

UniMotion: Unifying 3D Human Motion Synthesis and Understanding

Chuqiao Li¹ Julian Chibane^{1,2} Yannan He¹ Naama Pearl¹
Andreas Geiger¹ Gerard Pons-Moll^{1,2}

¹Tübingen AI Center, University of Tübingen, Germany, ²Max Planck Institute for Informatics, Saarland Informatics Campus, Germany
{chuqiao.li, yannan.he, naama.pearl, a.geiger, gerard.pons-moll}@uni-tuebingen.de,
jchibane@mpi-inf.mpg.de

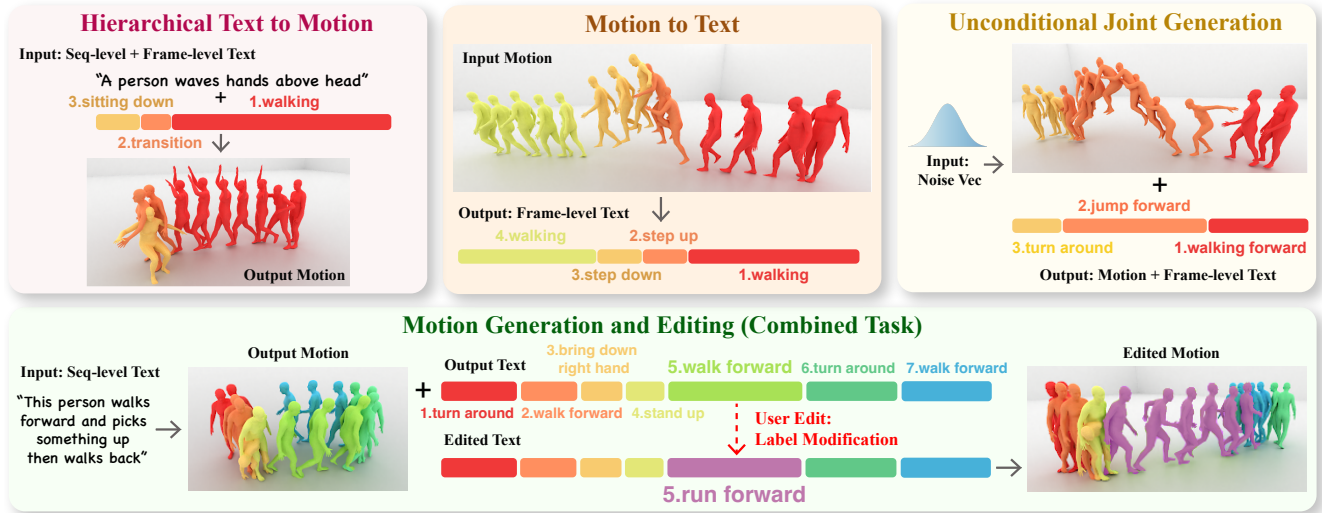


Figure 1. **Universality of UniMotion.** Our model can generate motion from compositional sequence- and frame-level text (**Hierarchical Text to Motion**), generate detailed per frame motion descriptions (**Motion to Text**), generate motion and accurate frame-level text descriptions from noise (**Unconditional Joint Generation**), amongst other use cases outlined in our experiments section. Tasks can be combined for a controllable generation: users can generate motion from a coarse sentence, our model additionally generates detailed text descriptions, which can be edited and used for regeneration, generating the desired edited motion (**Motion Generation and Editing**).

Abstract

We introduce *UniMotion*, the first unified multi-task human motion model capable of both flexible motion control and frame-level motion understanding. While existing works control avatar motion with global text conditioning, or with fine-grained per frame scripts, none can do both at once. In addition, none of the existing works can output frame-level text paired with the generated poses. In contrast, *UniMotion* allows to control motion with global text, or local frame-level text, or both at once, providing more flexible control for users. Importantly, *UniMotion* is the first model which by design outputs local text paired with the generated poses, allowing users to know what motion happens and when, which is necessary for a wide range of applications. We show *UniMotion* opens up new applications: 1.) hierarchical control, allowing users to specify motion at different levels of detail, 2.) obtaining motion text descriptions for existing MoCap data or youtube videos 3.) allowing for editability, generating motion from text and editing the motion via text edits. Moreover, *UniMotion* attains state-of-the-art results for the frame-level text-

to-motion task on the established HumanML3D dataset. The pre-trained model and code are available on our project page at <https://coral79.github.io/uni-motion/>.

1. Introduction

Human motion synthesis is important for gaming, robotics and AR/VR applications. In real-world scenarios, avatars need to be controlled at multiple levels of abstraction. Effective controllability requires that an avatar be capable of executing detailed local sub-tasks according to a timeline while simultaneously understanding the overall global objective. In addition, the synthesis model should be aware of what action happens and when – an essential feature of biological intelligence to react to the external world.

However, current motion synthesis methods focus on either global per sequence-level control, or local per frame-level control, but don't allow for both. This results in single-level conditioning, thereby **lacking hierarchical control**. Importantly, these models also lack fine-grained motion awareness, specifically, the ability to output motion descrip-

tions for each pose in the generated output motion sequence. The frame-level text-to-motion methods [1, 3, 30] provide detailed manipulation of individual frames. However, it can be impractical to specify the exact duration of each action in some situations, and ensuring overall semantic plausibility throughout the entire sequence remains challenging for these models. Conversely, the sequence-level text-to-motion methods [17, 34, 47] focus on achieving natural overall motion but struggles with fine-grained control. Furthermore, current models lack semantic awareness of the synthesized motion – there is no understanding of what action occurs when. Thus, they are **lacking motion understanding**, which is crucial for reacting to the external world and allows for action-specific editing in animation applications. While some works have made progress in this direction [17, 47] by predicting sequence-level text descriptions from motion, they fail to provide fine-grained frame-level text. Overall, despite their potential synergies, motion understanding and synthesis have been treated in isolation in the literature.

In this paper, we introduce UniMotion, the first unified multi-task model capable of both flexible motion control and frame-level motion understanding. UniMotion takes as input, global sequence level or local frame-level text inputs or human motion sequences, or any subsets thereof, or no input in case of unconditional generation. The output of our model is either fine grained, per pose text descriptions, or human motion sequences. This flexibility, allows us to train our model from different data sources. Moreover, by design, we unify tasks that are usually treated in separation by prior works, such as *Frame-Level Text-to-Motion*, *Sequence-Level Text-to-Motion* and *Motion-to-Text*, into a single simple unified model, trained a single time. Importantly, UniMotion’s flexibility also allows for novel tasks not previously considered by prior work like 1.) unconditional generation of human motion with corresponding frame-level text descriptions and 2.) generation of frame-level text from motion, providing granular, time-aware annotations (see Fig. 1 for an illustration our diverse tasks).

To accomplish this, our model utilizes a transformer architecture with temporal alignment between the motion and frame-level text. We further enhance this by diffusing the local text together with the poses, using different diffusion time variables for each, inspired by the approach in Uni-diffuser [2]. Specifically, the local text is tokenized and frame-wise aligned with the 3D poses, while the global text is injected as a global token. This design allows UniMotion to dynamically switch between global, local, or combined conditioning signals at test time, providing flexibility in motion generation and understanding. During training, we sample from all possible distributions (global and/or local conditioning, or unconditional), alternating between providing noise and signal to the model for each modality.

This method effectively teaches the model both unconditional and conditional distributions, equipping it with the ability to handle various inputs.

Real-world applicability. We demonstrate practical utility across various real-world scenarios:

2D Video Annotation: We annotate human motion extracted from YouTube videos with frame-level text, by feeding UniMotion with human pose estimation (HPE) results. This annotation can serve as close captions for the visually impaired. *4D Mocap Annotation:* We annotate human motion captures, e.g. obtained from IMUs, with frame-level text. This provides automated insights and descriptions into the captured motions, e.g. allowing for text search retrieval of motion sequences. *Hierarchical Control:* We provide examples of generating motion sequences with two levels of abstraction, specifying a general motion for arms via global text, and a fine-grained motion sequence for the rest of the body via local-level text. *Motion editing:* We show that UniMotion can be used for content creation, where controllability of the motion is important. Given a global text description, a user can generate an initial motion including a local-level text description. The user can then edit the motion as desired, editing the text segments and regenerating the motion. In summary, our **key contributions** are:

- **Unified Synthesis and Understanding:** We introduce UniMotion, the first unified probabilistic motion model allowing for sampling from the joint and all possible conditionals. It unifies tasks that are usually treated in separation by prior works, while also allowing for novel tasks not previously considered.
- **Results and Applications:** We show applicability to 2D Video Annotation, 4D Mocap Annotation, Hierarchical Control and Motion editing. Moreover UniMotion attains state-of-the-art results for the frame-level text-to-motion task on the established HumanML3D dataset. Code and models will be released upon acceptance.

2. Related Work

Conditional human motion synthesis. Synthesizing human motion has been a long-standing challenge. Recent studies in motion generation have shown notable progress in synthesizing movements conditioned on diverse modalities such as text [25, 26, 30, 31, 33, 34], music [20, 21], scenes [24, 35], and interactive objects [13, 19, 32, 37, 39, 40, 45, 46]. Recent years have witnessed substantial advancements in text-driven motion generation [8, 11, 12, 25, 33, 41]. Notably, diffusion-based generative models have emerged as potent tools, exhibiting impressive performance on leading benchmarks for text-to-motion tasks. Pioneering efforts such as MotionDiffuse [42], MDM [34], and FLAME [18] represent early applications of diffusion models to text-driven motion generation. Building upon

this foundation, MLD [6] further harnesses latent diffusion models, while ReMoDiffuse [43] integrates retrieval techniques into the motion generation pipeline. Recent MotionLCM [8] accelerates the sampling speed by adopting consistency model in motion latent space. Noteworthy, OmniControl [38] specializes in fine-grained spatial control of body joints.

Text-to-motion generation models. The current landscape of text-to-motion generation models can be categorized into two main streams of controllability: (a) global text-based control and (b) Fine-grained local text-based control. Among the former, MotionGPT [17] utilizes pre-trained language models and motion-specific vector quantized models to conceptualize human motion as a language. Similarly, AvatarGPT [47] proposes a top-down approach to address end-to-end motion planning and synthesis.

Conversely, research focusing on short, specific instructions presents another avenue. PriorMDM [30] introduces a two-stage method that synthesizes short motion sequences and their padded transitions. However, due to the lack of effective supervised learning, motions generated by such methods often exhibit artifacts, such as abrupt speed changes. FineMoGen [44] proposes diffusion-based motion generation and editing for fine-grained per-body part motion control, albeit requiring detailed per-body part instructions as input. Closely aligned with our work are methods enabling temporal control of motion, where the length of each motion segment can be controlled at the frame level. FlowMDM [3] demonstrates impressive results in seamless transitions between local motion segments, while STMC [26] proposes a hybrid method for spatial and temporal motion composition of multi-stream motion using off-the-shelf pre-trained motion models [34]. Notably, these methods do not condition on global text, resulting in a lack of awareness of the global motion context and less natural motion transitions.

Our method combines the advantages of both categories. It is the first method enabling the generation of human motion conditioned both at the abstract level with global text and at the detailed level with local texts.

Human motion understanding. Understanding the meaning of human motion has been a long-standing research topic, this has been approached by describing human motion with predefined action labels [7, 48], which have dominated this field for some time. However, these methods have obvious limitations, they are not appropriate to describe complex motion sequences. Recently, the text annotated motion datasets [5, 11, 27] have enabled the methods [12, 17, 41] that learn the mutual mapping between human motion sequences and natural language descriptions. While these works produce impressive, they fall short in generating accurate per-frame language descriptions. More recently, methods such as [9, 16] have achieved motion editing

based on more fine-grained conditions, such as per body part condition. However, they still lack the capability for temporal editing. UniMotion is the first approach that not only generates per-frame language descriptions but also allows for motion generation over specified time spans, thus advancing the understanding of human motion.

3. Preliminary: Motion Diffusion Model

We provide a brief overview of the Human Motion Diffusion Model (MDM) [34], which is designed for sequence-level text-to-motion synthesis. This model serves as a building block for our UniMotion, which extends its capabilities by (a) incorporating frame-level text input and (b) enabling the joint generation of both motion and text. MDM aims to synthesize human motion sequences, denoted as $\mathbf{x}^{1:N}$, where N is the length of the sequence. The synthesis process is guided by a sequence-level text condition c , meaning the entire motion sequence is described by a single text prompt. In cases of unconditioned motion generation, the condition is represented as $c = \emptyset$.

Diffusion is modeled as a Markov noising process, where $t = 0$ represents the timestep corresponding to the clean data and $t = T$ corresponds to the fully corrupted data. The samples generated during this process are denoted as $\{\mathbf{x}_t^{1:N}\}_{t=0}^T$, with $\mathbf{x}_0^{1:N}$ being drawn from the data distribution. The transition between steps is defined by:

$$q(\mathbf{x}_t^{1:N} | \mathbf{x}_{t-1}^{1:N}) = \mathcal{N}(\sqrt{\alpha_t} \mathbf{x}_{t-1}^{1:N}, (1 - \alpha_t) \mathbf{I}). \quad (1)$$

where $\alpha_t \in (0, 1)$ indicates the noise level, with $\alpha_i = 1 - \beta_i$ and β_i being the noise schedule. We drop the sequence length and use \mathbf{x}_t to denote the full sequence at noising step t for simplicity. The reverse diffusion process gradually denoises the noisy sequence \mathbf{x}_T , with the conditioned motion generation modeling the distribution $p(\mathbf{x}_0 | c)$. The denoised data is directly predicted using a model G , where $\hat{\mathbf{x}}_0 = G(\mathbf{x}_t, t, c)$ [29].

To adapt the diffusion model for human motion, we follow [11] to parameterize the human motion as a 263 dimensional vector. Due to its redundancy inherent in the motion representation, a simple training objective [15, 34] can be used, minimizing the expected distance between the original noisy motion \mathbf{x}_0 and the predicted motion $\hat{\mathbf{x}}_0$:

$$\mathcal{L}_{\text{simple}} = E_{\mathbf{x}_0 \sim q(\mathbf{x} | c), t \sim \mathcal{U}\{1, \dots, T\}} \|\mathbf{x}_0 - G(\mathbf{x}_t, t, c)\|_2^2. \quad (2)$$

4. UniMotion: Unifying Motion Synthesis and Understanding

In this section, we introduce UniMotion, a unified model for joint motion synthesis and understanding, including hierarchical control via text. UniMotion generates high-quality motion and text outputs, either from full noise or given conditional inputs such as frame-level text, sequence-level text,

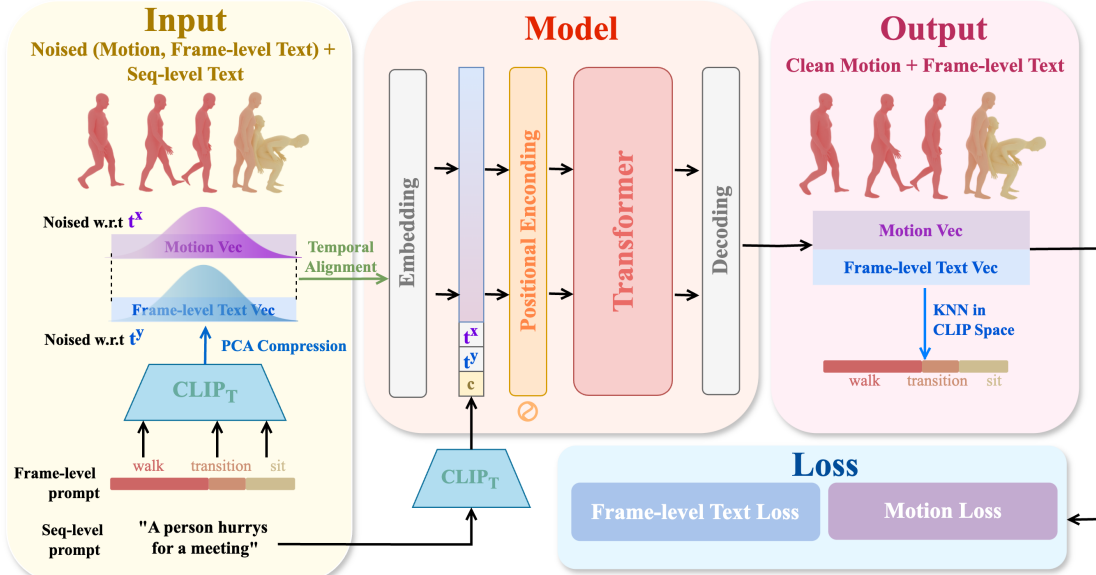


Figure 2. **Overview of UniMotion.** UniMotion is a transformer-based diffusion model (**Model**) that can be input conditioned on a) human motion, b) clip embedded frame-level text, or c) sequence-level text (**Input**) or any subsets thereof or none, and instead supplied with noise. At it's core it allows to diffuse motion and text individually, implemented via separate denoising timesteps t^x and t^y . After training with Frame-level text Losses and Motion losses (**Loss**), see Sec. 4.1. UniMotion can output clean, noise-free motion, and frame-level text descriptions explaining the generated motions. (**Output**)

a motion sequence, or any subsets thereof, (see Fig. 2) spanning a variety of applications treated in isolation by related works.

To achieve this, our model advances prior single-modality motion diffusion models (see Subsec. 4.1) to encompass multi-modal distributions, specifically motion, and fine-grained text. We combine motion sequence and fine-grained frame-level texts, maintaining the temporal alignment of these two modalities to enable temporal semantic awareness (see Sec. 4.2). Unlike previous works, our multi-modality diffusion process supports joint training across datasets with varying annotations (sequence-level and frame-level) (see Subsec. 4.3).

4.1. Multi-Modal Motion and Text Diffusion

UniMotion models the distribution of motion and text sequences in a unified framework, enabling synchronized generation and conditional sampling. Similar in spirit to [2], which jointly models 2D images and text, our approach extends this probabilistic modeling to temporal modalities. More concretely, a frame-level text sequence, $\mathbf{y}^{1:N}$ is denoted analogously to the motion sequence $\mathbf{x}^{1:N}$, where N denotes the sequence length and $\{\mathbf{y}_t^{1:N}\}_{t=0}^T$ are the noise samples created via Eq. 1. Similarly, we drop the notation of sequence length in the following for simplicity. With that, multi-modal diffusion can be achieved by extending G to $G_\theta(\mathbf{x}_{t^x}, \mathbf{y}_{t^y}; t^x, t^y)$ with the additional process and including the separately scheduled diffusion timesteps t^x, t^y

for motion and text respectively. These two timesteps are sampled independently from the same noise schedule sequence as [34]. By virtue of this formulation, the joint distribution $p(\mathbf{x}, \mathbf{y})$ can be sampled at inference time, starting the denoising process with $G_\theta(\mathbf{x}_T, \mathbf{y}_T; T, T)$, and the conditional $p(\mathbf{x}|\mathbf{y})$ by $G_\theta(\mathbf{x}_T, \mathbf{y}; T, 0)$ and analogously $p(\mathbf{y}|\mathbf{x})$. Specifically, we jointly train the model via

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}_0, \mathbf{y}_0), t^x, t^y} \mathbb{E}_{\mathbf{x}_{t^x}, \mathbf{y}_{t^y}} \|G_\theta(\mathbf{x}_{t^x}, \mathbf{y}_{t^y}; t^x, t^y, c) - (\mathbf{x}_0, \mathbf{y}_0)\|_2^2 \quad (3)$$

where θ are weights parametrizing G and \mathcal{U} is the discrete uniform distribution and expectation is taken over distributions: $(\mathbf{x}_0, \mathbf{y}_0) \sim p(\mathbf{x}, \mathbf{y})$, $t^x \sim \mathcal{U}\{0, \dots, T\}$, $t^y \sim \mathcal{U}\{0, \dots, T\}$, $\mathbf{x}_{t^x} \sim q(\mathbf{x}_{t^x}|\mathbf{x}_0)$, $\mathbf{y}_{t^y} \sim q(\mathbf{y}_{t^y}|\mathbf{y}_0)$. In practice, motion and text losses are computed separately using Mean Squared Error (MSE) on the output motion vectors and text embeddings.

4.2. Temporally aligned Text and Motion Encoding

We find that appropriate architectural integration of two modalities (text and motion) into the joint formulation is a key performance factor.

A simple integration is to treat motion and text as separate modalities and feed them independently into the Transformer. Inspired by UniDiffuser [2], our initial approach concatenated motion and text embeddings horizontally, placing N frames of motion vectors first, followed by N frames of text embeddings. However, we find that this lack of alignment negatively impacts performance.

| Method | Training Set | Input | Per-crop semantic correctness | | | Per-crop Realism | | Per-seq Realism | |
|------------|------------------|-------|-------------------------------|-------------------|-------------------|-------------------|-------------------------|-------------------|-------------------------|
| | | | R-Prec@3 \uparrow | M2T \uparrow | M2M \uparrow | FID \downarrow | Diversity \rightarrow | FID \downarrow | Diversity \rightarrow |
| GT | - | - | 0.735 ± 0.008 | 0.663 ± 0.000 | 1.000 ± 0.000 | 0.000 ± 0.000 | 1.375 ± 0.005 | 0.000 ± 0.000 | 1.391 ± 0.003 |
| TEACH | BABEL | f | 0.588 ± 0.007 | 0.623 ± 0.001 | 0.575 ± 0.000 | 0.155 ± 0.001 | 1.340 ± 0.003 | 0.304 ± 0.001 | 1.344 ± 0.003 |
| DoubleTake | BABEL | f | 0.544 ± 0.013 | 0.602 ± 0.002 | 0.560 ± 0.001 | 0.195 ± 0.002 | 1.332 ± 0.005 | 0.353 ± 0.002 | 1.337 ± 0.004 |
| STMC | HML | f | 0.528 ± 0.012 | 0.599 ± 0.000 | 0.616 ± 0.010 | 0.156 ± 0.000 | 1.358 ± 0.005 | 0.233 ± 0.000 | 1.362 ± 0.005 |
| FlowMDM | BABEL | f | 0.618 ± 0.007 | 0.631 ± 0.002 | 0.652 ± 0.001 | 0.101 ± 0.001 | 1.352 ± 0.006 | 0.211 ± 0.002 | 1.375 ± 0.005 |
| Ours | BABEL | f | 0.636 ± 0.017 | 0.633 ± 0.004 | 0.677 ± 0.002 | 0.087 ± 0.002 | 1.366 ± 0.009 | 0.180 ± 0.004 | 1.374 ± 0.002 |
| Ours | HML \cap BABEL | f | 0.668 ± 0.009 | 0.643 ± 0.002 | 0.698 ± 0.002 | 0.071 ± 0.001 | 1.372 ± 0.005 | 0.150 ± 0.001 | 1.378 ± 0.003 |
| Ours | HML \cap BABEL | f + s | 0.679 ± 0.006 | 0.644 ± 0.001 | 0.706 ± 0.002 | 0.066 ± 0.002 | 1.373 ± 0.009 | 0.133 ± 0.004 | 1.381 ± 0.006 |

Table 1. **Frame-Level Text-to-Motion evaluation.** *Per-crop* refers to text segment level evaluation. *Training Set* specifies the dataset used for training. *Input* specifies the type of text input. *f*: frame-level text, *s*: sequence-level text. *f+s* demonstrates that combining multi-level conditioning signals can enhance model performance in terms of semantic correspondence. The evaluation is repeated 10 times, and \pm indicates the 95% confidence intervals.

In contrast, in our setting, where motion and text sequences are inherently aligned, temporal alignment is the key factor. Instead of treating text as independent tokens, we extend each frame’s motion vector with its corresponding frame-level text annotation, ensuring direct alignment along the temporal dimension. This eliminates the need for the model to learn correspondences between word positions and motion frames.

However, this alone does not guarantee performance. We encode text into the space of CLIP [28] with a pertained model. Using the full encodings of pose and text as token creates issues. We hypothesize this is due to an excessive capacity spent on the high-dimensional text tokens. We solve this by projecting CLIP embeddings down to 50 dimensions via PCA [36] and find this improves performance drastically. To get back to text labels from embeddings after diffusion, we match the predicted clip embedding to our database of text labels to obtain the output text using the closest match.

4.3. Data Merging

The popular AMASS dataset [23] of natural human motion, represented by the SMPL body model [22] has recently been annotated in two efforts, namely BABEL [27] and HumanML3D [11]. While HumanML3D annotations consist of sequence-level text annotation, that is, a single text annotation for a motion clip, the BABEL annotations consist of frame-level annotations, assigning semantic label to the pose for each frame of the motion sequence. Instead of restricting to use one at a time, as in prior works, UniMotion is directly trained on both jointly, using sequence level HumanML3D annotations as condition c and frame level sequences as $y^{1:N}$.

A challenge however lies in that both datasets annotate different subsets of AMASS. A trivial solution is to consider overlapping annotations of motions. We denote our model trained with this scenario UniMotion *overlap*, and investigate the performance in experiments (see Sec. 5).

5. Experiments

In this section we investigate the benefit of hierarchical text at inference and training time (i.e. usage of frame-level and sequence-level text). We show the versatility of UniMotion’s unification of synthesis and understanding. Specifically, allowing for frame-level text to motion (Subsec. 5.1) and we for the first time show motion-to-frame-level text (Subsec. 5.2), including a real-world application scenario. Finally, UniMotion is the first model to show joint generation of motion along with frame-level understanding (Subsec. 5.2). In the ablation study, we show that the proposed multi-modality strongly improves generation quality compared to our backbone MDM [34]. (Subsec. 5.3)

Implementation Details.. We utilize a temporally aware transformer, similar to MDM [34]. Text inputs are encoded using pretrained CLIP, followed by PCA reduction. Our model is trained on a single A100, with training spanning approximately 40 hours. Please refer to 4.3 for details on training data.

Baselines. We compare our model to the publicly released works that are capable of frame-level text-to-motion generation: auto-regressive model **TEACH** [1], **DoubleTake** [30] based on diffusion sampling, **FlowMDM** [3], a diffusion model based on Blended Positional Encoding and **STMC** [26], a post-hoc test time method stitching individual predictions of MDM [34]. Note that neither Teach, FlowMDM nor STMC supports hierarchical training. Since STMC admits overlapping control signals we compare to it in terms of hierarchical control. Since no prior works allow for training on sequence and frame-level text input, models are either trained on BABEL(frame-level) or HumanML3D(sequence-level) data, as indicated in our result tables. Please refer to our supp. document for more details.

Evaluation Metrics. First, we introduce our **semantic metrics**, measuring how well the generated motions correspond to their text descriptions. **R-Precision** [11] assesses the accuracy of ranking the correct ground-truth text correspond-

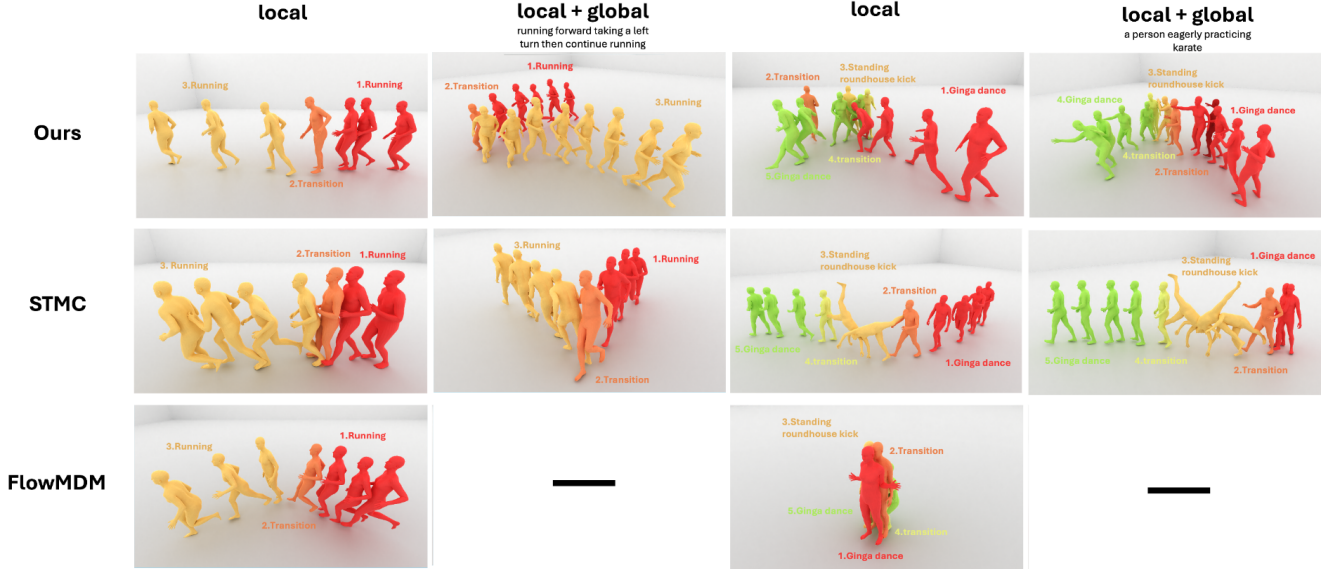


Figure 3. **Text2Motion qualitative results.** **Columns 1,3:** Local text is the input to our method and baselines STMC [26] (adapted) FlowMDM [3]. **Columns 2, 4:** Both local and global text are the input our method and STMC. Our model performs well regardless of the complexity of the local text, in contrast to STMC which fails to generate Ginga dance in columns 3 and 4 and performs walking instead. FlowMDM cannot be conditioned on both global+local text.

ing to a predicted motion at the top positions (Top-1, Top-2, and Top-3) within a set that includes 32 randomly sampled incorrect text matches. With **M2T** [26], we measure how well the per-crop motion matches their textual description, we calculate their cosine similarity in the joint-embedding space of TMR++ [4]. Similarly, the **M2M** [26] score is the cosine similarity between the generated and the ground-truth motion embeddings.

With our **realism metrics**, we measure how well the generated motion distributions fit the ground truth one. We utilize the **Frechet Inception Distance (FID)** [14] to measure the distribution distances and **Diversity** [11] computes the distributions variance, both in the TMR++ as the embedding space.

5.1. Frame-Level Text2Motion Results

We evaluate the Text2Motion task (Tab. 1), where we investigate the effects of frame-level and sequence-level training data. Qualitative analysis is presented in Fig. 3. When we train our model as FlowMDM (best performing prior work) on frame-level labels of all Babel annotations (Tab. 1, Ours *BABEL*) we observe our UniMotion to be consistently better but still roughly on par as expected since both models are using a backbone similar to MDM [34]. The slight improvement can be attributed to the temporal input alignment (see Sec.4.2) and the multi-timestep diffusion training (see Sec. 4.1). Next, we significantly reduce the training dataset size to the subset sequences annotated with both HML (frame-level text) and BABEL (sequence-level text) (cf. Tab. 1

Ours HML-BABEL f). Although one could expect a performance decrease, we find the opposite, a strong consistent performance increase in all metrics - suggesting the strong positive impact of multi-model training. Notably, this is the case although only frame level inputs are given for the evaluation and sequence-level inputs only enrich the models training data. Finally, we investigate the effect of adding sequence-level text into the model for evaluation (cf. Tab. 1 Ours HML-BABEL f + s), again showing a consistent improvement. In conclusion, the evaluation shows cross-modal generalization, consistently improving the results.

5.2. Applications

Please see these and further results in motion in the supplementary video.

Motion2Text. Here we show UniMotions capabilities of predicting frame-level text given human motion. This is a novel task, prior work is not able to do. We, therefore, restrict ourselves to qualitative evaluations. See Fig. 4, where we use UniMotion to annotate MoCap data and Youtube videos with motion descriptions.

Hierarchical Text2Motion. We show that UniMotion, although not directly trained for this task, shows generalization capabilities to compositional text conditioning, where global-text and local-text are giving different but complementary conditioning (see. Fig. 1).

Joint text and motion generation. UniMotion can jointly generate human motion and corresponding frame-level text, allowing users to not only generate motion but also to di-

| Method | Training Set | Input | FID ↓ | Diversity → | R-Prec@1 ↑ | R-Prec@2 ↑ | R-Prec@3 ↑ | M2T ↑ |
|--------|--------------|-------|-------------|-------------|-------------|-------------|-------------|-------------|
| GT | - | - | 0.000±0.000 | 1.391±0.003 | 0.699±0.014 | 0.834±0.011 | 0.878±0.005 | 0.748±0.000 |
| MDM | HML | s | 0.449±0.025 | 1.315±0.014 | 0.376±0.008 | 0.536±0.010 | 0.639±0.010 | 0.631±0.003 |
| Ours | HML-BABEL | f | 0.152±0.002 | 1.377±0.006 | 0.344±0.010 | 0.508±0.019 | 0.587±0.007 | 0.648±0.003 |
| Ours | HML-BABEL | s | 0.195±0.003 | 1.381±0.011 | 0.375±0.021 | 0.539±0.018 | 0.655±0.016 | 0.653±0.004 |
| Ours | HML-BABEL | f + s | 0.133±0.003 | 1.382±0.002 | 0.424±0.005 | 0.593±0.011 | 0.677±0.011 | 0.678±0.002 |

Table 2. **Ablation Study on Sequence-level Text-to-Motion generation.** In this table, we compare with our backbone model MDM[34] to study whether introducing multi-modality helps the motion generation performance. Symbols ↓, and → indicate that lower, or values closer to the ground truth (GT) are better, respectively. The evaluation is repeated 10 times, and ± indicates the 95% confidence interval.



Figure 4. **Motion-to-Text understanding of MoCap and YouTube data.** (a) Given an input MoCap sequence, we use UniMotion to predict frame-level local text. (b) We annotate human motion from YouTube videos with frame-level text. We lift 2D videos to 3D human motion via frame-by-frame pose estimators [10]. We visualize the SMPL human pose (Pink) overlaid on the YouTube videos frames. Then we run UniMotion to predict frame-level annotations (colored text descriptions below the frames). Annotations could serve as valuable audio close captions for the visually impaired.

rectly understand the generated sequence on a frame level. Prior work can not perform this task, see Fig. 5 for conditional joint generation and in Fig. 6 for unconditional generation.

Motion Editing for Content Creation. We show the application of UniMotion to content creation, where a user specifies a desired motion sequence via rough global text and obtains the motion sequence with a frame-level script. The user succeeds by editing the frame-level script and regenerates the motion to obtain the desired edits, see Fig. 1.

5.3. Ablation: Importance of Multi-Modality

In this section, we investigate the importance of our unification of multiple modalities.

Flexibility. As seen in previous experiments, this allows to generate high-quality motion and text outputs, either from full noise or given conditional inputs such as frame-level text, sequence-level text, a motion sequence, or any sub-

sets thereof, (see Fig. 2) spanning a variety of applications treated in isolation by related works.

Improved quality. Additionally, we ablate that the included multi-modality also allows for improved generation quality. For this, we compare our model trained on multi-modal against our backbone architecture MDM [34], which does not include frame-level text in output or input, nor is equipped with the flexible multi-modal diffusion.

Our model, used with the same sequence-level text input data (Table 2, input: s), as MDM, drastically improves MDM in terms of FID and diversity, but also improves or is on par in other metrics. Since the backbone transformer is the same, this shows the strength of the proposed multi-modal training. Notably, this effect is visible even though our training dataset is only a 30% subset of the MDM training dataset.

Combining sequence-level and frame-level text (Table 2, input: f+s) shows a further significant improvement, improving MDM in all metrics. This improvement does not

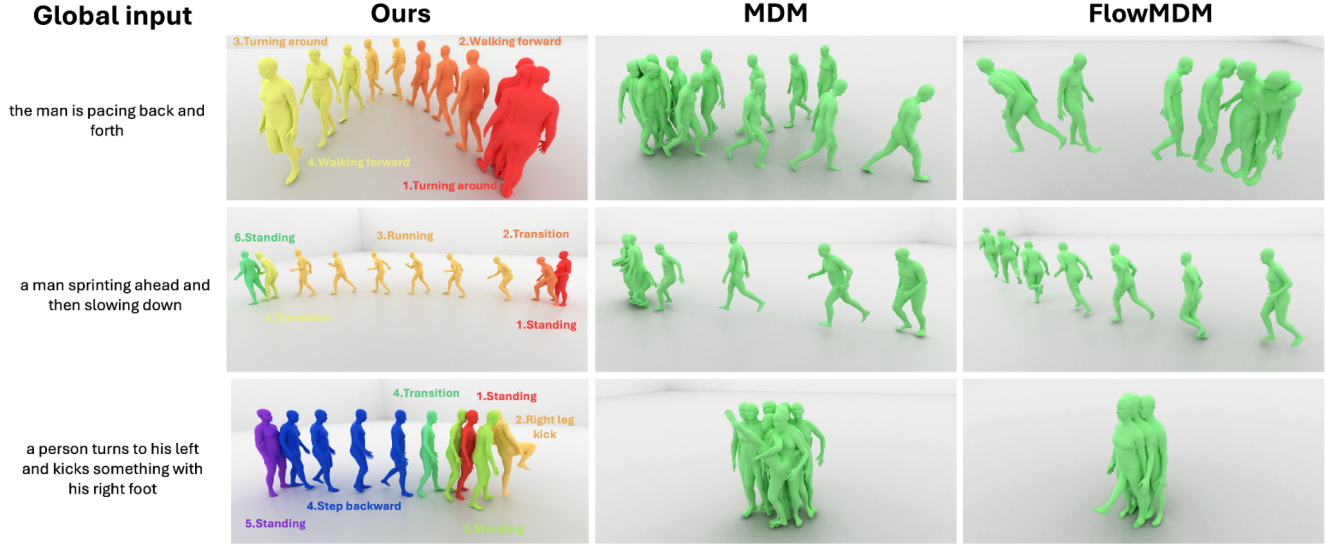


Figure 5. **Joint text and motion generation results.** Input to the models is only the global text shown on the left. We compare the generated motion of ours, MDM [34] and FlowMDM [3]. Our method jointly predicts the frame-level labels, so we can annotate subsequences, while MDM and FlowMDM can only generate the motion.



Figure 6. **Unconditional joint text and motion generation.** Our model, by design, generates poses aligned with local text.

stem from the addition of frame-level text input alone since, in isolation, frame-level labels do not achieve this quality (see Table 2, input: f). We find the interaction between frame-level and sequence-level inputs is the reason for the improvements. In conclusion, the proposed multi-modality is the key factor allowing for improved generation quality.

6. Conclusions

We introduced UniMotion, the first unified multi-task human motion model capable of both flexible motion control and frame-level motion understanding. Using a flexible multi-model diffusion scheme, UniMotion solves several tasks in a unified fashion. Specifically, it unifies tasks that are usually treated in separation by prior works, such as *Frame-Level Text-to-Motion*, *Sequence-Level Text-to-Motion* and *Motion-to-Text*, into a single simple unified model, trained a single time. Importantly, UniMotion’s flexibility also allows for novel tasks not previously considered by prior work like 1.) unconditional generation of human motion with corresponding frame-level text descrip-

tions and 2.) generation of frame-level text from motion, providing granular, time-aware annotations. We show UniMotion opens up new applications: 1.) hierarchical control, allowing users to specify motion at different levels of detail, 2.) obtaining motion text descriptions for existing Mo-Cap data or YouTube videos and 3.) allowing for editability, generating motion from text, and editing the motion via text edits. Moreover, UniMotion attained state-of-the-art results for the frame-level text-to-motion task on the established HumanML3D dataset showing the proposed multi-modality is the key factor allowing for improved generation quality.

Acknowledgments: Special thanks RVH and AVG members for the help and discussion. Prof. Gerard Pons-Moll and Prof. Andreas Geiger are members of the Machine Learning Cluster of Excellence, EXC number 2064/1 - Project number 390727645. Gerard Pons-moll is endowed by the Carl Zeiss Foundation. Andreas Geiger was supported by the ERC Starting Grant LEGO-3D (850533). Julian Chibane is a fellow of the Meta Research PhD Fellowship Program.

References

- [1] Nikos Athanasiou, Mathis Petrovich, Michael J. Black, and Gül Varol. TEACH: Temporal Action Compositions for 3D Humans. In *International Conference on 3D Vision (3DV)*, 2022. 2, 5
- [2] Fan Bao, Shen Nie, Kaiwen Xue, Chongxuan Li, Shi Pu, Yaole Wang, Gang Yue, Yue Cao, Hang Su, and Jun Zhu. One transformer fits all distributions in multi-modal diffusion at scale. In *Proceedings of the 40th International Conference on Machine Learning*. JMLR.org, 2023. 2, 4
- [3] German Barquero, Sergio Escalera, and Cristina Palmero. Flowmdm: Seamless human motion composition with blended positional encodings. *arXiv preprint arXiv:2402.15509*, 2024. 2, 3, 5, 6, 8
- [4] Léore Bensabath, Mathis Petrovich, and Gul Varol. A cross-dataset study for text-based 3d human motion retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1932–1940, 2024. 6
- [5] Zhongang Cai, Daxuan Ren, Ailing Zeng, Zhengyu Lin, Tao Yu, Wenjia Wang, Xiangyu Fan, Yang Gao, Yifan Yu, Liang Pan, Fangzhou Hong, Mingyuan Zhang, Chen Change Loy, Lei Yang, and Ziwei Liu. HuMMan: Multi-modal 4d human dataset for versatile sensing and modeling. In *17th European Conference on Computer Vision, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VII*, pages 557–577. Springer, 2022. 3
- [6] Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, and Gang Yu. Executing your commands via motion diffusion in latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18000–18010, 2023. 3
- [7] Hyung-gun Chi, Myoung Hoon Ha, Seunggeun Chi, Sang Wan Lee, Qixing Huang, and Karthik Ramani. Infogcn: Representation learning for human skeleton-based action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [8] Wenxun Dai, Ling-Hao Chen, Jingbo Wang, Jinpeng Liu, Bo Dai, and Yansong Tang. Motionlcm: Real-time controllable motion generation via latent consistency model. *arXiv preprint arXiv:2404.19759*, 2024. 2, 3
- [9] Purvi Goel, Kuan-Chieh Wang, C Karen Liu, and Kayvon Fatahalian. Iterative motion editing with natural language. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–9, 2024. 3
- [10] Shubham Goel, Georgios Pavlakos, Jathushan Rajasegaran, Angjoo Kanazawa, and Jitendra Malik. Humans in 4D: Reconstructing and tracking humans with transformers. In *ICCV*, 2023. 7
- [11] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and natural 3d human motions from text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 5, 6
- [12] 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. 2, 3
- [13] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael Black. Stochastic scene-aware motion prediction. In *Proceedings of the International Conference on Computer Vision 2021*, 2021. 2
- [14] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2017. 6
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 3
- [16] Yiming Huang, Weilin Wan, Yue Yang, Chris Callison-Burch, Mark Yatskar, and Lingjie Liu. Como: Controllable motion generation through language guided pose code editing, 2024. 3
- [17] Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. Motiongpt: Human motion as a foreign language. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 3
- [18] Jihoon Kim, Jiseob Kim, and Sungjoon Choi. Flame: Free-form language-based motion synthesis and editing. *arXiv preprint arXiv:2209.00349*, 2022. 2
- [19] Jiaman Li, Jiajun Wu, and C. Karen Liu. Object motion guided human motion synthesis. *ACM Transactions on Graphics*, 42(6), 2023. 2
- [20] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d dance generation with aist++, 2021. 2
- [21] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Learn to dance with aist++: Music conditioned 3d dance generation, 2021. 2
- [22] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34, 2015. 5
- [23] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS: archive of motion capture as surface shapes. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. 5
- [24] Aymen Mir, Xavier Puig, Angjoo Kanazawa, and Gerard Pons-Moll. Generating continual human motion in diverse 3d scenes. In *International Conference on 3D Vision (3DV)*, 2024. 2
- [25] Mathis Petrovich, Michael J. Black, and Gül Varol. TEMOS: Generating diverse human motions from textual descriptions. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [26] Mathis Petrovich, Or Litany, Umar Iqbal, Michael J Black, Gül Varol, Xue Bin Peng, and Davis Rempe. Stmc: Multi-track timeline control for text-driven 3d human motion generation. *arXiv preprint arXiv:2401.08559*, 2024. 2, 3, 5, 6

- [27] Abhinanda R. Punnakal, 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)*, 2021. 3, 5
- [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, 2021. 5
- [29] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 3
- [30] Yoni Shafir, Guy Tevet, Roy Kapon, and Amit Haim Bermano. Priormdm: Human motion diffusion as a generative prior. In *ICLR*, 2023. 2, 3, 5
- [31] Yi Shi, Jingbo Wang, Xuekun Jiang, and Bo Dai. Controllable motion diffusion model. *CoRR*, abs/2306.00416. 2
- [32] Sebastian Starke, He Zhang, Taku Komura, and Jun Saito. Neural state machine for character-scene interactions. *ACM Trans. Graph.*, 38(6), 2019. 2
- [33] Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. Motionclip: Exposing human motion generation to clip space. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [34] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *The Eleventh International Conference on Learning Representations*, 2023. 2, 3, 4, 5, 6, 7, 8
- [35] Jingbo Wang, Yu Rong, Jingyuan Liu, Sijie Yan, Dahua Lin, and Bo Dai. Towards diverse and natural scene-aware 3d human motion synthesis. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE. 2
- [36] Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, 2(1):37–52, 1987. Proceedings of the Multivariate Statistical Workshop for Geologists and Geochemists. 5
- [37] Zeqi Xiao, Tai Wang, Jingbo Wang, Jinkun Cao, Bo Dai, Dahua Lin, and Jiangmiao Pang. Unified human-scene interaction via prompted chain-of-contacts. *Arxiv*, 2023. 2
- [38] Yiming Xie, Varun Jampani, Lei Zhong, Deqing Sun, and Huaizu Jiang. Omnicontrol: Control any joint at any time for human motion generation. In *The Twelfth International Conference on Learning Representations*, 2024. 3
- [39] Sirui Xu, Zhengyuan Li, Yu-Xiong Wang, and Liang-Yan Gui. Interdiff: Generating 3d human-object interactions with physics-informed diffusion. In *ICCV*, 2023. 2
- [40] Hongwei Yi, Justus Thies, Michael J. Black, Xue Bin Peng, and Davis Rempe. Generating human interaction motions in scenes with text control. *arXiv:2404.10685*, 2024. 2
- [41] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 3
- [42] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. *arXiv preprint arXiv:2208.15001*, 2022. 2
- [43] Mingyuan Zhang, Xinying Guo, Liang Pan, Zhongang Cai, Fangzhou Hong, Huirong Li, Lei Yang, and Ziwei Liu. Remodiffuse: Retrieval-augmented motion diffusion model. *arXiv preprint arXiv:2304.01116*, 2023. 3
- [44] Mingyuan Zhang, Huirong Li, Zhongang Cai, Jiawei Ren, Lei Yang, and Ziwei Liu. Finemogen: Fine-grained spatio-temporal motion generation and editing. *NeurIPS*, 36, 2024. 3
- [45] Xiaohan Zhang, Bharat Lal Bhatnagar, Sebastian Starke, Vladimir Guzov, and Gerard Pons-Moll. Couch: Towards controllable human-chair interactions. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [46] Xiaohan Zhang, Bharat Lal Bhatnagar, Sebastian Starke, Ilya Petrov, Vladimir Guzov, Helisa Dhamo, Eduardo Pérez Pelitero, and Gerard Pons-Moll. Force: Dataset and method for intuitive physics guided human-object interaction. 2024. 2
- [47] Zixiang Zhou, Yu Wan, and Baoyuan Wang. Avatargpt: All-in-one framework for motion understanding, planning, generation and beyond. *CoRR*, abs/2311.16468, 2023. 2, 3
- [48] Wentao Zhu, Xiaoxuan Ma, Zhaoyang Liu, Libin Liu, Wayne Wu, and Yizhou Wang. Motionbert: A unified perspective on learning human motion representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 3