

Purposer: Putting Human Motion Generation in Context

Nicolas Ugrinovic¹ Thomas Lucas² Fabien Baradel² Philippe Weinzaepfel²
Grégory Rogez² Francesc Moreno-Noguer¹

¹Institut de Robòtica i Informàtica Industrial, CSIC-UPC, Barcelona, Spain

²NAVER LABS Europe

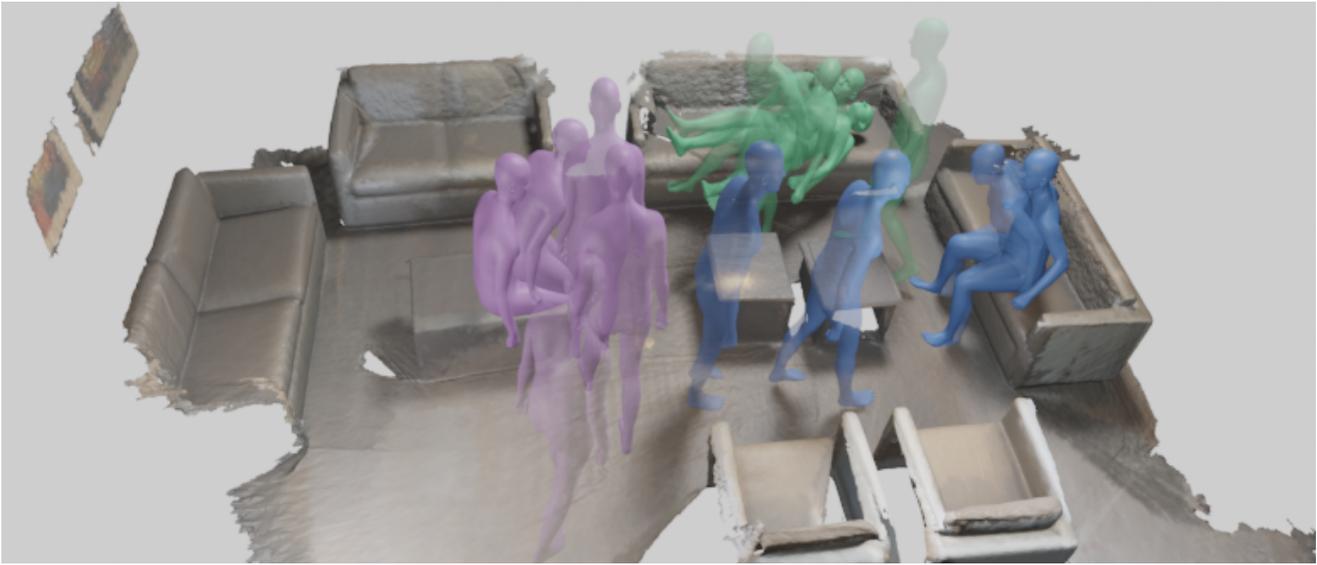


Figure 1. **An example of human motion generation in context.** We propose a method able to generate realistic-looking motions that interact with virtual scenes. In this example we take a scene from ScanNet [11]. The motion can be controlled with semantic action/object queries: here the human is first commanded ‘sit on table’, then ‘sit on couch’, and finally ‘lie on couch’. Purposer is a learning-based probabilistic model that can work efficiently with diverse types of conditioning.

Abstract

We present a novel method to generate human motion to populate 3D indoor scenes. It can be controlled with various combinations of conditioning signals such as a path in a scene, target poses, past motions, and scenes represented as 3D point clouds. State-of-the-art methods are either models specialized to one single setting, require vast amounts of high-quality and diverse training data, or are unconditional models that do not integrate scene or other contextual information. As a consequence, they have limited applicability and rely on costly training data. To address these limitations, we propose a new method, dubbed Purposer, based on neural discrete representation learning. Our model is capable of exploiting, in a flexible manner, different types of information already present in open access large-scale datasets such as AMASS. First, we encode unconditional human motion into a discrete latent space. Second, an auto-

regressive generative model, conditioned with key contextual information, either with prompting or additive tokens, and trained for next-step prediction in this space, synthesizes sequences of latent indices. We further design a novel conditioning block to handle future conditioning information in such a causal model by using a network with two branches to compute separate stacks of features. In this manner, Purposer can generate realistic motion sequences in diverse test scenes. Through exhaustive evaluation, we demonstrate that our multi-contextual solution outperforms existing specialized approaches for specific contextual information, both in terms of quality and diversity. Our model is trained with short sequences, but a byproduct of being able to use various conditioning signals is that at test time different combinations can be used to chain short sequences together and generate long motions within a context scene.

1. Introduction

Generating realistic and diverse human motion is a decades-old research problem [4, 5] but has gained traction in recent years [3, 27, 34, 35, 61]. In this work, we propose a learning-based model for motion generation that can be controlled using various forms of *contextual information* in order to navigate and interact with virtual scenes. In practice and as illustrated in Figure 1, human motion is strongly determined by several forms of context. Among them are: scene geometry, semantics of the surrounding objects, past motion and target actions and poses. So far, already established approaches have focused on narrow subsets of these. For instance, [27, 30] both condition on actions and past poses but do not consider scene or target goals. To extend their applicability to VR/AR and other potential areas, the generated motion needs to make sense for a given scene. This requires taking into account past motion [2, 7, 15, 56, 63] together with scene geometry [17, 52, 53]. However, human motion data in context is scarce; this hinders the development of powerful conditional models.

In the PROX [16] dataset, the amount of human motion data available together with detailed scene information is two orders of magnitude smaller than AMASS [29]. The lack of conditional data limits the expressivity of the models used and is not the regime in which recent deep learning methods excel. In that scarce data regime, existing scene-conditioned methods rely on test time optimization loops, which allow them to effectively take into account scene boundaries, but affects the realism of the generated motion [52, 53]. In contrast, we leverage the recent HUMANISE dataset [54] to learn scene interactions from the data. We also use unconditional data from where we learn a powerful motion prior.

We build our model on top of PoseGPT [27], itself based on neural discrete representation learning [51]. Thus, in our model, human motion is first mapped into an abstract discrete feature space, *without any conditioning*. Any human motion given as input can be represented as a *trajectory* in that discrete latent space, *i.e.*, a sequence of centroids. After this, motion is modeled in a *probabilistic* manner, directly in that latent space, by predicting latent trajectories in an auto-regressive manner. At this stage, various forms of contextual information can be used to condition the model and reduce prediction uncertainty. The latent trajectory is then mapped back into a continuous motion representation and latent trajectories are finally decoded into motion.

We propose a method that can take advantage of various combinations of contextual information. We account for three broad categories of contextual information, that can be combined together arbitrarily. First, we use the **scene geometry**. The scene is represented as a point cloud, encoded, and used to condition our generative model in latent space to exploit this information. Second, we use **past**

observations and future targets. A limitation in existing auto-regressive approaches is that they cannot easily be conditioned on time-dependent future information, because of their causal design. To remedy this, we propose a simple and flexible architecture that allows us to effectively condition our model on future trajectories or randomly selected future poses. Finally, we use **semantic information**. To achieve semantic control, we condition the second stage model on *target poses* which are generated using pairs of actions and object labels as targets as proposed in [64]. This offers semantic control over the generated sequences.

By combining this with conditioning on the past, we are able to chain multiple action/object targets together, which offers even more flexible semantic control and allows us to generate longer motion sequences, despite training on short-term sequences (HUMANISE). For instance, one can generate long sequences with multiple actions at different locations in the scene (*e.g.* conditioned on an interaction with nearby objects) while using a conditioning corresponding to locomotion to navigate (*i.e.*, move along a path in the scene from this first object to this second object).

In summary, we present a model capable of leveraging unconditional data together with combinations of contextual information and generate motion to populate virtual scenes. Our model (a) can leverage large amounts of unconditional data, (b) can adapt to various contexts and (c) offers fine control on model outputs.

We train our auto-encoder on large-scale unconditional data from the BABEL dataset [36] and our auto-regressive component with various combinations of conditioning signals on the HUMANISE dataset [54] and further fine-tune it on PROX [16]. To evaluate our approach, we measure sample quality and sample diversity, as well as our model’s generalization capabilities following practices established by existing work on uncontextualized motion generation [27, 34], inspired from the image generative modeling literature [6, 26, 31, 42]. We also evaluate the synthesized motion’s coherence with the scene, namely *physical plausibility*, using contact and non-collision scores [52, 62]. In this manner, we show that our proposed approach generates high-quality motions to populate virtual scenes. We provide video results and code at the [project page](#).

2. Related work

Human motion generation. The task of class-conditional human motion synthesis was first tackled assuming cyclic human actions such as walking [46, 50]. More recent work have focused on adapting generative models to action conditional 3D human motion generation [13, 34], and some approaches have explored conditioning on past poses [27, 30]. However, these methods do not condition on contextual information about the scene, which limits their applicability in practice. Another promising research av-

enue to control generated motion is to condition the model on high-level but detailed textual descriptions, as explored in [1, 3, 12, 24, 25, 35, 40, 47, 48, 57] or audio representations [22, 23]. While these approaches offer fine-grained control over the generated motions, they do not allow to generate motions in a given environment.

Scene interaction synthesis. It was not until recent years that the community focused its attention on estimating [10, 16, 49] and generating [18, 58, 62] human poses taking into account a 3D scene context. This was shown to improve both 3D pose and motion estimations [16, 28, 38, 41, 59]. Most recently, COINS [64] propose a framework that adds semantic control to this generation process. By augmenting the PROX dataset [16] with action-object paired labels and developing a specialized model, they generate semantically coherent poses. Building on their work, we go beyond static poses and propose a *motion* model that can be conditioned on action object pairs.

Object-conditioned human motion generation. One existing line of work focuses on conditioning motion generation on contextual or interaction information, be it nearby small [45, 55], medium [60] or dynamic [9] objects. In these cases, emphasis is given to one single object at a time. GOAL [45] and SAGA [55] focus on generating whole-body motions to match a final hand-grasping pose. COUCH [60] on the other hand focuses only on chairs, thus capturing human-chair interactions. By contrast, we focus on modeling more general interactions between human motion and an unconstrained number of objects within a scene.

Scene-conditioned human motion generation. Up until now, few works have fully studied scene-conditioned human motion generation [17, 44, 52, 53, 61]. Neural State Machine [44] generates different modes of motion that can be blended between different actions while interacting with the environment. This model allows excellent motion control while providing smooth transitions between modes. However, this method was designed for simple hand-crafted environments and relies on a deterministic model, limiting its ability to produce diverse motion and to model the full extent of human motion. Recently, [52] tackles the task of generating long-term motion given a 3D scene and start/end goal positions using a hierarchical framework that decomposes the task by synthesizing shorter motion sequences. This method relies on a post-optimization step to ensure smoothness, robust foot contact, and avoiding collisions with the scene. While effective, this optimization step reduces the naturality of the motion. SAMP [17] creates human motion conditioned on the action and a final target object, position and orientation in a stochastic manner. Given a starting position, and a target object (*e.g.* chair, sofa), they first estimate a goal position and orientation and then estimate a plausible path between the start and goal positions. Finally, they generate a sequence of poses with an auto-

regressive conditional Variational Auto-Encoder (cVAE). Wang *et al.* [53] propose a model composed of various networks each specialized on one sub-task: generating target start/end poses, path planning, and sequential human poses generation. They rely on cVAE networks conditioned on actions and on a generated path. However, they use the same optimization step as in [52] to reduce foot skating and scene penetration and thus suffer from a similar lack of naturality.

Reinforcement learning methods have also been used to tackle this problem [19, 39, 61, 65]. Zhang *et al.* [61] and Rempe *et al.* [39] mostly focus on generating realistic locomotion taking the scene topology into account. Hassan *et al.* [19] use a physical simulated character and imitation learning to generate diverse actions within an environment. Concurrent to our work, DIMOS [65] extends [61] to include more actions that interact with the environment.

Finally, [54] contributed a synthetic dataset, HUMANISE, that places a subset of motion capture (MoCap) sequences from AMASS [29] dataset in scenes from [11]. In this work, we leverage HUMANISE dataset to include scene context. This way, we are able to generate realistic motions of humans navigating and interacting in a scene. Our approach is most similar to [52, 53], but yield richer interactions and more realistic motions. We present a direct comparison to these two approaches in Section 4.

3. Purposer

We build our Purposer model on auto-regressive discrete-based generative models such as PoseGPT [27], T2M2 [14], T2M-GPT [57] or Bailando [43]. we detail how we propose to condition such causal methods, in particular in the case of future conditioning, in Section 3.2. We then discuss the different forms of contextual information that we consider for motion generation (Section 3.3). Finally, we detail our training setup in Section 3.4 and how various conditioning signals can be combined to generate long-term sequences while being trained on short ones (Section 3.5).

3.1. Background on discrete auto-regressive models

Discretization-based auto-regressive models proceed in two stages, see Figure 2: (a) an auto-encoder is learned to move from the continuous input space to a discrete latent space and vice-versa, (b) an auto-regressive model is learned in this discrete space, and can be fed to the decoder for obtaining the output in the desired space. We now give more background on these two stages.

Discrete motion auto-encoder. An auto-encoder is learned to compress motion sequences into discrete latent representations with neural discrete representation learning [51], see top row of Figure 2. Concretely, an encoder $E(\cdot)$, a quantizer $Q(\cdot)$ with a codebook and a decoder $D(\cdot)$ are trained such that the reconstruction error is minimized. A given

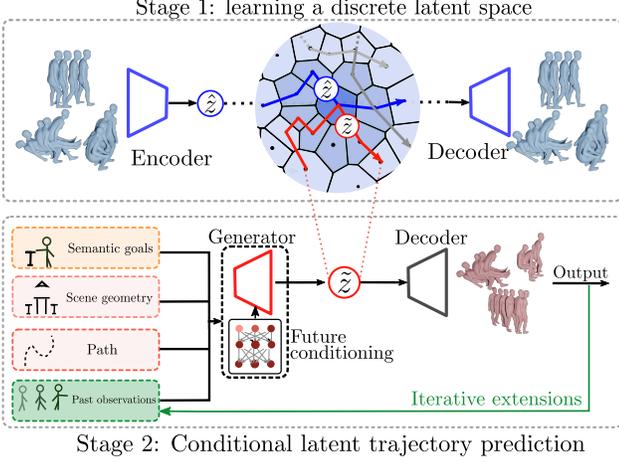


Figure 2. **Method Overview.** An auto-encoder is learned to compress human motion, without any context, into a discrete latent sequence space (top). A probabilistic model (bottom) is trained directly in that space, with three types of optional context: (a) scene geometry, (b) semantic goals, (c) observation of past motion.

motion sequence \mathbf{p} of length T can be represented by a discrete sequence of indices $\mathbf{z} = \{\mathbf{z}_1, \dots, \mathbf{z}_{T'}\}$ of length T' , by computing $Q(E(\mathbf{p}))$. Conversely, any sequence of discrete latent indices can be decoded into a motion sequence by forwarding it to the decoder D . Note that here, we use T' instead of T as the sequence in pose space \mathbf{p} can be downsampled when converting to the latent discrete space \mathbf{z} and then upsampled again by using the decoder D . To allow conditioning on past observations, a causal encoder is used, such that for any $t \leq T'$, $\hat{\mathbf{z}}_t$ is a function of $\{\mathbf{p}_1, \dots, \mathbf{p}_{\lfloor t \cdot T/T' \rfloor}\}$ only. In this work, we rely on the discrete motion auto-encoder from PoseGPT [27] that further uses product quantization for better leveraging the discrete space.

Auto-regressive prediction. The auto-regressive model can then be learned directly in the frozen discretized latent space. Any motion sequence \mathbf{p} of length T can be represented as $\mathbf{z} = \{\mathbf{z}_1, \dots, \mathbf{z}_{T'}\}$. To generate trajectories in the latent space, an auto-regressive model $G(\cdot)$ can be trained to predict the next index as the successful GPT [8] family in natural language processing, *i.e.*, by maximizing:

$$p_G(\mathbf{z}) = p(\mathbf{z}_1) \prod_{t=2}^{T'} p(\mathbf{z}_t | \mathbf{z}_1, \dots, \mathbf{z}_{t-1}). \quad (1)$$

To obtain motion samples, latent sequences are sampled from p_G and decoded using the decoder $D(\cdot)$. Such auto-regressive models can elegantly be conditioned on past motion when using a causal encoder.

3.2. Future conditioning in auto-regressive models

While conditioning an auto-regressive model such as a GPT with a *sequence-wide* information, *i.e.*, a fixed context across the full sequence (*e.g.*, static scene information),

can be easily implemented, it is not straightforward to condition on *future* information (*e.g.*, a target pose or a path).

Sequence-wide conditioning. Some types of conditioning are valid for the full sequence – for instance static scene information, a sequence duration T , or a constant action label. In that case, given an input sequence $(\mathbf{z}_1, \dots, \mathbf{z}_{T'})$ embedded into features $\mathbf{h} = (\mathbf{h}_1, \dots, \mathbf{h}_{T'})$ and some conditioning signal \mathbf{c} represented by a feature vector \mathbf{h}_c , conditioning the auto-regressive model can be done simply by *prompting*, *i.e.*, adding \mathbf{h}_c as an extra token at the start of the input sequence, as commonly done *e.g.* in language models [32]:

$$\tilde{\mathbf{h}}_{\text{prompt}} = (\mathbf{h}_c, \mathbf{h}_1, \dots, \mathbf{h}_{T'}). \quad (2)$$

Another solution is to inject it into all input tokens:

$$\tilde{\mathbf{h}}_{\text{feat}} = (\mathbf{h}_1 \oplus \mathbf{h}_c, \dots, \mathbf{h}_{T'} \oplus \mathbf{h}_c), \quad (3)$$

where the \oplus operation denotes any operator that combines the two features, such as concatenation or sum.

Conditioning with causal masking. Let us now consider a time-dependent conditioning $\mathbf{c}_1, \dots, \mathbf{c}_{T'}$, *i.e.*, an information from the ‘future’; for instance, the path to be followed defined as set of locations that varies with t . Conditioning the input features directly as in Equation 3 remains possible, by replacing \mathbf{h}_c by \mathbf{h}_{c_t} at each timestep t . However, because of the causal masking, an auto-regressive model predicts \mathbf{z}_i only from past information; therefore with the conditioning in Equation 3 only past context is available when predicting a timestep i . The model then has to predict the future path rather than use the available information, which will deteriorate output quality. This is probably why most methods are limited to a sequence-wide context, *i.e.*, a single action label for [27] or a text prompt represented with CLIP features [57]. We now detail our proposed solution to circumvent this issue.

Future conditioning. To implement future conditioning, we propose to use a network with two branches to compute two stacks of features – a causal one responsible for the prediction of the next timestep and a non-causal one that can propagate information about the conditioning at all timesteps – and inject the non-causal one into the causal one. More precisely, given an input token sequence $(\mathbf{z}_1, \dots, \mathbf{z}_{T'})$ and a conditioning sequence $(\mathbf{c}_1, \dots, \mathbf{c}_{T'})$, both are embedded into feature sequences \mathbf{h}^0 and \mathbf{g}^0 , respectively.

As in standard auto-regressive models, a stack of L causal layers $\mathbf{f}_c^1, \dots, \mathbf{f}_c^L$ is used to compute features $\mathbf{h}^1, \dots, \mathbf{h}^L$ such that for any l and any t , \mathbf{h}_t^l is a function of $\mathbf{z}_1, \dots, \mathbf{z}_{t-1}$ only. In addition, a second stack of non-causal layers $\mathbf{f}^1, \dots, \mathbf{f}^L$ is used to process \mathbf{g}^0 , and for any $1 \leq t \leq T'$:

$$\begin{cases} \mathbf{h}_t^l = \mathbf{f}_c^l(\mathbf{h}_1^{l-1}, \dots, \mathbf{h}_{t-1}^{l-1}), \\ \mathbf{g}_t^l = \mathbf{f}^l(\mathbf{g}_1^{l-1}, \dots, \mathbf{g}_t^{l-1}, \dots, \mathbf{g}_{T'}^{l-1}), \\ \tilde{\mathbf{h}}_t^l = \mathbf{h}_t^l + \mathbf{g}_t^l. \end{cases} \quad (4)$$

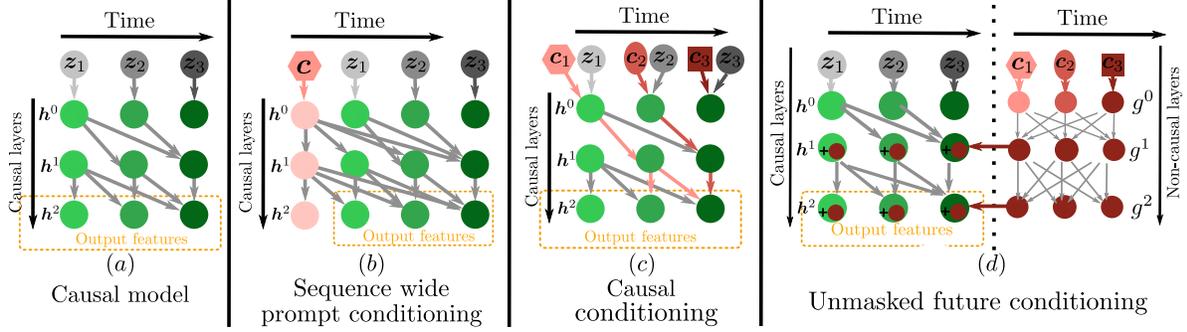


Figure 3. **Ways of conditioning an auto-regressive model.** (a): an auto-regressive model without conditioning is based on causal attention. (b): by adding a prompt token c_0 to the sequence, sequence-wide conditioning can be added. (c): for time-dependent conditioning $c_1, \dots, c_{T'}$, features could be combined but the model will be unaware of the future conditioning when predicting a given timestep. (d): we include future conditioning by making a non-causal network to process the time-varying conditioning, and combine their future with the standard causal generative model.

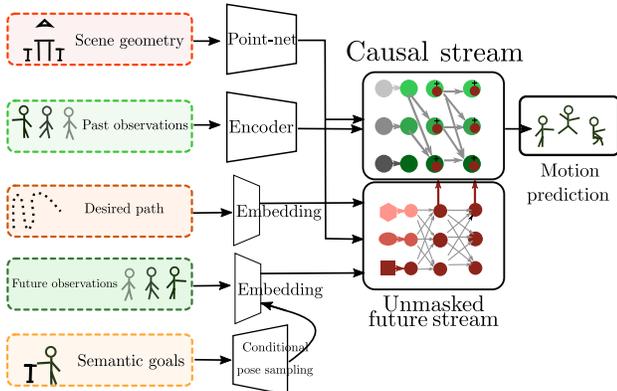


Figure 4. **Different conditionings used in Purposer:** high-level view of how different motion contexts are implemented.

With this construction, h_t^L is a function of $\mathbf{z}_1, \dots, \mathbf{z}_{t-1}$ only and can be used to predict \mathbf{z}_t , see Figure 3 for an illustration. Additionally, any feature h_t^l is a function of all conditioning signals $(c_1, \dots, c_{T'})$, and increasingly complex conditioning features g^1, \dots, g^L can be learned. This construction circumvents the causal masking, and the standard auto-regressive architecture does not need other modifications.

3.3. Motion generation with context

We now describe the different types of context we consider to put human motion generation in context, namely scene geometry, past observations or future target trajectories, and semantic information such as action labels or target human-object interaction.

Conditioning on scene geometry. To condition motions on a scene, we represent the geometry of the scene, given as an input point cloud, with a constant feature embedding \mathbf{h}_c and condition with prompting. A scene point cloud $\mathcal{S} = \{\mathbf{s}_i | i = 1, \dots, N_s\}$ is embedded using PointNet [37] following [52]. The output is then projected with a learn-

able linear layer \mathbf{W}_s : $\mathbf{c}_s = \mathbf{W}_s \cdot \text{PointNet}(\mathcal{S})$, with \mathbf{c}_s the vector containing the scene information to condition p_G with prompting.

Conditioning on past observations. An auto-regressive architecture naturally allows conditioning on observed past motion, as long as the latent sequence is produced by a causal encoder. More precisely, if a past motion $\mathbf{p}_1, \dots, \mathbf{p}_t$ of arbitrary length is observed as context, it can be encoded by the encoder into $E(\mathbf{p}_1, \dots, \mathbf{p}_t) = \mathbf{z}_1, \dots, \mathbf{z}_{t'}$. A future human motion of length T can be sampled from our model conditioned on the observation:

$$p_G(\mathbf{z} | \mathbf{z}_1, \dots, \mathbf{z}_{t'}) = \prod_{l=t'+1}^{T'+t'} p(\mathbf{z}_l | \mathbf{z}_1, \dots, \mathbf{z}_{l-1}). \quad (5)$$

Thus by design, an auto-regressive model has the flexibility to be conditioned on past motion without change and/or retraining.

Conditioning on trajectories and target poses. In addition to past observations, our model can also be conditioned on future target poses, or trajectories. For this type of conditioning, the *causal conditioning* approach taken by PoseGPT does not allow the future path or targets to be observed when generating a given timestep. However, by using the *future conditioning* from Section 3.2, we can condition on trajectories or arbitrarily chosen future poses. Note that while this would allow us to condition the model on a future pose in arbitrary time steps in the future, for our purposes we chose the last pose in the sequence. Let $((x_1, y_1), \dots, (x_{T'}, y_{T'}))$ be a 2D trajectory on a bird’s eye view of the scene, \mathbf{p}_{t_j} some future poses, and \mathbf{W}_p and \mathbf{W}_f two learnable linear layers. We can write the path and future poses conditions as:

$$\mathbf{c}_p = (\mathbf{W}_p \cdot (x_1, y_1), \dots, \mathbf{W}_p \cdot (x_{T'}, y_{T'})) \quad (6)$$

$$\mathbf{c}_f = (\mathbf{W}_f \cdot \mathbf{p}_t)_{t \in \mathcal{T}} \quad (7)$$

where \mathcal{T} is a set of time steps; c_p and c_f are then used as input to compute h^0 in Equation 4. Note that conditioning on the final target pose is a special case of this, which could also be achieved with a simpler conditioning such as an additional token.

Conditioning on semantic information. We consider two semantic contexts: (a) *action labels* and (b) a *target final human-object interaction*.

- *Action labels:* We use sequence-wide conditioning to include action information into the model.

Action labels are embedded and projected into c_a and p_G is conditioned on the result by inputting c_a to the unmasked stream.

- *Target human-object interactions:* We also integrate the possibility to control the interaction of the motion with specific objects by conditioning our motion samples on *pairs* of actions and objects, $\{(a_t, o_t)\}$, *e.g.* (lie, couch). To achieve that, we decouple the problem into (a) generating a *static* pose $\mathbf{p}_{(a,o)}$ conditionally on (a, o) , following the pioneering approach of [64], and (b) conditioning the motion model on this target pose, which is embedded into $\mathbf{c}_{(a,o)} = \mathbf{W}_{ao} \cdot \mathbf{p}_{(a,o)}$, where \mathbf{W}_{ao} is a learnable linear layer. Then, a latent sequence is sampled using this pose as target, and the latent sequence is decoded into a human motion with the decoder $D(\cdot)$. At train time, we condition the model on final target poses extracted from the data rather than generated.

Finally, based on the different types of conditioning we have introduced above we can re-write Equation (1) with the following conditioning:

$$p_G(\mathbf{z}_i | \mathbf{c}_i) = \prod_{t=1}^{T'} p(\mathbf{z}_t | \{\mathbf{z}_i\}_{i < t}; \mathbf{c}_s, \mathbf{c}_a, \mathbf{c}_{(a,o),t}, \mathbf{c}_p, \mathbf{c}_f). \quad (8)$$

In Figure 4, we recap how the different contextual information is embedded. To generate motion, latent variables are iteratively sampled and added to the conditioning sequence before being decoded.

3.4. Training setup

We directly use the discrete auto-encoder from [27] trained on BABEL [36]. To train the auto-regressive model we use the scene-conditioned data from HUMANISE, a synthetic dataset composed of a subset of BABEL motions placed in ScanNet [11] scenes. It contains 19.6K human motion sequences in 643 3D scenes and consists of actions that are commonly performed when interacting with a scene. We train each of our generative networks with the relevant inputs as conditioning. Finally, to evaluate on real-world but smaller-scale data, we fine-tune the model on PROX [16], which consists of 100K frames with pseudo ground truth. It captures dynamic sequences of 20 subjects in 12 scenes

at 30fps. Implementation details are included in the supplementary material.

3.5. Generating long-term sequences

While our model is trained on short sequences, taking advantage of its autoregressive nature, and using various sets of conditioning allows us to generate long-term motions that are coherent with a virtual scene by chaining short-term motions. Specifically, we can define an ‘object interaction’ configuration to generate motions that interact with scene objects, *e.g.* sit/lie down, by giving the proper conditioning as input. We can then use a ‘locomotion’ configuration to generate walking motion sequences to navigate the scene and connect between different interaction motions that are far apart in the scene, if needed. Each of these configurations is conditioned with relevant information. The ‘object interaction’ model is mainly conditioned with a target pose that encodes a correct interaction to guide the motion to reach that pose at the end of the sequence. Other than the target pose, this model is configured with the rest of the relevant conditioning: action labels, scene geometry and past observations. The ‘locomotion’ model is mainly guided with a path, *e.g.* the shortest path between the two locations where the two motion interactions happen in the scene, *e.g.* from the A* algorithm. Although, this path could be set in any other different manner if needed. The ‘locomotion’ configuration does not need the use of a target pose but instead only the desired (x, y) position is specified. For chaining any two consecutive sequences together we take the last n poses from the first sequence and use those to condition the generation of the next sequence. Concretely, we use $n = 2$ for our results as the downsample factor from the pose space \mathbf{p} to the latent space \mathbf{z} is 2. Thus, the minimum number of conditioning poses is $n = 2$ (Section 3.1). For more detail about the exact conditioning both of these configurations please refer to Table 1 in Section 4.2.

4. Experiments

Given that without any conditioning our method boils down to a standard discrete auto-regressive model such as [27], we focus our experiments on different conditioning scenarios. We aim to apply our model to populate virtual scenes. We thus focus our experiments on scene-conditioned motion generation, with various types of context on top. After introducing the datasets and metrics in Section 4.1, we present several ablations in Section 4.2. Finally, we compare our method to the state of the art in Section 4.3.

Conditionings. We clearly state in each table the conditioning that was used. In addition to scene conditioning either at the feature-level (F) or with prompting, *i.e.*, an extra token (T), we consider other conditioning forms: first pose (first), target pose (target/P) or target position (target/XY),

action label and path. Action labels conditioning is used in all cases unless noted otherwise.

4.1. Datasets and metrics

Datasets. We perform most experiments on the HUMANISE dataset, following the official splits [54] with 16.5K motions in 543 scenes for training and 3.1K motions in 100 scenes for testing. We also experiment on the PROX dataset with the standard splits (8 training scenes and 4 scenes for testing) as in [52, 53], and rely on the improved fittings provided in [59] and action labels from [64].

Physical plausibility. We evaluate the physical plausibility of generated interactions using the non-collision metric proposed in [62], which measures how much the generated human mesh interpenetrates the mesh of the scene, and the contact score proposed in [52] which is complementary as it ensures that the motion actually makes contact with the object – non-collision alone would be maximized by standing away from everything. For the contact scores, we follow [54] and use a threshold of 0.02, except for some tables that are clearly specified where we follow [52, 53] use 0.01.

Diversity metrics. To evaluate the diversity of generated samples, we follow common practices in the literature [21, 52, 54] and report the average pairwise distance (APD) metric. This metric measures the average L_2 distance between all pairs of motion of K samples computed with exactly the same input information. When evaluating on HUMANISE, we follow the practice of [54]. When comparing on PROX, we follow [21]. In both cases $K = 20$. Additionally, following [21], we measure APD for a specific set of 61 markers (APD mark.) extracted from the body meshes. As advocated for in [27], and given that all components of our model are likelihood based, we also report likelihood-based metrics for the generator G.

Quality metrics. To measure the quality of sampled motion, we compute the Fréchet distance score [20]. For comparison with existing work, we compute the FD metric with a VPoser [33] model. We denote this metric by FD_{static} as this model only takes individual frames as input. Please also find qualitative video results in the supplementary for a complementary perspective.

4.2. Ablation of design choices

In Table 1, we first ablate the impact of the different conditioning information used by our model: future stream, first and target poses, scene point-cloud, and path. Note that ‘future stream’ is not a type of information but a novel component of our model that conditions the information in a special way. The first two rows correspond to our baseline, namely, PoseGPT [27].

First pose. In Table 1, rows 1 and 2, we observe that using the first observed pose as conditioning significantly improves the non-collision and contact metrics. This is ex-

stream	conditioning				NLL↓	APD↑ mark.	phys. plausibility↑	
	first	scene	target	path			non-coll.	contact
X	X	X	X	X	0.86	4.83	55.73	93.93
X	✓	X	X	X	0.86	4.09	69.56	92.68
X	✓	F	X	X	0.96	4.08	69.10	93.10
X	✓	T	X	X	0.87	3.91	70.19	92.79
X	✓	T	P	X	0.62	3.05	71.64	91.86
✓	✓	T	XY	✓	0.70	5.91	71.24	92.59
✓	✓	T	P	✓	0.48	3.02	71.76	94.15
✓	✓	T	P	X	0.42	3.13	73.28	94.29

Table 1. **Ablation study** on HUMANISE [54] with action label conditioning and without post-processing optimization. XY means that it uses target position instead of target pose (P).

pected as it guides the motion in a correct direction where less collisions are likely to occur, whether the model is conditioned on the scene or not. Thus adding a first pose, which can be obtained from past observed motion, improves scene interaction and motion control while retaining generation diversity.

Scene. Using scene information helps to improve the quality of the generated motion by forcing it to better fit the given scene. However, this depends on how we input this information to our model. As seen in Table 1 row 3, if done at the feature level (F), *i.e.*, by concatenation with the input embeddings, the model’s performance does not improve, or even deteriorates, both in terms of next index prediction negative likelihood (NLL) and non-collision. On the other hand, if we introduce this information with token prompting (T) (row 4), we both maintain the generation quality (NLL) and improve penetration (69.56% vs. 70.19%).

Target pose. To guide motion towards human-scene interaction, we design our model to take as input a target conditioning pose. These target poses can be taken from the ground-truth data or sampled at test time, given action object pairs. As observed in Table 1 (row 5), using a target pose gives an important increase of 1.46% in the non-collision score while also improving the generation quality (NLL). Contact score is slightly reduced, possibly as a consequence of reducing scene penetration. Furthermore, APD decreased but this is expected, as the generated motion is more constrained and therefore there is less variability in the outputs.

Future Stream. The future conditioning block proposed in Section 3.3 allows conditioning on time-step dependent future signal and is a key component of our method. The last three rows of Table 1 show that variants using this component outperforms non-‘future stream’ counterparts. We hypothesize that this is because (a) it enables the use of all available path information and (b) it is a more flexible way to condition on the rest of features, *e.g.*, target pose and scene. Using the time-dependent path as input (row 6) yields an improvement in non-collision score (0.70%), maintains contact score, and substantially improves NLL.

We refer to this setting as the *locomotion model*. The

optim	NLL↓	APD↑ (all)	FD _{static} ↓	phys. plausibility↑ non-coll.	contact
×	2.69	2.06	30.11	95.23	99.98
✓	2.69	3.02	29.76	99.24	99.96

Table 2. **Impact of the optimization step** on PROX, performed with first/last pose cond., scene prompt, but without action and path information to match the conditions of [52]. Contact score threshold of 0.01.

Name	conditioning act path	APD↑ (all)	FD _{static} ↓	phys. plausibility↑ non-coll.	contact
Wang <i>et al.</i> [52] Purposer	×	×	0.00	—	99.35
	×	×	3.02	29.76	99.24
Wang <i>et al.</i> [53] Purposer	✓	✓	2.78	111.65	99.35
	✓	✓	2.58	29.84	99.89

Table 3. **Comparison on PROX** with Wang *et al.* [52] and Wang *et al.* [53]. Contact score threshold is 0.01. Results use first and last pose conditioning to match the compared SOTA. Results are refined by an optimization step.

last row presents the performance when conditioned on a target pose, but no path. It obtains the best performance in motion quality, non-collision, and contact, and is referred to as the *object interaction* model. Note that all configurations in Table 1 can generate any action present in the dataset; thus, we evaluate them all together.

Impact of post-processing optimization. Previous methods [52, 53] leverage a post-processing step that optimizes the generated motions to avoid collision and favor contacts. We measure the impact of such post-processing optimization on the PROX dataset and later compare to other SOTA methods that apply the same post-processing, see Table 2. This optimization step slightly improves the physical plausibility scores, but produces less natural and stiffer motion. This is best appreciated in the qualitative results.

Qualitative results and long-term generation. In Figure 5, we present qualitative results that illustrate the impact of incorporating target pose conditioning (top row) to enable scene interaction. Our model is able to interact with the same object when initialized with random initial body locations and orientations. Additionally, we present the effect of the path conditioning (bottom row) for walking actions. Our model is capable of generating realistic locomotion in random directions and with arbitrarily chosen paths. In the supplementary video, we also show examples of long-term motions combining short-term motions for object interaction and for locomotion, see Section 3.5.

4.3. Comparison to the state of the art

We compare our method to the state of the art [52, 53] on PROX in Table 3. Different sets of metrics and conditionings have been reported in the literature. To make a fair comparison, we match the conditioning with each method

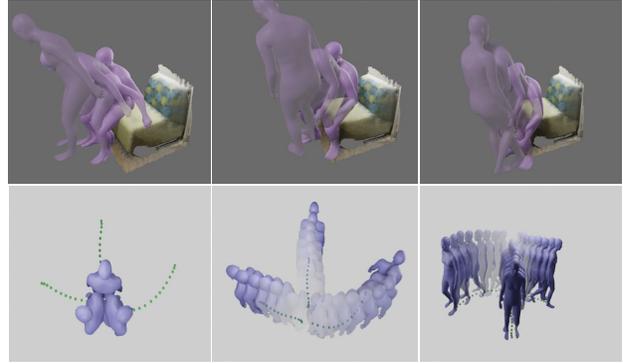


Figure 5. **Effect of target pose and path conditioning.** **Upper row:** examples of object interaction. Here we use the same object with different and random initial body position and orientation. **Lower row:** demonstration of the effects of path conditioning: we can define the final position and trajectory given a common starting point. The green dots represent the conditioning path.

being compared.

We observe that our model provides a good trade-off between diversity (APD), quality (FD), and physical plausibility. In particular, non-collision scores for Purposer does not vary substantially from the competing approaches while our model consistently has the highest contact score, which indicates a rich interaction with the scene. Simultaneously, a low FD score is achieved, which is a measure of realistic generations. To provide a more comprehensive understanding of how our method compares to these baselines, we have included qualitative video results in the supplementary material.

4.4. Limitations

Since our method is purely kinematic, *i.e.*, it does not take into account physics constraints, interpenetrations with the scene may occur in some cases. Furthermore, HUMANISE is a synthetic dataset and provides short motion clips that may not comprehensively capture all the subtleties associated with interacting with objects. We leave the use of more realistic datasets such as SAMP [17] for future work.

5. Conclusion

In this work, we introduce Purposer, a novel approach grounded in neural discrete representation for generating human motion within 3D virtual scenes. Purposer can generate realistic motions while also capturing human-object interactions. This is an important step forward due to its potential applications in indoor activity simulation and synthetic data creation, among others. Experiments show that Purposer consistently improves upon the baselines. Additionally, our model can be controlled semantically, generalizes to a variety of new scenes, and generates long-term motions even if only very short sequences are present in the training data.

Acknowledgements: This work is partly supported by the project MoHuCo PID2020-120049RB-I00 funded by MCIN/AEI/10.13039/501100011033.

References

- [1] Hyemin Ahn, Timothy Ha, Yunho Choi, Hwiyeon Yoo, and Songhwa Oh. Text2action: Generative adversarial synthesis from language to action. In *ICRA*, 2018.
- [2] Emre Aksan, Manuel Kaufmann, and Otmar Hilliges. Structured prediction helps 3d human motion modelling. In *ICCV*, 2019.
- [3] Nikos Athanasiou, Mathis Petrovich, Michael J. Black, and Güll Varol. TEACH: Temporal Action Compositions for 3D Humans. In *3DV*, 2022.
- [4] Norman Badler. Temporal scene analysis: Conceptual descriptions of object movements. In *PhD thesis, University of Toronto*, 1975.
- [5] Norman I. Badler, Cary B. Phillips, and Bonnie Lynn Webber. Simulating humans: Computer graphics animation and control. In *Oxford University Press*, 1993.
- [6] Shane Barratt and Rishi Sharma. A note on the inception score. *arXiv preprint arXiv:1801.01973*, 2018.
- [7] Emad Barsoum, John Kender, and Zicheng Liu. Hp-gan: Probabilistic 3d human motion prediction via gan. In *CVPRW*, 2018.
- [8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 2020.
- [9] Sammy Christen, Muhammed Kocabas, Emre Aksan, Jemin Hwangbo, Jie Song, and Otmar Hilliges. D-grasp: Physically plausible dynamic grasp synthesis for hand-object interactions. In *CVPR*, 2022.
- [10] Enric Corona, Albert Pumarola, Guillem Alenyà, and Francesc Moreno-Noguer. Context-aware human motion prediction. In *CVPR*, 2020.
- [11] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017.
- [12] Anindita Ghosh, Noshaba Cheema, Cennet Oguz, Christian Theobalt, and Philipp Slusallek. Synthesis of compositional animations from textual descriptions. In *CVPR*, 2021.
- [13] Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun, and Li Cheng. Action2motion: Conditioned generation of 3d human motions. In *ACMMM*, 2020.
- [14] Chuan Guo, Xinxin Zuo, Sen Wang, and Li Cheng. Tm2t: Stochastic and tokenized modeling for the reciprocal generation of 3d human motions and texts. In *ECCV*, 2022.
- [15] Ikhsanul Habibie, Daniel Holden, Jonathan Schwarz, Joe Yearsley, and Taku Komura. A recurrent variational autoencoder for human motion synthesis. In *BMVC*, 2017.
- [16] Mohamed Hassan, Vasileios Choutas, Dimitrios Tzionas, and Michael J. Black. Resolving 3D human pose ambiguities with 3D scene constraints. In *ICCV*, 2019.
- [17] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael Black. Stochastic scene-aware motion prediction. In *ICCV*, 2021.
- [18] Mohamed Hassan, Partha Ghosh, Joachim Tesch, Dimitrios Tzionas, and Michael J. Black. Populating 3D scenes by learning human-scene interaction. In *CVPR*, 2021.
- [19] Mohamed Hassan, Yunrong Guo, Tingwu Wang, Michael Black, Sanja Fidler, and Xue Bin Peng. Synthesizing physical character-scene interactions. In *SIGGRAPH*, 2023.
- [20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a nash equilibrium. In *NeurIPS*, 2017.
- [21] Siyuan Huang, Zan Wang, Puhao Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song-Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes. *CVPR*, 2023.
- [22] Hsin-Ying Lee, Xiaodong Yang, Ming-Yu Liu, Ting-Chun Wang, Yu-Ding Lu, Ming-Hsuan Yang, and Jan Kautz. Dancing to music. *NeurIPS*, 2019.
- [23] Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. Learn to dance with aist++: Music conditioned 3D dance generation. *arXiv preprint arXiv:2101.08779*, 2021.
- [24] Angela S. Lin, Lemeng Wu, Rodolfo Corona, Kevin Tai, Qixing Huang, and Raymond J. Mooney. Generating animated videos of human activities from natural language descriptions. In *Visually Grounded Interaction and Language Workshop at NeurIPS*, 2018.
- [25] Xiao Lin and Mohamed R Amer. Human motion modeling using dv-gans. *arXiv preprint arXiv:1804.10652*, 2018.
- [26] Thomas Lucas, Konstantin Shmelkov, Karteek Alahari, Cordelia Schmid, and Jakob Verbeek. Adaptive density estimation for generative models. In *NeurIPS*, 2019.
- [27] Thomas Lucas, Fabien Baradel, Philippe Weinzaepfel, and Grégory Rogez. PoseGPT: Quantization-based 3D Human Motion Generation and Forecasting. In *ECCV*, 2022.
- [28] Zhengyi Luo, Shun Iwase, Ye Yuan, and Kris Kitani. Embodied scene-aware human pose estimation. In *NeurIPS*, 2022.
- [29] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. Amass: Archive of motion capture as surface shapes. In *ICCV*, 2019.
- [30] Wei Mao, Miaomiao Liu, and Mathieu Salzmann. Weakly-supervised action transition learning for stochastic human motion prediction. In *CVPR*, 2022.
- [31] Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models. In *ICML*, 2020.
- [32] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [33] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3d hands, face, and body from a single image. In *CVPR*, 2019.

- [34] Mathis Petrovich, Michael J Black, and Gül Varol. Action-conditioned 3d human motion synthesis with transformer vae. In *ICCV*, 2021.
- [35] Mathis Petrovich, Michael J Black, and Gül Varol. TEMOS: Generating diverse human motions from textual descriptions. In *ECCV*, 2022.
- [36] Abhinanda R. Punnakal, Arjun Chandrasekaran, Nikos Athanasiou, Alejandra Quiros-Ramirez, and Michael J. Black. BABEL: Bodies, action and behavior with english labels. In *CVPR*, 2021.
- [37] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *CVPR*, 2017.
- [38] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J. Guibas. Humor: 3d human motion model for robust pose estimation. *ICCV*, 2021.
- [39] Davis Rempe, Zhengyi Luo, Xue Bin Peng, Ye Yuan, Kris Kitani, Karsten Kreis, Sanja Fidler, and Or Litany. Trace and pace: Controllable pedestrian animation via guided trajectory diffusion. In *CVPR*, 2023.
- [40] Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H. Bermano. Human motion diffusion as a generative prior, 2023.
- [41] Soshi Shimada, Vladislav Golyanik, Zhi Li, Patrick Pérez, Weipeng Xu, and Christian Theobalt. Hulc: 3d human motion capture with pose manifold sampling and dense contact guidance. In *ECCV*, 2022.
- [42] Konstantin Shmelkov, Cordelia Schmid, and Karteek Alahari. How good is my gan? In *ECCV*, 2018.
- [43] Li Siyao, Weijiang Yu, Tianpei Gu, Chunze Lin, Quan Wang, Chen Qian, Chen Change Loy, and Ziwei Liu. Bailando: 3d dance generation by actor-critic gpt with choreographic memory. In *CVPR*, 2022.
- [44] Sebastian Starke, He Zhang, Taku Komura, and Ju Saito. Neural state machine for character-scene interactions. In *ACM ToG*, 2019.
- [45] Omid Taheri, Vasileios Choutas, Michael J. Black, and Dimitrios Tzionas. GOAL: Generating 4D whole-body motion for hand-object grasping. In *CVPR*, 2022.
- [46] Graham W Taylor, Geoffrey E Hinton, and Sam Roweis. Modeling human motion using binary latent variables. *NeurIPS*, 2006.
- [47] Guy Tevet, Brian Gordon, Amir Hertz, Amit H. Bermano, and Daniel Cohen-Or. Motionclip: Exposing human motion generation to clip space. *arXiv preprint arXiv:2203.08063*, 2022.
- [48] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *ICLR*, 2023.
- [49] Nicolas Ugrinovic, Adria Ruiz, Antonio Agudo, Alberto Sanfeliu, and Francesc Moreno-Noguer. Body size and depth disambiguation in multi-person reconstruction from single images. In *2021 International Conference on 3D Vision (3DV)*, pages 53–63. IEEE, 2021.
- [50] Raquel Urtasun, David J Fleet, and Neil D Lawrence. Modeling human locomotion with topologically constrained latent variable models. In *Workshop on Human Motion*, 2007.
- [51] Aaron van den Oord, Vinyals Oriol, and Koray Kavukcuoglu. Neural discrete representation learning. In *ICML*, 2018.
- [52] Jiashun Wang, Huazhe Xu, Jingwei Xu, Sifei Liu, and Xiaolong Wang. Synthesizing long-term 3d human motion and interaction in 3d scenes. In *CVPR*, 2021.
- [53] Jingbo Wang, Yu Rong, Jingyuan Liu, Sijie Yan, Dahua Lin, and Bo Dai. Towards diverse and natural scene-aware 3d human motion synthesis. In *CVPR*, 2022.
- [54] Zan Wang, Yixin Chen, Tengyu Liu, Yixin Zhu, Wei Liang, and Siyuan Huang. HUMANISE: Language-conditioned human motion generation in 3d scenes. In *NeurIPS*, 2022.
- [55] Yan Wu, Jiahao Wang, Yan Zhang, Siwei Zhang, Otmar Hilliges, Fisher Yu, and Siyu Tang. Saga: Stochastic whole-body grasping with contact. In *ECCV*, 2022.
- [56] Ye Yuan and Kris Kitani. Dlow: Diversifying latent flows for diverse human motion prediction. In *ECCV*, 2020.
- [57] Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Shaoli Huang, Yong Zhang, Hongwei Zhao, Hongtao Lu, and Xi Shen. T2m-gpt: Generating human motion from textual descriptions with discrete representations. *arXiv preprint arXiv:2301.06052*, 2023.
- [58] Siwei Zhang, Yan Zhang, Qianli Ma, Michael J. Black, and Siyu Tang. PLACE: Proximity learning of articulation and contact in 3D environments. In *3DV*, 2020.
- [59] Siwei Zhang, Yan Zhang, Federica Bogo, Pollefeys Marc, and Siyu Tang. Learning motion priors for 4d human body capture in 3d scenes. In *ICCV*, 2021.
- [60] Xiaohan Zhang, Bharat Lal Bhatnagar, Sebastian Starke, Vladimir Guzov, and Gerard Pons-Moll. Couch: Towards controllable human-chair interactions. In *ECCV*, 2022.
- [61] Yan Zhang and Siyu Tang. The wanderings of odysseus in 3d scenes. In *CVPR*, 2022.
- [62] Yan Zhang, Mohamed Hassan, Heiko Neumann, Michael J. Black, and Siyu Tang. Generating 3d people in scenes without people. In *CVPR*, 2020.
- [63] Yan Zhang, Michael J Black, and Siyu Tang. We are more than our joints: Predicting how 3D bodies move. In *CVPR*, 2021.
- [64] Kaifeng Zhao, Shaofei Wang, Yan Zhang, Thabo Beeler, , and Siyu Tang. Compositional human-scene interaction synthesis with semantic control. In *ECCV*, 2022.
- [65] Kaifeng Zhao, Yan Zhang, Shaofei Wang, Thabo Beeler, , and Siyu Tang. Synthesizing diverse human motions in 3d indoor scenes. In *ICCV*, 2023.