermore, following Müller et al. (2023), we run both optimization stages for at most 1000 steps. We use the Adam optimizer (Kingma & Ba, 2014) with learning rate 0.01. For other implementation details such as prompts, the list of coarse regions in each setting, and additional differences between the 1- and 2-person cases, please refer to § 6.

**Metrics.** As is standard in the pose estimation literature, we report Procrustes-aligned Mean Per Joint Position Error (PA-MPJPE) in millimeters. This metric finds the best alignment between the estimated and ground-truth pose before computing the joint error. In the two-person setting, we focus on the *joint* PA-MPJPE, as this evaluation incorporates the relative translation and orientation of the two people. See § 7.2 for the per-person PA-MPJPE.

We also include the percentage of correct contact points (PCC) metric introduced by (Müller et al., 2023). This metric captures the fraction of ground-truth contact pairs that are accurately predicted. For a given radius  $r$ , a pair is classified as “in contact” if the two regions are both within the specified radius. We use the set of fine-grained regions defined in Fieraru et al. (2020) to compute PCC. The metric is averaged over  $r \in 0, 5, 10, 15, \dots, 95$  mm. Please note that since these regions are defined on the SMPL-X mesh topology, we convert the regression baselines– BEV and HMR2– from the SMPL mesh topology to SMPL-X to compute this metric. Please see § 7.1 for more details on the regions and on the mesh conversion.

#### 4.1 TWO-PERSON POSE REFINEMENT

**Datasets** We evaluate on three datasets, and our dataset processing largely follows (Müller et al., 2023). **Hi4D** (Yin et al., 2023) is a motion capture dataset of pairs of people interacting. Each sequence has a subset of frames marked as contact frames, and we take every fifth contact frame. We use the images from a single camera, resulting in roughly 247 images. **Flickr Close Interactions 3D (FlickrCI3D)** (Fieraru et al., 2020) is a collection of Flickr images of multiple people in close interaction. The dataset includes manual annotations of the contact maps between pairs of people. (Müller et al., 2023) used these contact maps to create pseudo-ground truth 3D meshes and curated a version of the test set to exclude noisy annotations, which has roughly 1403 images. **CHI3D** (Fieraru et al., 2020) is a motion capture dataset of pairs of people interacting. We present results on the validation set. There are 126 different sequences, each of which has a single designated “contact frame.” Each frame is captured from 4 cameras, so there are roughly 504 images in this set.

To develop our method, we experimented on the validation sets of FlickrCI3D and Hi4D, and a sample of the training set from CHI3D. For our experiments, we can compute the PCC on FlickrCI3D and CHI3D, which have annotated ground-truth contact maps. Since all baselines also use BEV for initialization, we exclude images where BEV fails to detect one of the subjects in the interaction pair.

**Baselines** We compare our estimated poses to the following:

Table 2: **Two-person PCC**. Percent of correct contact points (PCC) for five different radii  $r$  in mm. **Bold** indicates the best score without contact supervision in each column. At the ground-truth contact points, our method brings the meshes closer together than the baselines.

	PCC <sub>↑</sub> @ $r$ on FlickrCI3D					PCC <sub>↑</sub> @ $r$ on CHI3D				
	5	10	15	20	25	5	10	15	20	25
<i>Without contact supervision</i>										
BEV (Sun et al., 2022)	3.6	6.3	10.8	17.1	28.6	5.8	17.4	32.5	47.3	61.9
Heuristic	14.6	33.9	49.3	60.8	70.3	11.1	28.0	45.3	55.3	64.4
ProsePose	<b>15.6</b>	<b>39.9</b>	<b>57.1</b>	<b>67.9</b>	<b>75.8</b>	<b>13.5</b>	<b>35.2</b>	<b>52.5</b>	<b>61.3</b>	<b>68.4</b>
<i>With contact supervision</i>										
BUDDI (Müller et al., 2023)	18.5	44.2	61.8	73.1	80.8	15.7	39.4	57.1	68.8	78.0

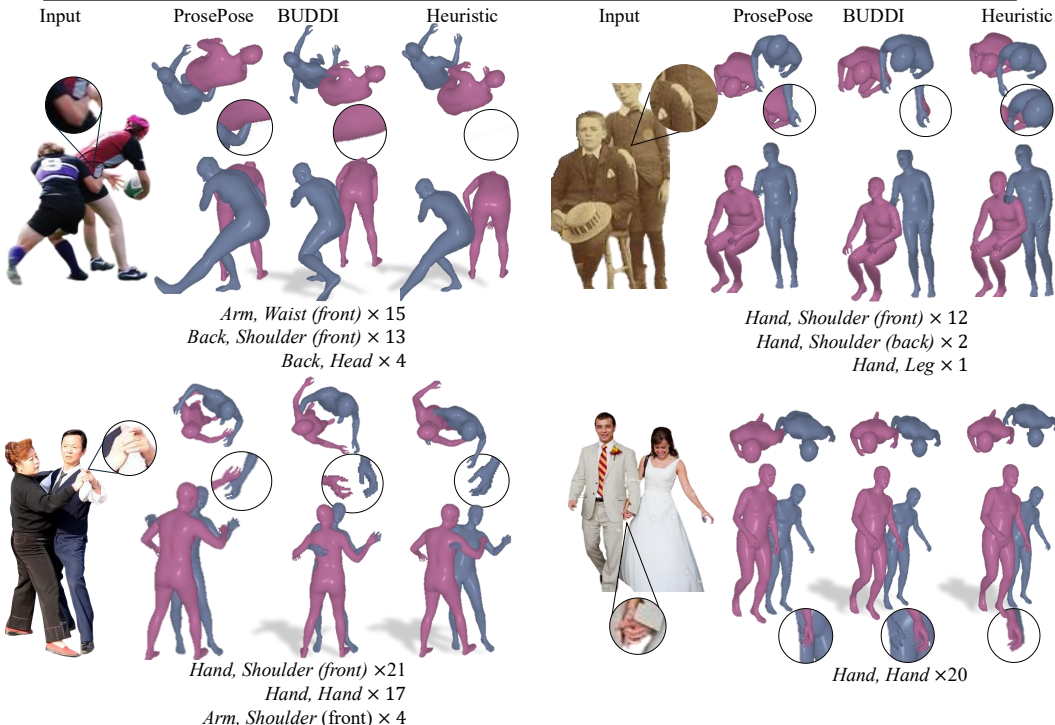


Figure 4: **Two-person results** We show qualitative results from ProsePose , BUDDI (Müller et al., 2023), and the contact heuristic. Under each example, we show the top 3 constraints predicted by GPT4-V and the number of times each constraint was predicted across all 20 samples. Our method correctly reconstructs people in a variety of interactions, and the predicted constraints generally align with the interaction type in each example.

- **BEV (Sun et al., 2022)** Multi-person 3D pose estimation method. Uses relative depth to reason about spatial placement of individuals in the scene. ProsePose , Heuristic, and BUDDI use BEV to initialize pose estimates.
- **Heuristic** A contact heuristic which includes the auxiliary losses in Section 3.4 as well as a term that minimizes the minimum distance between the two meshes. Introduced by (Müller et al., 2023). We use their hyperparameters for this heuristic. Please note, this baseline is used as the default when the number of empty constraint sets is at least the threshold  $t$ .
- **BUDDI (Müller et al., 2023)** This method uses a learned diffusion prior to constrain the optimization. We stress that BUDDI requires a large amount of annotated training data on pairs of interacting bodies, which is not used in our method.

**Quantitative Results** Table 1 provides quantitative results on the three datasets.

Across datasets, ProsePose consistently improves over the strongest baseline, **Heuristic**.

On the Hi4D dataset, ProsePose reduces 85% of the gap in PA-MPJPE between **Heuristic** and the fully supervised **BUDDI**. On the FlickrCI3D and CHI3D datasets, ProsePose narrows the gap in the average PCC between **Heuristic** and **BUDDI** by more than one-



third. (While ProsePose achieves a better PA-MPJPE than **BUDDI** on FlickrCI3D, for this dataset, we rely primarily on PCC since PA-MPJPE is computed on *pseudo-ground-truth* fits.)

On CHI3D, ProsePose outperforms **Heuristic** but underperforms **BEV** in terms of PA-MPJPE. We find that on the subset of images where we do not default to the heuristic (i.e. on images where GPT4-V predicts enough non-empty constraint sets), the PA-MPJPE for ProsePose and BEV is 86 and 87, respectively. In other words, in the cases where our method is actually used, the joint error is slightly less than that of BEV. As a result, we can attribute the worse overall error to the poorer performance of the heuristic. Overall, our method improves over the other methods that do not use 3D supervision in terms of both joint error and PCC. Table 2 shows the PCC for each method at various radii. The results show that ProsePose brings the meshes closer together at the correct contact points. On both the FlickrCI3D and CHI3D datasets, ProsePose outperforms the other baselines that do not use contact supervision.

Next, we ablate important aspects of ProsePose. In Figure 5, we show that averaging the loss over several samples from the LMM improves performance, mitigating the effect of LMM hallucination. Table 3 presents an ablation of all the losses involved in our optimization on the Hi4D validation set.  $\mathcal{L}_{LMM}$  and  $\mathcal{L}_{2D}$  have the greatest impact, indicating that our LMM-based loss is crucial for the large improvement in joint error.

**Qualitative Results** Figure 4 shows examples

of reconstructions from ProsePose, **Heuristic**, and **BUDDI**. Below each of our predictions, we list the most common constraints predicted by GPT4-V for the image. The predicted constraints correctly capture the semantics of each interaction. For instance, it is inherent that in tango, one person’s arm should touch the other’s back. In a rugby tackle, a player’s arms are usually wrapped around the other player. Using these constraints, ProsePose correctly reconstructs a variety of interactions, such as tackling, dancing, and holding hands. In contrast, the heuristic struggles to accurately position individuals and/or predict limb placements, often resulting in awkward distances.

## 4.2 ONE-PERSON POSE REFINEMENT

**Datasets** Next, we evaluate ProsePose on a single-person setting. For this setting, we evaluate on MOYO (Tripathi et al., 2023), a motion capture dataset with videos of a single person performing various yoga poses. The dataset provides views from multiple different cameras. We pick a single camera that shows the side view for evaluation. For each video, we take single frame from the middle as it generally shows the main pose. There is no official test set, and the official validation set consists of only 16 poses. Therefore, we created our own split by picking 79 arbitrary examples from the training set to form our validation set. We then combine the remaining examples in the training set with the official validation set to form our test set. In total, our test set is composed of 76 examples. Since this dataset does not have annotated region contact pairs, we compute the pseudo-ground-truth contact maps using the Euclidean and geodesic distance following Muller et al. (2021).

**Baselines** We compare against the following baselines:

- **HMR2** (Goel et al., 2023) State-of-the-art pose regression method. We use this baseline to initialize our pose estimates for optimization.
- **HMR2+opt** Optimization procedure that is identical to our method without  $\mathcal{L}_{LMM}$ . This method is the default when the number of empty constraint sets is at least the threshold  $t$ .

Both the quantitative and qualitative results echo the trends discussed in the 2-person setting. Table 4 provides the quantitative results. The PCC metrics show that our LMM loss improves the predicted

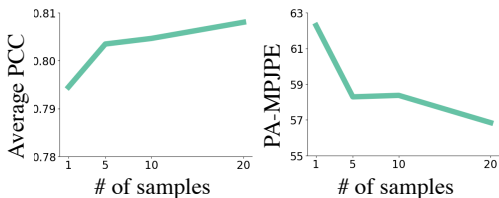


Figure 5: **More samples improve pose estimation.** On the FlickrCI3D validation set, taking more samples from the LMM and averaging the resulting loss functions improves joint PA-MPJPE (left) and average PCC (right).

	PA-MPJPE $_{\downarrow}$
All Losses	81
w/o. $\mathcal{L}_{LMM}$	138
w/o. $\mathcal{L}_{GMM}$	85
w/o. $\mathcal{L}_{\beta}$	91
w/o. $\mathcal{L}_{\theta}$	84
w/o. $\mathcal{L}_{2D}$	130
w/o. $\mathcal{L}_P$	78

Table 3: **Ablations on Hi4D.** Joint PA-MPJPE (lower is better). We evaluate the impact of each loss in our optimization on the Hi4D by removing one loss at a time. For all experiments, we use the same settings. The set of cases where we default to the baseline (Heuristic) is also kept the same.

Table 4: **One-person Results.** PA-MPJPE (lower is better) and Avg. PCC (higher is better). Our method captures ground-truth contacts better than the baseline methods, as shown by the PCC.

	PA-MPJPE $\downarrow$	PCC $\uparrow$	PCC $\uparrow$ @ $r$				
			5	10	15	20	25
HMR2 (Goel et al., 2023)	84	83.0	34.2	55.2	69.5	78.4	83.9
HMR2+opt	<b>81</b>	85.2	47.7	65.5	74.6	80.9	86.2
ProsePose	82	<b>87.8</b>	<b>54.2</b>	<b>73.8</b>	<b>81.4</b>	<b>86.5</b>	<b>91.3</b>

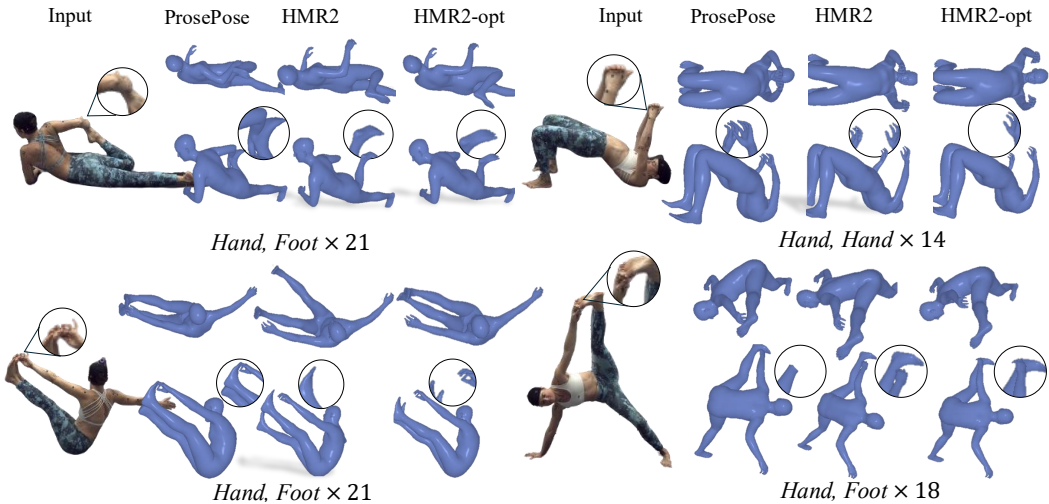


Figure 6: **Single-person results** We show qualitative results from ProsePose, HMR2 (Goel et al., 2023), and HMR2-opt on complex yoga poses. Each example also shows the constraints that are predicted by the LMM at least  $f = 10$  times (and are thus used to compute  $\mathcal{L}_{LMM}$ ) with their counts. ProsePose correctly identifies self-contact points and optimizes the poses to respect these contacts.

self-contact in complex yoga poses relative to the two baselines. Figure 6 provides a qualitative comparison of poses predicted by ProsePose versus the two baselines. Below each of our predictions, we list the corresponding constraints predicted by GPT4-V. In each case, the predicted constraint captures the correct self-contact, which is reflected in the final pose estimates. With the addition of the semantically guided loss, ProsePose effectively refines the pose to ensure proper contact between hand-foot or hand-hand, an important detail consistently overlooked by the baselines.

### 4.3 LIMITATIONS

While ProsePose consistently improves contact across settings and datasets, it has some limitations. First, though we mitigate it through averaging, LMM hallucination of incorrect constraints may lead to an unexpected output. Second, when taking the minimum loss across the possible chiralities of limbs, the pose initialization may lead to a suboptimal choice. We show in § 7.3 examples of failure cases like these. We also note that the LMM may be biased toward poses common in certain cultures due to its training data. In addition, we find that GPT4-V performs worse with some of the camera angles in the MOYO dataset (e.g. frontal or aerial), perhaps because in photos yoga poses are most often captured from a side view.

## 5 CONCLUSION

We present ProsePose, a zero-shot framework for refining 3D pose estimates to capture touch accurately using the implicit semantic knowledge of poses in LMMs. Our key novelty is that we generate structured pose descriptions from LMMs and convert them into loss functions used to optimize the pose. Since ProsePose does not require training, we eliminate the need for the expensive contact annotations used in prior work to train priors for contact estimation. Our framework applies in principle to an arbitrary number of people, and our experiments show in both one-person and two-person settings, ProsePose improves over previous zero-shot baselines. More broadly, this work provides evidence that LMMs are promising tools for 3D pose estimation, which may have implications beyond touch.

## REFERENCES

- 540  
541  
542 OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni  
543 Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor  
544 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum,  
545 Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg  
546 Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage,  
547 Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory  
548 Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen,  
549 Jason Chen, Mark Chen, Benjamin Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave  
550 Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,  
551 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty  
552 Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Sim'on Posada Fishman, Juston Forte,  
553 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel  
554 Goh, Raphael Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene,  
555 Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He,  
556 Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,  
557 Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,  
558 Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn,  
559 Heewoo Jun, Tomer Kaftan, Lukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish  
560 Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik  
561 Kirchner, Jamie Ryan Kiros, Matthew Knight, Daniel Kokotajlo, Lukasz Kondraciuk, Andrew  
562 Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai  
563 Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin,  
564 Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Adeola Makanju,  
565 Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer,  
566 Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake  
567 McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela  
568 Mishkin, Vinnie Monaco, Evan Morikawa, Daniel P. Mossing, Tong Mu, Mira Murati, Oleg Murk,  
569 David M'ely, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo,  
570 Hyeonwoo Noh, Ouyang Long, Cullen O'Keefe, Jakub W. Pachocki, Alex Paino, Joe Palermo,  
571 Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alexandre Passos,  
572 Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov,  
573 Henrique Pondé de Oliveira Pinto, Michael Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly  
574 Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya  
575 Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri  
576 Roussez, Nick Ryder, Mario D. Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather  
577 Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica  
578 Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin,  
579 Katarina Slama, Ian Sohl, Benjamin D. Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski  
580 Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas A. Tezak, Madeleine Thompson, Phil  
581 Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan  
582 Felipe Cer'on Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright,  
583 Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila  
584 Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter,  
585 Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao,  
586 Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang,  
587 Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph.  
588 Gpt-4 technical report. 2023. URL [https://api.semanticscholar.org/CorpusID:  
589 257532815](https://api.semanticscholar.org/CorpusID:257532815). 2, 6
- 587 Anurag Arnab, Carl Doersch, and Andrew Zisserman. Exploiting temporal context for 3d human  
588 pose estimation in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
589 Pattern Recognition*, pp. 3395–3404, 2019. 2
- 590  
591 Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J  
592 Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In  
593 *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October  
11-14, 2016, Proceedings, Part V 14*, pp. 561–578. Springer, 2016. 6

- 594 Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh. Openpose: Realtime multi-person  
595 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis and Machine*  
596 *Intelligence*, 2019. 6
- 597
- 598 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang,  
599 Boyang Albert Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose  
600 vision-language models with instruction tuning. *ArXiv*, abs/2305.06500, 2023. URL <https://api.semanticscholar.org/CorpusID:258615266>. 2, 4
- 601
- 602 Delmas, Ginger and Weinzaepfel, Philippe and Lucas, Thomas and Moreno-Noguer, Francesc and  
603 Rogez, Grégory. PoseScript: 3D Human Poses from Natural Language. In *ECCV*, 2022. 3
- 604
- 605 Delmas, Ginger and Weinzaepfel, Philippe and Moreno-Noguer, Francesc and Rogez, Grégory.  
606 PoseFix: Correcting 3D Human Poses with Natural Language. In *ICCV*, 2023. 3
- 607
- 608 Yao Feng, Jing Lin, Sai Kumar Dwivedi, Yu Sun, Priyanka Patel, and Michael J. Black. Posegpt:  
609 Chatting about 3d human pose. *ArXiv*, abs/2311.18836, 2023. URL <https://api.semanticscholar.org/CorpusID:265506071>. 3
- 610
- 611 Mihai Fieraru, Mihai Zanfir, Elisabeta Oneata, Alin-Ionut Popa, Vlad Olaru, and Cristian Sminchis-  
612 escu. Three-dimensional reconstruction of human interactions. In *Proceedings of the IEEE/CVF*  
613 *Conference on Computer Vision and Pattern Recognition*, pp. 7214–7223, 2020. 2, 4, 7, 3
- 614
- 615 Mihai Fieraru, Mihai Zanfir, Teodor Szente, Eduard Bazavan, Vlad Olaru, and Cristian Sminchis-  
616 escu. Remips: Physically consistent 3d reconstruction of multiple interacting people under weak  
617 supervision. *Advances in Neural Information Processing Systems*, 34:19385–19397, 2021. 2
- 618
- 619 Shubham Goel, Georgios Pavlakos, Jathushan Rajasegaran, Angjoo Kanazawa\*, and Jitendra Malik\*.  
620 Humans in 4D: Reconstructing and tracking humans with transformers. In *International Conference*  
621 *on Computer Vision (ICCV)*, 2023. 1, 2, 6, 9, 10
- 622
- 623 Peng Guan, Alexander Weiss, Alexandru O Balan, and Michael J Black. Estimating human shape  
624 and pose from a single image. In *2009 IEEE 12th International Conference on Computer Vision*,  
pp. 1381–1388. IEEE, 2009. 2
- 625
- 626 Riza Alp Guler and Iasonas Kokkinos. Holopose: Holistic 3d human reconstruction in-the-wild.  
627 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
628 10884–10894, 2019. 2
- 629
- 630 Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating  
631 diverse and natural 3d human motions from text. In *Proceedings of the IEEE/CVF Conference on*  
632 *Computer Vision and Pattern Recognition (CVPR)*, pp. 5152–5161, June 2022. 3
- 633
- 634 Buzhen Huang, Chen Li, Chongyang Xu, Liang Pan, Yangang Wang, and Gim Hee Lee. Closely  
635 interactive human reconstruction with proxemics and physics-guided adaption. In *Proceedings of*  
636 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1011–1021,  
June 2024. 2
- 637
- 638 Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. Motiongpt: Human motion as a  
639 foreign language. *arXiv preprint arXiv:2306.14795*, 2023. 3
- 640
- 641 Wen Jiang, Nikos Kolotouros, Georgios Pavlakos, Xiaowei Zhou, and Kostas Daniilidis. Coherent re-  
642 construction of multiple humans from a single image. In *Proceedings of the IEEE/CVF Conference*  
643 *on Computer Vision and Pattern Recognition*, pp. 5579–5588, 2020. 2
- 644
- 645 Hanbyul Joo, Natalia Neverova, and Andrea Vedaldi. Exemplar fine-tuning for 3d human model  
646 fitting towards in-the-wild 3d human pose estimation. In *2021 International Conference on 3D*  
647 *Vision (3DV)*, pp. 42–52. IEEE, 2021. 2
- 648
- 649 Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of  
human shape and pose. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

- 648 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization.  
649 *CoRR*, abs/1412.6980, 2014. URL [https://api.semanticscholar.org/CorpusID:  
650 6628106](https://api.semanticscholar.org/CorpusID:6628106). 7
- 651
- 652 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
653 Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick.  
654 Segment anything. *arXiv:2304.02643*, 2023. 6
- 655 Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to recon-  
656 struct 3d human pose and shape via model-fitting in the loop. In *Proceedings of the IEEE/CVF  
657 international conference on computer vision*, pp. 2252–2261, 2019. 2
- 658
- 659 Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J Black, and Peter V Gehler.  
660 Unite the people: Closing the loop between 3d and 2d human representations. In *Proceedings of  
661 the IEEE conference on computer vision and pattern recognition*, pp. 6050–6059, 2017. 2
- 662
- 663 Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong  
664 Bing. Mitigating object hallucinations in large vision-language models through visual contrastive  
665 decoding. *arXiv preprint arXiv:2311.16922*, 2023. URL [https://arxiv.org/abs/2311.  
666 16922](https://arxiv.org/abs/2311.16922). 4
- 667 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object  
668 hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023. 4
- 669
- 670 Han Liang, Wenqian Zhang, Wenxuan Li, Jingyi Yu, and Lan Xu. Intergen: Diffusion-based multi-  
671 human motion generation under complex interactions. *arXiv preprint arXiv:2304.05684*, 2023.  
672 3
- 673
- 674 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*,  
675 2023. 2, 4
- 676 Yujie Lu, Dongfu Jiang, Wenhu Chen, William Wang, Yejin Choi, and Bill Yuchen Lin. Wild-  
677 vision arena: Benchmarking multimodal llms in the wild, February 2024. URL [https:  
678 //huggingface.co/spaces/WildVision/vision-arena/](https://huggingface.co/spaces/WildVision/vision-arena/). 6
- 679
- 680 Lea Muller, Ahmed AA Osman, Siyu Tang, Chun-Hao P Huang, and Michael J Black. On self-contact  
681 and human pose. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
682 Recognition*, pp. 9990–9999, 2021. 2, 9
- 683
- 684 Lea Müller, Vickie Ye, Georgios Pavlakos, Michael Black, and Angjoo Kanazawa. Generative  
685 proxemics: A prior for 3d social interaction from images. *arXiv preprint arXiv:2306.09337*, 2023.  
686 2, 6, 7, 8, 3, 4
- 687 Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios  
688 Tzionas, and Michael J. Black. Expressive body capture: 3D hands, face, and body from a single  
689 image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp.  
690 10975–10985, 2019a. 2, 4, 6
- 691
- 692 Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios  
693 Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single  
694 image. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
695 pp. 10975–10985, 2019b. 2, 6
- 696
- 697 Matthias Plappert, Christian Mandery, and Tamim Asfour. The KIT motion-language dataset. *Big  
698 Data*, 4(4):236–252, dec 2016. doi: 10.1089/big.2016.0028. URL [http://dx.doi.org/10.  
699 1089/big.2016.0028](http://dx.doi.org/10.1089/big.2016.0028). 3
- 700
- 701 Abhinanda R. Punnakkal, Arjun Chandrasekaran, Nikos Athanasiou, Alejandra Quiros-Ramirez,  
and Michael J. Black. BABEL: Bodies, action and behavior with english labels. In *Proceedings  
IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 722–731, June 2021. 3

- 702 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
703 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
704 models from natural language supervision. In *International conference on machine learning*, pp.  
705 8748–8763. PMLR, 2021. 1
- 706 Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J Guibas.  
707 Humor: 3d human motion model for robust pose estimation. In *Proceedings of the IEEE/CVF*  
708 *international conference on computer vision*, pp. 11488–11499, 2021. 2
- 709 Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermano. Human motion diffusion as a  
710 generative prior. *arXiv preprint arXiv:2303.01418*, 2023. 3
- 711 Yu Sun, Qian Bao, Wu Liu, Yili Fu, Michael J Black, and Tao Mei. Monocular, one-stage, regression  
712 of multiple 3d people. In *Proceedings of the IEEE/CVF international conference on computer*  
713 *vision*, pp. 11179–11188, 2021. 2
- 714 Yu Sun, Wu Liu, Qian Bao, Yili Fu, Tao Mei, and Michael J Black. Putting people in their place:  
715 Monocular regression of 3d people in depth. In *Proceedings of the IEEE/CVF Conference on*  
716 *Computer Vision and Pattern Recognition*, pp. 13243–13252, 2022. 2, 6, 7, 8, 4
- 717 Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human  
718 motion diffusion model. In *The Eleventh International Conference on Learning Representations*,  
719 2023. URL <https://openreview.net/forum?id=SJ1kSyO2jwu>. 3
- 720 Shashank Tripathi, Lea Müller, Chun-Hao P. Huang, Taheri Omid, Michael J. Black, and Dimitrios  
721 Tzionas. 3D human pose estimation via intuitive physics. In *Conference on Computer Vision and*  
722 *Pattern Recognition (CVPR)*, pp. 4713–4725, 2023. URL <https://ipman.is.tue.mpg.de>.  
723 9
- 724 X Wang, J Wei, D Schuurmans, Q Le, E Chi, S Narang, A Chowdhery, and D Zhou. Self-consistency  
725 improves chain of thought reasoning in language models. *arxiv. Preprint posted online March*, 21:  
726 10–48550, 2022. 5
- 727 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc  
728 Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models.  
729 *ArXiv*, abs/2201.11903, 2022. URL [https://api.semanticscholar.org/CorpusID:  
730 246411621](https://api.semanticscholar.org/CorpusID:246411621). 5
- 731 Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong  
732 Wang. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the*  
733 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18134–18144,  
734 June 2022a. 1
- 735 Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. ViTPose: Simple vision transformer  
736 baselines for human pose estimation. In *Advances in Neural Information Processing Systems*,  
737 2022b. 6
- 738 Yifei Yin, Chen Guo, Manuel Kaufmann, Juan José Zárate, Jie Song, and Otmar Hilliges. Hi4d:  
739 4d instance segmentation of close human interaction. *2023 IEEE/CVF Conference on Com-*  
740 *puter Vision and Pattern Recognition (CVPR)*, pp. 17016–17027, 2023. URL [https://api.  
741 semanticscholar.org/CorpusID:257766362](https://api.semanticscholar.org/CorpusID:257766362). 7
- 742 Andrei Zanfir, Elisabeta Marinoiu, and Cristian Sminchisescu. Monocular 3d pose and shape  
743 estimation of multiple people in natural scenes-the importance of multiple scene constraints. In  
744 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2148–2157,  
745 2018. 2
- 746 Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Shaoli Huang, Yong Zhang, Hongwei Zhao,  
747 Hongtao Lu, and Xi Shen. T2m-gpt: Generating human motion from textual descriptions with  
748 discrete representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
749 *Pattern Recognition (CVPR)*, 2023. 3
- 750  
751  
752  
753  
754  
755

## APPENDIX FOR POSE PRIORS FROM LANGUAGE MODELS

In this appendix, we provide additional details about our method (Section 6), details about metrics (Section 7.1), additional quantitative results (Section 7.2), examples of failure cases (Section 7.3), experiments with a different LMM (Section 7.4), and more qualitative comparisons (Section 7.5). We also provide a video overview of the method and qualitative results: <https://drive.google.com/file/d/1blaLnALiOd4C-au8GW61CtThsolWeLf3/view?usp=sharing>.

## 6 ADDITIONAL METHOD DETAILS

### 6.1 LMM PROMPTS

The box below contains our prompt for the two-person experiments.

You are a helpful assistant. You follow all directions correctly and precisely.  
 For each image, identify all pairs of body parts of Person 1 and Person 2 that are touching.  
 Write all of these in a Markdown table where the first column is "Person 1 Body Part" and the second column is "Person 2 Body Part".  
 You can pick which is Person 1 and which is Person 2.  
 The list of possible body parts is: head, neck, chest, stomach, waist (back), waist (front), back, shoulder (back), shoulder (front), arm, hand, leg, foot, butt.  
 Do not include left/right.  
 List ALL pairs you are confident about.  
 If you are not confident about any pairs, output an empty table.  
 Carefully write your reasoning first, and then write the Markdown table.

The box below contains our prompt for the one-person experiment.

You are a helpful assistant. You answer all questions carefully and correctly.  
 Identify which body parts of the yogi are touching each other in this image (if any).  
 Write each pair in a Markdown table with two columns.  
 Each body part MUST be from this list:  
 head, back, shoulder, arm, hand, leg, foot, stomach, butt, ground  
 Do not write "left" or "right".  
 Describe and name the yoga pose, and then write the Markdown table.  
 Note that the pose may differ from the standard version, so pay close attention.  
 Only list a part if you're certain about it.

In each setting, the prompt is given as the "system prompt" to the GPT-4 API, and the only other message given as input contains the input image with the "high" detail setting.

### 6.2 COARSE REGIONS

Figure 7 illustrates the coarse regions referenced in the prompt in our two-person experiments. Figure 8 illustrates the coarse regions referenced in the prompt in our one-person experiments. In the one-person case, the prompt does not mention the "chest," "neck," or "waist" regions, since they tend to be less important for contacts in yoga poses, and the front/back shoulders are merged into one region, since the distinction tends to be less important for contacts in yoga poses.

### 6.3 CONVERTING CONSTRAINTS TO LOSSES IN 1 VS. 2 PERSON CASES

Our implementation of the conversion from constraints output by the LMM to loss functions differs slightly between the two-person and one-person cases.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

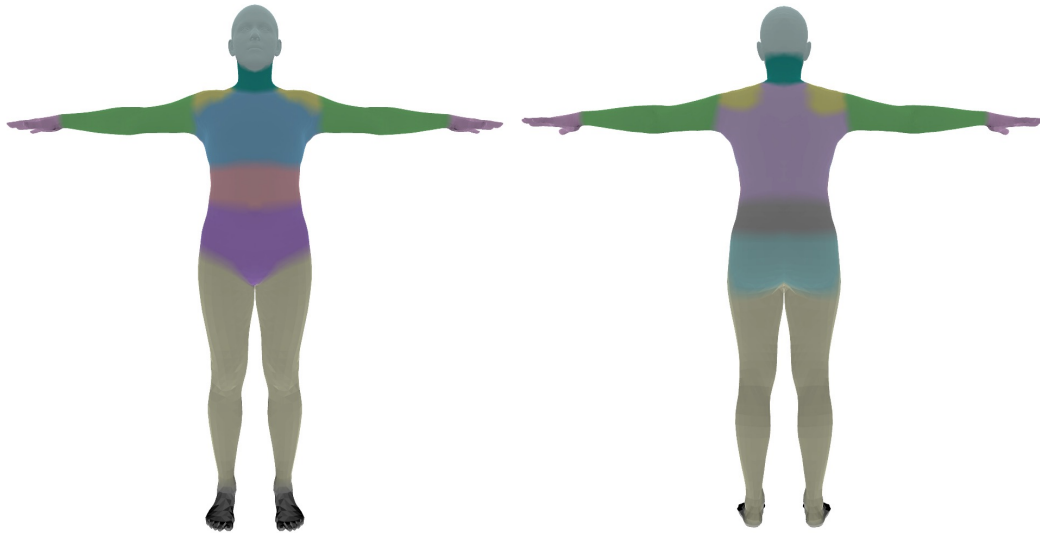


Figure 7: Color-coded coarse regions in the two-person prompt: head, neck, chest, stomach, waist (back), waist (front), back, shoulder (back), shoulder (front), arm, hand, leg, foot, butt. Note that some of these regions overlap. For instance, the “back” includes the “waist (back)” and “shoulder (back)” regions as a subset.

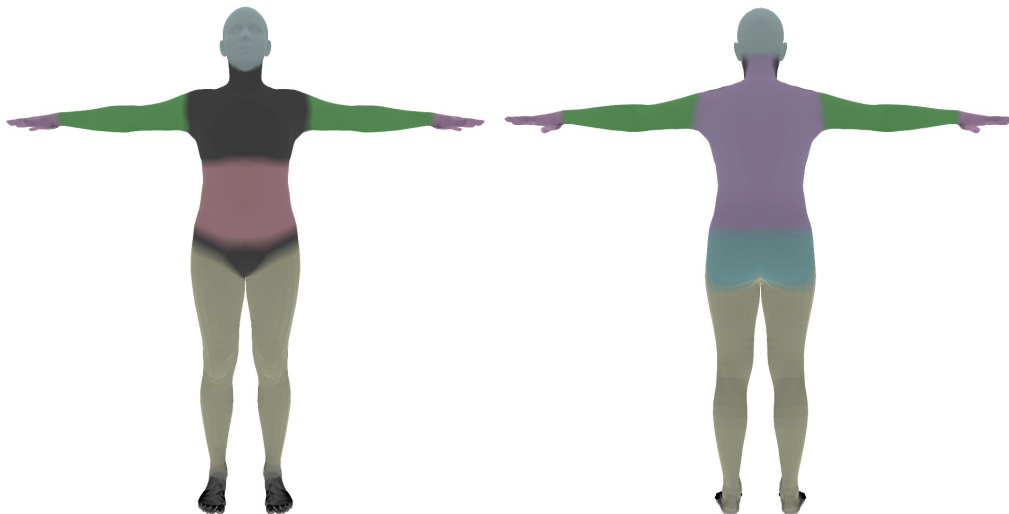


Figure 8: Color-coded coarse regions in the one-person prompt: head, stomach, back, shoulder, arm, hand, leg, foot, butt. Note that the “chest,” “neck,” and “waist (front)” regions are not covered by the regions in the prompt, since they tend to have less importance for contacts in yoga poses.



### 6.3.1 TWO-PERSON

Since we ask the VLM not to differentiate between “left” and “right” limbs, when there should be a constraint on both limbs (e.g. both hands), taking the minimum distance independently for each constraint pair may lead to a constraint on only one limb. Consequently, if the same body part (e.g. “hand”) is mentioned in at least two separate rows of the table output by the LMM (without any “left” or “right” prefix), we enforce that both the left and right limbs of this type must participate in the loss.

We also handle some variations in how the LMM references body parts. First, we check for the following terms in addition to the coarse regions named in the prompt: left hand, right hand, left arm, right arm, left foot, right foot, left leg, right leg, left shoulder, right shoulder, left shoulder (front), right shoulder (front), left shoulder (back), right shoulder (back), waist. “waist” corresponds to the union of “waist (front)” and “waist (back).” Each of these terms is mapped to the corresponding set of fine-grained regions, similar to the coarse regions shown in Figure 7. As stated in Section 3.3 of the main paper, if a “left” or “right” part is explicitly named by the LMM’s output, this part of the coarse region is directly used without considering the other part.

Second, we find there are some cases where the LMM expresses uncertainty between regions using a delimiter like “/” (e.g. “hand / arm”). So we split each entry in the Markdown Table’s output by the delimiter “/” and we compute the loss for each possible region that is listed; we then sum all of these losses.

### 6.3.2 ONE-PERSON

In the one-person experiment, we do not make use of the constraints involving the “ground” that the LMM outputs. Similar to the two-person case, the code for converting the LMM’s output to a loss function checks for the following terms in addition to the body regions listed in the prompt: left hand, right hand, left arm, right arm, left foot, right foot, left leg, right leg, left shoulder, right shoulder, left shoulder (front), right shoulder (front), left shoulder (back), right shoulder (back), waist. Each of these terms is mapped to the corresponding set of fine-grained regions, similar to the coarse regions shown in Figure 7.

## 6.4 BOUNDING BOXES AND CROPPING

As stated in Section 3 of the main paper, we take bounding boxes of the subjects of interest as input and use them to crop the image in order to isolate the person/people of interest when prompting the LMM. For FlickrCI3D, we use the ground-truth bounding boxes of the two subjects of interest. For the other datasets, we use keypoints detected by ViTPose/OpenPose to create the bounding boxes. For the single-person MOYO dataset, we manually check that the bounding boxes from the keypoints and the selected HMR2 outputs correspond to the correct person in the image. We note that the baseline HMR2+opt also benefits from this manual checking, since HMR2+opt also depends on the HMR2 outputs and accurate keypoints.

## 7 EXPERIMENTS

### 7.1 PCC CALCULATION

Figure 9 illustrates the 75 fine-grained regions used for PCC calculation, which are the same as those used in Fieraru et al. (2020). We opted to compute PCC on the fine-grained regions rather than on the coarse ones since prior work uses the fine-grained regions Müller et al. (2023) and since we want to measure contact correctness at a finer granularity (e.g. upper vs. lower thigh vs. knee). Since the regressors BEV and HMR2 use the SMPL mesh while the fine-grained regions are defined on the SMPL-X mesh, we use a matrix  $M \in \mathbb{R}^{\text{num\_vertices\_smplx} \times \text{num\_vertices\_smpl}}$  to convert the SMPL meshes to SMPL-X in order to compute PCC.

### 7.2 PER-PERSON PA-MPJPE

Table 5 shows the per-person PA-MPJPE for each of the datasets used in our two-person experiments.

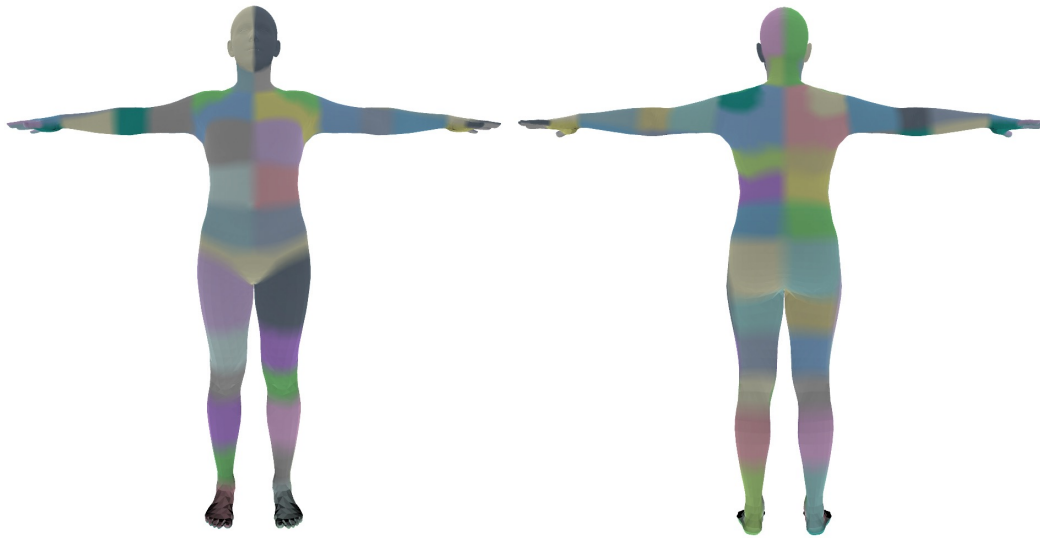


Figure 9: Color-coded 75 fine-grained regions used for PCC calculation

Table 5: **Two-person Results.** Per-person PA-MPJPE (lower is better). For FlickrCI3D, PA-MPJPE is computed using the pseudo-ground-truth fits.

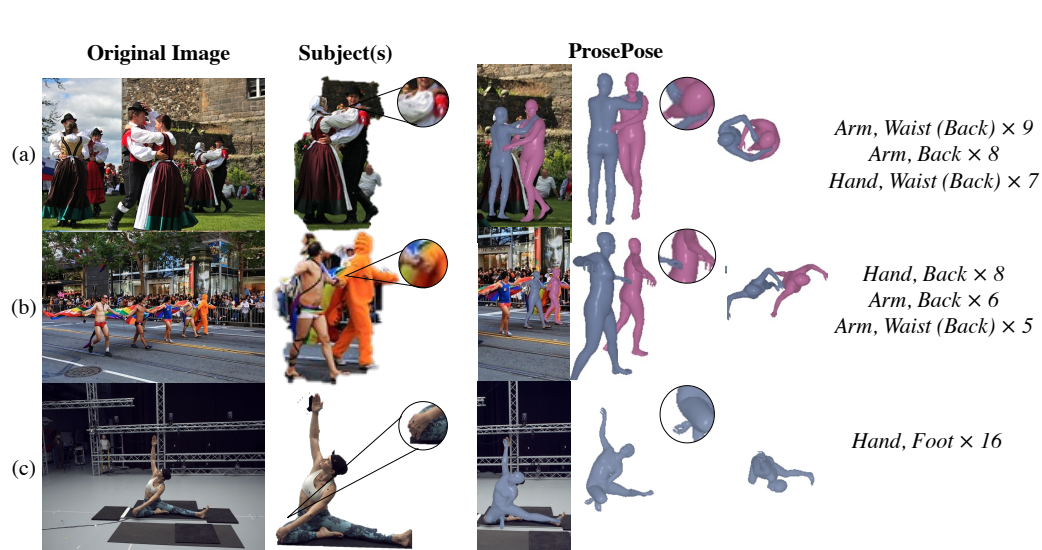
	Hi4D PA-MPJPE $\downarrow$	FlickrCI3D PA-MPJPE $\downarrow$	CHI3D PA-MPJPE $\downarrow$
<i>Without contact supervision</i>			
BEV Sun et al. (2022)	76	71	51
Heuristic	65	31	48
ProsePose	65	31	49
<i>With contact supervision</i>			
BUDDI Müller et al. (2023)	70	43	47

### 7.3 FAILURE CASES

Figure 10 shows examples of two types of LingoPose failures: (1) incorrect chirality (example a) and (2) hallucination (examples b and c). In example (a), the top constraints are correct but without the chirality specified. The optimization then brings both hands of one person to roughly the same point on the other person’s waist, rather than positioning one hand on each hip. Similarly, both hands of the other person are positioned on the same shoulder of the first person. Examples (b) and (c) both show cases of hallucination. In example (b), the hand is predicted to touch the back rather than the hand. In example (c), the hand is predicted to touch the foot rather than the leg. Interestingly, in the yoga example, GPT4-V correctly predicts the name of the yoga pose in all 20 samples (“Parivrtta Janu Sirsasana”). However, it outputs a constraint between a hand and a foot, which is true in the standard form of this pose but not in the displayed form of the pose. Consequently, the optimization brings the left hand closer to the right foot than to the right knee.

### 7.4 DIFFERENT MULTIMODAL MODEL

In this section, we evaluate ProsePose when using a different LMM. We use LLaVA-NeXT 34B (i.e. LLaVA v1.6) Liu et al. (2023) as the LMM. We find that the model does not perform well in directly generating the table of constraints from the image. This is presumably a result of a weaker language model in LLaVA compared to GPT4. Therefore, we instead generate a caption from the LMM, and we feed the caption alone to GPT4 in order to convert it into a table of constraints. We evaluated a few different prompts on the validation sets and chose the prompts with the best performance therein. For the two-person experiments, we use the following prompt for LLaVA:



990 **Figure 10: Failure cases** We show examples in which ProsePose fails to output a semantically correct pose.  
991 The constraints shown are the top 3 constraints (or the total number of constraints, whichever is smaller)  
992 that meet the threshold  $f$  along with their counts ( $f = 1$  for two-person experiments and  $f = 10$  for the one-person  
993 experiment).

994  
995 Describe the pose of the two people.  
996

997 We then use the following prompt with GPT4 to rewrite the caption so that it does not mention left  
998 and right to refer to limbs, since we find that the LMM is not reliably correct in doing so:  
999

1000 Rewrite the caption below so that it doesn't mention "left" or "right" to describe any hand,  
1001 arm, foot, or leg. The revised caption should otherwise be identical. Write only the revised  
1002 caption and no other text.  
1003

1004 We then use the following prompt with GPT4 to create the formatted table.  
1005

1006  
1007 You are a helpful assistant. You will follow ALL rules and directions entirely and precisely.  
1008 Given a description of Person 1 and Person 2 who are physically in contact with each other,  
1009 create a Markdown table with the columns "Person 1 Body Part" and "Person 2 Body Part",  
1010 listing the body parts of the two people that are guaranteed to be in contact with each other,  
1011 from the following list. ALL body parts that you list must be from this list. You can choose  
1012 which person is Person 1 and which is Person 2. Body parts: "chest", "stomach", "waist  
1013 (front)", "waist (back)", "shoulder (front)", "shoulder (back)", "back", "hand", "arm", "foot",  
1014 "leg", "head", "neck", "butt" Note that "back" includes the entire area of the back.  
1015 Include all contact points that are directly implied by the description, not just those that  
1016 are explicitly mentioned. If there are no contact points between these body parts that the  
1017 description implicitly or explicitly implies, your table should contain only the column names  
1018 and no other rows.  
1019 First, write your reasoning. Then write the Markdown table.

1020 For the one-person case, we use the following prompt for LLaVA:  
1021

1022 Describe the person's pose.  
1023  
1024

1025 We use the same prompt as above to rewrite the caption. We then use the following prompt to create  
the formatted table:

Table 6: **LLaVA Results**. Err denotes Joint PA-MPJPE for the two-person datasets (Hi4D, FlickrCI3D, CHI3D) and PA-MPJPE for MOYO. Lower is better for Err, and higher is better for Avg. PCC. **Bold** indicates best method without contact supervision in each column.

	Hi4D	FlickrCI3D		CHI3D		MOYO	
	Err $\downarrow$	Err $\downarrow$	PCC $\uparrow$	Err $\downarrow$	PCC $\uparrow$	Err $\downarrow$	PCC $\uparrow$
Heuristic	116	67	77.8	105	74.1	–	–
HMR2+opt	–	–	–	–	–	81	85.2
GPT4-V	93	58	79.9	100	75.8	82	87.8
LLaVA+GPT4	95	60	79.7	101	75.2	82	85.2

You are a helpful assistant. You will follow ALL rules and directions entirely and precisely. Given a description of a yoga pose, create a Markdown table with the columns "Body Part 1" and "Body Part 2", listing the body parts of the person that are guaranteed to be in contact with each other, from the following list. ALL body parts that you list must be from this list. Body parts: "head", "back", "shoulder", "arm", "hand", "leg", "foot", "stomach", "butt", "ground" Note that "back" includes the entire area of the back. Include all contact points that are directly implied by the description, not just those that are explicitly mentioned. If there are no contact points between these body parts that the description implicitly or explicitly implies, your table should contain only the column names and no other rows. First, write your reasoning. Then write the Markdown table.

We use the `gpt-4-0125-preview` version of GPT4 via the OpenAI API (we obtained better results using this model than `gpt-4-1106-preview`). The latency of this approach is much higher than the single-stage approach used with GPT4-V, since we must feed each caption individually to the OpenAI API. Therefore, we set  $N = 5$  for these experiments. Since we change  $N$ , we also need to select appropriate thresholds  $f$  and  $t$ . As in the experiments with GPT4-V, we set  $t = N$  for all datasets except CHI3D. For CHI3D, we find on the validation set that  $t = 2$  works better than  $t = 1$ , so we set  $t = 2$ . As in the experiments with GPT4-V, we set  $f = 1$  for the 2-person datasets, and we set  $f = 3$  for MOYO, to approximate the ratio  $f/N$  used in the GPT4-V experiments. Finally, when converting the constraint pairs to loss functions, we found that on a small number of examples, the pipeline produced a large number of constraints, leading to very slow loss functions. Therefore, we discarded loss functions that are longer than 10000 characters.

Table 6 shows the results. On the 2-person datasets, the LLaVA+GPT4 approach performs better than the contact heuristic but not as well as GPT4-V. This is in line with holistic multimodal evaluations that indicate that GPT4-V performs better than LLaVA Lu et al. (2024). On the 1-person yoga dataset, the performance of LLaVA+GPT4 is comparable with that of the baseline (HMR2+opt). The reason that LLaVA performs worse than GPT4-V in this setting may be that LLaVA does not have enough training data on yoga to provide useful constraints.

## 7.5 ADDITIONAL QUALITATIVE RESULTS

Figures 11, 12, 13, and 14 show additional, randomly selected examples from the multi-person FlickrCI3D test set. Figures 15, 16, 17, and 18 show the same examples comparing ProsePose with the pseudo-ground truth fits. Figures 19, 20, and 21 show additional, randomly selected examples from the Hi4D test set. Figures 22 and 23 show additional, randomly selected examples from the CHI3D validation set (which we use as the test set following Müller et al. (2023)). Figures 24 and 25 show additional, randomly selected examples from the 1-person yoga MOYO test set.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

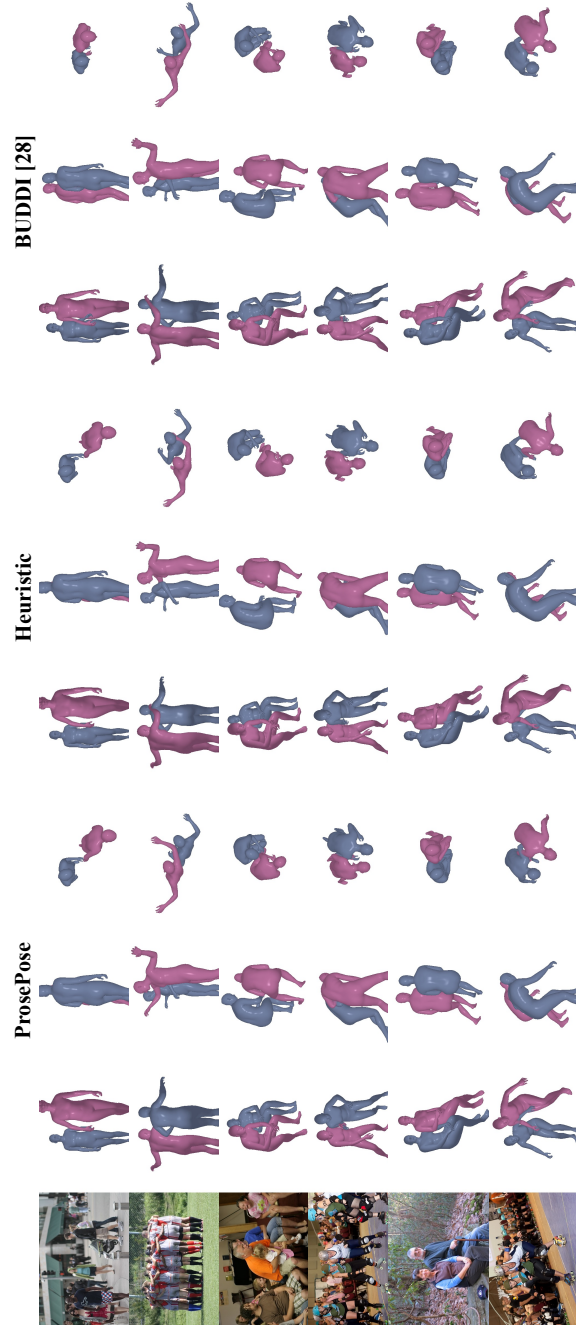


Figure 11: Non-curated examples from the FlickrCI3D test set. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

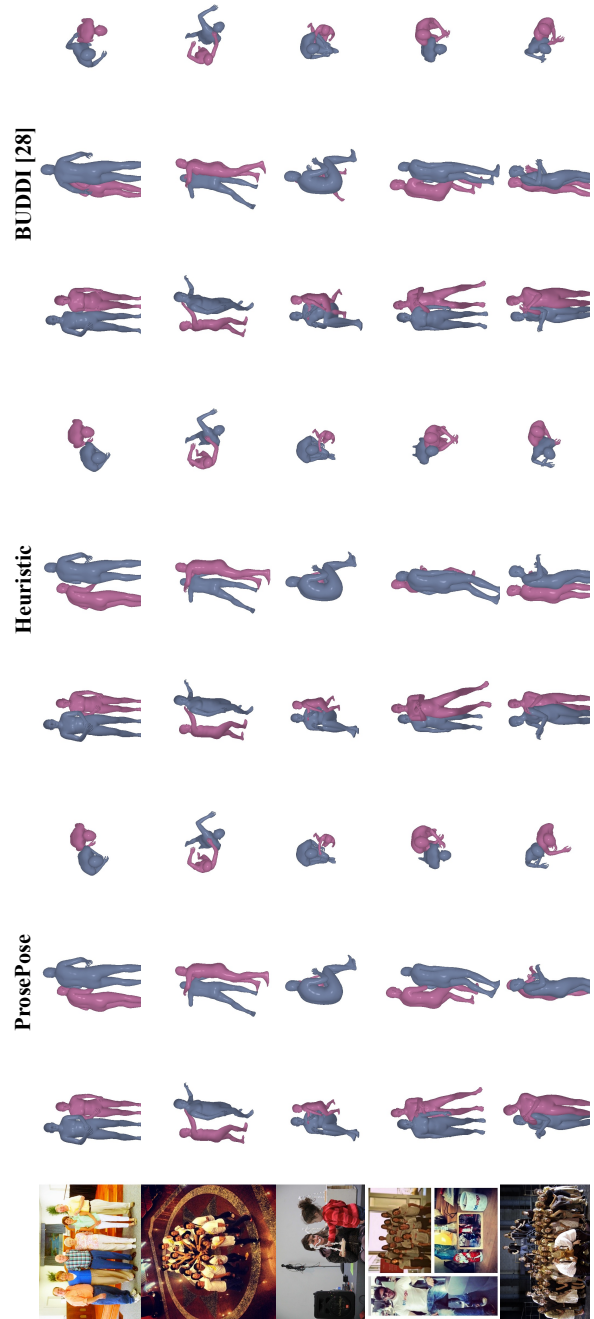


Figure 12: Non-curated examples from the FlickrCI3D test set. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

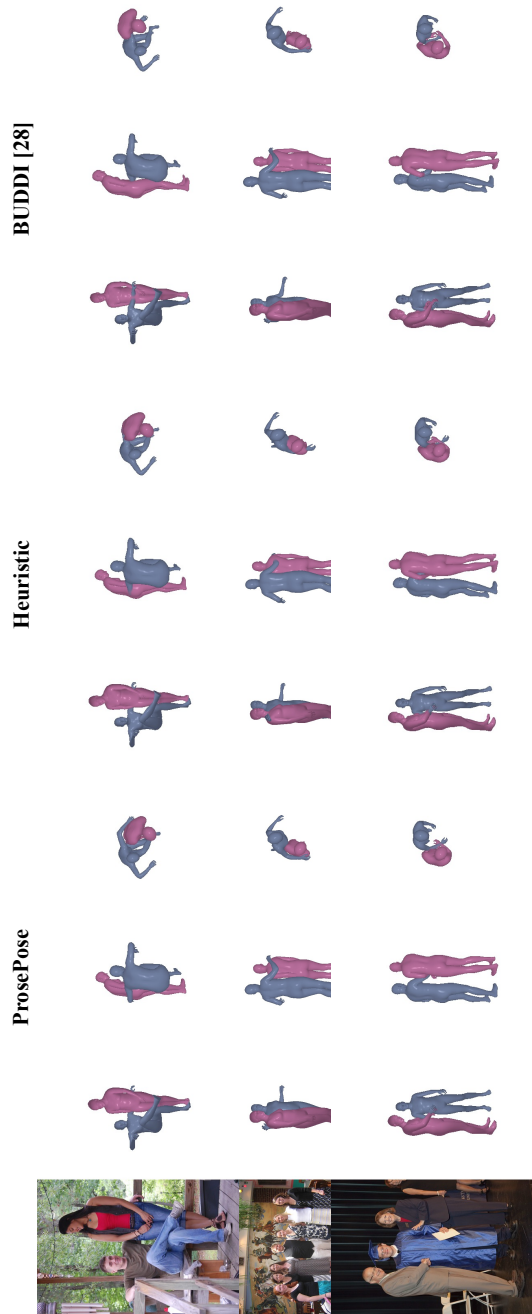


Figure 13: Non-curated examples from the FlickrCI3D test set. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

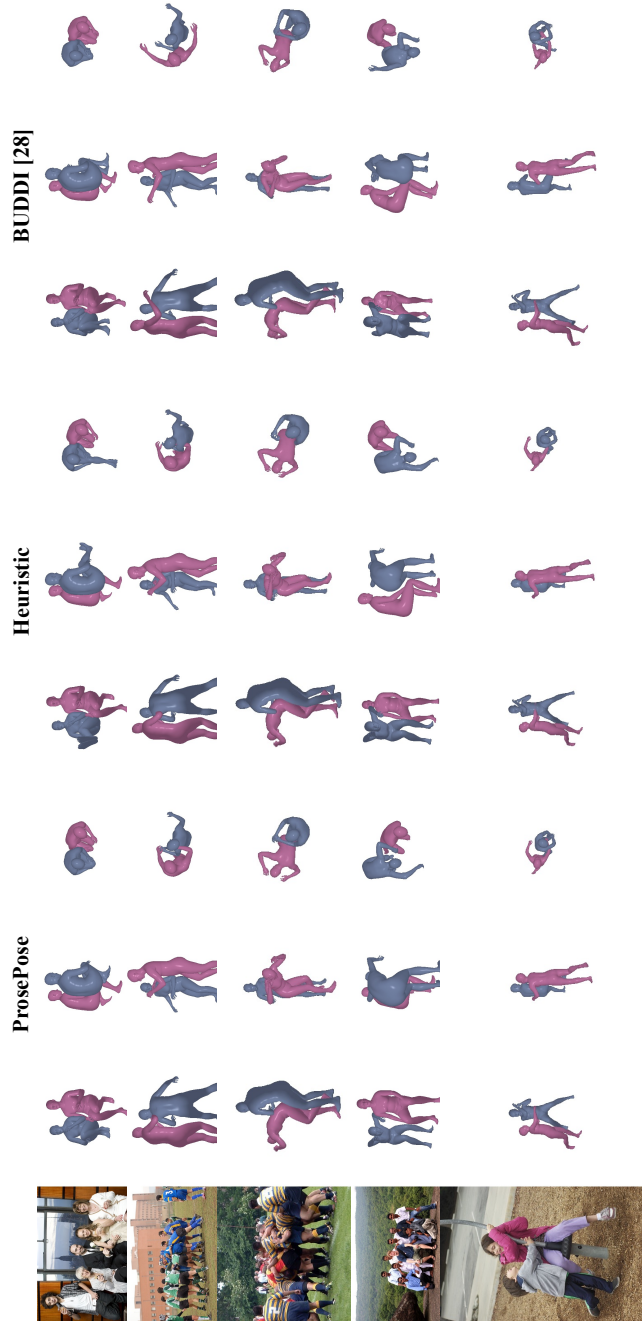


Figure 14: Non-curated examples from the FlickrCI3D test set. They are randomly selected from the examples for which there is at least one non-empty constraint set.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

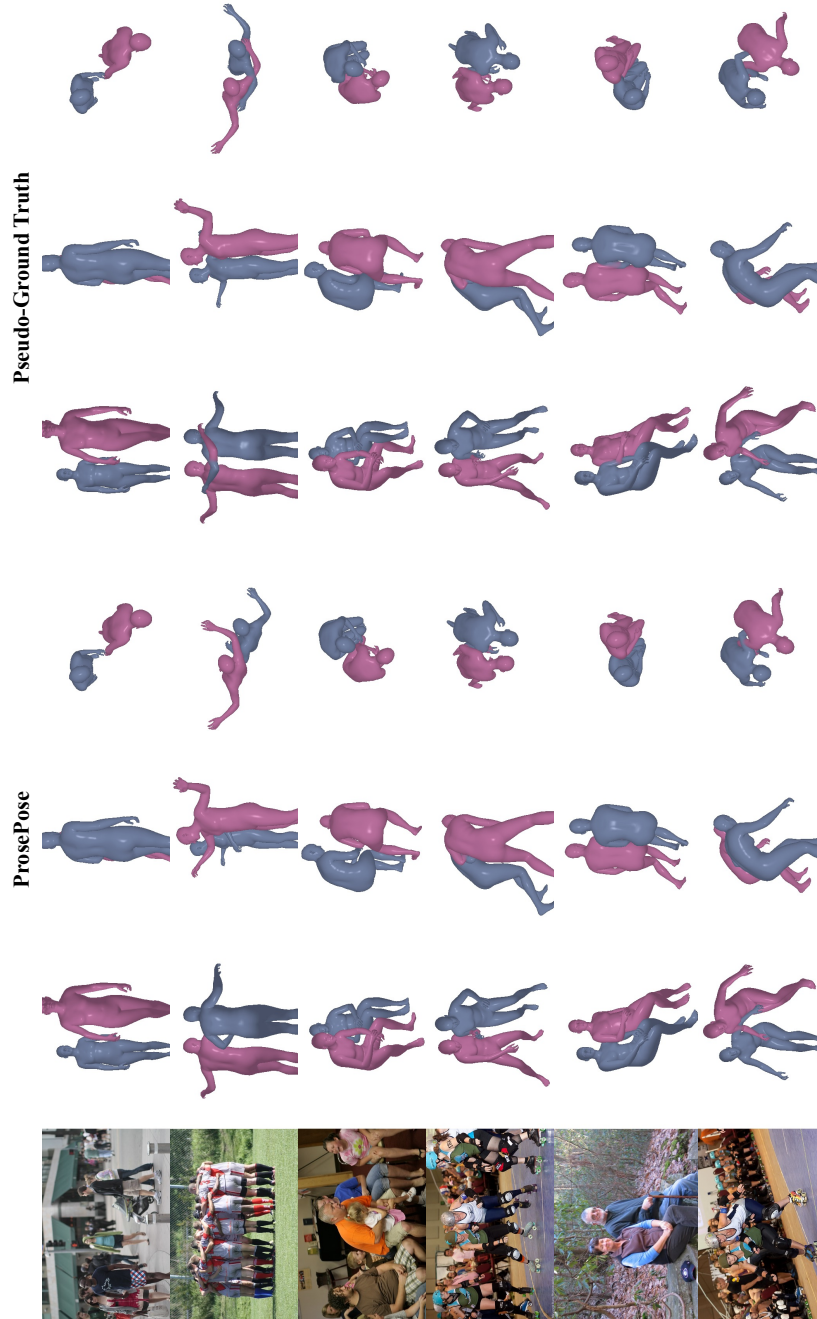


Figure 15: Non-curated examples from the FlickrCI3D test set, comparing ProsePose with the pseudo-ground truth fits. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

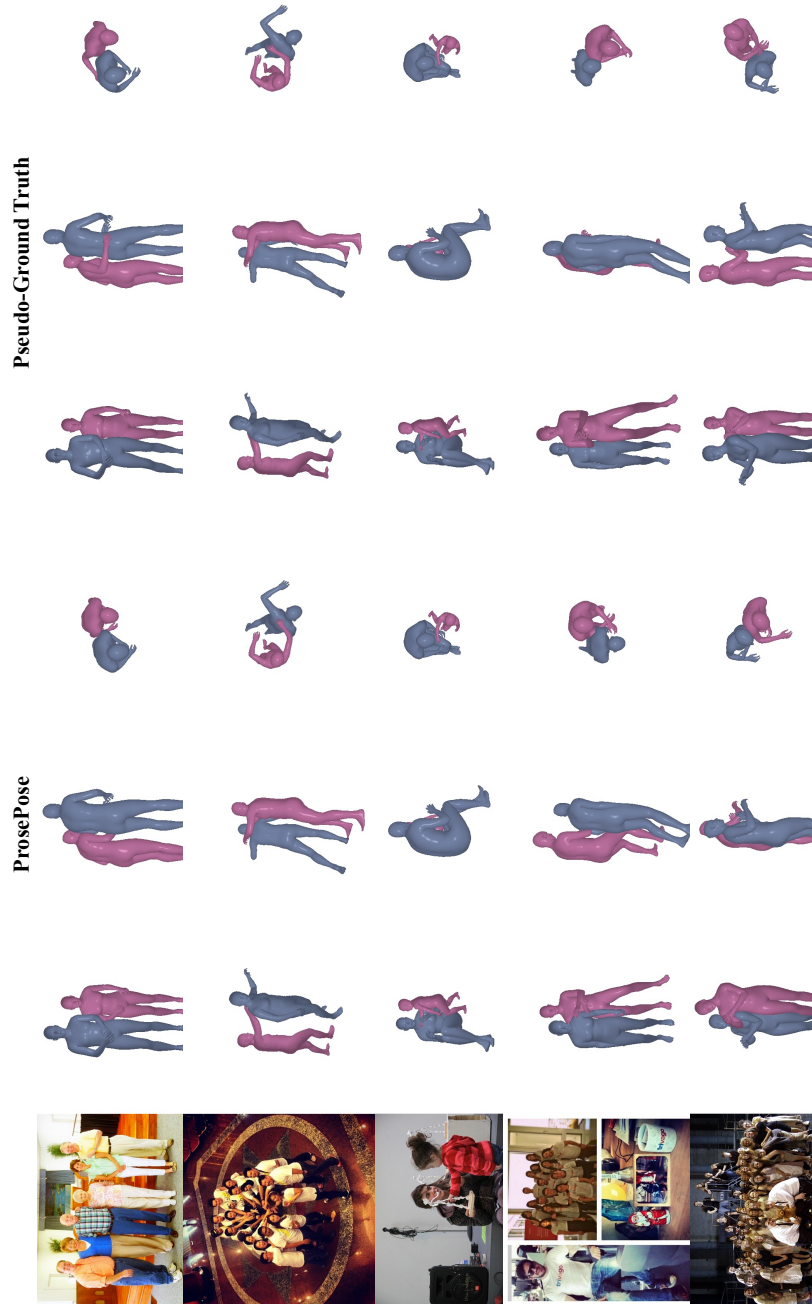


Figure 16: Non-curated examples from the FlickrCI3D test set, comparing ProsePose with the pseudo-ground truth fits. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

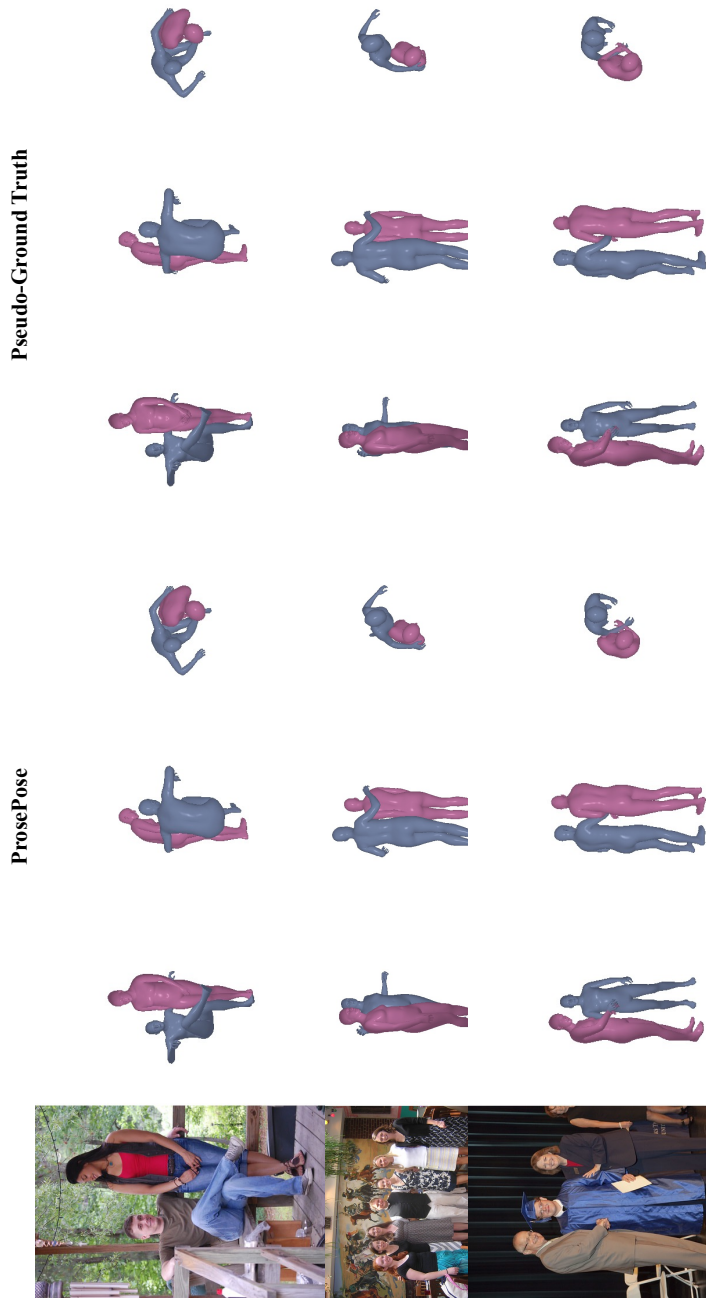


Figure 17: Non-curated examples from the FlickrCI3D test set, comparing ProsePose with the pseudo-ground truth fits. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

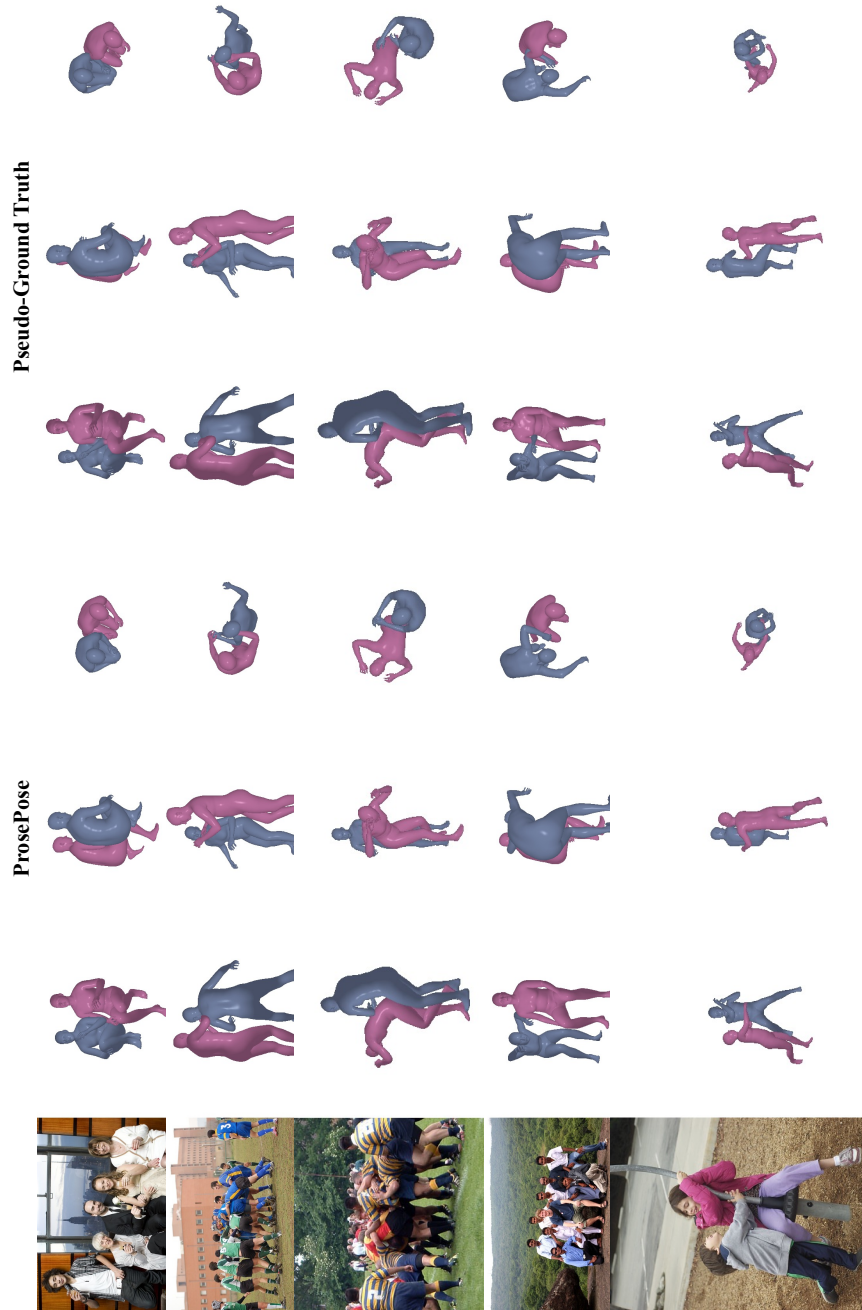


Figure 18: Non-curated examples from the FlickrCI3D test set, comparing ProsePose with the pseudo-ground truth fits. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

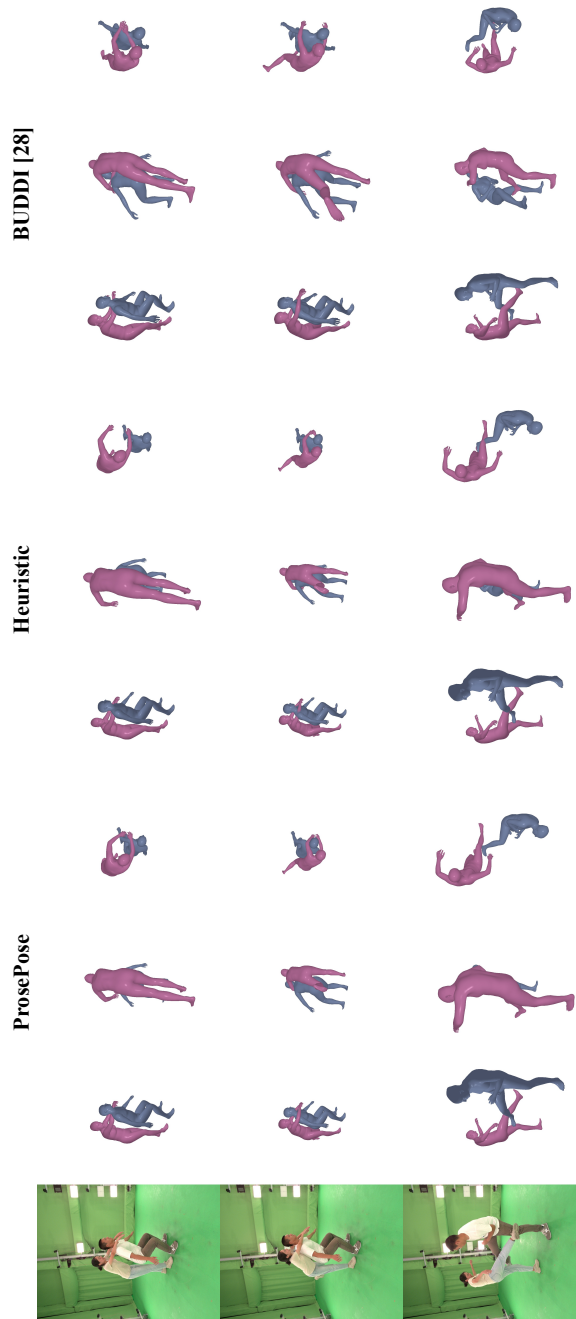


Figure 19: Non-curated examples from the Hi4D test set. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

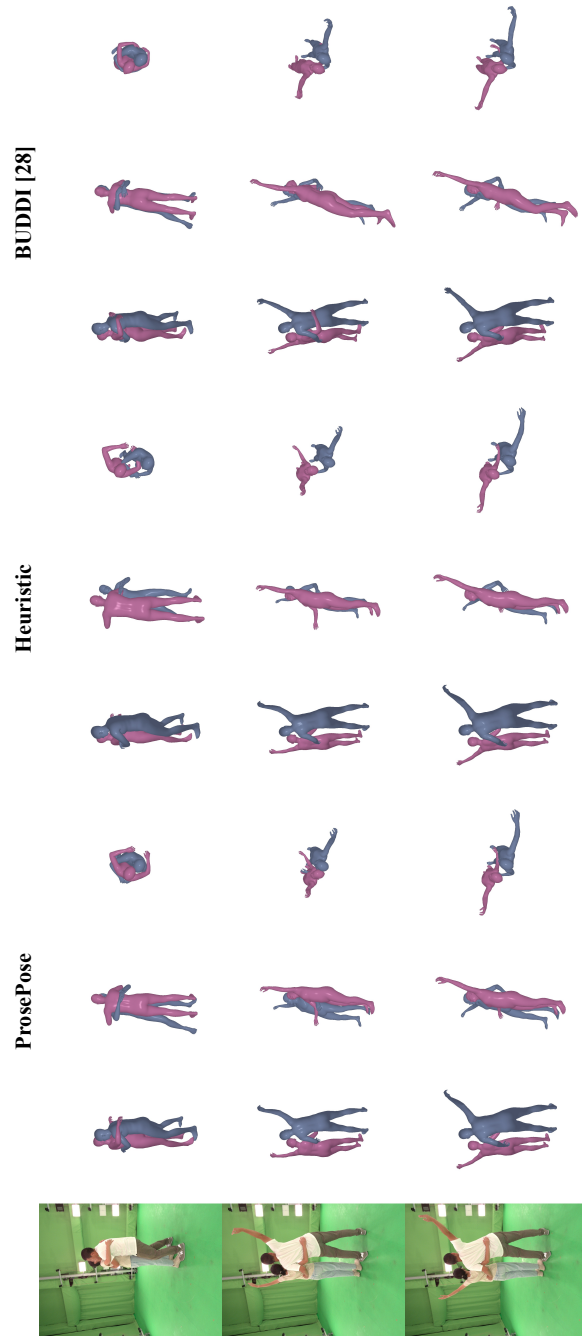


Figure 20: Non-curated examples from the Hi4D test set. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673

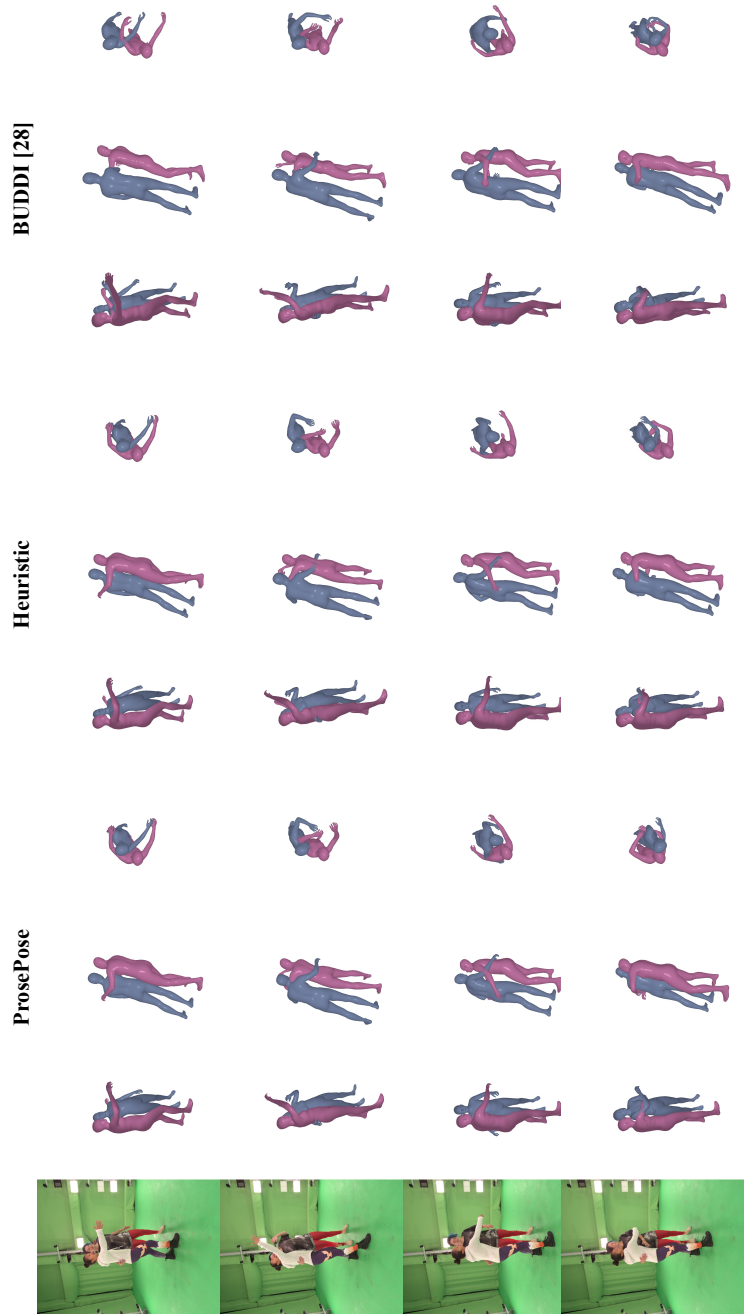


Figure 21: Non-curated examples from the Hi4D test set. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727

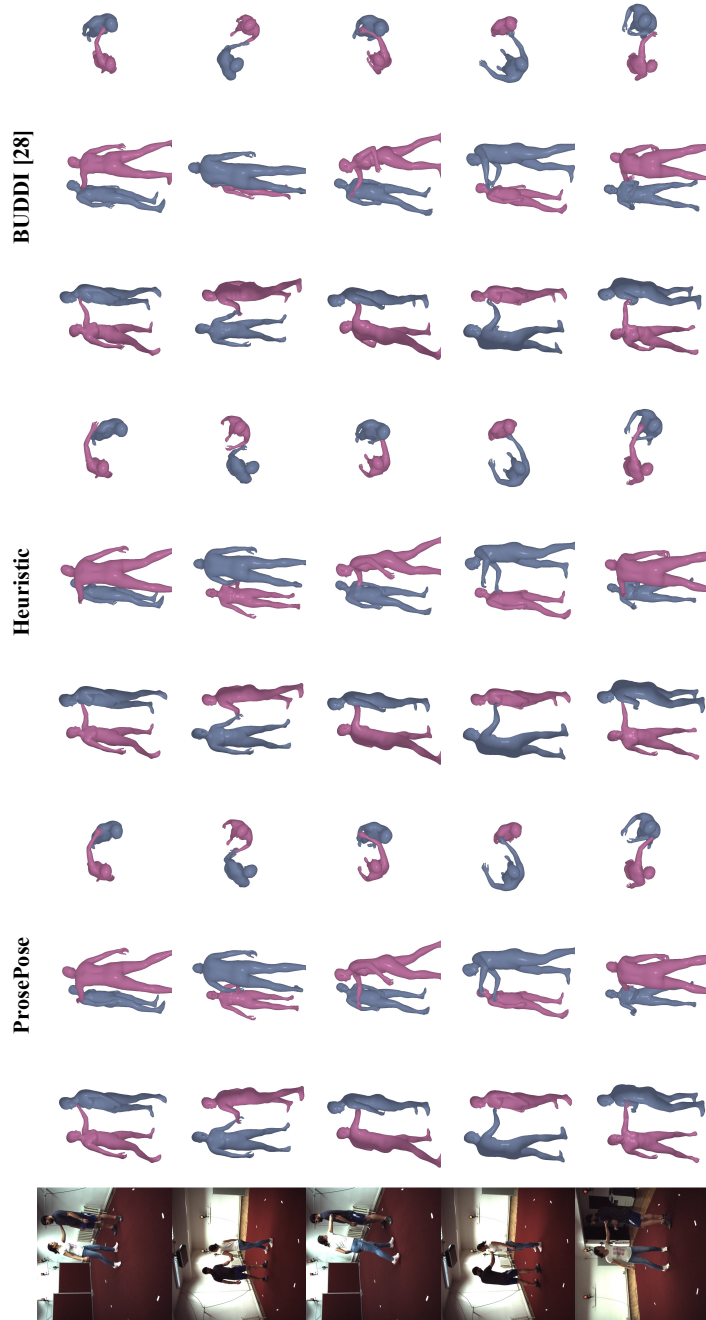


Figure 22: Non-curated examples from the CHI3D validation set (which we use as the test set). They are randomly selected from the examples for which there are at least nineteen non-empty constraint sets (since we set  $t = 2$  for CHI3D).



1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

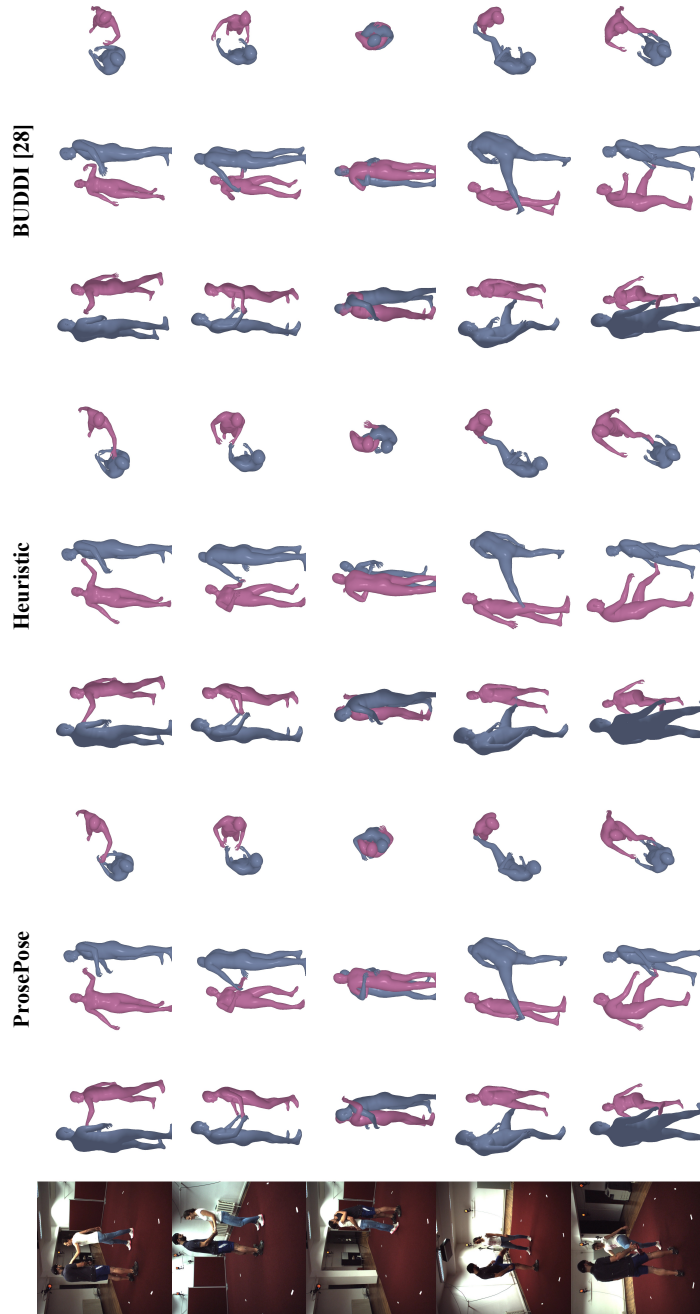


Figure 23: Non-curated examples from the CHI3D validation set (which we use as the test set). They are randomly selected from the examples for which there are at least nineteen non-empty constraint sets (since we set  $t = 2$  for CHI3D).

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835



Figure 24: Non-curated examples from the MOYO test set. They are randomly selected from the examples for which there is at least one non-empty constraint set.

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

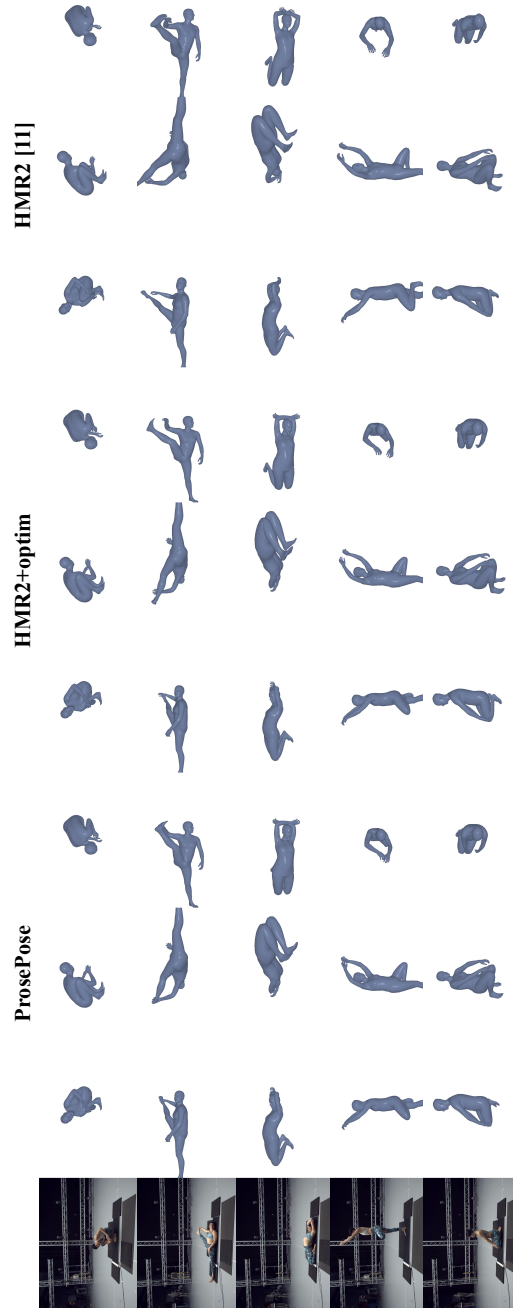


Figure 25: Non-curated examples from the MOYO test set. They are randomly selected from the examples for which there is at least one non-empty constraint set.