VERSATILE MOTION-LANGUAGE MODELS FOR MULTI-TURN INTERACTIVE AGENTS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031 032

033

Paper under double-blind review

ABSTRACT

Recent advancements in large language models (LLMs) have greatly enhanced their ability to generate natural and contextually relevant text, making AI interactions more human-like. However, generating and understanding interactive human-like motion, where two individuals engage in coordinated movements, remains a challenge due to the complexity of modeling these coordinated interactions. Furthermore, a versatile model is required to handle diverse interactive scenarios, such as chat systems that follow user instructions or adapt to their assigned role while adjusting interaction dynamics. To tackle this problem, we introduce VIM, short for the Versatile Interactive Motion language model, which integrates both language and motion modalities to effectively understand, generate, and control interactive motions in multi-turn conversational contexts. To address the scarcity of multi-turn interactive motion data, we introduce a synthetic dataset called $INTER-MT^2$; where we utilize pre-trained models to create diverse instructional datasets with interactive motion. Our approach first trains a motion tokenizer that encodes interactive motions into residual discrete tokens. In the pretraining stage, the model learns to align motion and text representations with these discrete tokens. During the instruction fine-tuning stage, VIM adapts to multi-turn conversations using INTER- MT^2 . We evaluate the versatility of our method across motion-related tasks-motion-to-text, text-to-motion, reaction generation, motion editing, and reasoning about motion sequences. The results highlight VIM's versatility and effectiveness in handling complex interactive motion synthesis.

1 INTRODUCTION

Agents that reflect how humans communicate and interact with each other through motion have the potential to revolutionize our interaction with technology. By capturing the subtleties of human communication, including gestures, expressions, and interactive behaviors, these agents can offer more intuitive and natural interfaces. This holistic understanding enables technology to adjust its responses and behaviors based on the user's physical motions and situational context, leading to more personalized and engaging interactions. Such capabilities are crucial for enhancing support across various domains, including robotics, virtual humans, entertainment, and more.

041 Recent advancements in large language models (LLMs) (Dubey et al., 2024; Team et al., 2024; Yang 042 et al., 2024) have demonstrated significant potential in generating human-like text and understanding 043 complex linguistic interactions. They have even extended their capability to multi-modal contexts, 044 successfully integrating various input sources such as images, speech, and videos (Ge et al., 2024; Liu et al., 2024; Chen et al., 2024b; Tang et al., 2024; Shu et al., 2023). Building upon these developments, there is a growing interest in incorporating human (or robot) motion as a new modality (Jiang 046 et al., 2024; Chen et al., 2024a), leading to the emergence of the "motion-language models" (MLM). 047 However, existing approaches (Zhang et al., 2023; Guo et al., 2024a; 2022; Zhang et al., 2024d; Cai 048 et al., 2024) often focus on single tasks, such as text-to-motion or motion-to-text translation, and consider only single motions without interactions. This limitation hinders the agents' ability to handle scenarios involving multi-agents, complex interactions, and multi-turn conversations. 051

Beyond modeling the motion of a single person, interactive motion, where two individuals partic ipate in interactions, allows the model to learn about social behavior. Modeling such interactions requires versatility to effectively control interactions, allowing users to provide instructions, assign



064

065

066

067

068 069

071

073 074



Figure 1: We introduce VIM, a versatile interactive motion language model. Left: reaction generation, motion reasoning, and generation. Right top: text-to-motion generation. Right bottom: Motion editing and motion understanding.

roles, or modify behaviors. In this paper, we aim to build a unified yet versatile motion-language
 model designed to generate, control, and comprehend sophisticated interactive motions.

One of the primary challenges in developing these interactive agents is the lack of multi-turn interactive motion data. Datasets containing motions of two individuals interacting with each other, along with multi-turn conversational instructions, are scarce and challenging to collect. This makes it difficult for models to learn the nuances of interactive motions and multi-turn dynamics.

To address this, we present a new synthesized dataset called INTER-MT², which includes various instructions about the interactive motions, in a multi-turn conversational format. We utilize a large language model to produce diverse instructions with motion captions and a diffusion-based text-tomotion model to generate corresponding motions. We expect that leveraging such foundation models to construct training data allows the model to generalize more effectively with prior knowledge.

087 Building upon our synthesized dataset, we present VIM, a Versatile Interactive Motion-language model designed for multi-turn conversations involving interactive motions. We pursue the versatility of VIM through a unified architecture that can simultaneously input and output both motion and text modalities. Based on the pre-trained LLMs, our training process can be divided into three stages: 090 (1) training of the interaction motion tokenizer, (2) pre-training for motion and text representation 091 alignment, and (3) instruction tuning with our synthesized dataset, INTER-MT², to handle more 092 complex and multi-turn instructions. This enables VIM to effectively comprehend, generate, and control interactive motions, as illustrated in Figure 1. To assess VIM's capabilities, we introduce 094 new evaluation protocols that evaluate its performance across various motion-related tasks. This in-095 clude editing motions and reasoning about motion sequences based on contextual cues, highlighting 096 its versatility and effectiveness in complex motion interaction scenarios.

In summary, the main contributions of this paper are threefold: (1) We propose VIM that can simultaneously process and generate both motion and text modalities, along with a three-stage training pipeline consisting of motion tokenizer training, pre-training for modality alignment, and instruction tuning. (2) We present INTER-MT², a multi-turn interactive motion-text dataset, to address the lack of multi-turn interactive motion data. (3) We introduced a new evaluation protocol to evaluate the performance of motion-language models on complex motion interaction scenarios.

103 104 105

- 2 RELATED WORK
- **Human Motion Modeling & Control** Advancements in human motion modeling have driven significant progress in both motion generation and control. Diffusion-based methods, such as MDM

108 (Tevet et al., 2023), FG-T2M (Wang et al., 2023), and MotionDiffuse (Zhang et al., 2024a) ex-109 cel at synthesizing realistic human motions from the text. Transformer-based models with vector 110 quantization, such as TM2T (Guo et al., 2022) and T2M-GPT (Zhang et al., 2023), effectively cap-111 ture complex motion patterns. MoMASK (Guo et al., 2024a) improves motion granularity with 112 residual tokenizers. For motion editing, some approaches focus on style transfer (Aberman et al., 2020; Guo et al., 2024b) or specific body part modifications (Zhang et al., 2024a; Kim et al., 2023). 113 FineMoGEN (Zhang et al., 2024c) offers fine-grained motion synthesis based on user instructions. 114 MEOs (Goel et al. (2024)) use captions and large language models to identify frames and body parts 115 to edit, while MotionFix (Athanasiou et al. (2024)) conditions diffusion models on both source mo-116 tion and edit text for seamless motion edits. However, these models usually target single tasks (e.g., 117 text-to-motion, or motion editing) and lack versatility in handling input and output of both motion 118 and text simultaneously in a unified architecture. 119

120

Motion Language Model Recent developments in motion language models have aimed to achieve 121 versatility across various motion-related tasks. MotionGPT (Jiang et al., 2023) demonstrates versa-122 tility in motion comprehension and generation based on a unified framework. MotionChain (Jiang 123 et al., 2024) introduces a multi-turn conversational system for interpreting and generating motions 124 within dialogue contexts, including image inputs. Zhou et al. (2024) introduces AvatarGPT integrat-125 ing motion generation and planning ability in motion large language model. Some studies, like Chen 126 et al. (2024a), expand modalities into speech, music, and videos but focus primarily on comprehen-127 sion rather than generation. Zhang et al. (2024b) proposed unified models for generating motion from various input modalities. M³-GPT, from Luo et al. (2024), models speech, music, text, and 128 motion interchangeably. However, modeling interactive motions in versatile large models remains 129 under-explored. While some efforts, such as Wu et al. (2024), address this direction, they often lack 130 multi-turn interactions and complex reasoning abilities. Our work addresses this gap with a model 131 trained on our synthesized INTER- MT^2 dataset, enabling the understanding and generation of in-132 teractive motions in multi-turn conversations with advanced reasoning capabilities. This approach 133 facilitates more nuanced, context-aware motion generation in complex interactive behaviors. 134

134 135

Human-Human Interactive Motion Modeling Modeling human-human interactions has gar-136 nered increasing attention in recent research. Several multi-person interaction datasets (Ng et al., 137 2020; Fieraru et al., 2020; Yin et al., 2023) have been developed, and recent efforts like Inter-X (Xu 138 et al., 2024a) and InterHuman (Liang et al., 2024) have collected interactive motions paired with tex-139 tual descriptions for text-based motion control. In text-to-motion tasks, InterGEN (Xu et al., 2024a) 140 introduces a diffusion-based model with spatial constraint loss. PriorMDM (Shafir et al., 2024) 141 leverages pre-trained motion diffusion models with slim communication blocks. For reaction gener-142 ation, ReMoS Ghosh et al. (2023) synthesizes reactive motion using spatio-temporal cross-attention, while ReGenNet Xu et al. (2024b) employs a transformer-based model with distance-based inter-143 action loss to predict human reactions. While existing models have advanced interactive motion 144 modeling, they lack versatility and focus on specific tasks, failing to capture complex multi-turn 145 dynamics. To address this, we introduce INTER-MT², enabling agents to generate sophisticated 146 motions, respond to instructions, adapt roles, and adjust behaviors based on context.

- 147 148
- 149 150

3 INTER-MT²: INTERACTIVE MULTI-TURN MOTION-TEXT DATASET

Current datasets (Yin et al., 2023; Liang et al., 2024; Xu et al., 2024a) for modeling interactive motions lack sufficient diversity in instructions and do not include multi-turn conversations. To address this gap, we introduce INTER-MT²: INTERactive MUTI-Turn Motion-Text dataset. This dataset covers a variety of interactive motion scenarios with multi-turn conversations, diverse instructions, and spatiotemporally aligned motions between two individuals.

Collecting diverse instructional data for interactive motions and multi-turn conversation samples
 poses significant challenges. The first challenge is obtaining instruction and conversational data
 that encompass complex reasoning and generation capabilities. Leveraging pre-trained foundation
 models to generate a broad range of instructions can enrich the dataset with diverse and intricate
 examples. The second challenge is acquiring motion data that aligns with these text instructions.
 Functional approaches using rule-based methods may struggle to maintain spatial and temporal con straints in complex interactive scenarios, while retrieval-based methods are limited by dependence



Figure 2: Overview of synthetic data generation for multi-turn conversations with interactive motions. (a) Motion captions and instructions are generated using GPT-4 based on interactions between two characters, followed by (b) the corresponding motion being synthesized using the InterGEN.

179 on existing cases and lack diversity. Alternatively, generative approaches using pre-trained models show promise in producing diverse, complex text-to-motion sequences, offering more flexibility for 181 modeling interactive motions.

182 We utilize the Inter-X (Xu et al., 2024a) and InterHuman (Liang et al., 2024) datasets as the foun-183 dation for building our datasets. We further employ GPT-40 (OpenAI, 2024) to generate motion captions and conversational instructions for a variety of tasks, such as motion editing, reasoning, 185 and story generation, enhancing the model's versatility. Motion captions are sourced from these base datasets or generated by large language models (LLMs). We utilize the state-of-the-art textto-motion diffusion model, InterGEN Liang et al. (2024), to generate corresponding motions that 187 align with the generated caption from LLMs. Our data collection pipeline, shown in Figure 2, com-188 prises 82K samples of multi-turn conversational data involving interactive motions, including 96K 189 of synthesized interactive motions and 56K motions from the source dataset. 190

191 192 193

199

200

175

176

177 178

4 VIM: VERSATILE INTERACTIVE MOTION-LANGUAGE MODEL

In this section, we introduce VIM, a versatile interactive motion language model, designed to incor-194 porate multi-turn conversations considering both language and interactive motion as input or output 195 modality. First, we will explain the underlying philosophy behind our design choices for the model 196 architectures, followed by a detailed description of the training methodologies. Then, we introduce 197 advanced interactive motion tasks in multi-turn conversations. 198

4.1 NOTATIONS

201 Formally, we denote interactive motion from two individual a and b as $\{\mathbf{m}_a, \mathbf{m}_b\}$, following non-202 canonical representation from Liang et al. (2024) based on SMPL-X structure (Pavlakos et al., 2019). 203 Each timestep of the motion $\mathbf{m}^i = [\mathbf{j}_g^p, \mathbf{j}_g^v, \mathbf{j}^r, \mathbf{c}^f]$ is composed of global joint positions $\mathbf{j}_g^p \in \mathbb{R}^{3N_j}$, 204 global joint velocities $\mathbf{j}_q^v \in \mathbb{R}^{3N_j}$, 6D representation of local rotations $\mathbf{j}^r \in \mathbb{R}^{6N_j}$, with the number 205 of joints N_i , and binary ground contact features $\mathbf{c}^f \in \mathbb{R}^4$. We aim to train a motion language model 206 p_{θ} which can model texts and motions for both inputs and outputs. We define input as an instruction 207 or previous context and output as an answer, with the template below

208 209

210

211 where X_u and X_a are composed of both a mixture of text modalities and motion modalities.

- 212
- 213 4.2 ARCHITECTURE 214
- Our architecture for modeling and generating interactive motions consists of three primary com-215 ponents: encoders (tokenizers), a large language model block, and decoders. This design allows



Figure 3: Method Overview. Stage 1 involves training a motion tokenizer that encodes and decodes interactive motion data. In Stage 2, we pre-train the model by integrating motion and text data, allowing it to learn the alignment between text and motion. Stage 3 focuses on Instruction Tuning, fine-tuning the model to follow instructions and improve its responsiveness to conversational cues.

for the integration of both motion and text data within a unified framework. For motion, we use a residual vector quantized variational auto-encoder (RQ-VAE (Lee et al., 2022; Guo et al., 2024a)) as a tokenizer. Vector quantized variational auto-encoders (Van Den Oord et al., 2017) are effective, but their quantization causes information loss and reduces reconstruction quality, which is critical for accurately modeling interactive motions. The motion encoder \mathcal{E}_M applies 2D convolutions to motion features along the time axis, converting motion pairs $\mathbf{m}_a, \mathbf{m}_b$ into latent vectors $\{\mathbf{z}_a^{1:L}, \mathbf{z}_b^{1:L}\} = \mathcal{E}_M(\{\mathbf{m}_a^{1:M}, \mathbf{m}_b^{1:M}\})$, where *M* is a motion length and L = M/l with down-sample rate *l*. Then the latent vectors \mathbf{z} are quantized by RQ-VAE as an ordered *D* discrete codes:

$$\mathcal{Q}(\mathbf{z}^{i};\mathcal{C},D) = (k_{1}^{i},\cdots,k_{D}^{i}) \in [K]^{D}$$

$$\tag{1}$$

where C is the codebook, K = |C|, and k_d^i is code of z at timestep *i* and depth *d*. These discrete codes form a motion vocabulary. For text, we use a standard text tokenizer to process textual instructions and descriptions into tokens.

 $\mathcal{T}_{\mathbf{r}}$

Tokens are then proceeded to the language model block, which serves as the central processing unit. In this work, we have utilized the LLaMA-3.1-8B (Dubey et al., 2024) model as a base model. We integrate motion tokens with text tokens using a unified vocabulary, which combines both the text and motion vocabularies into one, with special tokens added to mark the start and end of the motion sequences. This shared token space enables the model to efficiently process and generate both modalities for motion-related tasks. Interactive motion is represented as $X_m = \{k_{1:D}^{1;a}, k_{1:D}^{1;b}, \dots, k_{1:D}^{L;a}, k_{1:D}^{L;b}\}$, where X_m is a sequence of motion represented in unified vocabulary and $k_{1:D}^{i;a} \in [K]^D$ is the *i*-th token of motion *a*.

Finally, the decoding stage reverses the encoding process. For motion, the decoder \mathcal{D}_M projects the quantized features $\hat{\mathbf{z}}^i = \sum_{d=1}^{D} \mathbf{e}(k_d^i)$, where \mathbf{e} is codebook embedding, back into motion sequences using 1D convolution. Text decoding follows standard language model decoding procedures.

4.3 TRAINING

Motion Tokenizer The first stage is to train a motion tokenizer composed of an encoder, decoder, and quantizer. We followed the original objective functions from Lee et al. (2022) to train this model, minimizing the reconstruction loss, the codebook loss to align the encoder's outputs with the codebook, and the commitment loss to ensure encoder consistency. Once the encoder and decoder are optimized, we maintain this model freezed during the rest of the training stage.

268

261

262

232

233

234

235 236

245

Pre-training Strategy In the pre-training stage, we train a pre-trained large language model to align motion representations with textual representations. We design tasks including motion-to-

270 text, text-to-motion, motion prediction, and reaction generation to train the model, leveraging paired 271 datasets like Inter-X Xu et al. (2024a) and InterHuman Liang et al. (2024). Using the template from 272 Jiang et al. (2023), we create input sequences y from motion sequences X_m and the corresponding 273 motion caption. Since both motion tokens and text tokens are discrete, we train our model with the general language modeling next-token prediction objective: $\mathcal{L} = -\log \sum^{T} p_{\theta}(y_i|y_{<i})$, where T is 274 275 the length of the multi-modal sequence and i only counts when the text token appears at position 276 *i*. To improve training efficiency, we train the LLM using a low-rank adaptor (LoRA) (Hu et al., 2022), similar to Ge et al. (2024). We then merged the LoRA parameters to the LLM backbone for 278 further training. Furthermore, due to a limited number of interactive motion data, we also leverage larger single motion-text datasets from Motion-X (Lin et al., 2024). This offers prior knowledge of 279 how individual motions are described in language, enhancing the model's ability to align motions 280 with corresponding textual descriptions. 281

Instruction-tunning with INTER-MT² Data In this stage, we aim to enhance the model's ability to follow a wide range of instructions presented in a conversational format. We utilize INTER-MT² dataset combined with single-turn data from prior interactive motion datasets (Xu et al., 2024a; Liang et al., 2024). We follow instruction templates from MotionGPT (Jiang et al., 2023), to format the input and output for single-turn data. The multi-modal sequence y consists of user instructions and corresponding responses, formatted as $y = (X_u^1, X_a^1, X_u^2, X_a^2, \cdots)$. The training objective remains the same as in the pre-training stage, with user instructions omitted in the loss function.

289 290 291

4.4 Advanced Downstream Interactive Motion Tasks

292 After training, our model can perform complex reasoning and generate interactive motions within 293 multi-turn conversations. To verify this, we introduce two additional tasks requiring advanced capabilities: motion reasoning and editing. Motion reasoning involves predicting past or future events, or 294 reasoning about current motions, based on prior conversational data. This task requires the model to 295 understand the context of the conversation, interpret how the given motion fits within that context, 296 and adjust its reasoning accordingly. In the motion editing task, we focus on altering a person's 297 persona or shifting scenarios, such as emotions or relationship dynamics. This adds complexity, as 298 changes to one person's behavior affect the other's motion. The model must edit the target motion 299 while maintaining contextual coherence, requiring a deep understanding of social dynamics. 300

301 302

303

5 EXPERIMENTS

In our experiments, we evaluated VIM's ability to generate detailed motion-based chat responses, requiring complex reasoning about interactive motions, alongside traditional motion-related tasks. We focused on two main questions: first, whether the model can reason effectively about interactive motions, such as refining motions in editing tasks or generating contextually accurate narratives in motion reasoning. Second, we evaluated whether the training with INTER-MT² dataset improves the model's performance in text-to-motion, motion-to-text, and reaction generation tasks.

309 310

311

5.1 EVALUATION TASKS AND BASELINES

312 Motion Reasoning Motion reasoning involves predicting past or future events or interpreting cur-313 rent motions using prior conversational context. We utilize powerful LLMs, i.e., GPT-40 (OpenAI, 314 2024) to assess the content alignment, naturalness, and logical coherence of the generated textual re-315 sponses. Content alignment evaluates how accurately the text reflects the given motion data, logical 316 coherence checks the consistency and reasoning accuracy of inferences made about past or future 317 events, and naturalness evaluates the fluency of generated texts, with rating each metric on a 10-point 318 scale. In addition, we utilized linguistic metrics like Rouge-L (Lin (2004)), METEOR (Banerjee & Lavie (2005)), and MAUVE (Pillutla et al. (2021)), to quantitatively assess the relevance, accuracy, 319 and naturalness of the generated responses compared with labeled texts in INTER-MT² test dataset. 320 We present the results on motion reasoning in §5.2. 321

322

Motion Editing The goal of motion editing is to modify a reference motion based on user instructions. We conducted within-subject user studies to compare edited motion samples, with participants

Methods	LI	LM-Assiste	d	Linguistic Metrics			
	Coh. \uparrow	Align. ↑	Nat. ↑	ROUGE-L	METEOR	MAUVE	
two-stage approach							
TM2T + LLaMA-3.1-8B	3.852	3.050	6.348	0.158	0.226	0.009	
TM2T + GPT-40	<u>4.266</u>	<u>3.455</u>	<u>6.790</u>	0.162	0.227	0.019	
unified approach							
MotionGPT*	1.855	1.303	3.574	0.113	0.096	0.005	
MotionGPT $_{I}^{*}$	3.690	3.160	5.291	0.207	0.218	0.417	
VIM w/o INTER-MT ²	2.770	2.141	4.968	0.155	0.145	0.004	
VIM (Ours)	5.252	4.511	6.981	0.239	0.260	0.794	

Table 1: Evaluation on Motion Reasoning task with INTER-MT² test set. Coh., Align., and Nat. denote logical coherence, content alignment, and naturalness, respectively. **Bold** indicates best performance and <u>underline</u> denotes the second best performance.

337 338

327 328

339 rating them on three metrics: content similarity, instruction alignment, and motion quality, using a 340 5-point Likert scale, following Goel et al. (2024). Content similarity evaluates whether the edited 341 motion preserves the original meaning of the source motion, while instruction alignment assesses 342 how accurately the edited motion follows the given command. The study involved 30 participants, 343 each evaluating five samples from a set of 30 randomly selected test samples. Participants evaluated four baselines and our method, viewing randomly shuffled motion outputs side by side and providing 344 feedback on all metrics. We also employed data-driven metrics, such as Frechet Inception Distance 345 (FID) and mean per joint position error (MPJPE) in meters, to evaluate the quality of the generated 346 edited motion against the labeled motions in the INTER-MT² test set, following Goel et al. (2024). 347 The results on motion editing is shown in §5.3. 348

349 **Traditional Motion Relevant Tasks** For standard motion-related tasks, we evaluated the proposed 350 method in three traditional motion-relevant tasks in interactive motions: motion-to-text, text-to-351 motion, and reaction generation, in the union of the test set in InterHuman (Liang et al., 2024) 352 and Inter-X (Xu et al., 2024a) datasets. To evaluate the text-motion matching score, we report the 353 retrieval precision based on the feature space of retrieval models (Petrovich et al. (2023)). This 354 evaluates the accuracy of matching between texts and motions using Top 3 retrieval accuracy with 355 a fixed batch of 32. The quality of motion was measured by Frechet Inception Distance (FID), 356 which measures the distance of feature distribution between motion data and generated motion. In 357 addition, we measure the mean per joint position error (MPJPE) in meters to evaluate the accuracy 358 of the reaction motion. The results on motion-related tasks are shown in §5.4.

360 **Baselines** Since our interactive multi-turn scenarios and tasks are novel, there is no exact comparison method. We consider and compare reasonable baselines that handle both motion and texts as 361 input and output. We first employ a two-stage approach using off-the-shelf methods. For the motion 362 reasoning task, we convert all motions into text descriptions using the state-of-the-art motion-to-text 363 method TM2T (Guo et al., 2022). These textual descriptions are then used for text-based reason-364 ing with large language models such as GPT-40 (OpenAI, 2024) and LLaMA-3.1-8B (Dubey et al., 2024). For the motion editing task, we combine the text converted by TM2T with the editing text 366 command, and generate the edited motion using the off-the-shelf text-to-motion method InterGen 367 (Liang et al., 2024). As baselines for the *unified approach*, we leverage MotionGPT (Jiang et al., 368 2023) framework with the following configurations: (1) MotionGPT*: a modified MotionGPT fine-369 tuned on interaction data with instruction templates; (2) Motion GPT_I^* : Motion GPT^* enhanced with 370 INTER-MT² dataset; (3) VIM w/o INTER-MT²: our method fine-tuned with instruction templates 371 from MotionGPT, but without INTER- MT^2 data.

372

359

373 5.2 MOTION REASONING

374

In the motion reasoning task, conversations about two interactive motions are examined to assess the model's ability to deduce past or future events and comprehend the motivations driving the motions. The experimental results in Table 1 demonstrate that our unified model, VIM, significantly outperforms two-stage approaches across all LLM-assisted and linguistic metrics. Specifi-



Figure 4: Generated samples for interactive motion reasoning task. This example shows how VIM explains behaviors and their motivations, demonstrating a deeper understanding of scenarios by incorporating context from prior interactions.

Table 2: Data-driven evaluation in motion editing in INTER-MT² test set.



407 cally, VIM achieves improvements with performance increases exceeding 1.9 points in logical co-408 herence, 1.1 points in content alignment, and nearly 0.2 points in naturalness compared to the best 409 two-stage model. The baseline models trained solely on text-motion pair datasets, such as VIM w/o 410 INTER-MT² and MotionGPT^{*}, show limited reasoning capabilities. Although MotionGPT^{*}_I, which 411 incorporates INTER-MT² datasets, exhibits improved performance compared to baselines trained 412 without INTER-MT², it still does not match the effectiveness of the two-stage approaches.

413 The improved performance of our unified model, VIM, over two-stage approaches, appears to result 414 from two key factors: error accumulation and interpretation ambiguity. First, in two-stage models, 415 errors can accumulate; if the motion captioning model generates incorrect motion captions, those 416 mistakes carry over to the second stage, reducing content alignment and coherence. In contrast, VIM's unified architecture integrates motion encoding and reasoning in a single framework, mini-417 mizing error propagation. Second, interpreting motions is not always straightforward, with multiple 418 ways to understand and describe the same motion. In two-stage models, mapping the motion to a 419 single caption for the second stage can lead to more contextually accurate reasoning of the given sce-420 narios or contexts. Our unified model, however, is built to recognize these varied interpretations and 421 generate reasoning that is more contextually accurate. Figure 4 showcases it's ability to dynamically 422 adjust interpretations and responses by incorporating context from previous conversations. 423

424 425 5.3 MOTION EDITING

We aim to validate the hypothesis that people will perceive the generated edited interactive motion from the proposed method to be more content-consistent, instruction-aligned, and better quality, through user subject studies. To analyze the results, we conducted a repeated-measures multivariate analysis of variance on the rated measures. We observed that methods significantly affect the user's

430 431

392

393

394 395

396

397

398

399

400

401

402

403

404

¹We plotted the difference in a post hoc pairwise comparison of the proposed method only. We denote * as 0.01 , ** as <math>p < 0.01, and *** as p < 0.001. The error bars represent 95% confidence intervals.



Figure 6: Generated samples for interactive motion editing. The proposed method excels in capturing nuances, outperforming alternatives in content similarity and instruction alignment.

Mathada	M2T	T2N	1	Reaction Gen.	
Methods	R Top3 ↑	R Top3 ↑	$FID\downarrow$	$MPJPE \downarrow$	$FID\downarrow$
Real	0.867	0.869	0.00	-	0.00
MotionGPT*	0.494	0.328	0.123	3.444	0.355
Motion GPT_I^*	0.503	0.331	0.118	1.436	0.380
VIM w/o INTER-MT ²	0.894	0.561	0.082	0.984	0.031
VIM (Ours)	0.901	0.568	0.059	0.691	0.019

Table 3: Comparisons for three motion-related tasks on Inter-X and InterHuman datasets. M2T denotes motion-to-text, T2M for text-to-motion, and Reaction Gen. for reaction generation.

perception of all dimensions; F(4) = 4.591, p = 0.002, $\eta^2 = 0.137$ for content similarity, F(4) = 7.134, p = 0.000, $\eta^2 = 0.197$ for instruction alignment, and F(4) = 4.781, p = 0.001, $\eta^2 = 0.142$ for motion quality, with all $\alpha = 0.05$. The estimated marginal mean of the rated score is reported in Figure 5. The results show that the proposed method had better instruction alignment, quality, and content consistency across other baselines with significant differences.

During post hoc pairwise comparisons, we identified a significant difference, with our proposed method outperforming the two-stage model (TM2T (Guo et al., 2022) with InterGEN (Liang et al., 2024)) in content similarity (p = 0.017) and instruction alignment (p = 0.010). The two-stage model showed lower content similarity due to motion-to-text conversion errors, leading to unin-tended motions, whereas our unified framework avoids such error accumulation. Additionally, the two-stage model struggled with instruction alignment because InterGEN was trained to generate motions from captions, limiting its ability to adapt to varying personas or contexts. In contrast, our method, trained on diverse instructions and tasks, demonstrated superior reasoning and adaptability, resulting in more accurate motion generation based on instructions and source motions.

In post hoc pairwise comparisons with MotionGPT^T, we observed significant differences, with our proposed method performing better in content similarity (p = 0.005), instruction alignment (p < 0.0005), and motion quality (p = 0.009). This suggests that the VQ-VAE-based tokenizer and conditional generation model negatively impacted performance. Additionally, compared to VIM w/o INTER-MT², there were significant differences in content similarity (p = 0.010) and instruction alignment (p = 0.001), indicating that without INTER-MT² data, the model struggles to control motion based on context and instructions. We also evaluated the proposed method using data-driven metrics, including FID and MPJPE, as shown in Table 2. The proposed method outper-forms the baselines on both measures, which is consistent with the results from user studies. Figure 6 illustrates the generated edited motions based on the source motion and instruction.

5.4 TRADITIONAL MOTION RELATED TASKS

The results in Table 3 support our hypothesis that utilizing the INTER-MT² dataset enhances the model's performance in traditional motion tasks like motion-to-text (M2T), text-to-motion (T2M),

Table 4:	Ablation	Studies	on	motion	tokenizer.
----------	----------	---------	----	--------	------------

Mathada	Reasoning		Editing		M2T	2T T2M		I Reaction Gen.		
Methous	Coh. \uparrow	Align. ↑	Nat.↑	$MPJPE \downarrow$	FID↓	R Top3 ↑	R Top3↑	$FID\downarrow$	$MPJPE \downarrow$	$FID\downarrow$
VIM-VQ	5.004	4.256	6.915	0.892	0.128	0.861	0.601	0.101	1.109	0.055
VIM (Ours)	5.252	4.511	6.981	0.758	0.064	0.901	0.568	0.059	0.691	0.019

493 and reaction generation. The first row ("Real") shows retrieval accuracy, and FID scores from the 494 dataset labels. Note that both VIM w/o INTER-MT² and MotionGPT* were trained on all of these 495 tasks for fair comparison. Comparing the VIM w/o INTER-MT² to the version trained with INTER-496 MT^2 ("Ours"), we see improvements across all tasks. In M2T, Top-3 retrieval accuracy rose from 497 0.894 to 0.901. For T2M, Top-3 retrieval accuracy increased from 0.561 to 0.568, with FID dropping 498 from 0.082 to 0.059, indicating better motion generation. In reaction generation, MPJPE dropped 499 from 0.984 to 0.691, and FID from 0.031 to 0.019, confirming that multi-turn datasets improve mo-500 tion comprehension and generation. Using the INTER-MT² dataset provides diverse, context-rich examples, helping the model learn more nuanced relationships between text and motion. Addi-501 tionally, incorporating INTER-MT² in MotionGPT^{*}, denoted as MotionGPT^{*}_I, improved retrieval 502 precision accuracy for M2T and T2M tasks, and joint position error in reaction generation. 503

505 5.5 Ablation Studies on Motion Tokenizer

506 We conducted ablation studies comparing the VQ-VAE-based model with our RQ-VAE-based ap-507 proach, as shown in Table 4. The RQ-VAE-based motion tokenizer outperformed the VQ-VAE 508 model in motion reasoning tasks, achieving higher scores in coherence, alignment, and naturalness. 509 This improvement is attributed to reduced information loss, allowing our model to capture finer mo-510 tion details while also enhancing its motion-to-text retrieval precision. For generation and editing 511 tasks, the VQ-VAE model achieved slightly better text-to-motion retrieval accuracy but performed 512 worse in FID and MPJPE across editing, reaction generation, and T2M tasks, indicating degraded 513 motion quality and less precise motion details. In contrast, our approach reduced MPJPE by 0.055 for reaction generation, preserving joint dynamics and producing more realistic and natural motions. 514 VQ-VAE's limitations are especially problematic for modeling interactive motions, where precise 515 relative positioning is crucial, making its information loss and reconstruction quality more evident. 516

517 518

519

504

6 CONCLUSION AND DISCUSSIONS

Conclusion In this paper, we introduced VIM, a versatile motion-language model designed to
 model, understand, and reason about interactive motions. We outlined its architecture and provided
 detailed training strategies to create a unified framework integrating large language models with
 interactive motion modality. To enhance the model's reasoning capabilities and versatility, we pre sented a specialized dataset, INTER-MT², which incorporates a variety of reasoning tasks set within
 multi-turn conversations centered on interactive motions. Our experiments demonstrated VIM's
 ability to effectively follow instructions, edit motions, and reason about interactive motions.

527

Limitations and Impact Statement There are several limitations that warrant attention. First, 528 the model's expressiveness remains limited when handling complex or previously unseen actions, 529 indicating a need for further diverse motion source data in its ability to generalize across diverse 530 motion scenarios. Second, the sequence length becomes excessively long as we flatten the resid-531 ual motion tokens, which can impact efficiency and computational resources. Leveraging addi-532 tional transformer models to predict the residual token can reduce this work. Lastly, our method 533 faces challenges in personalization and interpretability, as motion is inherently ambiguous and users 534 may interpret the same motion in different ways. Addressing this issue will require incorporating more tailored approaches that adapt to individual user preferences and expectations through further human-in-the-loop feedback and refinement processes. In terms of broader impact, VIM opens up 536 new possibilities for interactive motion modeling and understanding in AI, potentially benefiting 537 fields like robotics, virtual environments, and human-computer interaction. However, careful con-538 sideration of ethical concerns, such as misinterpretation of motions or unintended behavioral biases, is crucial as the model evolves.

540 REFERENCES

558

559

560

561

- Kfir Aberman, Yijia Weng, Dani Lischinski, Daniel Cohen-Or, and Baoquan Chen. Unpaired motion
 style transfer from video to animation. *ACM Transactions on Graphics (TOG)*, 39(4):64–1, 2020.
- Nikos Athanasiou, Alpár Ceske, Markos Diomataris, Michael J Black, and Gül Varol. MotionFix: Text-driven 3d human motion editing. *arXiv preprint arXiv:2408.00712*, 2024.
- Satanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proc. of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pp. 65–72, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/W05-0909.
- Zhongang Cai, Jianping Jiang, Zhongfei Qing, Xinying Guo, Mingyuan Zhang, Zhengyu Lin, Haiyi
 Mei, Chen Wei, Ruisi Wang, Wanqi Yin, Liang Pan, Xiangyu Fan, Han Du, Peng Gao, Zhitao
 Yang, Yang Gao, Jiaqi Li, Tianxiang Ren, Yukun Wei, Xiaogang Wang, Chen Change Loy, Lei
 Yang, and Ziwei Liu. Digital life project: Autonomous 3d characters with social intelligence.
 In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.
 557 582–592, June 2024.
 - Ling-Hao Chen, Shunlin Lu, Ailing Zeng, Hao Zhang, Benyou Wang, Ruimao Zhang, and Lei Zhang. Motionllm: Understanding human behaviors from human motions and videos. *arxiv:2405.20340*, 2024a.
- Wei Chen, Lin Li, Yongqi Yang, Bin Wen, Fan Yang, Tingting Gao, Yu Wu, and Long Chen. CoMM:
 A coherent interleaved image-text dataset for multimodal understanding and generation. *arXiv* preprint arXiv:2406.10462, 2024b.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Mihai Fieraru, Mihai Zanfir, Elisabeta Oneata, Alin-Ionut Popa, Vlad Olaru, and Cristian Smin chisescu. Three-dimensional reconstruction of human interactions. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7214–7223, 2020.
- Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying
 Shan. SEED-X: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024.
- Anindita Ghosh, Rishabh Dabral, Vladislav Golyanik, Christian Theobalt, and Philipp Slusallek. Re Mos: Reactive 3d motion synthesis for two-person interactions. *arXiv preprint arXiv:2311.17057*, 2023.
- Purvi Goel, Kuan-Chieh Wang, C. Karen Liu, and Kayvon Fatahalian. Iterative motion editing with natural language. New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400705250. doi: 10.1145/3641519.3657447. URL https://doi.org/10.1145/3641519.3657447.
- 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 *Proc. of the European Conference on Computer Vision (ECCV)*, pp. 580–597. Springer, 2022.
- Chuan Guo, Yuxuan Mu, Muhammad Gohar Javed, Sen Wang, and Li Cheng. Momask: Generative masked modeling of 3d human motions. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1900–1910, 2024a.
- 591 Chuan Guo, Yuxuan Mu, Xinxin Zuo, Peng Dai, Youliang Yan, Juwei Lu, and Li Cheng. Generative
 592 human motion stylization in latent space. In *Proc. of the Twelfth International Conference on* 593 *Learning Representations (ICLR)*, 2024b. URL https://openreview.net/forum?id=
 daEqXJ0yZo.

- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *Proc. of the International Conference on Learning Representations (ICLR)*, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.
- Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. MotionGPT: Human motion as a foreign language. *Proc. of the Advances in Neural Information Processing Systems (NEURIPS)*, 36:20067–20079, 2023.
- Biao Jiang, Xin Chen, Chi Zhang, Fukun Yin, Zhuoyuan Li, Gang Yu, and Jiayuan Fan.
 MotionChain: Conversational motion controllers via multimodal prompts. *arXiv preprint arXiv:2404.01700*, 2024.
- Jihoon Kim, Jiseob Kim, and Sungjoon Choi. Flame: Free-form language-based motion synthesis
 & editing. In *Proc. of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 37, pp. 8255–8263, 2023.
- Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Autoregressive image
 generation using residual quantization. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11523–11532, 2022.
- Han Liang, Wenqian Zhang, Wenxuan Li, Jingyi Yu, and Lan Xu. Intergen: Diffusion-based multi human motion generation under complex interactions. *International Journal of Computer Vision*,
 pp. 1–21, 2024.
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/W04-1013.
- Jing Lin, Ailing Zeng, Shunlin Lu, Yuanhao Cai, Ruimao Zhang, Haoqian Wang, and Lei Zhang.
 Motion-X: A large-scale 3d expressive whole-body human motion dataset. *Proc. of the Advances in Neural Information Processing Systems (NEURIPS)*, 36, 2024.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Proc. of the Advances in neural information processing systems (NEURIPS)*, 36, 2024.
- Mingshuang Luo, Ruibing Hou, Hong Chang, Zimo Liu, Yaowei Wang, and Shiguang Shan. M3
 GPT: An advanced multimodal, multitask framework for motion comprehension and generation. arXiv preprint arXiv:2405.16273, 2024.
- Evonne Ng, Donglai Xiang, Hanbyul Joo, and Kristen Grauman. You2me: Inferring body pose in
 egocentric video via first and second person interactions. In *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9890–9900, 2020.
- 631 OpenAI. Hello gpt-40. 2024. URL https://openai.com/index/hello-gpt-40/.
- Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3D hands, face, and body from a
 single image. In *Proc.of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*,
 pp. 10975–10985, 2019.
- Mathis Petrovich, Michael J. Black, and Gül Varol. TMR: Text-to-motion retrieval using contrastive
 3D human motion synthesis. In *Proc. of the International Conference on Computer Vision (ICCV)*, 2023.
- Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. MAUVE: Measuring the gap between neural text and human text using divergence frontiers. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Proc. of the Advances in Neural Information Processing Systems (NEURIPS)*, 2021. URL https://openreview.net/forum?id=Tqx7nJp7PR.
- Yoni Shafir, Guy Tevet, Roy Kapon, and Amit Haim Bermano. Human motion diffusion as a generative prior. In *Proc. of the Twelfth International Conference on Learning Representations (ICLR)*, 2024. URL https://openreview.net/forum?id=dTpbEdN9kr.

664

665

666

677

681

684

688

689

690

691

692

693

694

- 648 Yu Shu, Siwei Dong, Guangyao Chen, Wenhao Huang, Ruihua Zhang, Daochen Shi, Qiqi Xiang, 649 and Yemin Shi. Llasm: Large language and speech model. arXiv preprint arXiv:2308.15930, 650 2023. 651
- Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, 652 and Chao Zhang. SALMONN: Towards generic hearing abilities for large language models. In 653 Proc. of the Twelfth International Conference on Learning Representations (ICLR), 2024. URL 654 https://openreview.net/forum?id=14rn7HpKVk. 655
- 656 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 658 2: Improving open language models at a practical size. arXiv preprint arXiv:2408.00118, 2024. 659
- Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. 660 Human motion diffusion model. In Proc. of the Eleventh International Conference on 661 Learning Representations (ICLR), 2023. URL https://openreview.net/forum?id= 662 SJ1kSy02jwu. 663
 - Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. Proc. of the Advances in neural information processing systems (NEURIPS), 30, 2017.
- 667 Yin Wang, Zhiying Leng, Frederick WB Li, Shun-Cheng Wu, and Xiaohui Liang. Fg-t2m: Finegrained text-driven human motion generation via diffusion model. In Proc. of the IEEE/CVF 668 International Conference on Computer Vision (CVPR), pp. 22035–22044, 2023. 669
- 670 Qi Wu, Yubo Zhao, Yifan Wang, Yu-Wing Tai, and Chi-Keung Tang. Motionllm: Multimodal 671 motion-language learning with large language models. arXiv preprint arXiv:2405.17013, 2024. 672
- 673 Liang Xu, Xintao Lv, Yichao Yan, Xin Jin, Shuwen Wu, Congsheng Xu, Yifan Liu, Yizhou Zhou, 674 Fengyun Rao, Xingdong Sheng, et al. Inter-x: Towards versatile human-human interaction anal-675 ysis. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition 676 (CVPR), pp. 22260–22271, 2024a.
- Liang Xu, Yizhou Zhou, Yichao Yan, Xin Jin, Wenhan Zhu, Fengyun Rao, Xiaokang Yang, and 678 Wenjun Zeng. ReGenNet: Towards human action-reaction synthesis. In Proc. of the IEEE/CVF 679 Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1759–1769, June 2024b. 680
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, 682 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. arXiv preprint 683 arXiv:2407.10671, 2024.
- Yifei Yin, Chen Guo, Manuel Kaufmann, Juan Jose Zarate, Jie Song, and Otmar Hilliges. Hi4d: 685 4d instance segmentation of close human interaction. In Proc. of the IEEE/CVF Conference on 686 Computer Vision and Pattern Recognition (CVPR), pp. 17016–17027, 2023. 687
 - Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Yong Zhang, Hongwei Zhao, Hongtao Lu, Xi Shen, and Ying Shan. Generating human motion from textual descriptions with discrete representations. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR), pp. 14730-14740, 2023.
 - Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2024a.
- 696 Mingyuan Zhang, Daisheng Jin, Chenyang Gu, Fangzhou Hong, Zhongang Cai, Jingfang Huang, 697 Chongzhi Zhang, Xinying Guo, Lei Yang, Ying He, et al. Large motion model for unified multi-698 modal motion generation. arXiv preprint arXiv:2404.01284, 2024b. 699
- Mingyuan Zhang, Huirong Li, Zhongang Cai, Jiawei Ren, Lei Yang, and Ziwei Liu. Finemogen: 700 Fine-grained spatio-temporal motion generation and editing. In Proc. of the Advances in Neural 701 Information Processing Systems (NEURIPS), volume 36, 2024c.

702 703	Yaqi Zhang, Di Huang, Bin Liu, Shixiang Tang, Yan Lu, Lu Chen, Lei Bai, Qi Chu, Nenghai Yu, and Wanli Ouyang. Motiongpt: Finetuned llms are general-purpose motion generators. In <i>Proc.</i>
704	of the AAAI Conference on Artificial Intelligence (AAAI), volume 38, pp. 7368–7376, 2024d.
705	Ziviang Zhou Vu Wan and Paousan Wang AvatarCDT: All in one framework for motion up
706	derstanding planning generation and beyond. In <i>Proceedings of the IEFE/CVF Conference on</i>
707	Computer Vision and Pattern Recognition (CVPR), pp. 1357–1366. June 2024.
708	
709	
710	
711	
712	
713	
714	
715	
717	
718	
719	
720	
721	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	