# Multi-Track Timeline Control for Text-Driven 3D Human Motion Generation

Anonymous CVPR submission

Paper ID 9



Figure 1. **Multi-track timeline control:** We introduce a new problem setting for text-driven motion synthesis, where the input consists of parallel tracks allowing simultaneous actions, as well as continuous temporal intervals enabling sequential actions. A long and complex motion can be generated (top) given the structured input of multiple simple textual descriptions, each corresponding to a temporal interval (bottom).

#### Abstract

Recent advances in generative modeling have led to promis-001 002 ing progress on synthesizing 3D human motion from text, with 003 methods that can generate character animations from short 004 prompts and specified durations. However, using a single text 005 prompt as input lacks the fine-grained control needed by animators, such as composing multiple actions and defining precise 006 007 durations for parts of the motion. To address this, we intro-008 duce the new problem of timeline control for text-driven motion 009 synthesis, which provides an intuitive, yet fine-grained, input interface for users. Instead of a single prompt, users can specify 010 011 a multi-track timeline of multiple prompts organized in temporal intervals that may overlap. This enables specifying the 012 013 exact timings of each action and composing multiple actions 014 in sequence or at overlapping intervals. To generate compos-015 ite animations from a multi-track timeline, we propose a new 016 test-time denoising method. This method can be integrated with any pre-trained motion diffusion model to synthesize realistic 017 motions that accurately reflect the timeline. At every step of 018 019 denoising, our method processes each timeline interval (text prompt) individually, subsequently aggregating the predictions 020 021 with consideration for the specific body parts engaged in each 022 action. Experimental comparisons and ablations validate that 023 our method produces realistic motions that respect the semantics 024 and timing of given text prompts.

## 1. Introduction

Motivated by applications in video games, entertainment, and virtual avatar creation, recent work has demonstrated substantial progress in learning to generate 3D human motion [27, 37, 44, 60]. Generating motions from text descriptions is of particular interest; it has the potential to democratize animation with a natural language interface that is intuitive for beginner and expert users alike. To this end, several methods have been proposed that synthesize reasonable character animations given a single text prompt and fixed duration as input [38, 53, 65].

While these methods are a promising first step towards faster 035 and more accessible animation interfaces, they lack the precise 036 control that is crucial for many animators. Consider the input 037 prompt (see Fig. 2d): "A human walks in a circle clockwise, 038 then sits, simultaneously raising their right hand towards the 039 end of the walk, the hand raising halts midway through the 040 sitting action." Due to a lack of representative training data, 041 prior work struggles with such complex text prompts [38, 53]. 042 Namely, the prompt includes *temporal* composition [4] where 043 multiple actions are performed in sequence (e.g., walking 044 then sitting), along with spatial composition [5] where several 045 actions are performed simultaneously with differing body parts 046 (e.g., walking *while* raising hand). Furthermore, such lengthy 047 prompts quickly become unwieldy for the user and, despite 048 their detailed descriptions, are still ambiguous with respect to 049 the timing and duration of the constituent actions. 050

025 026

027

028

029

030

031

032

033

071

072

073



Figure 2. **Text-driven motion synthesis tasks:** Our framework generalizes (a) traditional *text-to-motion synthesis* given one text and one duration, (b) *temporal composition* given a sequence of texts for non-overlapping intervals, and (c) *spatial composition* given a set of texts for a single interval. (d) *Multi-track timeline control* uses a set of texts for arbitrary intervals, allowing fine-grained control over the timings of several complex actions.

051 To improve controllability, we propose the new problem of 052 multi-track timeline control for text-driven 3D human motion synthesis. In this task, the user provides a structured and 053 054 intuitive timeline as input (Fig. 1), which contains several (potentially overlapping) temporal intervals. Each interval 055 056 corresponds to a precise textual description of a motion. As shown in Fig. 2d, the complex example prompt discussed 057 058 earlier becomes simple to specify within the timeline, and 059 allows animators to control the timing of each action. Such a timeline interface is already common in animation and video 060 editing software, and is analogous to control interfaces that 061 062 have recently emerged from the text-to-image community [64], e.g., image generation from a segmentation mask. 063

Multi-track timeline control for text-driven motion synthesis
is a generalization of several motion synthesis tasks, and
therefore brings many challenges. In particular, the multi-track
timeline input can achieve (see Fig. 2):

- *Text-to-motion synthesis* [18, 38] specifying a single interval
   (i.e., duration) with one textual description,
  - *Temporal composition* [4, 66] a sequence of textual descriptions corresponding to non-overlapping intervals,
  - Spatial (body-part) composition [5] a set of text prompts performed simultaneously with differing body parts.

074 Solving this task is difficult due to the lack of training data con-075 taining complex compositions and long durations. For example, a timeline-controlled model must handle the multi-track input 076 077 containing several prompts, rather than a single text description. 078 Moreover, the model must account for both spatial and temporal compositions to ensure seamless transitions, unlike prior work 079 080 that has addressed each of these individually. The timeline also 081 relaxes the assumption of a limited duration (<10 sec) made 082 by many recent text-to-motion approaches [11, 53, 65].

To address these challenges, we introduce a method for 083 084 Spatio-Temporal Motion Collage (STMC). Our method copes with the lack of appropriate training data by operating at test 085 time, leveraging a pre-trained motion diffusion model such 086 087 as off-the-shelf MDM [53] or MotionDiffuse [65]. At each denoising step, STMC first applies the diffusion model on 088 089 each text prompt in the timeline independently to predict a 090 denoised motion for the corresponding intervals. Our key

insight is to stitch together such independent generations 091 in both space and time before continuing to denoise. For 092 spatial compositions, automatic body part associations [5] 093 allow coherently concatenating predictions together. Score 094 arithmetic [66] is used to ensure smooth transitions for temporal 095 compositions. To further improve the performance of STMC, 096 we introduce MDM-SMPL, which makes several improvements 097 to prior motion diffusion models [53], including directly using 098 the SMPL [34] body representation. 099

The performance of STMC on timeline control for 100 text-driven motion synthesis is verified through comprehensive 101 comparisons and a perceptual user study. In summary, the 102 central contribution of this work consists of: (i) the new problem 103 of multi-track timeline control for text-driven 3D human motion 104 synthesis, and (ii) a novel test-time technique, STMC, that 105 effectively structures the denoising process to ensure faithful 106 execution of all prompts in a timeline. As a side contribution, 107 (iii) we upgrade MDM to directly support the SMPL body 108 representation instead of skeletons, and reduce runtime through 109 fewer denoising steps. Code will be released upon publication. 110

### 2. Related Work

Human motion synthesis. A large body of work in both vi-112 sion and graphics has been dedicated to generating 3D hu-113 man motions [70]. This generation process can be uncon-114 ditional [36, 56] or conditioned on actions [10, 17, 37], mu-115 sic [32, 50, 52, 57], speech [3, 69], goals [30, 51, 60], previous 116 motion [13, 15, 44, 62] (i.e., future motion prediction), sce-117 nes/objects [21, 31, 58, 59], and text [1, 2, 11, 16, 19, 29, 53, 118 65]. Technical approaches vary from early statistical models 119 [8, 15] to modern generative models like VAEs [20, 37, 38], 120 GANs [6, 12, 49, 61], normalizing flows [22, 57], and diffu-121 sion [11, 29, 30, 60, 68]. Our work is most related to recent 122 text-conditioned diffusion models [53, 65], however we solve 123 a new problem where the model is conditioned on a timeline 124 containing several text inputs instead of a single prompt. 125

Motion composition. Due to the lack of training data, a par-126 ticular challenge for action and text-conditioned motion gen-127 eration is to synthesize compositional motions. Several works 128 [4, 41, 66] focus on generating motions from a sequence of text 129 prompts and durations, i.e., *temporal* compositions. TEACH [4] 130 autoregressively generates one motion (per text prompt) at a 131 time, conditioning the next motion in the sequence with the 132 previous one. EMS [41] proposes a two-stage approach, by 133 first generating each action separately and then merging them 134 through a subsequent network. Diffusion models EDGE [54] 135 and PriorMDM [48] ensure consistency between adjacent mo-136 tions by enforcing temporal constraints at transitions. Our ap-137 proach to temporal composition is based on DiffCollage [66], 138 which stitches motions (or images) together throughout the de-139 noising process via score arithmetic at overlapping transitions. 140

Other work generates motions from a set of texts to be 141 executed at the same time, i.e., *spatial* (body-part) composition. 142

215

233

SINC [5] labels ground truth motion capture (mocap) sequences 143 with corresponding body parts by prompting GPT-3 [9]. 144 These labels are used to create a synthetic dataset of motions 145 146 stitched together from mocap sequences with compatible 147 body parts, thereby improving performance of VAE-based 3D motion generation methods [38] for spatial composition. 148 MotionDiffuse [65] proposes a noise interpolation method to 149 control different body part motions separately. Our approach, 150 151 STMC, takes inspiration from SINC [5] by using body part labels to stitch motions together during test-time denoising. 152 153 Overall, our problem of timeline-conditioned generation generalizes temporal and spatial composition, and STMC must 154 tackle both issues simultaneously, unlike most prior work. 155

Controllable motion diffusion. Following success in im-156 age [43, 46, 47], video [26], and 3D [33, 40, 63] domains, 157 158 diffusion has become a useful approach to generate highquality 3D human motions [3, 28, 54], especially from text 159 160 inputs [11, 14, 53, 65]. Some works focus on improving the controllability of motion diffusion models, e.g., by enabling 161 162 temporal [48, 66] and spatial [65] composition of text prompts. 163 Other controls such as following specific keyframe poses, joint trajectories, and waypoints have also been achieved using a mix 164 of test-time diffusion guidance [28, 30, 45], in-painting [48, 53], 165 and direct conditioning [60]. We focus on making text-to-166 motion generation more controllable by handling several text 167 168 prompts in a fine-grained timeline format through a composi-169 tional denoising process.

#### **3. Human Motion Synthesis from Timelines**

We first formulate the new problem setup of multi-track
timeline control (Sec. 3.1), then propose a motion denoising
strategy to handle timeline inputs (Sec. 3.2 and Sec. 3.3), and
finally summarize our improved diffusion model (Sec. 3.4).

#### **175 3.1. Timeline Control Problem Formulation**

176 **Inputs**. As illustrated in Fig. 1, the multi-track timeline enables 177 users to define multiple intervals, each linked to a natural language prompt describing the desired human motion. For the 178 *i*th prompt in the timeline, we represent its temporal interval as 179  $[a_i, b_i]$  and the corresponding prompt as  $C_i$ . The intervals are 180 181 arranged in a multi-track layout on the timeline, allowing for overlaps. Both the duration of each interval and of the overall 182 timeline are variable, and users can add an arbitrary number of 183 tracks (rows) to the timeline (although, in practice, a character 184 can most often perform a handful of actions simultaneously). 185

Outputs. The goal is to generate a 3D human motion that fol-186 187 lows all the text instructions at the specified intervals. A human motion  $\boldsymbol{x}$  lasting N timesteps is represented as a sequence of 188 pose vectors  $\boldsymbol{x} = (\boldsymbol{x}^1, ..., \boldsymbol{x}^N)$  with each pose  $\boldsymbol{x}^i \in \mathbb{R}^d$ . Several 189 recent works [53, 65] use the pose representation from Guo 190 et al. [18] with d=263, which contains root velocities along 191 192 with local joint positions, rotations, and velocities. Other pose 193 representations like SMPL [34] can also be used (see Sec. 3.4).

#### **3.2. Background: Motion Diffusion Models**

Our generation method (Sec. 3.3) leverages a pre-trained motion 195 diffusion model such as MDM [53] or MotionDiffuse [65] 196 trained on single text prompts, which we briefly review here. 197 These methods follow a denoising diffusion scheme and 198 synthesize animations through iterative denoising of a noisy 199 pose sequence. Given a clean motion  $x_0$ , a Gaussian diffusion 200 process is employed to corrupt the data to be approximately 201  $\mathcal{N}(\mathbf{0},\mathbf{I})$ . Each step of this process is given by: 202

$$q(\boldsymbol{x}_t | \boldsymbol{x}_{t-1}) = \mathcal{N}(\boldsymbol{x}_t; \sqrt{1 - \beta_t} \boldsymbol{x}_{t-1}, \beta_t \mathbf{I})$$
(1) 203

with  $\beta_t$  defined by the noise schedule. Note the denoising step 204 *t* is not to be confused with the temporal timestep *i*, which 205 indexes the sequence of poses in the motion. In practice, one 206 can make sampling  $x_t$  easier by using the reparameterization 207 trick  $x_t = \sqrt{\overline{\alpha}_t} x_0 + \sqrt{1 - \overline{\alpha}_t} \epsilon$ , where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \alpha_t = 1 - \beta_t$ , 208 and  $\overline{\alpha}_t = \prod_{s=0}^{t} \alpha_s$ . 209

Sampling from a diffusion model requires reversing this210process to recover a clean motion from random noise. While211 $q(x_{t-1}|x_t)$  is hard to compute, the probability conditioned on212 $x_0$  is tractable [25]:213

$$q(x_{t-1}|x_t,x_0) = \mathcal{N}(x_{t-1};\mu_t(x_t,x_0),\Sigma_t)$$
, (2) 214

where

CVPR 2024 Submission #9. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

$$\mu_t(\boldsymbol{x}_t, \boldsymbol{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \boldsymbol{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \boldsymbol{x}_0 \qquad (3) \qquad \textbf{216}$$

$$\boldsymbol{\Sigma}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \mathbf{I} . \tag{4}$$

Since  $x_t$  is known at sampling time, we approximate the reverse distribution by training a denoising model  $\hat{x}_{\theta}(x_t, t, C)$  219 to estimate  $x_0$ , where C is the text conditioning. This model is trained with the simplified loss function as in Ho et al. [25] 221 (i.e., without the t-dependent factor): 222

$$\mathcal{L} = \mathbb{E}_{\boldsymbol{\epsilon}, t, \boldsymbol{x}_0, C} \| \hat{\boldsymbol{x}}_{\boldsymbol{\theta}}(\boldsymbol{x}_t, t, C) - \boldsymbol{x}_0 \|_2^2$$
(5) 223

with  $x_0$  and C sampled from a dataset of motion-text pairs, 224 step t sampled uniformly, and noise  $\epsilon \sim \mathcal{N}(\mathbf{0},\mathbf{I})$  used to corrupt 225 the ground truth motion. To enable classifier-free guidance [24] 226 at sampling time, the text conditioning C is dropped with 227 some probability at each training iteration. At test time, the 228 sampling (reverse) process starts from random noise and 229 denoises iteratively for T steps to obtain a clean 3D human 230 motion. At each denoising step, the model is conditioned on 231 the single input text prompt (e.g., Fig. 2a). 232

## 3.3. STMC: Spatio-Temporal Motion Collage

STMC operates only at test time, enabling an off-the-shelf, pretrained denoising model to generate motion conditioned on a multi-track timeline. At *every* denoising step, our method takes as input the current noisy motion  $x_t$  encapsulating the entire 237

#### CVPR 2024 Submission #9. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 3. **Overview of STMC:** Before denoising, the multi-track timeline is first (**a**) partitioned into relevant body parts per text (using LLM-based labeling [5]) to create body part timelines, which are then (**b**) extended to overlap, leading to the transition intervals used for temporal stitching *per body part* with DiffCollage [66]. (**c**) At each denoising step, motions for each prompt are denoised independently before being combined based on the body-part timelines. The composite motion is re-noised by sampling  $x_{t-1}$  from  $\mathcal{N}(\mu_t(x_t, \hat{x}_0), \Sigma_t)$  (as in Eq. (2)) before being passed to the next step.

238 timeline and outputs a corresponding clean motion  $\hat{x}_0$ . As shown in Fig. 3c, STMC uses the denoising model to indepen-239 dently predict a clean motion crop corresponding to each of the 240 241 input text prompts. These predictions are stitched together spatially using body part annotations for each text prompt (Fig. 3a), 242 243 and stitched in time to ensure the clean motion smoothly spans the entire timeline (Fig. 3b). This final composite motion be-244 245 comes the output of the current step  $\hat{x}_0$ , which is used to sample  $x_{t-1}$  with Eq. (2) and continue the denoising process. To enable 246 body part stitching, STMC assumes the denoiser operates on 247 248 explicit poses [53, 65], rather than in a latent space [11].

249 Motion cropping and denoising. The input  $x_t$  at denoising 250 step t extends over the duration of the entire timeline. As shown in Fig. 3c, we first temporally split the input into motion "crops" 251 to separately denoise each text prompt. For each interval  $[a_j, b_j]$ , 252 the motion is cropped in time to  $x_t^{a_j:b_j} = x_t[a_j:b_j]$ . The crop, along with the text prompt  $C_j$ , is given to the denoising model 253 254 to predict a corresponding clean motion crop  $\hat{x}_0^{a_j:b_j}$ . Denoising 255 each text prompt independently gives high-quality motion from 256 pre-trained models since each prompt typically contains a single 257 action and the interval duration is reasonably short (<10 sec). 258

259 Two or more text prompts in the timeline may overlap in time, meaning the predicted clean crops will also overlap. As 260 a concrete example, suppose the crops for "walking in a circle" 261 262 and "raising right hand" are overlapping, as in Fig. 3. In this case, it is not clear which of the two generated motions should 263 264 be assigned to the overlapping region. To construct a motion that 265 matches both prompts, we need the leg motion from "walking in a circle" and the right arm motion from "raising right hand". 266 We therefore stitch together outputs from overlapping prompts 267 268 based on automatically labeled body parts, as detailed next.

Spatial (body-part) stitching. Spatial stitching follows
 SINC [5], which proposed to combine compatible body-part

motions from mocap sequences through simple concatenation. 271 While SINC applies stitching only once, STMC does so at every 272 step of denoising, encouraging a more coherent composition 273 of movements by allowing the denoiser to correct any artifacts. 274 This is possible because the denoiser outputs explicit human 275 poses (i.e., we know which indices correspond to arms, legs, 276 etc. within the pose vector), so we can extract body-part mo-277 tions from separate crops and spatially combine them to obtain 278 a composite motion. To achieve this, we first pre-process the 279 input timeline to assign a text prompt to each body part at every 280 timestep, thereby creating a separate motion timeline for every 281 body part (see Fig. 3a): left arm, right arm, torso, legs and head. 282

As shown in Fig. 3a, each text prompt in the multi-track 283 timeline is first annotated with a set of body parts involved in the 284 motion. This can be done automatically by querying GPT-3 [9] 285 as in SINC, or directly given by the user for additional creative 286 control. Then, each text prompt is assigned to its annotated 287 body parts within the corresponding time interval, which 288 assumes that body parts at overlapping intervals are compatible 289 (e.g., if a prompt is annotated with "legs", then no other prompt 290 should involve legs throughout its entire interval). To fill in the 291 remainder of the body-part timelines where body parts have not 292 been annotated to a text prompt, heuristics similar to SINC are 293 used. Please see the Appendix B and the Fig. A.1 for full details. 294 Finally, during the denoising step (Fig. 3c), each crop  $x_t^{a_j:b_j}$  is 295 split into separated body-part motions and concatenated together 296 as specified by the body-part timelines to obtain the output  $\hat{x}_0$ . 297 Temporal stitching. Because the motion crops are denoised 298 independently, simple temporal concatenation of body-part mo-299 tions from different text prompts will cause abrupt transitions. 300 To mitigate these potential artifacts, we apply DiffCollage [66] 301 to each body-part motion. As shown in Fig. 3b, instead of 302 directly denoising  $x_t^{a_j:b_j}$  for each text prompt, we denoise an 303 expanded time interval  $[a_i - l, b_i + l]$ , where l is the desired over-304

305 lap length between adjacent motion crops (e.g., fixed to 0.25sec). Concretely, for the temporal transition between prompts j306 and k, we have  $\hat{x}_0^{a_j-l:b_j+l}$  and  $\hat{x}_0^{a_k-l:b_k+l}$  after denoising. We 307 then unconditionally denoise a small (0.5 sec) crop of motion 308 centered on the overlap between j and k to obtain  $\hat{x}_0^{\text{uncond}}$ . The 309 final predicted motion spanning intervals  $\boldsymbol{j}$  and  $\boldsymbol{k}$  is computed as 310  $\hat{x}_0 = \hat{x}_0^{a_j - l:b_j + l} + \hat{x}_0^{a_k - l:b_k + l} - \hat{x}_0^{\text{uncould}}$ , as depicted in Fig. 3c. 311 This equation derives from a factor graph representation of the 312 problem, as detailed in DiffCollage [66]. 313

## 314 3.4. SMPL Support for Motion Diffusion Model

315 While STMC works well with off-the-shelf models [53, 65] 316 (see Sec. 4), we propose several practical improvements to 317 MDM [53] to further enhance results. Our model, MDM-SMPL, 318 employs a skinned human body SMPL [34]: we use SMPL pose 319 parameters instead of the joint rotation features in the original pose representation of Guo et al. [18]. In contrast to models that 320 use the joint position outputs from the pose representation of 321 322 [18], this SMPL-based representation avoids the need for expensive test-time optimization [7, 71] to fit the generated motion 323 on a SMPL body. Moreover, the local joint rotations in SMPL, 324 325 which are relative to parents in the kinematic tree, are more amenable to body-part stitching than root-relative joint positions. 326 327 This is because any change to a joint rotation is propagated to all 328 children in the kinematic tree, unlike root-relative joint positions which may not be coherent when simply concatenated together. 329 Additional improvements include lowering the number of 330 diffusion steps to T=100 from 1000 to substantially speed up 331 332 sampling, and various architectural changes. We provide more details on MDM-SMPL in Appendix D. 333

## **334 4. Experiments**

We first present the data (Sec. 4.1) and the evaluation protocols (Sec. 4.2) used in the experiments. We then show comparisons with baselines quantitatively (Sec. 4.3) and with a perceptual study (Sec. 4.4), followed by qualitative results (Sec. 4.5). We conclude with a discussion of the limitations (Sec. 4.6).

#### **340 4.1. Datasets**

341 **HumanML3D** [18] is a text-motion dataset that provides textual descriptions for a subset of the AMASS [35] and Human-342 343 Act12 [17] motion capture datasets. It consists of 44970 text annotations for 14616 motions. This dataset is used to train 344 all diffusion models used in our experiments. For MDM [53] 345 and MotionDiffuse [65], we use publicly available models pre-346 347 trained on the released version of HumanML3D with the original motion representation from Guo et al. [18]. Consequently, 348 these methods require test-time optimization to obtain SMPL 349 350 pose parameter outputs. For training our MDM-SMPL diffusion model, which is designed to directly generate SMPL pose 351 352 parameters, we re-process the dataset and exclude the Human-353 Act12 subset as SMPL poses are not available for this dataset.

Multi-track timeline (MTT) dataset. To properly evaluate our 354 new task, we introduce a new challenging dataset of 500 multi-355 track timelines. Each timeline in the dataset is automatically 356 constructed and contains three prompts on a two-track timeline 357 (e.g., Fig. 2d). To construct these timelines, we first manually 358 collect a set of 60 texts covering a diverse set of "atomic" ac-359 tions (e.g., "punch with the right hand", "jump forward", "run 360 backwards", see Appendix C for the full list), and annotate the 361 involved body parts for each text. To serve as ground truth for 362 computing evaluation metrics (Sec. 4.2), we also select motion 363 samples from AMASS that correspond to each text. Based on 364 the atomic texts, we automatically generate timelines containing 365 three prompts and two tracks (rows). For each timeline, the first 366 track is filled with two consecutive prompts sampled from the 367 set of texts and given randomized durations. A third random 368 text with complementary body-part annotations is then placed 369 in the second track at a random location in time. 370

The main reasons for restricting the evaluation to three 371 prompts are (i) to keep the cognitive load for users low in the 372 perceptual study, subsequently increasing the reliability of the 373 results, and (ii) to construct a minimal setup where we can fairly 374 compare against baselines in a controlled setting, eliminating 375 confounding factors such as the number of prompts. Though 376 these timelines contain only three prompts, they already pose 377 a significant challenge (see Sec. 4.3). Examples of timelines 378 in the dataset are provided in Fig. A.2 and qualitative results 379 beyond three prompts can be found in the supplementary video. 380

### 4.2. Evaluation Metrics

Given the novelty of the task, identifying relevant metrics to 382 evaluate different methods is crucial. Instead of relying on a 383 single metric, we disentangle the evaluation of semantic cor-384 rectness (how faithful individual motion crops are to the textual 385 descriptions) from that of realism (e.g., temporal smoothness). 386 Semantic metrics. Firstly, we evaluate the alignment between 387 the generated motion and the text description within the speci-388 fied intervals on the timeline, which we term "per-crop semantic 389 correctness". To assess this, we utilize the recent text-to-motion 390 retrieval model TMR [39]. Similar to how CLIP [42] functions 391 for images and texts, TMR provides a joint embedding space 392 that can be used to determine the similarity between a text and 393 motion. Using TMR, we encode each atomic text prompt and 394 corresponding motion from our MTT dataset to obtain ground 395 truth text and motion embeddings, respectively. Each generated 396 motion crop is also embedded and the TMR-Score, a measure of 397 cosine similarity ranging from 0 to 1, is calculated between the 398 generated motion embedding and the ground truth. We report 399 both motion-to-text similarity by comparing against the ground 400 truth text embedding (TMR-Score M2T) and motion-to-motion 401 similarity against the ground truth motion embedding (TMR-402 Score M2M). Such embedding similarity measures are akin to 403 BERT-Score [67] for text-text, CLIP-Score [23] for image-text, 404 and more recently TEMOS-Score [4] for motion-motion similar-405

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

ity. Since TMR is trained contrastively, its retrieval performance
is better than TEMOS [38] which only trains with positive pairs,
leading to our decision to instead use TMR-Score. Moreover,
its embedding space is optimized with cosine similarity, making
the values potentially more calibrated across samples.

411 Ideally, the *TMR-Score M2T* between a generated motion 412 crop and the corresponding input text prompt should surpass 413 those of other texts. Hence, we also measure motion-to-text 414 retrieval metrics (as in [18]) including the frequency of the 415 correct text prompt being in the top-1 (R@I) and top-3 (R@3) 416 retrieved texts from the entire set of atomic texts.

417 **Realism metrics**. Secondly, we evaluate the realism of the gen-418 erated motions, which includes transitioning smoothly between actions. While the Frechet Inception Distance (FID) between 419 420 generated and ground truth motion in a learned feature space 421 (e.g., TMR) is a common metric for quality, the embedding 422 space of TMR is not trained on motions that are longer than 10 423 sec, and may therefore be unreliable for longer motions. Hence, 424 we follow DiffCollage [66] and compute the FID+ to evaluate 425 transitions. The FID+ metric measures FID based on 5 random 5-second motion crops from each timeline-conditioned motion 426 generation. Following TEACH [4], we also measure the tran-427 428 sition distance as the Euclidean distance (in cm) between the 429 poses in two consecutive frames around the transition time. We 430 choose to compute this distance in the local coordinate system of the body to more effectively capture transitions for individual 431 body parts, rather than being dominated by global motion. This 432 metric is sensitive to abrupt pose changes, and a motion should 433 434 not have high transition distance to remain realistic.

435 Perceptual study. Since no quantitative metric can fully cap-436 ture the subtleties of human motion, we also conduct perceptual studies, where human raters on Amazon Mechanical Turk judge 437 438 the quality of the generated motions [55]. To compare two 439 generation methods, raters are presented with two videos of 440 generated motions side-by-side rendered on a skeleton. The 441 multi-track timeline is also visible with an animated bar that 442 progresses along the timeline as the videos play. Users are 443 asked which motion is *more realistic* and which one is *better at* following the text in the timeline; they may choose one of the 444 two motions or mark "no preference". The studies presented in 445 Sec. 4.3 are performed on a set of 100 motions with multiple 446 447 raters judging each pair. The preference for each video is deter-448 mined by a majority vote from all raters. Responses are filtered for quality by using three "warmup" questions at the start of 449 each 15-question survey along with two "honeypot" examples 450 451 with objectively correct answers. The honeypot examples test a 452 rater's understanding of the task: one example shows a motion with obviously severe limb stretching (realism understanding 453 test) and the other displays a motion generated from a different 454 timeline than the one displayed (timeline understanding test). If 455 456 a rater fails to answer either of these questions correctly, all of 457 their responses are discarded.

#### **4.3.** Quantitative Comparison with Baselines

We apply our STMC test-time approach on the pretrained 459 diffusion models of MotionDiffuse [65], MDM [53], and 460 MDM-SMPL (ours). For each denoiser, we establish several 461 strong baselines by repurposing existing methods to the 462 timeline-conditioned generation task for comparison. Results 463 are shown in Tab. 1. Next to each method, the table indicates 464 how many tracks the input timelines have (#tracks) and how 465 many text prompts can be contained in a track (#crops). Next, 466 we introduce each baseline and analyze results. 467

Single-text input [53, 65] baseline. The simplest approach 468 to condition motion diffusion on a timeline is to convert the 469 timeline into a single text description, which aligns with the 470 model's training input format (e.g., Fig. 2a). Given that our 471 timeline dataset is consistently comprised of three motions (A, 472 B, and C), we formulate single-text prompts as follows: "A and 473 then B while C". While timing information can be included 474 in the prompt, e.g., "A for 4 seconds", this is out-of-distribution 475 for models trained on HumanML3D, leading to worse results. 476 This method parallels the baseline strategies of SINC [5] for 477 spatial composition and TEACH [4] for temporal composition. 478

As shown for each denoiser in Tab. 1, this approach is ineffective for both semantic correctness metrics and realism. Since these models cannot generate motions longer than 10 sec and there is no timing information in the prompt, for this experiment, outputs are limited to a maximum duration of 10 sec and semantic correctness metrics are reported over the entire duration of the motion rather than per-crop. The poor performance is a result of the models not being trained on the types of complex compositional prompts that result from collapsing the timeline to a single text description.

**DiffCollage [66] baseline**. Instead of converting the multi-track timeline into a single prompt, one can collapse it into a single track timeline containing a series of consecutive text prompts, i.e., transform the problem to be one of temporal composition. DiffCollage can then be used to temporally compose the sequence of actions. For example, the timeline in Fig. 2d would be split into ["walking in a circle," "walking in a circle while raising the right hand," "sitting down while raising the right hand," the timeline, this splitting preserves the timings (*#crops*) in the timeline.

While the DiffCollage baseline generally produces smooth 499 transitions and reasonable FID scores, the semantic accuracy 500 is consistently worse than STMC. This is due to the complex 501 spatial compositions within the prompts after collapsing the time-502 line into a single track, which models trained on HumanML3D 503 struggle with. In contrast, STMC uses body-part stitching 504 throughout denoising to compose actions from simpler prompts. 505 SINC [5] baseline. Rather than performing body-part stitching 506 iteratively at every denoising step, an alternative approach is 507 to stitch body motions together only once after all crops have 508 finished the entire denoising process. This is most similar to 509 SINC and forms the basis for two baselines that accept the full 510

Method	Input type		Per-crop semantic correctness				Realism	
	#tracks	#crops	R@1↑	R@3↑	TMR-S M2T	Score↑ M2M	FID $\downarrow$	Transition distance $\downarrow$
Ground truth	-	-	55.0	73.3	0.748	1.000	0.000	1.5
MotionDiffuse [65]	Single	Single	10.9	21.3	0.558	0.546	0.621	1.9
DiffCollage	Single	Multi	22.6	43.3	0.633	0.612	0.532	4.6
SINC w/o Lerp	Multi	Multi	23.8	45.9	0.656	0.630	0.554	<u>3.8</u>
SINC w/ Lerp	//	//	24.9	46.7	0.663	0.632	0.552	1.0
STMC (ours)	//	"	24.8	46.7	0.660	0.632	0.531	1.5
MDM [53]	Single	Single	9.5	19.7	0.556	0.549	0.666	2.5
DiffCollage	Single	Multi	24.9	42.3	0.636	0.623	0.600	2.2
SINC w/o Lerp	Multi	Multi	21.5	41.8	0.629	0.626	0.638	<u>10.2</u>
SINC w/ Lerp	//	//	23.3	43.1	0.634	0.628	0.630	2.8
STMC (ours)	"	"	25.1	46.0	0.641	0.633	0.606	2.4
MDM-SMPL	Single	Single	12.1	23.5	0.573	0.578	0.484	1.8
DiffCollage	Single	Multi	29.1	49.7	0.675	0.656	0.446	1.2
SINC w/o Lerp	Multi	Multi	32.3	50.5	0.676	0.667	0.463	4.2
SINC w/ Lerp	//	//	31.8	51.0	0.679	0.668	0.457	1.2
STMC (ours)	//	//	30.5	50.9	0.675	0.665	0.459	0.9

Table 1. **Quantitative baseline comparison**: Our method STMC is compared to several strong baselines when using three different denoising models. The single-text and DiffCollage baselines struggle to handle complex compositional prompts that results from collapsing the timeline down to a single track. The SINC baselines produce reasonable semantic accuracy by denoising prompts independently as in STMC, but cause abrupt or unnatural transitions with higher transition distance (underlined) or FID.



Figure 4. **Perception study results:** Our STMC method is preferred over baselines by human raters for both motion realism and semantic accuracy. (Left) Comparison against the strong SINC with Lerp baseline. (Right) Comparison against the DiffCollage baseline. MDM [53] is used as the denoiser in these experiments.

multi-track timeline as input, similar to STMC.

512 SINC w/o Lerp concatenates body part motions at the end of denoising without considering temporal transitions. As 513 a result, transitions tend to be abrupt as evidenced by high 514 transition distances in Tab. 1 and occasional "teleporting" limbs 515 516 in qualitative results. To mitigate this, SINC w/ Lerp employs linear interpolation (lerp) at transitions for smoother results, 517 518 similar to the approach in TEACH [4]. Though this leads 519 to smoothness at transitions, FID scores tend to be slightly higher than STMC. The cause is obvious qualitatively, where 520 521 the generated motion often appears mechanical and unnatural, 522 sometimes resulting in foot sliding. Despite issues with motion 523 quality, these SINC baselines effectively capture the semantics 524 of each motion crop since crops are denoised independently.

525 Analysis of the results. Our method STMC consistently per-

forms effectively across both semantic and realism metrics, 526 unlike baselines that tend to sacrifice performance in one cate-527 gory for the other. For example, DiffCollage achieves the best 528 FID using MDM, but its inability to handle spatial compositions 529 results in worse semantics than STMC across all models. Addi-530 tionally, SINC baselines perform best in terms of semantics for 531 MotionDiffuse and MDM-SMPL, but result in abrupt or unnatu-532 ral transitions with FID or transition distance that is often higher 533 than STMC. Such transitions are also readily apparent in quali-534 tative results (see supplementary video). It is also notable that 535 using MDM-SMPL with STMC performs on par with MDM 536 and MotionDiffuse, while enabling direct SMPL output and 537 significantly reducing (by  $10 \times$ ) the number of diffusion steps. 538 Fewer steps, combined with pre-computing text embeddings, 539 enable sampling MDM-SMPL in less than 5 seconds on average. 540 This is a substantial improvement over MDM, which takes 4 541 minutes to generate motions followed by 8 min of optimization 542 to obtain SMPL poses, on average. 543

While the performance of STMC is promising, the semantic544metrics for ground truth motions indicate room for improvement.545As discussed in Sec. 4.6, STMC is currently limited by the546pre-trained diffusion model that it leverages for each motion547crop; we expect improvements in these models to also boost548STMC. An additional experiment on varying the overlap length549for temporal stitching can be found in Appendix E.550

#### 4.4. Perceptual Study

We perform two separate user studies to compare STMC to 552 SINC with Lerp and DiffCollage when using MDM. Fig. 4 553

#### CVPR 2024 Submission #9. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 5. **Qualitative results:** We visualize the results of STMC with MDM-SMPL on several input timelines and color the bodies depending on their location in the timeline. We see that STMC is capable of generating realistic motions, which capture the semantics of the given text prompts with the desired timing and duration. In (a) and (c), STMC generates motions that precisely follow the instructions, controlling a single arm while still performing another action. The accurate timing of intervals is demonstrated in (b) where the arms are still up in the air when transitioning from "walking" to "jumping", which is difficult to achieve with alternative methods. In (c) and (d), we observe that STMC is capable of generating compositions that were not present in the ground truth data, such as "walking backwards while eating" or "walking while playing violin".

shows results of both studies, measuring human preference for 554 motion realism and semantic accuracy. On the left, STMC is 555 556 preferred or similar to SINC 66% of the time for realism and 62% of the time for semantic accuracy, with 4.2 raters judging 557 558 each video on average after filtering bad responses. Compared to DiffCollage on the right, our method is preferred or similar 559 560 68% of the time for realism and 70% for semantic accuracy, with 561 2.8 raters judging each video after filtering. This demonstrates 562 that STMC improves the motion in ways that are discernible by 563 humans but may not be fully captured in quantitative metrics.

## **564 4.5. Qualitative Results**

We visualize motions generated by STMC with MDM-SMPL 565 566 in Figure 5, given multi-track timelines as input from our MTT dataset. The coloring follows the input text, prioritizing the 567 newest prompt when there is an overlap across tracks. These 568 results show that STMC is capable of generating realistic 569 motions for complex multi-prompt timelines, which follow 570 the timing and duration of the given intervals. Please see the 571 572 caption for full analysis of these examples, and we refer to 573 the supplementary video for additional qualitative results and comparison to generated motions from baseline methods. 574

### 575 4.6. Limitations

576 While STMC expands the capabilities of pre-trained motion 577 diffusion models to take a multi-track timeline as input, it is also limited by the models that it relies on. For example, 578 our proposed body-part stitching process produces spatially 579 composed motions throughout denoising that the off-the-shelf 580 models are not trained to robustly handle. One potential 581 direction to ameliorate this is a more sophisticated stitching 582 "schedule" where body parts are not combined until later in the 583 denoising process instead of at every step. STMC also inherits 584 the limitations of SINC, e.g., restricting overlapping motions 585 to have compatible body part combinations. 586

## 5. Conclusion

587

In this work, we proposed the new problem of multi-track time-588 line control for text-driven 3D human motion generation. The 589 timeline input gives users fine-grained control over the timing 590 and duration of actions, while still maintaining the simplicity 591 of natural language. We tackled this challenging problem 592 using a new test-time denoising process called spatio-temporal 593 motion collage (STMC), which enables pre-trained diffusion 594 models to handle the spatial and temporal compositions present 595 in timelines. Finally, extensive quantitative and qualitative 596 evaluation demonstrated the advantage of STMC over strong 597 baseline methods and its ability to generate realistic motions 598 that are faithful to a multi-track timeline from the user. 599

605

607

608

609

610

611

612

613

614

618

622

623

624

625

626

627

628

629

- 601 [1] Hyemin Ahn, Timothy Ha, Yunho Choi, Hwiyeon Yoo, and 602 Songhwai Oh. Text2Action: Generative adversarial synthesis 603 from language to action. In International Conference on Robotics 604 and Automation (ICRA), 2018. 2
- [2] Chaitanya Ahuja and Louis-Philippe Morency. Language2Pose: 606 Natural language grounded pose forecasting. In International Conference on 3D Vision (3DV), 2019. 2
  - [3] Simon Alexanderson, Rajmund Nagy, Jonas Beskow, and Gustav Eje Henter. Listen, denoise, action! audio-driven motion synthesis with diffusion models. ACM Transactions on Graphics (TOG), 2023. 2, 3
    - [4] Nikos Athanasiou, Mathis Petrovich, Michael J. Black, and Gül Varol. TEACH: Temporal action composition for 3D humans. In International Conference on 3D Vision (3DV), 2022. 1, 2, 5, 6, 7
- 615 [5] Nikos Athanasiou, Mathis Petrovich, Michael J. Black, and Gül Varol. SINC: Spatial composition of 3D human motions for 616 617 simultaneous action generation. In International Conference on Computer Vision (ICCV), 2023. 1, 2, 3, 4, 6
- [6] Emad Barsoum, John Kender, and Zicheng Liu. HP-GAN: Prob-619 620 abilistic 3D human motion prediction via GAN. In Computer Vision and Pattern Recognition Workshops (CVPRW), 2018. 2 621
  - [7] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J. Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In European Conference on Computer Vision (ECCV), 2016. 5
  - [8] Richard Bowden. Learning statistical models of human motion. In Computer Vision and Pattern Recognition (CVPR), Workshop on Human Modeling, Analysis and Synthesis, 2000. 2
- 630 [9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, 631 Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav 632 Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel 633 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, 634 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, 635 Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam 636 637 McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 638 Language models are few-shot learners. In Neural Information Processing Systems (NeurIPS), 2020. 3, 4 639
- 640 [10] Pablo Cervantes, Yusuke Sekikawa, Ikuro Sato, and Koichi 641 Shinoda. Implicit neural representations for variable length 642 human motion generation. In European Conference on Computer 643 Vision (ECCV), 2022. 2
- 644 [11] Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao 645 Chen, Jingyi Yu, and Gang Yu. Executing your commands via 646 motion diffusion in latent space. In Computer Vision and Pattern 647 Recognition (CVPR), 2023. 2, 3, 4
- 648 [12] B. Chopin, N. Otberdout, M. Daoudi, and A. Bartolo. Human 649 motion prediction using manifold-aware wasserstein gan. In 650 International Conference on Automatic Face and Gesture 651 Recognition, 2021. 2
- 652 [13] Enric Corona, Albert Pumarola, Guillem Alenyà, and Francesc 653 Moreno-Noguer. Context-aware human motion prediction. In 654 Computer Vision and Pattern Recognition (CVPR), 2020. 2
- 655 [14] Rishabh Dabral, Muhammad Hamza Mughal, Vladislav 656 Golyanik, and Christian Theobalt. MoFusion: A framework for

denoising-diffusion-based motion synthesis. In Computer Vision and Pattern Recognition (CVPR), 2023. 3

- [15] Aphrodite Galata, Neil Johnson, and David Hogg. Learning variable length markov models of behaviour. Computer Vision and Image Understanding (CVIU), 2001. 2
- Anindita Ghosh, Noshaba Cheema, Cennet Oguz, Christian [16] Theobalt, and Philipp Slusallek. Synthesis of compositional animations from textual descriptions. In International Conference on Computer Vision (ICCV), 2021. 2
- [17] Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun Gong, and Li Cheng. Action2Motion: Conditioned generation of 3D human motions. In ACM International Conference on Multimedia (ACMMM), 2020. 2, 5
- [18] 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 Computer Vision and Pattern Recognition (CVPR), 2022. 2, 3, 5, 6
- [19] 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 European Conference on Computer Vision (ECCV), 2022. 2
- [20] Ikhsanul Habibie, Daniel Holden, Jonathan Schwarz, Joe Yearsley, and Taku Komura. A recurrent variational autoencoder for human motion synthesis. In British Machine Vision Conference (BMVC), 2017. 2
- [21] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael J Black. Stochastic scene-aware motion prediction. In Computer Vision and Pattern Recognition (CVPR), 2021. 2
- [22] Gustav Eje Henter, Simon Alexanderson, and Jonas Beskow. MoGlow: Probabilistic and controllable motion synthesis using normalising flows. ACM Transactions on Graphics (TOG), 2020.
- [23] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: A reference-free evaluation metric for image captioning. In EMNLP, 2021. 5
- [24] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. arXiv:2207.12598, 2022. 3
- [25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Neural Information Processing Systems (NeurIPS), 2020. 3
- [26] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. arXiv:2204.03458, 2022. 3
- [27] Daniel Holden, Jun Saito, and Taku Komura. A deep learning framework for character motion synthesis and editing. ACM Transactions on Graphics (TOG), 2016. 1
- [28] Siyuan Huang, Zan Wang, Puhao Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song-Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023. 3
- [29] Peng Jin, Yang Wu, Yanbo Fan, Zhongqian Sun, Yang Wei, 709 and Li Yuan. Act as you wish: Fine-grained control of motion 710 diffusion model with hierarchical semantic graphs. In Neural 711 Information Processing Systems (NeurIPS), 2023. 2 712

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

777

778

779

780

781

782

783

784

785 786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818 819

820

821

822

823

824

825

826

827

- [30] Korrawe Karunratanakul, Konpat Preechakul, Supasorn
  Suwajanakorn, and Siyu Tang. Gmd: Controllable human
  motion synthesis via guided diffusion models. In *International Conference on Computer Vision (ICCV)*, 2023. 2, 3
- [31] Nilesh Kulkarni, Davis Rempe, Kyle Genova, Abhijit Kundu,
  Justin Johnson, David Fouhey, and Leonidas Guibas. NIFTY:
  Neural object interaction fields for guided human motion
  synthesis. *arXiv:2307.07511*, 2023. 2
- [32] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa.
  AI choreographer: Music conditioned 3D dance generation
  with AIST++. In *International Conference on Computer Vision*(ICCV), 2021. 2
- [33] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa,
  Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu
  Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3D
  content creation. In *Computer Vision and Pattern Recognition*(*CVPR*), 2023. 3
- [34] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard PonsMoll, and Michael J. Black. SMPL: A skinned multi-person linear
  model. *ACM Transactions on Graphics (TOG)*, 2015. 2, 3, 5
- [35] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard
  Pons-Moll, and Michael J. Black. AMASS: Archive of motion
  capture as surface shapes. In *International Conference on Computer Vision (ICCV)*, 2019. 5
- [36] Dirk Ormoneit, Michael J. Black, Trevor Hastie, and Hedvig
   Kjellström. Representing cyclic human motion using functional
   analysis. *Image and Vision Computing*, 2005. 2
- [37] Mathis Petrovich, Michael J. Black, and Gül Varol. Actionconditioned 3D human motion synthesis with transformer VAE.
  In *International Conference on Computer Vision (ICCV)*, 2021.
  1, 2
- [38] Mathis Petrovich, Michael J. Black, and Gül Varol. TEMOS: Generating diverse human motions from textual descriptions. In *European Conference on Computer Vision (ECCV)*, 2022. 1, 2, 3, 6
- [39] Mathis Petrovich, Michael J. Black, and Gül Varol. TMR: Textto-motion retrieval using contrastive 3D human motion synthesis. In *International Conference on Computer Vision (ICCV)*, 2023. 5
- [40] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall.
  Dreamfusion: Text-to-3D using 2D diffusion. In *International Conference on Learning Representations (ICLR)*, 2022. 3
- [41] Yijun Qian, Jack Urbanek, Alexander G. Hauptmann, and
  Jungdam Won. Breaking the limits of text-conditioned 3D
  motion synthesis with elaborative descriptions. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [42] 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 (ICML)*, 2021. 5
- [43] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray,
  Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever.
  Zero-shot text-to-image generation. In *International Conference on Machine Learning (ICML)*, 2021. 3
- [44] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang,
  Srinath Sridhar, and Leonidas J. Guibas. HuMoR: 3D human
  motion model for robust pose estimation. In *International Conference on Computer Vision (ICCV)*, 2021. 1, 2

- [45] Davis Rempe, Zhengyi Luo, Xue Bin Peng, Ye Yuan, Kris Kitani, Karsten Kreis, Sanja Fidler, and Or Litany. Trace and pace: Controllable pedestrian animation via guided trajectory diffusion. In *Conference on Computer Vision and Pattern Recognition* (*CVPR*), 2023. 3
  775
- [46] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [47] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *Neural Information Processing Systems (NeurIPS)*, 2022. 3
- [48] Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermano. Human motion diffusion as a generative prior. *arXiv:2303.01418*, 2023. 2, 3
- [49] Ayumi Shiobara and Makoto Murakami. Human motion generation using wasserstein GAN. In *International Conference* on Digital Signal Processing (ICDSP), 2021. 2
- [50] Jiangxin Sun, Chunyu Wang, Huang Hu, Hanjiang Lai, Zhi Jin, and Jian-Fang Hu. You never stop dancing: Non-freezing dance generation via bank-constrained manifold projection. In *Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [51] Omid Taheri, Vasileios Choutas, Michael J. Black, and Dimitrios Tzionas. GOAL: Generating 4D whole-body motion for handobject grasping. In *Computer Vision and Pattern Recognition* (CVPR), 2022. 2
- [52] Taoran Tang, Jia Jia, and Hanyang Mao. Dance with melody: An lstm-autoencoder approach to music-oriented dance synthesis. In ACM International Conference on Multimedia (ACMMM), 2018. 2
- [53] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *International Conference on Learning Representations* (*ICLR*), 2023. 1, 2, 3, 4, 5, 6, 7
- [54] Jonathan Tseng, Rodrigo Castellon, and C Karen Liu. EDGE: Editable dance generation from music. In *Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 3
- [55] Amazon Mechanical Turk. Amazon mechanical turk, 2023. 6
- [56] Raquel Urtasun, David J. Fleet, and Neil D. Lawrence. Modeling human locomotion with topologically constrained latent variable models. In *Human Motion – Understanding, Modeling, Capture and Animation*, 2007. 2
- [57] Guillermo Valle-Pérez, Gustav Eje Henter, Jonas Beskow, André Holzapfel, Pierre-Yves Oudeyer, and Simon Alexanderson. Transflower: probabilistic autoregressive dance generation with multimodal attention. ACM Transactions on Graphics (TOG), 2021. 2
- [58] Jiashun Wang, Huazhe Xu, Jingwei Xu, Sifei Liu, and Xiaolong Wang. Synthesizing long-term 3D human motion and interaction in 3d scenes. In *Computer Vision and Pattern Recognition* (*CVPR*), 2021. 2
- [59] Zan Wang, Yixin Chen, Tengyu Liu, Yixin Zhu, Wei Liang, and Siyuan Huang. HUMANISE: Language-conditioned human motion generation in 3d scenes. In *Neural Information Processing Systems (NeurIPS)*, 2022. 2

- [60] Yiming Xie, Varun Jampani, Lei Zhong, Deqing Sun, and Huaizu
  Jiang. Omnicontrol: Control any joint at any time for human
  motion generation. *arXiv:2310.08580*, 2023. 1, 2, 3
- [61] Liang Xu, Ziyang Song, Dongliang Wang, Jing Su, Zhicheng
  Fang, Chenjing Ding, Weihao Gan, Yichao Yan, Xin Jin,
  Xiaokang Yang, Wenjun Zeng, and Wei Wu. ActFormer: A
  gan-based transformer towards general action-conditioned 3d
  human motion generation. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [62] Ye Yuan and Kris Kitani. DLow: Diversifying latent flows for
   diverse human motion prediction. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [63] Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic,
  Or Litany, Sanja Fidler, and Karsten Kreis. Lion: Latent point
  diffusion models for 3d shape generation. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 3
- [64] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding
  conditional control to text-to-image diffusion models. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- 848 [65] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. MotionDiffuse:
  850 Text-driven human motion generation with diffusion model.
  851 arXiv:2208.15001, 2022. 1, 2, 3, 4, 5, 6, 7
- [66] Qinsheng Zhang, Jiaming Song, Xun Huang, Yongxin Chen,
  and Ming yu Liu. Diffcollage: Parallel generation of large
  content with diffusion models. In *Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 3, 4, 5, 6
- [67] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger,
  and Yoav Artzi. BERTScore: Evaluating text generation with
  BERT. 2020. 5
- [68] Zihan Zhang, Richard Liu, Kfir Aberman, and Rana Hanocka.
  TEDi: Temporally-entangled diffusion for long-term motion synthesis. *arXiv:2307.15042*, 2023. 2
- [69] Lingting Zhu, Xian Liu, Xuanyu Liu, Rui Qian, Ziwei Liu, and
  Lequan Yu. Taming diffusion models for audio-driven co-speech
  gesture generation. In *Computer Vision and Pattern Recognition*(*CVPR*), 2023. 2
- [70] Wentao Zhu, Xiaoxuan Ma, Dongwoo Ro, Hai Ci, Jinlu Zhang,
  Jiaxin Shi, Feng Gao, Qi Tian, and Yizhou Wang. Human motion
  generation: A survey. *arXiv:2307.10894*, 2023. 2
- [71] Xinxin Zuo, Sen Wang, Jiangbin Zheng, Weiwei Yu, Minglun
  Gong, Ruigang Yang, and Li Cheng. Sparsefusion: Dynamic
  human avatar modeling from sparse rgbd images. *IEEE Transactions on Multimedia*, 2021. 5