

MDMP: Multi-modal Diffusion for supervised Motion Predictions with uncertainty

Anonymous CVPR submission

Paper ID 6

Abstract

001 This paper introduces a Multi-modal Diffusion model for
 002 Motion Prediction (MDMP) that integrates and synchronizes
 003 skeletal data and textual descriptions of actions to
 004 generate refined long-term motion predictions with quantifiable
 005 uncertainty. Existing methods for motion forecasting or motion
 006 generation rely solely on either prior motions or text prompts,
 007 facing limitations with precision or control, particularly over
 008 extended durations. The multi-modal nature of our approach
 009 enhances the contextual understanding of human motion, while
 010 our graph-based transformer framework effectively captures both
 011 spatial and temporal motion dynamics. As a result, our model
 012 consistently outperforms existing generative techniques in
 013 accurately predicting long-term motions. Additionally, by
 014 leveraging diffusion models' ability to capture different modes
 015 of prediction, we estimate uncertainty, significantly improving
 016 spatial awareness in human-robot interactions by incorporating
 017 zones of presence with varying confidence levels.
 018

019 1. Introduction

020 Through collaboration and assistance, robots could significantly
 021 augment human capabilities across diverse sectors, including
 022 smart manufacturing, healthcare, agriculture, construction and
 023 many others. Indeed, they can complement the critical and
 024 adaptive decision-making skills of human workers with higher
 025 precision and consistency in repetitive tasks. However, one
 026 challenge prohibiting human-robot collaboration is the safety
 027 of workers in the presence of robots. To act safely and
 028 effectively together, continuous knowledge of future human
 029 motion and location in the common workspace with a measure
 030 of uncertainty is pivotal. This real-time awareness allows
 031 robots to adjust their trajectories to avoid collision and
 032 perform precise collaborative tasks [5, 25, 58].
 033

034 Humans can predict future events based on their self-constructed
 035 models of physical and socio-cultural systems.

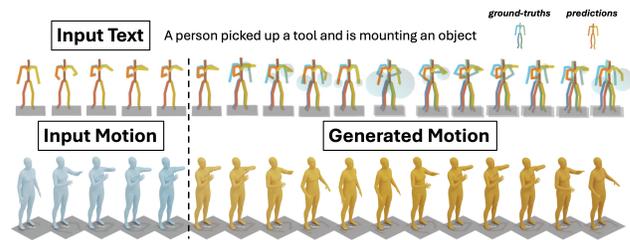


Figure 1. MDMP integrates skeletal motion and text to generate long-term motion predictions with uncertainty zones, shown in both skeletal and 3D human mesh formats.

This skill, developed from childhood through observation and active participation in society, enables them to anticipate others' movements. Researchers are now trying to transfer this capability often referred as "Theory of mind" [10] to machines by training them to learn similar motion estimation tasks. Current methodologies fall into two main categories: Human Motion Forecasting (see Section 2.2) and Human Motion Generation (see Section 2.3). While the former uses only a short input sequence of skeletal motion to predict its future trajectory, the latter relies exclusively on textual prompts to generate motion sequences.

Despite advancements in text-to-motion models, challenges remain in controlling generation due to the expansive action space a simple prompt can describe, which may not always align with human expectations or behavior. Moreover, while some text-to-motion methods have been adapted to perform tasks like motion editing or motion prediction by conditioning their generative process on motion data during sampling, our study demonstrates that, since they are only fed with textual prompts during training, our method consistently outperforms them in terms of accuracy metrics.

Conversely, motion prediction using past sequences is a long-standing challenge that has achieved high accuracy over short-term predictions but struggles with long-term predictions. Even for humans, predicting someone's immediate future movement based on past motions is feasible, but beyond one or two seconds, the multitude of possibilities makes it nearly impossible without context. However, knowledge of the intended action provides a rough idea of

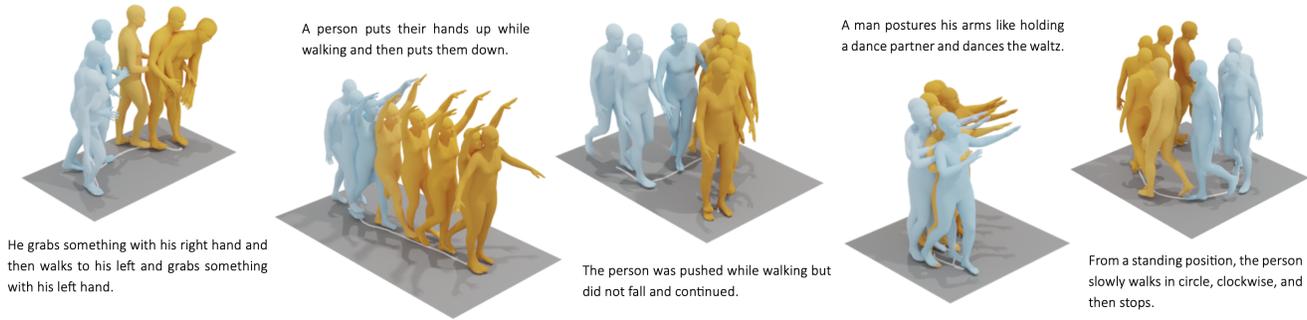


Figure 2. **SMPL [34] Meshes of MDMP Predicted motions of different scenarios.** The text descriptions vertically associated to the motions as well as the blue frames are the inputs of the model. The orange frames are the predictions, darker colors indicate later frames.

065 future positions, as the contextual information of the action
066 guides intuition.

067 In Human Robot Collaboration (HRC), there is a crucial
068 need for longer-term predictions to coordinate precise in-
069 teractive tasks, avoid collisions, and maintain efficient tra-
070 jectory planning. As a result, our method uniquely com-
071 bines and synchronizes textual and skeletal data to gener-
072 ate precise, longer-term predictions. Indeed, this integra-
073 tion allows for a richer, more contextually aware generation
074 of motion predictions. To the best of our knowledge, this
075 model is the first to be trained on a combination of both
076 types of inputs to leverage context in motion.

077 In this work, inspired by MDM [50] (Motion Diffu-
078 sion Model) and LTD [37] (Learning Trajectory Depend-
079 encies), we propose a transformer-based diffusion model with
080 a Graph Convolutional Encoder optimized for the spatio-
081 temporal dynamics of motion data. A key design element is
082 the use of learnable graph connectivity, as introduced by
083 Mao et al. [37], to more effectively capture joint depen-
084 dencies. Additionally, our Multi-modal Diffusion Model
085 for Motion Predictions (MDMP) harnesses the stochastic
086 nature of diffusion models to predict presence zones with
087 varying confidence levels. This uncertainty measure is par-
088 ticularly crucial for long-term motion predictions, where
089 uncertainty grows over time. By offering a spatial under-
090 standing of human presence, our model significantly en-
091 hances collision avoidance, improving safety and real-time
092 interaction in dynamic collaboration scenarios.

093 We summarize the contributions as follows: 1) A novel
094 multi-modal diffusion model trained on both textual and
095 skeletal data for precise long-term motion predictions. 2)
096 An uncertainty estimation method to significantly enhance
097 spatial awareness and safety in HRC scenarios. 3) A graph-
098 based transformer capturing spatial-temporal dynamics ef-
099 fectively. 4) A comprehensive validation of uncertainty es-
100 timation, with an open-source implementation.

2. Related Work 101

In this section, we review key works that inform our ap-
proach. We cover Diffusion Generative Models, Human
Motion Forecasting, and Human Motion Generation, high-
lighting the advancements and limitations in each area as
they relate to our method. 102
103
104
105
106

2.1. Diffusion Generative Models 107

Diffusion models [18, 46, 47] are neural generative models
based on a stochastic diffusion process as modeled in Ther-
modynamics. The training process involves two phases:
forward and backward. The forward process takes observed
samples x and progressively adds Gaussian noise until the
original information is completely obscured. In contrast,
the backward or reverse process employs a neural model
that learns to denoise a sample from pure noise back to
the original data distribution $p(x)$, hence the term Denoising
Diffusion Probabilistic Models [18]. DDPMs have gained
prominence in generative modeling, initially demonstrating
excellent performance in image generation, and later in con-
ditioned generation [9] and latent text representation [41]
using CLIP [44]. Recently, diffusion models have also
been applied to various generation tasks, such as text-to-
speech [43], text-to-sound [56], and text-to-video [19]. 108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123

While diffusion models excel in performance, a signifi-
cant trade-off is the lengthy inference time required for the
reverse process, which is impractical for real-time applica-
tions. However, many work such as DDIM [48] and Consis-
tency Models [49] tackles that issue and trade off computa-
tion for sample quality. Nichol et al. [40] found that instead
of fixing variances of distributions modeling the progres-
sively denoised data as a hyperparameter [18], learning it
would improve log-likelihood, forcing generative models to
capture all data distribution modes, and enable faster sam-
pling with minimal quality loss. Considering the paramount
importance of efficiency in HRC, we follow Nichol et al.’s
approach by learning variances and leverage the different
modes as a factor for uncertainty. Our method demonstrates 124
125
126
127
128
129
130
131
132
133
134
135
136
137

138 better performance with just 50 time steps instead of 1000,
139 achieving over 20 times faster inference.

140 2.2. Human Motion Forecasting

141 Human Motion Forecasting aims to predict future full-body
142 motion trajectories in 3D space based on past observations
143 from motion capture data or real-time Human Pose Esti-
144 mation methods. This task is formulated as a sequence-to-
145 sequence problem, using past motion segments to predict
146 future motion. Deep learning methods have shown notable
147 results due to their ability to learn motion patterns and un-
148 derstand spatio-temporal relationships. Early methods em-
149 ployed RNNs [11, 20, 26, 32, 38], then CNNs [31, 57] and
150 GANs [8, 13, 17, 22, 27, 53, 60] but either accumulated er-
151 rors led to unrealistic predictions or faced limitations due
152 to prefixed kinematic dependencies between body joints.
153 GCNs have proven effective for the task [7, 29, 30, 33, 37,
154 59], considering that the human skeleton can be effectively
155 modeled as a graph. Transformer-based models, leveraging
156 self-attention [51] for long-range dependencies, have also
157 been adopted [2, 4, 39, 54]. Considering the efficiency and
158 accuracy of the previously mentioned methods, our denois-
159 ing model leverages GCNs to encode joint features due to
160 their effectiveness in capturing spatial patterns, and a Trans-
161 former backbone in the latent space to address the temporal
162 nature of motion data. However, since none of these meth-
163 ods can learn contextual information from the data they are
164 fed, they tend to diverge for durations beyond one second.

165 2.3. Human Motion Generation

166 Instead of predicting future motion based on past sequences,
167 some generative methods are conditioned on natural lan-
168 guage [1, 42] to overcome this short-term issue. This ap-
169 proach faces other challenges such as the vast variability of
170 possible motions corresponding to the same label. How-
171 ever, Text2Motion has garnered significant interest and var-
172 ied successful approaches. TEMOS [42] and T2M [15] em-
173 ploy a VAE to map text prompts to a latent space distri-
174 bution of language and motion. MotionGPT [21] furthers
175 this by proposing a unified motion-language framework.
176 MDM [50] proved that diffusion models are a better candi-
177 date for human motion generation, as they can retain the
178 formation of the original motion sequence and thus allows
179 them to easily apply more constraints during the denoising
180 process. Then, LDM [6] performed the Diffusion in the la-
181 tent space and MoMask [16] leveraged Masked Transform-
182 ers.

183 By fixing some parts of a motion sequence and filling
184 in the gaps, some of these Text2Motion baselines such as
185 MDM [50], MotionGPT [21] and MoMask [16] propose a
186 form of "motion editing" by forcing their models to gener-
187 ate motions with preserved original data. Unlike these
188 methods, which only edit motions during sampling, our ap-

proach trains the model with both textual prompts and mo-
tion sequence conditioning to learn contextuality and guide
generation towards precise predictions. While these models
are compared on diversity and multi-modality metrics, our
goal is to minimize the distance between predictions and
ground-truth for accurate predictions in HRC.

3. Methodology

We now explain the architecture of our proposed MDMP
in detail. For an overview, please refer to Figure 3. As
part of the Diffusion Process MDMP progressively denoises
a motion sample conditioned by the input motion through
masking. Our architecture employs a GCN encoder to cap-
ture spatial joint features. We encode text prompts using
CLIP followed by a linear layer; the textual embedding
 c and the noise time-step t are projected to the same di-
mensional latent space by separate feed-forward networks.
These features, summed with a sinusoidal positional em-
bedding, are fed into a Transformer encoder-only back-
bone [51]. The backbone output is projected back to the
original motion dimensions via a GCN decoder. Our model
is trained both conditionally and unconditionally on text, by
randomly masking 10% of the text embeddings. This ap-
proach balances diversity and text-fidelity during sampling.

Our method uses the building blocks of MDM [50],
but with three key differences: (1) a denoising model that
includes variance learning to increase log-likelihood and
perform uncertainty estimates, (2) the GCN encoder with
learnable graph connectivity, and (3) a learning framework
that incorporates contextuality by synchronizing skeletal in-
puts with initial textual inputs.

3.1. Problem Formulation

A motion sample can be represented by a temporal skele-
ton sequence $X = \{p^i\}_{i=1}^N$ of length N where a frame p_i
denotes a pose that can be modeled using different joint
feature representations depending on the dataset (see Sec-
tion 4.1). The simplest form that any representation can
easily revert to without any loss of information is the joints'
position in 3D space where $p_i = \{x(1)_i, \dots, x(J)_i\}$ with
joints $x(j)_i \in \mathbb{R}^{M=3}$ and J being the total number of joints.
Some parameterizations use rotation matrices ($M = 9$),
angle-axis ($M = 4$), or quaternion ($M = 4$) to represent
each joint, some also include information such as angular
and/or linear velocity.

3.2. The Variational Diffusion Process

A Diffusion model can be described as a Markovian Hierar-
chical Variational Auto-Encoder [35] with a constant latent
dimension. During training, we draw X_0 from the data dis-
tribution, and at each time step t , the fixed encoder adds lin-
ear Gaussian noise centered around the output of the previ-
ous latent sample X_{t-1} until its distribution becomes a stan-

189
190
191
192
193
194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

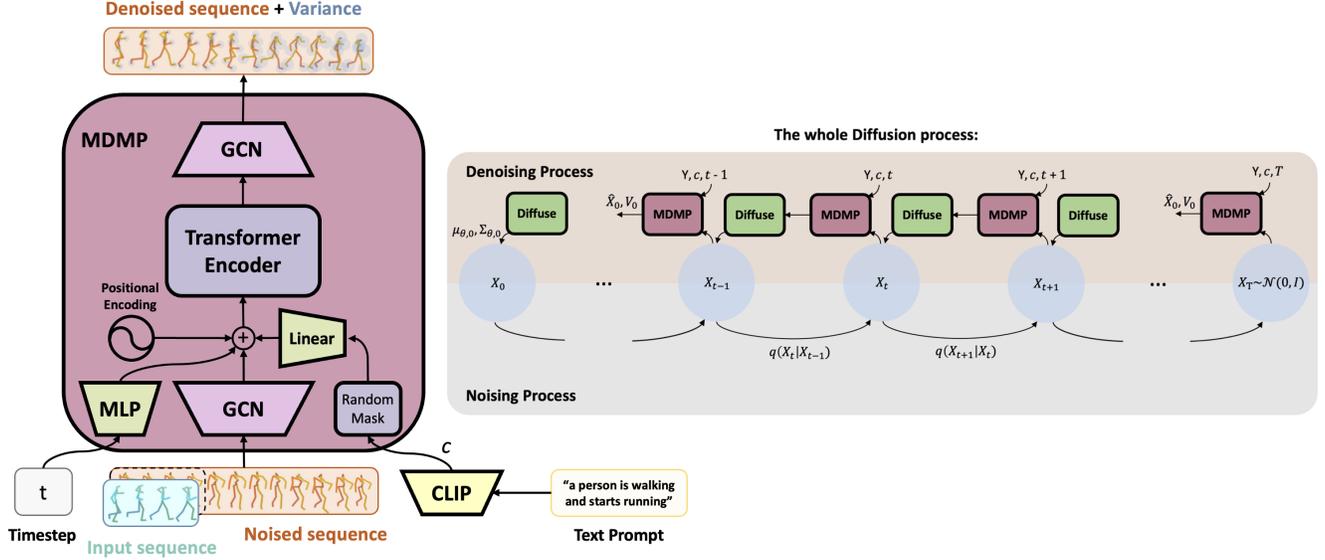


Figure 3. **(Left) Architecture of MDMP.** The denoising model takes as input a motion sample $X_t = \{p_i^i\}_{i=1}^N$ from the previous latent distribution, the diffusion time step t and the conditioning parameters: $Y = \{p^i\}_{i=1}^n$ with $n < N$ the motion input sequence and c the textual embedding encoded by CLIP [44]. At each time step, MDMP outputs a prediction of the final motion \hat{X}_0 along with V_0 , the variance of each predicted joint feature. **(Right) Overview of the Diffusion Process.** On top is the denoising Process, where the Sampling starts from $t = T$ and recursively calls MDMP and uses the output along with X_t to diffuse back to X_{t-1} by calculating $\mu_{\theta,t}$ and $\Sigma_{\theta,t}$.

239 dard Gaussian at the final time step T . Hence, the Gaussian
240 encoder is parameterized with mean $\mu_t(X_t) = \sqrt{\alpha_t}X_{t-1}$
241 and variance $\Sigma_t(X_t) = (1 - \alpha_t)I$:

$$242 \quad q(X_t|X_{t-1}) = \mathcal{N}(X_t; \sqrt{\alpha_t}X_{t-1}, (1 - \alpha_t)I) \quad (1)$$

243

$$244 \quad q(X_{1:T}|X_0) = \prod_{t=1}^T q(X_t|X_{t-1}) \quad (2)$$

245 Inspired by Nichol et al. [40], we use a cosine scheduler
246 for β_t and $\alpha_t = 1 - \beta_t$ such that $\beta_t, \alpha_t \in [0, 1]$. α_t is slowly
247 decreasing, so that for $T = 1000$, α_t is small enough to say
248 that $X_T \sim \mathcal{N}(0, I)$.

249 Then, during both training and inference, we use MDMP
250 (see Fig 3) as the decoder—conditioned at each step by the
251 previously mentioned inputs Y and c —to progressively de-
252 noise X_t from a standard Gaussian. Instead of predicting
253 the noise ϵ_0 as formulated in DDPM [18], we follow [45]
254 and [50] and predict the signal itself along with its vari-
255 ance: $\hat{X}_0, V_0 = \text{MDMP}(X_t, t, Y, c)$

256 Then we use this prediction \hat{X}_0 along with the current
257 X_t to diffuse back to the posterior mean:

$$258 \quad \mu_{\theta,t-1} = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})X_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{X}_0}{1 - \bar{\alpha}_t} \quad (3)$$

259

$$260 \quad \text{with } \bar{\alpha}_t = \prod_{s=1}^t \alpha_s. \quad (4)$$

261 We use the simple objective from [18] to train our model:

$$262 \quad L_{\text{simple}} = \mathbb{E}_{X_0 \sim q(X_0|c, Y), t \sim [1, T]} [\|X_0 - \hat{X}_0\|^2] \quad (5)$$

263 One subtlety is that L_{simple} provides no learning signal for
264 variances, as Ho et al. [18] chose to fix the variance rather
265 than learn it. However, in our framework, we leverage
266 learned variances to generate presence zones with varying
267 confidence levels to help ensure safety in HRC scenarios.

268 3.3. Learning the Variances of the Denoising process

269 To learn the reverse process variances, our model outputs a
270 vector V_0 of the same shape as \hat{X}_0 , and—following Nichol et
271 al. [40]—we parameterize the variance as an interpolation
272 between β_t and $\tilde{\beta}_t$ in the log domain by turning this output
273 V_0 into $\Sigma_{\theta,t}$ as follows:

$$274 \quad \Sigma_{\theta,t} = \exp(V_0 \log \beta_t + (1 - V_0) \log \tilde{\beta}_t) \quad (6)$$

$$275 \quad \text{with } \tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t. \quad (7)$$

277 Then, we leverage the reparameterization trick $x_t =$
278 $\bar{\alpha}_t x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$ with $\epsilon \sim \mathcal{N}(0, I)$ to sample from an
279 arbitrary step of the forward noising process and estimate
280 the variational lower bound (VLB). As mentioned ear-
281 lier, the diffusion model can be thought of as a VAE [23]
282 where q represents the encoder and $p_{\theta}(x_{t-1}|x_t) =$
283 $\mathcal{N}(x_{t-1}; \mu_{\theta,t}, \Sigma_{\theta,t})$ is the decoder, so we can write:

$$284 \quad L_{\text{VLB}} := L_0 + L_1 + \dots + L_{T-1} + L_T \quad (8)$$

$$285 \quad L_0 := -\log p_{\theta}(x_0|x_1) \quad (9)$$

$$286 \quad L_{t-1} := D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \| p_{\theta}(x_{t-1}|x_t)) \quad (10)$$

$$287 \quad L_T := D_{\text{KL}}(q(x_T|x_0) \| p(x_T)) \quad (11)$$

291 Finally, with $q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)$
 292 we estimate L_{t-1} and approximate L_{VLB} with the expecta-
 293 tion $\mathbb{E}_{t, X_0, \epsilon}[L_{t-1}]$.

294 Since L_{simple} does not depend on $\Sigma_{\theta, t}$, we define a new
 295 hybrid objective: $L_{\text{hybrid}} = L_{\text{simple}} + \lambda L_{\text{VLB}}$

296 Conversely to Nichol et al. [40], we apply a clamping on
 297 V_0 to prevent NaN values during the calculation of L_{VLB} .

298 3.4. Encoding the joint features with GCNs

299 To encode the spatial pose features, we leverage GCNs [52].
 300 Instead of relying on a predefined sparse graph, we fol-
 301 low Mao et al. [37] and learn the graph connectivity dur-
 302 ing training, thus essentially learning the dependencies be-
 303 tween the different joint trajectories. To this end, we use a
 304 fully-connected graph with N nodes, N being the length
 305 of the predicted sequence. The strength of the edges in
 306 this graph is represented by the weighted adjacency matrix
 307 $A \in \mathbb{R}^{N \times N}$. The graph convolutional encoder/decoder then
 308 takes as input a matrix $H^{(\text{in})} \in \mathbb{R}^{N \times F}$, where in our case
 309 F is the number of body joint features. Given the input a
 310 matrix $H^{(\text{in})}$, the adjacency matrix A and a set of trainable
 311 weights $W \in \mathbb{R}^{F \times \hat{F}}$, a graph convolutional layer outputs a
 312 matrix of the form: $H^{(\text{out})} = AH^{(\text{in})}W$. All operations are
 313 differentiable with respect to both the adjacency matrix A
 314 and the weight matrix W , which allows training on both.

315 3.5. Predicting Uncertainty

316 To derive an effective uncertainty index for each joint pre-
 317 diction over time, we explore three different approaches
 318 which we evaluate and compare in Section 4.4:

- 319 • **Mode Divergence:** This approach measures the variabil-
 320 ity between multiple motion sequences generated from
 321 the same input. We compute several predictions in par-
 322 allel, calculate the standard deviation of these sequences,
 323 and use this as the uncertainty index.
- 324 • **Denosing Fluctuations:** Here, we measure the fluctua-
 325 tions during the denosing process as an uncertainty indi-
 326 cator. As illustrated in Figure 1 which tracks the evolution
 327 of the x-coordinate of key joints (head, hands, feet) from
 328 random noise to the final prediction, earlier steps are very
 329 noisy and progressively converge with more or less stab-
 330 ility. Significant fluctuations in the last 20 timesteps are
 331 used as an indicator of uncertainty.
- 332 • **Predicted Variance:** The final approach uses the learned
 333 variance of the predicted distribution of each motion se-
 334 quence $\Sigma_{\theta, 0}$ as the uncertainty factor.

335 Both the second and third methods produce outputs in
 336 the same dimensions as the model, including predictions
 337 for root height, root angular and linear velocity, as well as
 338 joint positions and velocities in the local reference of the
 339 root. To calculate a single uncertainty index for each joint
 340 at each timestep, we average all features associated to the
 341 same joint.

4. Experiments and Results 342

343 In this section, we present the experimental setup and eval-
 344 uation of our proposed model. We describe the dataset used
 345 for training and testing, outlines and explains the choice of
 346 metrics used for accuracy and uncertainty, and provide de-
 347 tails on our model’s implementation. Our comprehensive
 348 quantitative and qualitative evaluation, includes compari-
 349 son with state-of-the-art Text2Motion baselines that pro-
 350 pose Motion-Editing re-implemented for a fair comparison
 351 with similar conditioning, analysis of uncertainty param-
 352 eters, and an ablation study to assess the effects of motion-
 353 text fusion and architectural design choices.

4.1. Dataset 354

355 To train and evaluate our model, we use the Hu-
 356 manML3D [15] dataset, which is the largest and most
 357 diverse collection of scripted human motions. It com-
 358 bines motion sequences from the HumanAct12 [14] and
 359 AMASS [36] datasets, processed to standardize the mo-
 360 tions to 20 FPS with a maximum length of 10 seconds per
 361 sequence. HumanML3D comprises 14,616 motions with
 362 44,970 descriptions, covering 5,371 distinct words, totaling
 363 28.59 hours of motion data with an average length of 7.1s
 364 and three textual descriptions per sequence. The dataset is
 365 split for training and evaluation. For evaluation, we filter
 366 the set to include only motions longer than 3s, allowing us
 367 to condition the models on 2.5s of motion and predict at
 368 least 0.5s into the future up until more than 5s for longer
 369 recorded motions. After filtering, the evaluation set con-
 370 tains 4,328 out of the initial 4,646 motion sequences.

4.2. Metrics for Accuracy and Uncertainty 371

372 To evaluate and compare the accuracy of our model we use
 373 the *Mean Per Joint Position Error (MPJPE)* on 3D joint
 374 coordinates, which is the most widely used metric for eval-
 375 uating 3D pose errors. This metric calculates the average
 376 L2-norm across different joints between the prediction and
 377 ground truth. Since HumanML3D [15] pose representation
 378 contains 263 redundant features per body frames including
 379 joint positions, velocities and rotations we use a transfor-
 380 mation process (described in the Appendix) to obtain the
 381 3D joint positions in order to both calculate the *MPJPE* and
 382 visualize the predicted sequences.

383 To further validate our method we have also added some
 384 more metrics in the C.4 Section of the Appendix. First of
 385 all, we have re-trained MDM [50] with skeleton data as an
 386 input for a direct comparison, demonstrating the efficacy
 387 of our architecture. Secondly, one issue with the *MPJPE*
 388 is that it is biased towards one “ground-truth sequence”
 389 and thus heavily penalizes frequency or phase shifts com-
 390 mon in longer-term predictions, leading to misleadingly
 391 large errors even if motions remain qualitatively realis-
 392 tic. Hence we compared our method with baselines on the

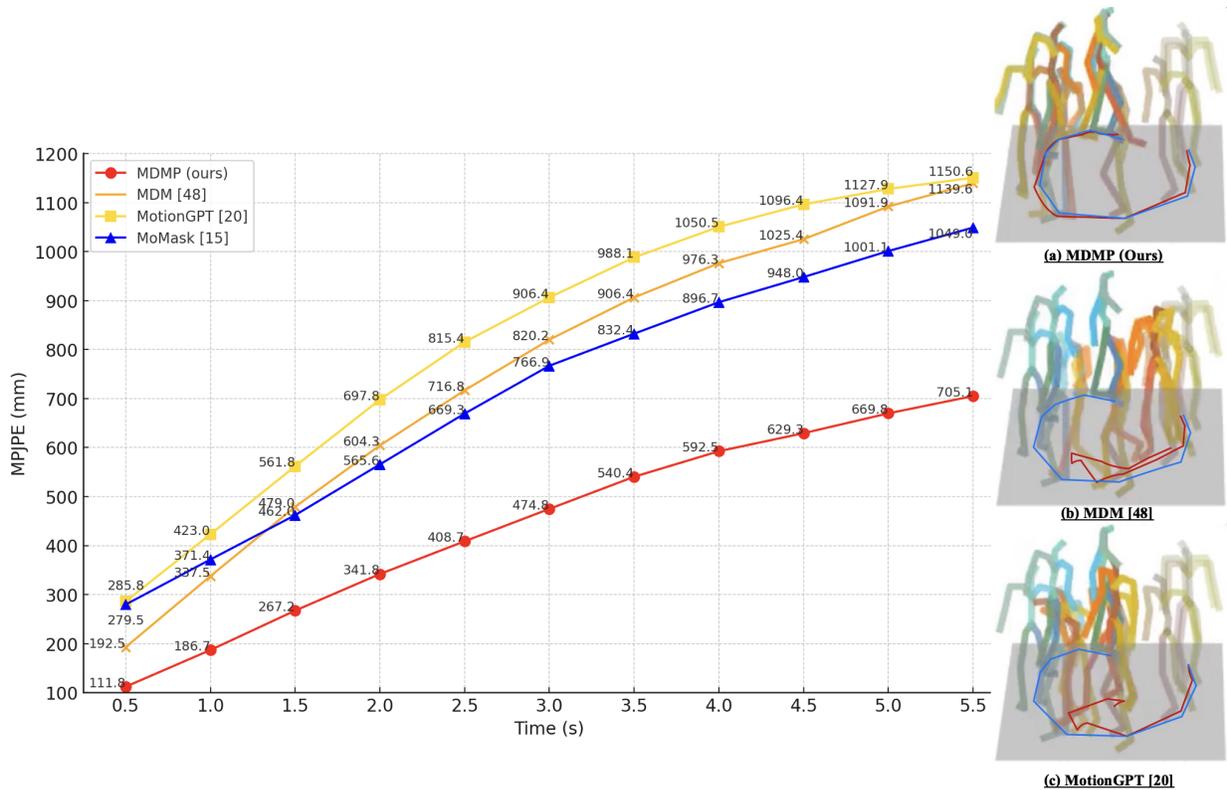


Figure 4. **(Left) Temporal evolution of error in predictions.** Quantitative Results on HumanML3D over *MPJPE* [mm]. **(Right) 3D Plots of Motion Predictions (orange) vs Ground truth (blue).** Motion Sequence example associated to textual prompt: “*from a standing position, the person slowly walks in circle, clockwise, then stops*”. Paler shades represents earlier frames.

393 NPSS [12] metric which measures similarity in frequency
 394 spectra rather than absolute frame-by-frame error, making
 395 it better suited to assess the quality of long-term predic-
 396 tions by capturing perceptually relevant motion coherence.
 397 Finally, we report results using metrics proposed by Guo
 398 et al. [15], such as *Frechet Inception Distance (FID)*, *R-*
 399 *Precision*, and *Multimodality*. However, these metrics pri-
 400 marily assess motion quality, semantic alignment with text-
 401 ual input, and variability rather than precise spatial accu-
 402 racy. Additionally, they depend on pretrained feature ex-
 403 tractors not tailored to motion-conditioned predictions.

404 To evaluate and compare our uncertainty indices, we
 405 use sparsification plots, a common approach for assess-
 406 ing how well estimated uncertainty aligns with true er-
 407 rors [3, 24, 28, 55]. In our implementation, we compute
 408 multiple motion sequences and rank each joint’s uncer-
 409 tainty. By progressively removing the joints with the high-
 410 est uncertainty and summing the remaining error, we obtain
 411 the sparsification curve. The ideal reference, known as the
 412 “Oracle”, is based on ranking joints by their true errors. A
 413 well-performing and reliable uncertainty index should pro-
 414 duce a curve that decreases monotonically and closely fol-
 415 lows the oracle.

4.3. Implementation Details

416 Our models were trained on an *NVIDIA Titan V* GPU over
 417 1.7 days and on *NVIDIA Tesla V100* GPU over 1.2 days
 418 with a batch size of 64. We used 8 layers of the Trans-
 419 former Encoder with 4 multi-head attention for each, sep-
 420 arated by a GeLU activation function and a dropout value
 421 of 0.1. The GCN layer encodes the joint features from X_t
 422 into a latent dimension of 1024 when learning variances
 423 and 512 without learning variances. 1024 corresponds to
 424 the concatenation of the joint features of \hat{X}_0 [512] and V_0
 425 [512]. To encode the text, we use a frozen CLIP-ViT-B/32
 426 model. Each model was trained for 600K steps, after which
 427 a checkpoint was chosen that minimizes the *MPJPE* metric
 428 to be reported. Our generative process is conditioned by a
 429 motion input sequence of 50 frames which represents 2.5
 430 seconds at 20FPS. We also set $\lambda = 0.001$ to prevent L_{VLB}
 431 from overwhelming L_{simple} . We evaluate our models with
 432 guidance-scale $\mu = 2.5$ but as discussed in the Motion &
 433 Text ablation study Section 4.4 this can be adapted for spe-
 434 cific applications (eg. short/long-term predictions).
 435

436 To evaluate the effectiveness of our multimodal fu-
 437 sion approach, we compare against state-of-the-art Mo-
 438 tion Editing baselines MoMask [16], MotionGPT [21] and
 439 MDM [50] which are all trained on HumanML3D [15].
 439

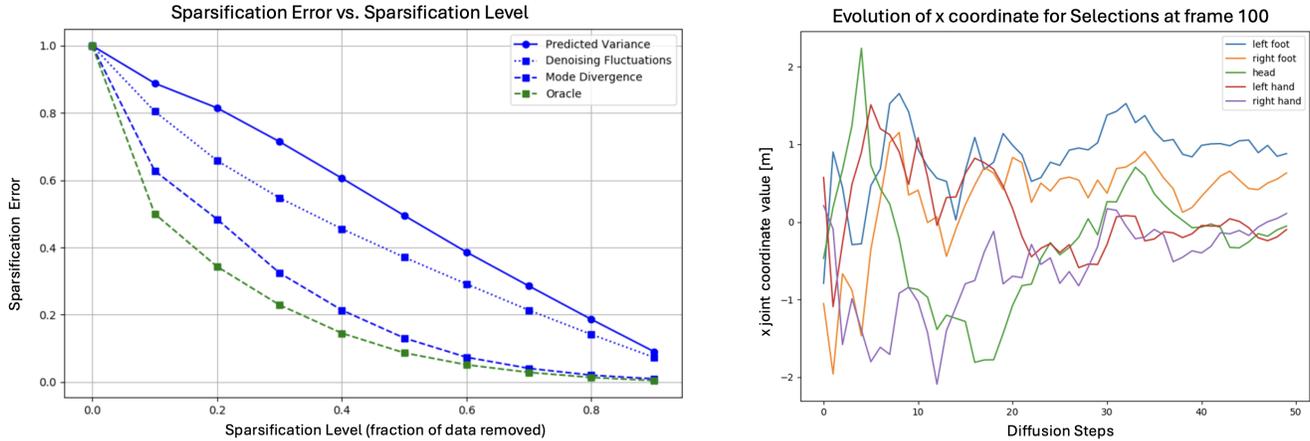


Figure 5. (Left) **Sparsification Error Plot**. Quantitative Results of the uncertainty parameters: The Mode Divergence index closely follows the Oracle curve, indicating the strongest alignment between uncertainty estimates and true errors. (Right) **Joint Position Evolution over the Denoising Process**. The position is progressively denoised until it converges to its final prediction. The fluctuations are used as a parameter for uncertainty.

440 We implement their pretrained versions (open-sourced)
 441 and compare on the entire test set of HumanML3D using
 442 *MPJPE*. Conversely to Motion Editing, to ensure a fair
 443 comparison setting we conditioned each baseline with only
 444 the same motion prequel sequence of 50 frames and compared
 445 the rest of the predicted sequence to the ground-truth.

446 4.4. Quantitative & Qualitative Results

447 **Model Accuracy Evaluation over *MPJPE***: Unlike the im-
 448 plemented baselines MoMask [16], MotionGPT [21] and
 449 MDM [50] that treat motion data as a masked input during
 450 sampling, our model is trained to leverage it as an additional
 451 supervision signal, which we find to lead to significantly en-
 452 hances performance, especially over longer sequences. In
 453 Fig. 4 the temporal evolution chart shows that our model
 454 outperforms these baselines in accuracy, with consistently
 455 lower *MPJPE* values over time and a more gradual increase
 456 in error. These results are also demonstrated qualitatively
 457 in the 3D plots Fig. 4 (Right) (see Appendix & Video for
 458 more examples) where our predictions align more closely
 459 with the ground truth, especially towards the end of the se-
 460 quence. Indeed, both baselines’ outputs fail to follow the
 461 indicated trajectory (projection of the root joint in the XZ-
 462 plane) whereas our model follows the “circle”, almost align-
 463 ing with the ground truth on the last frame.

464 **Uncertainty Parameters Evaluation**: The results of
 465 our comparison study between the different uncertainty in-
 466 dices are presented in Fig. 5 (Left). The Sparsification Error
 467 plot (explained in 4.2) shows that the best-performing index
 468 is the Mode Divergence, which closely follows the Oracle
 469 curve, indicating a strong alignment between uncertainty
 470 estimates and true errors. These results are also demon-
 471 strated qualitatively in the video as well as in the Additional
 472 Experimental Results (Appendix) where we visualize the

evolution of the zones of presence with varying confidence
 evolution of the zones of presence with varying confidence
 levels based on the different uncertainty indices. For clarity
 and visibility, we limit the uncertainty visualization to the
 “end-effector” joints—specifically the head, hands, and
 feet—since these are the most critical in human-robot col-
 laboration, and visualizing uncertainty for all joints would
 create overly cluttered visuals. We calculate the mean un-
 certainty across the x, y, and z coordinates for each key
 joint, using this value as the radius of the sphere represent-
 ing uncertainty around the end-effectors.

Uncertainty Results Interpretation: Although the De-
 noising Fluctuations and Predicted Variance methods show
 a general decline in their sparsification curves, the effect is
 less pronounced, suggesting these indices are less reliable
 for uncertainty estimation. The learned variance is sup-
 posed to generally follow the same trends as the original
 fixed schedule, consistently decreasing during denoising to
 reduce stochasticity. However, its effectiveness as an uncer-
 tainty factor is somewhat limited, as the final value, while
 still meaningful, becomes slightly less informative. Simi-
 larly, the instability of fluctuations diminishes their reli-
 ability. In contrast, the Mode Divergence factor consistently
 rises over time, aligning with the increasing error, making
 it the most robust and dependable indicator (see video and
 Appendix for visual confirmation in 3D plots).

Ablation Study - Motion and Text Effects: To evaluate
 the relevance of our multi-modal contribution, we perform
 an ablation study, presented in Table 1, comparing our stan-
 dard approach to one where models are fed with either mo-
 tion or textual inputs exclusively. Firstly, this study clearly
 confirms that combining both types of inputs results in sig-
 nificantly higher prediction accuracy. Secondly, the study
 indicates that our model relies more heavily on motion in-
 put sequences than textual prompts. Notably, it performs

Time (seconds)	0.5s	1s	1.5s	2s	2.5s	3s	3.5s	4s	4.5s	5s	5.5s
Ours with motion & text	111.8	186.7	267.2	341.8	408.7	474.8	540.4	592.5	629.3	669.8	705.1
MDM [50] with motion & text	192.5	337.5	479.0	604.3	716.8	820.2	906.4	976.3	1025.4	1091.9	1139.6
Ours with text no motion	254.8	418.6	609.5	796.8	972.2	1105.1	1253.3	1383.8	1526.1	1624.8	1679.8
MDM [50] with text no motion	237.9	362.6	482.9	595.5	687.8	783.2	871.6	965.3	1039.6	1085.2	1143.8
Ours with motion no text	100.2	186.9	271.9	358.9	445.7	528.6	608.7	677.6	739.1	810.8	902.0
MDM [50] with motion no text	406.1	614.5	852.3	1079.3	1288.6	1503.5	1684.8	1871.9	2001.6	2187.3	2332.0

Table 1. Ablation study: *MPJPE* (mm) to assess Motion and Text Effects

Time (seconds)	0.5s	1s	1.5s	2s	2.5s	3s	3.5s	4s	4.5s	5s	5.5s
Encoder/Decoder: Linear	118.6	205.5	298.8	385.5	472.0	551.3	629.7	692.0	741.0	791.9	852.1
Encoder/Decoder: GCN	111.8	186.7	267.2	341.8	408.7	474.8	540.4	592.5	629.3	669.8	705.1
Learning the Variance: False	86.3	163.0	250.1	332.4	409.3	485.4	560.5	622.6	676.8	729.5	775.5
Learning the Variance: True	111.8	186.7	267.2	341.8	408.7	474.8	540.4	592.5	629.3	669.8	705.1
Diffusion Steps: 1000	104.8	192.2	280.6	360.9	438.3	482.1	553.6	617.2	653.5	702.3	745.8
Diffusion Steps: 50	111.8	186.7	267.2	341.8	408.7	474.8	540.4	592.5	629.3	669.8	705.1

Table 2. Ablation study: *MPJPE* (mm) to evaluate Architectural Design and Parameter Choice

507 slightly better without text for very short-term predictions.
508 This means that our model could be used in a HRC setting
509 for continuous operation between different actions, even
510 without specific action context. This capability is presum-
511 ably not possible with Text2Motion models, which perform
512 poorly without text, as the study shows. Finally, the study
513 confirms that textual information is most useful for longer-
514 term predictions where the stochasticity and variability of
515 potential scenarios are much higher.

516 **Ablation Study - Architectural Design and Parame-**
517 **ter Choice:** To assess our architectural contributions, we
518 conduct a deeper analysis with additional ablation studies
519 presented in Table 2. In the first study, we retrain our model
520 with both the encoder and decoder composed of simple linear
521 layers, as in MDM [50]. The study confirms that learn-
522 able graph connectivity improves the understanding of hu-
523 man joint trajectory dependencies, especially for long-term
524 predictions. The second study evaluates our architectural
525 design that learns the variance of the motion sample distri-
526 bution. Although learning variances allow diffusion mod-
527 els to capture more data distribution modes with we lever-
528 age for uncertainty estimates, our study shows that it only
529 enhance accuracy over long-term predictions. In the final
530 study, inspired by Nichol et al. [40], we significantly reduce
531 the number of diffusion steps from 1000 to 50 which con-
532 siderably improves the computational efficiency-pivotal for
533 real-time Human-Robot Collaboration-and resulted in im-
534 proved accuracy over time.

5. Conclusion and Limitations

We present MDMP, a multimodal diffusion model that
learns contextuality from synchronized tracked motion se-
quences and associated textual prompts, enabling it to pre-
dict human motion over significantly longer terms than its
predecessors. Our model not only generates accurate long-
term predictions but also provides uncertainty estimates,
enhancing our predictions with presence zones of vary-
ing confidence levels. This uncertainty analysis was vali-
dated through a study, demonstrating the model’s capability
to offer spatial awareness, which is crucial for enhancing
safety in dynamic human-robot collaboration. Our method
demonstrates superior results over extended durations with
adapted computational time, making it well-suited for en-
suring safety in Human-Robot collaborative workspaces.

A limitation of this work is the reliance on textual de-
scriptions of actions, which can be a burden for real-time
Human-Robot Collaboration, as not every action is scripted
in advance. Currently, we use CLIP to embed these textual
descriptions into guidance vectors for our model. An inter-
esting future direction is to replace these descriptions with
images or videos captured in real time within the robotics
workspace. Since current Human Motion Forecasting meth-
ods already rely on human motion tracking data, often ob-
tained using RGB/RGB-D cameras, the necessary material
is typically already present in the workspace. Given that
CLIP leverages a shared multimodal latent space between
text and images, this approach could provide similar guid-
ance while being far less restrictive, making it more practi-
cal for dynamic and unsupervised HRC environments.

565

References

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

- [1] Chaitanya Ahuja and Louis-Philippe Morency. Language2pose: Natural language grounded pose forecasting. In *International Conference on 3D Vision (3DV)*, pages 719–728. IEEE, 2019. 3
- [2] Emre Aksan, Manuel Kaufmann, Peng Cao, and Otmar Hilliges. A spatio-temporal transformer for 3d human motion prediction. In *International Conference on 3D Vision (3DV)*, pages 565–574. IEEE, 2021. 3
- [3] Oisín M. Aodha, Arsalan Humayun, Marc Pollefeys, and Gabriel J. Brostow. Learning a confidence measure for optical flow. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(5):1107–1120, 2013. 6
- [4] Yujun Cai, Lin Huang, Yiwei Wang, Tat-Jen Cham, Jianfei Cai, Junsong Yuan, Jun Liu, Xu Yang, Yiheng Zhu, Xiaohui Shen, Ding Liu, Jing Liu, and Nadia Magnenat-Thalmann. Learning progressive joint propagation for human motion prediction. In *European Conference on Computer Vision (ECCV)*, pages 226–242, 2020. 3
- [5] Angelo Caregnato-Neto, Luciano Cavalcante Siebert, Arkady Zgonnikov, Marcos Ricardo Omena de Albuquerque Maximo, and Rubens Junqueira Magalhães Afonso. Armchair: integrated inverse reinforcement learning and model predictive control for human-robot collaboration, 2024. 1
- [6] Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, and Gang Yu. Executing your commands via motion diffusion in latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18000–18010, 2023. 3
- [7] Qiongjie Cui, Huaijiang Sun, and Fei Yang. Learning dynamic relationships for 3d human motion prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6519–6527, 2020. 3
- [8] Qiongjie Cui, Huaijiang Sun, Yue Kong, Xiaoqian Zhang, and Yanmeng Li. Efficient human motion prediction using temporal convolutional generative adversarial network. *Information Sciences*, 545:427–447, 2021. 3
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems*, pages 8780–8794, 2021. 2
- [10] Martin J. Doherty. *Theory of Mind: How Children Understand Others’ Thoughts and Feelings*. Psychology Press, Hove, England, 2009. 1
- [11] Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for human dynamics. In *International Conference on Computer Vision (ICCV)*, 2015. 3
- [12] Anand Gopalakrishnan, Ankur Mali, Dan Kifer, C. Lee Giles, and Alexander G. Ororbia. A neural temporal model for human motion prediction, 2019. 6
- [13] Liang-Yan Gui, Yu-Xiong Wang, Xiaodan Liang, and Jose M. F. Moura. Adversarial geometry-aware human motion prediction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 786–803, 2018. 3
- [14] 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*, pages 2021–2029, 2020. 5
- [15] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5152–5161, 2022. 3, 5, 6
- [16] Chuan Guo, Yuxuan Mu, Muhammad Gohar Javed, Sen Wang, and Li Cheng. Momask: Generative masked modeling of 3d human motions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3, 6, 7
- [17] Alejandro Hernandez, Jurgen Gall, and Francesc Moreno-Noguer. Human motion prediction via spatio-temporal inpainting. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 7134–7143, 2019. 3
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, pages 6840–6851, 2020. 2, 4
- [19] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv:2204.03458*, 2022. 2
- [20] Ashesh Jain, Amir R. Zamir, Silvio Savarese, and Ashutosh Saxena. Structural-rnn: Deep learning on spatio-temporal graphs. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 3
- [21] Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. Motiongpt: Human motion as a foreign language. *Advances in Neural Information Processing Systems (NeurIPS)*, 37, 2023. 3, 6, 7
- [22] Qihong Ke, Mohammed Bennamoun, Hossein Rahmani, Senjian An, Ferdous Sohel, and Farid Boussaid. Learning latent global network for skeleton-based action prediction. *IEEE Transactions on Image Processing*, 29:959–970, 2019. 3
- [23] Diederik P Kingma and Max Welling. Auto-encoding

- 669 variational bayes. *arXiv preprint arXiv:1312.6114*,
670 2013. 4 722
- 671 [24] Christian Kondermann, Rudolf Mester, and Christoph
672 Garbe. A statistical confidence measure for opti- 723
673 cal flows. In *Pattern Recognition*, pages 290–301.
674 Springer Berlin Heidelberg, Berlin, Heidelberg, 2008.
675 6 724
- 676 [25] Aadi Kothari, Tony Tohme, Xiaotong Zhang, and Ka- 725
677 mal Youcef-Toumi. Enhanced human-robot collabo- 726
678 ration using constrained probabilistic human-motion
679 prediction, 2023. 1 727
- 680 [26] Hsu kuang Chiu, Ehsan Adeli, Borui Wang, De-An
681 Huang, and Juan Carlos Niebles. Action-agnostic hu- 728
682 man pose forecasting. In *IEEE Winter Conference on* 729
683 *Applications of Computer Vision (WACV)*, 2019. 3 730
684 [27] Jogendra Nath Kundu, Maharshi Gor, and
685 R. Venkatesh Babu. Bihmp-gan: Bidirectional
686 3d human motion prediction gan. In *Proceedings of* 731
687 *the AAAI Conference on Artificial Intelligence*, pages
688 8553–8560, 2019. 3 732
- 689 [28] Jan Kybic and Clemens Nieuwenhuis. Bootstrap opti- 733
690 cal flow confidence and uncertainty measure. *Com- 734*
691 *puter Vision and Image Understanding*, 115(10):
692 1449–1462, 2011. 6 735
- 693 [29] Fanjia Li, Aichun Zhu, Yonggang Xu, Ran Cui,
694 and Gang Hua. Multi-stream and enhanced spatial- 736
695 temporal graph convolution network for skeleton- 737
696 based action recognition. *IEEE Access*, 8:97757–
697 97770, 2020. 3 738
- 698 [30] Maosen Li, Siheng Chen, Yangheng Zhao, Ya Zhang,
699 Yanfeng Wang, and Qi Tian. Dynamic multiscale
700 graph neural networks for 3d skeleton based human
701 motion prediction. In *Proceedings of the IEEE Con- 739*
702 *ference on Computer Vision and Pattern Recognition*
703 *(CVPR)*, pages 214–223, 2020. 3 740
- 704 [31] Xiaoli Liu, Jianqin Yin, Jin Liu, Pengxiang Ding, Jun
705 Liu, and Huaping Liub. Trajectorycnn: A new spatio- 741
706 temporal feature learning network for human motion
707 prediction. *IEEE Transactions on Circuits and Sys- 742*
708 *tems for Video Technology*, 2020. 3 743
- 709 [32] Zhenguang Liu, Shuang Wu, Shuyuan Jin, Qi Liu,
710 Shijian Lu, Roger Zimmermann, and Li Cheng. To- 744
711 wards natural and accurate future motion prediction
712 of humans and animals. In *Conference on Computer 745*
713 *Vision and Pattern Recognition (CVPR)*, 2019. 3 746
- 714 [33] Ziyu Liu, Hongwen Zhang, Zhenghao Chen, Zhiyong
715 Wang, and Wanli Ouyang. Disentangling and unifying
716 graph convolutions for skeleton-based action recogni- 747
717 tion. In *Proceedings of the IEEE/CVF Conference on* 748
718 *Computer Vision and Pattern Recognition*, pages 143–
719 152, 2020. 3 749
- 720 [34] Matthew Loper, Naureen Mahmood, Javier Romero,
721 Gerard Pons-Moll, and Michael J. Black. *SMPL: A* 750
Skinned Multi-Person Linear Model. Association for
Computing Machinery, New York, NY, USA, 1 edition,
2023. 2 751
- [35] Calvin Luo. Understanding diffusion models: A uni- 752
fied perspective. *arXiv preprint arXiv:2208.11970*,
2022. 3 753
- [36] Naureen Mahmood, Nima Ghorbani, Nikolaus F. 754
Troje, Gerard Pons-Moll, and Michael J. Black. 755
Amass: Archive of motion capture as surface shapes.
In *International Conference on Computer Vision*,
pages 5442–5451, 2019. 5 756
- [37] Wei Mao, Miaomiao Liu, Mathieu Salzmann, and
Hongdong Li. Learning trajectory dependencies for
human motion prediction. In *International Confer- 757*
ence on Computer Vision (ICCV), pages 9488–9496,
2019. 2, 3, 5 758
- [38] Julieta Martinez, Michael J. Black, and Javier
Romero. On human motion prediction using recurrent
neural networks, 2017. 3 759
- [39] Angel Martinez-Gonzalez, Michael Villamizar, and
Jean-Marc Odobez. Pose transformers (potr): Human
motion prediction with non-autoregressive transfor-
mers. In *International Conference on Computer Vision* 760
Workshops (ICCVW), pages 2276–2284, 2021. 3 761
- [40] Alex Nichol and Prafulla Dhariwal. Improved denois- 762
ing diffusion probabilistic models. In *International* 763
Conference on Machine Learning, 2021. 2, 4, 5, 8 764
- [41] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh,
Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya
Sutskever, and Mark Chen. Glide: Towards photore-
alistic image generation and editing with text-guided
diffusion models. *arXiv preprint arXiv:2112.10741*,
2021. 2 765
- [42] Mathis Petrovich, Michael J. Black, and Gul Varol.
Temos: Generating diverse human motions from tex-
tual descriptions. In *European Conference on Com- 766*
puter Vision (ECCV), 2022. 3 767
- [43] V. Popov, I. Vovk, V. Gogoryan, T. Sadekova, and M.
Kudinov. Grad-tts: A diffusion probabilistic model for
text-to-speech. In *International Conference on Ma- 768*
chine Learning, pages 8599–8608. PMLR, 2021. 2 769
- [44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-
try, Amanda Askell, Pamela Mishkin, Jack Clark,
et al. Learning transferable visual models from natural
language supervision. In *International Conference on* 770
Machine Learning, pages 8748–8763. PMLR, 2021.
2, 4 771
- [45] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol,
Casey Chu, and Mark Chen. Hierarchical text con-
ditional image generation with clip latents. *arXiv* 772
preprint arXiv:2204.06125, 2022. 4 773

- 774 [46] Jascha Sohl-Dickstein, Eric Weiss, Niru Mah- 826
775 eswaranathan, and Surya Ganguli. Deep unsupervised 827
776 learning using nonequilibrium thermodynamics. In *Inter- 828*
777 *national Conference on Machine Learning*, pages 829
778 2256–2265. PMLR, 2015. 2
- 779 [47] Jiaming Song, Chenlin Meng, and Stefano Ermon. 830
780 Denoising diffusion implicit models. *arXiv preprint 831*
781 *arXiv:2010.02502*, 2020. 2 832
- 782 [48] Jiaming Song, Chenlin Meng, and Stefano Er- 833
783 mon. Denoising diffusion implicit models. 834
784 *arXiv:2010.02502*, 2020. 2
- 785 [49] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya 835
786 Sutskever. Consistency models. In *Proceedings of the 836*
787 *40th International Conference on Machine Learning*. 837
788 JMLR.org, 2023. 2 838
- 789 [50] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, 839
790 Daniel Cohen-or, and Amit Haim Bermano. Hu- 840
791 man motion diffusion model. In *The Eleventh In- 841*
792 *ternational Conference on Learning Representations 842*
793 *(ICLR)*, 2023. 2, 3, 4, 5, 6, 7, 8 843
- 794 [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob 844
795 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz 845
796 Kaiser, and Illia Polosukhin. Attention is all you need. 846
797 In *Advances in Neural Information Processing Sys- 847*
798 *tems*, 2017. 3 848
- 799 [52] Petar Veličković, Guillem Cucurull, Arantxa 849
800 Casanova, Adriana Romero, Pietro Lio, and Yoshua 850
801 Bengio. Graph attention networks. In *International 851*
802 *Conference on Learning Representations (ICLR)*, 852
803 2018. 5 853
- 804 [53] Dong Wang, Yuan Yuan, and Qi Wang. Early action 854
805 prediction with generative adversarial networks. *IEEE 855*
806 *Access*, 7:35795–35804, 2019. 3 856
- 807 [54] Jiashun Wang, Huazhe Xu, Medhini Narasimhan, and 857
808 Xiaolong Wang. Multi-person 3d motion prediction 858
809 with multi-range transformers. In *Advances in Neu- 859*
810 *ral Information Processing Systems (NeurIPS)*, pages 860
811 6036–6049, 2021. 3 861
- 812 [55] Alexander S. Wannenwetsch, Margret Keuper, and 862
813 Stefan Roth. Probflow: Joint optical flow and uncer- 863
814 tainty estimation. In *IEEE International Conference 864*
815 *on Computer Vision (ICCV)*, 2017. 6 865
- 816 [56] Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, 866
817 Chao Weng, Yuexian Zou, and Dong Yu. Diffsound: 867
818 Discrete diffusion model for text-to-sound generation. 868
819 *IEEE/ACM Transactions on Audio, Speech, and Lan- 869*
820 *guage Processing*, 31:1720–1733, 2023. 2 870
- 821 [57] Hao Yang, Chunfeng Yuan, Li Zhang, Yunda Sun, 871
822 Weiming Hu, and Stephen J. Maybank. Sta-cnn: Con- 872
823 volutional spatial-temporal attention learning for ac- 873
824 tion recognition. *IEEE Transactions on Image Pro- 874*
825 *cessing*, 29:5783–5793, 2020. 3 875
- [58] Dianhao Zhang, Mien Van, Pantelis Sopasakis, and 826
Seán McLoone. An nmpc-ecbf framework for dy- 827
namic motion planning and execution in vision-based 828
human-robot collaboration, 2023. 1 829
- [59] Xikun Zhang, Chang Xu, and Dacheng Tao. Con- 830
text aware graph convolution for skeleton-based action 831
recognition. In *2020 IEEE/CVF Conference on Com- 832*
puter Vision and Pattern Recognition (CVPR), pages 833
14321–14330, 2020. 3 834
- [60] Tianhang Zheng, Sheng Liu, Changyou Chen, Jun- 835
song Yuan, Baochun Li, and Kui Ren. To- 836
wards understanding the adversarial vulnerability of 837
skeleton-based action recognition. *arXiv preprint 838*
arXiv:2005.07151, 2020. 3 839