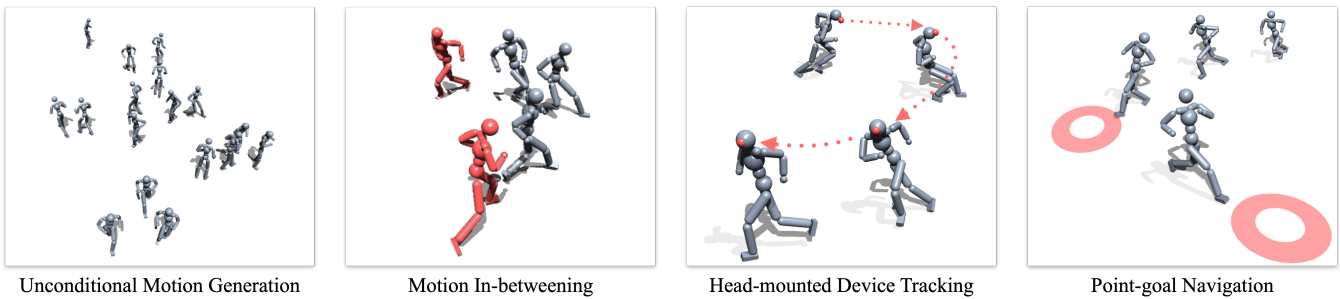# Versatile Physics-based Character Control with Hybrid Latent Representation

Jinseok Bae ⓘD, Jungdam Won ⓘD, Donggeun Lim ⓘD, Inwoo Hwang ⓘD, Young Min Kim[†] ⓘD

Seoul National University, South Korea

Unconditional Motion Generation    Motion In-betweening    Head-mounted Device Tracking    Point-goal Navigation

**Figure 1:** *Our hybrid latent representation serves as a versatile motion prior such that the agent can robustly adapt to a range of challenging downstream tasks.*

## Abstract

*We present a versatile latent representation that enables physically simulated character to efficiently utilize motion priors. To build a powerful motion embedding that is shared across multiple tasks, the physics controller should employ rich latent space that is easily explored and capable of generating high-quality motion. We propose integrating continuous and discrete latent representations to build a versatile motion prior that can be adapted to a wide range of challenging control tasks. Specifically, we build a discrete latent model to capture distinctive posterior distribution without collapse, and simultaneously augment the sampled vector with the continuous residuals to generate high-quality, smooth motion without jittering. We further incorporate Residual Vector Quantization, which not only maximizes the capacity of the discrete motion prior, but also efficiently abstracts the action space during the task learning phase. We demonstrate that our agent can produce diverse yet smooth motions simply by traversing the learned motion prior through unconditional motion generation. Furthermore, our model robustly satisfies sparse goal conditions with highly expressive natural motions, including head-mounted device tracking and motion in-betweening at irregular intervals, which could not be achieved with existing latent representations.*

**CCS Concepts**

*• Computing methodologies → Physical simulation; Reinforcement learning; Motion processing; Motion capture;*

## 1. Introduction

Recent advances in physics-based character control have shown remarkable progress in generating realistic human motions using simulated agents and reinforcement learning [YMG*23, BWL*23, HGW*23]. Hierarchical control strategies [MHG*18, MTA*20, KFPM21] offer a scalable approach to obtain and utilize a general motion prior, instead of manually selecting reference data tailored for each scenario. A hierarchical controller is trained in two phases: imitation learning and task learning. In the imitation learning phase, an agent constructs a latent embedding to serve as a scalable motion prior. During the task learning phase, the agent explores this latent space to discover optimal skills for the given scenario. Constructing a rich and structured latent space during the first phase is critical to enhancing training efficiency and motion quality during the task learning phase.

A standard method extracts the latent space of motion through the bottleneck layer of an autoencoder. For generative tasks, we

---

[†] Corresponding Author

often employ variational inference [WGH22, YSCL22, LCM*23], such that the latent representation matches a pre-defined probabilistic distribution, which is usually Gaussian. However, variational frameworks frequently suffer from posterior collapse [LTGN19], failing to capture details from the original dataset. Moreover, when combined with the hierarchical control, searching a latent vector within the continuous embedding space may be biased towards fulfilling the task reward and results in awkward motions that deviate from the original dataset. As a remedy, several works build the latent space to be a structured set of discrete vectors [ZZLH23, YSZ*23, GMJ*23]. However, action sequences selected from the discrete latent space are temporally discontinuous, leading to artifacts that can significantly degrade motion quality. Furthermore, these models face a critical trade-off: while a larger codebook increases the imitation policy's capacity, it complicates the exploration during the task learning phase.

We address these issues by integrating continuous and discrete latent space, namely *hybrid* latent representation. Unlike the continuous latent model, our model uses a discrete policy as a high-level strategy, ensuring that the generated motion thoroughly captures the diverse motion distribution in the dataset. At the same time, our agent consistently produces smooth motions, enhancing visual quality compared to the agent with discrete latent representation. We attribute this to our novel policy architecture, which augments discrete motion prior with the temporally coherent motion prior from the continuous latent space.

Additionally, we improve the effectiveness of the discrete latent representation using Residual Vector Quantization (RVQ) [LKK*22, ZLO*21]. Our experiments show that RVQ improves the latent model in both the imitation and task learning phases. During the imitation learning phase, RVQ significantly enhances the capacity of the imitation policy compared to a standard vector quantization layer. By employing the *quantizer dropout* technique [ZLO*21] of RVQ, our method allows an agent to efficiently reduce the training complexity of the task learning phase by constraining the action space. We additionally suggest a strategy to increase controllability in task learning by balancing the exploration and the exploitation of learned motion prior. While some works use RVQ for motion generation [GMJ*23, YSZ*23], their approaches are not suited for challenging control tasks like motion tracking with extremely sparse targets, *e.g.* head-mounted device tracking. Their models encode sequences of kinematic states over a temporal window, which limits precise control in highly dynamic scenarios where the agent must respond quickly to the goal. In contrast, we use RVQ to learn discrete embeddings of short-term motion dynamics in neighboring timesteps. The enhanced temporal precision enables our latent model to create a control policy that exhibits robust responsiveness to environmental changes.

In our experiments, we demonstrate the versatility of our model across various downstream tasks. Our agent consistently produces high-quality motion while discovering optimal motion combinations from the latent space. The model outperforms existing approaches in challenging scenarios with sparse goals, such as motion in-betweening over arbitrary time intervals. Furthermore, our agent generates highly natural motion without explicit reinforcement for motion style. We also show that our novel latent representation en-

ables the agent to easily solve tasks with simple reward designs, as seen in examples like head-mounted device tracking and point-goal navigation.

We summarize our contributions below.

- *Improved stability and quality*. Our hybrid representation offers integrated latent space that significantly enhances training stability and motion quality for downstream tasks.
- *Enhanced code usage and training*. We utilize Residual Vector Quantization (RVQ) to maximize the usage of a limited number of discrete latent codes while reducing the training complexity during the task learning phase.
- *Versatility*. We showcase our agent's capabilities across various downstream tasks, demonstrating that the proposed latent embedding can discover optimal motion combinations in challenging unseen target scenarios.
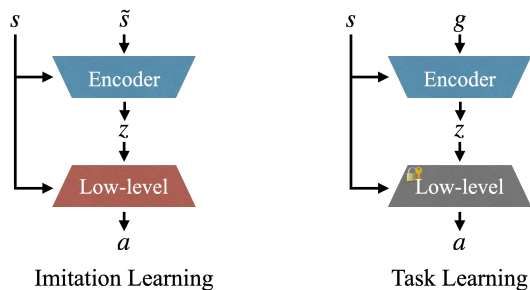
## 2. Related Works

### 2.1. Character Animation

Advancements in the data-driven approach demonstrated that virtual characters can learn natural, human-like motion from large datasets. With a sufficient number of samples from extensive motion datasets [IPOS13, MGT*19, HYNP20, LZL*24b], neural networks can understand patterns in human motion. Recent studies have demonstrated the ability to reconstruct complex interactions involving humans and objects [PXW*23, DD24, XWWG24, ZZS*24, BXP*22], humans and other humans [LZL*24a, STKB23, XZY*24], and humans within environments [JZL*24, WCL*22, ZZW*23, YTB*24].

On the other hand, physics-based control generates a sequence of actions relying on physics simulators to produce physically plausible, therefore highly realistic motion. Pioneering works generated high-quality motion of virtual characters by training imitation policies using deep reinforcement learning (DRL) [PALVdP18, XK21, PMA*21]. The approaches are further extended to generate complex physical interactions between agents [YKK*23] or between agents and objects [BWL*23, WLZ*23, HGW*23]. However, training DRL controllers from scratch for each scenario is not sample-efficient, making the process computationally demanding.

A stream of research suggests enhancing the efficiency of controllers by employing reusable skill prior that is built from motion datasets [MHG*18, MTA*20, KFPM21, WGH21]. These methods often use a hierarchical strategy. First, a low-level controller is trained, which extracts a latent space that encapsulates reference motions. Then, the high-level policy explores the latent space and outputs the latent vectors, which are then decoded into physical actions by the pre-trained low-level controller. A line of works have employed continuous latent space [LZCVDP20, WGH22, YMG*23, LCM*23, PGH*22, TKG*23, YSCL22, XXA*23], while some recent works [ZZLH23, YSZ*23] have leveraged discrete latent representation. However, it is not trivial to obtain stable yet rich latent space which allows structured exploration to adapt to novel scenarios with natural motion. We further examine the limitations of these approaches in Sec. 3

**Figure 2:** *Overview of the hierarchical control. (Left) At the first stage of imitation learning, an encoder and the low-level policy are jointly trained to learn latent space for a prior. (Right) During the task learning phase, high-level policy is trained to output proper latent action for the pretrained low-level policy.*

## 2.2. Latent Representation

Variational Autoencoders (VAE) [KW13] have been foundational for generative modeling, particularly in handling complex data distributions with continuous latent vectors. While VAEs have been successfully applied across various domains [CKD*15, GPB*18, TRG*22], they often face stability issues during training. This instability largely stems from the need to balance the trade-off between reconstruction fidelity and the regularization imposed by KL divergence [BVV*15, HMP*17]. As a partial remedy, one needs to adjust weights with proper scheduling, which requires complex tuning.

Building on the principles of VAE, Vector Quantized Variational Autoencoders (VQ-VAE) [VDOV*17, RVdOV19] introduce a significant evolution in the field of latent representations. By converting continuous data into distinct vectors, VQ-VAEs stabilize training and enhance output clarity, making them ideal for tasks like video and audio synthesis [GvdOL*19, YZAS21]. Recent studies have further advanced network architecture for VQ-VAE, mainly focusing on increasing the code usage [ZLO*21, YLH*23, YLK*21, LKK*22] and the efficiency of the quantization process [MMAT23, YLG*23]. However, VQ-VAEs have limitations, such as overfitting to small datasets and challenges in increasing quantization levels without compromising model efficiency. Additionally, when discrete codes from a VQ-VAE structure are used for hierarchical control [ZZLH23], increasing the codebook size directly enlarges the action space, which undesirably raises the training complexity.

## 3. Preliminaries

In this section, we briefly explain the base framework for hierarchical control (Sec. 3.1). Then, we review two types of motion priors, namely *continuous* (Sec. 3.2) and *discrete* (3.3) priors.

### 3.1. Hierarchical Control

Hierarchical control enables an agent to reuse pretrained motion priors in various types of downstream tasks. The training pipeline can be split into two stages: the first stage to train the low-level controller $\pi_{\text{low}}$, and the second stage to train the high-level controller $\pi_{\text{high}}$. In the first stage, $\pi_{\text{low}}$ is trained to output appropriate actions $a$ given a goal state $\tilde{s}$ sampled from reference motion clips. Here, the policy learning employs an encoder-decoder structure such that the agent can learn a compact yet expressive representation $z$ for each pair of $(s, \tilde{s})$. Specifically, the encoder maps the requested transition $(s, \tilde{s})$ to a latent embedding $z$, while $\pi_{\text{low}}$ decodes $s$ and $z$ to output action $a$. Once trained, the parameters of the decoder are frozen.

At the second stage, the framework trains a new encoder, the high-level policy $\pi_{\text{high}}$, that outputs desirable latent vector $z$ given $s$ and the task-specific goal $g$. Therefore, the pretrained latent space must faithfully represent the dataset, while allowing $\pi_{\text{high}}$ to safely explore the latent space. We illustrate the overall process of the hierarchical control in Figure 2.

### 3.2. Continuous Prior

Regularization technique can further improve a basic encoder-decoder structure (Figure 2) with comprehensive latent space. Previous works [WGH22, YSCL22, LCM*23] suggest a regularization technique inspired by the VAE structure [KW13]. Specifically, these methods allow the encoder to output a posterior distribution $q_\phi$, and sample a latent vector $z$ from $q_\phi$. Conventionally, the *Kullback-Leibler* (KL) divergence loss $\mathcal{L}_{\text{KL}}$ regularizes $q_\phi$ to match a desired continuous distribution $p_\theta$, which is either a standard normal distribution $\mathcal{N}(0, I)$ or a distribution parameterized by an additional prior network $f_\theta$ that is conditioned to $s$. We summarize the behavior of $\pi_{\text{low}}$ and $\mathcal{L}_{\text{KL}}$ in the first stage as

$$a \sim \pi_{\text{low}}(a|s, z), \text{ where } z \sim q_\phi(z|s, \tilde{s}),$$
$$\mathcal{L}_{\text{KL}} = \text{KL}(q_\phi(z|s, \tilde{s})||p_\theta(z|s)). \tag{1}$$

During the high-level policy training, a previous work [LCM*23] empirically suggests outputting a residual value $u$ relative to the mean of the prior distribution $\mu_p = f_\theta(s)$ produces better results than directly outputting $z$ on downstream tasks. Then the latent action $z$ from $\pi_{\text{high}}$ in the second stage can be written as

$$z = \mu_p + u, \text{ where } u \sim \pi_{\text{high}}(u|s, g). \tag{2}$$

We denote this type of method as a *continuous* prior. Although the *continuous* prior allows the agent to generate smooth trajectories, the agent may experience posterior collapse, which makes the latent representation $z$ less informative. Furthermore, the model requires extensive hyperparameter tuning, such as noise level for proper exploration, to stabilize the high-level policy training. Otherwise, we find that the agent with *continuous* prior easily generates unnatural motion that deviates from the valid motion space in the embedded prior, especially when the reward for the downstream task does not consider motion quality (*e.g.*, *point-goal navigation* that employs a reward for the final position only).

### 3.3. Discrete Prior

Previous works [ZZLH23, YSZ*23] suggest the vector quantization approach [VDOV*17] to overcome drawbacks of *continuous* motion prior. We categorize these models as *discrete* models and

explain the behavior of their motion prior based on Neural Categorical Prior (NCP) [ZZLH23]. The posterior network $f_\phi$ initially outputs a continuous latent vector $z$, and then quantizes it as $\bar{z}$ by selecting the nearest code $e^*$ from the codebook $\mathcal{B}$. This can be viewed as a mapping process that forces the quantizer to commit to a specific vector in the codebook for each input. The commitment loss $\mathcal{L}_{\text{commit}}$ in the VQ-VAE structure [VDOV*17] attaches the latent vectors to a code while updating the codebook parameters:

$$e^* = \operatorname*{argmin}_{e \sim \mathcal{B}} \|e - z\|^2, \text{ where } z = f_\phi(s, \tilde{s}),$$
$$\mathcal{L}_{\text{commit}} = \|z - \text{sg}(e^*)\|^2 + \|\text{sg}(z) - e^*\|^2, \tag{3}$$

where $\text{sg}(\cdot)$ denotes the stop-gradient function. Note that vector quantization assumes a uniform distribution as a prior distribution, so the KL divergence term becomes constant. Therefore, the training objective can alleviate KL divergence loss, which significantly accelerates the process without time-consuming hyperparameter tuning. Finally, the *straight-through* trick [VDOV*17] enables back-propagation through the quantized latent $\bar{z}$, which can be written as

$$\bar{z} = z + \text{sg}(e^* - z). \tag{4}$$

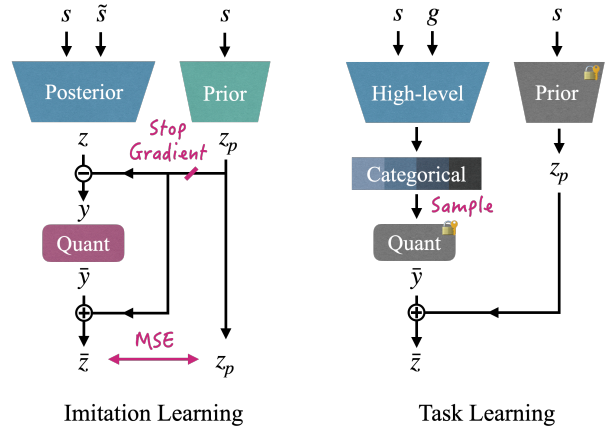Given $\bar{z}$ and $s$, the low-level controller $\pi_{\text{low}}$ predicts the action.

At the task learning phase, the *discrete* model treats the high-level policy $\pi_{\text{high}}$ as a discrete policy, which chooses the index of a desirable embedding $e$ from the pretrained codebook $\mathcal{B}$. Since $\pi_{\text{high}}$ explores a finite number of discrete actions, the model prevents $\pi_{\text{high}}$ from deviating largely from the pretrained motion prior. However, the discrete choice of the latent code triggers high-frequency jittering in motion since selected codes can be temporally discontinuous. Similarly, the motion generated from random prior sampling often encounters jerky movements. To employ a better prior distribution, *discrete* latent model may allocate an additional training stage for the prior network and use it as a regularizer for downstream tasks. However, we find that the prior network cannot imitate the dataset as much as the posterior network, especially when the motion dataset is large. The reported training curve from NCP [ZZLH23] also supports that the performances in the tasks are nearly the same with or without additional prior network.

## 4. Method

In this section, we detail our novel architecture, which contributes to smooth motion trajectories through efficient training across a variety of downstream tasks. We refer to the motion prior derived from our model as *hybrid* motion prior, as it benefits from both *continuous* and *discrete* priors. We first describe the *hybrid* motion prior learned from the imitation phase (Sec. 4.1), followed by a method to enhance *discrete* prior using Residual Vector Quantization (Sec. 4.2). Lastly, we elaborate on the improved training procedure for the imitation learning and the task learning phases (Sec. 4.3).

### 4.1. Hybrid Prior using Marginal Quantization (MQ)

As shown in Figure 3, the prior network outputs continuous latent vector $z_p$ to maintain the expressivity and temporal smoothness,



**Figure 3:** *Network architecture for hybrid motion prior. (Left) During the imitation learning phase, posterior and prior network output latent vector $z$ and $z_p$ given the current state $s$ and target state $\tilde{s}$. Next, the network quantizes the marginal value of $y = z - z_p$ into $\bar{y}$, then reconstructs the final latent representation $\bar{z}$ by biasing the quantized value with $z_p$. Once trained, the parameters of the codebook and prior network are frozen. (Right) In the task learning phase, high-level policy outputs categorical distribution, and then samples the index for the pretrained codebook. The latent $\bar{z}$ is then reconstructed with the selected $\bar{y}$ and continuous vector $z_p$.*

while the posterior and high-level policy networks achieve stability by generating categorical outputs that are discretized. We introduce Marginal Quantization (MQ), a method that quantizes only the difference between the posterior and prior latent vectors. The process begins by computing the margin $y = z - \text{sg}(z_p)$, where $z$ and $z_p$ are the latent vectors from the posterior and prior networks, respectively. The margin $y$ is then quantized to $\bar{y}$, and the final latent vector produced is $\bar{z} = \bar{y} + \text{sg}(z_p)$. The attached quantization potentially avoids posterior collapses and encourages the posterior network to output informative continuous latent $z$ to achieve accurate imitation. Nonetheless, training the final value of $\bar{z}$ observes the original gradient calculated for the continuous output $z$ as the subtracted value of $\text{sg}(z_p)$ is added back at the end. This is an extension of the *straight-through* technique, originally developed to enable back-propagation during vector quantization [VDOV*17]. Attaching the discrete vector $\bar{z}$ to the continuous vector $z_p$ from the prior network, the final vector $\bar{z}$ stays temporally coherent throughout an episode during the task learning phase.

Additionally, we incorporate a margin-minimizing loss $\mathcal{L}_{\text{mm}}$ to align the posterior with the prior latent space, defined as $\mathcal{L}_{\text{mm}} = \|y\|^2 = \|\bar{z} - z_p\|^2$. Unlike $\mathcal{L}_{\text{KL}}$ in the *continuous* model, we empirically find that the overall performance is not sensitive to the weight multiplied by the additional loss $\mathcal{L}_{\text{mm}}$. We summarize the behavior of our imitation policy in Algorithm 1.

### 4.2. Enhanced Code Usage for Discrete Prior

As we develop versatile motion prior that can adapt to various challenging tasks, it is crucial to allow efficient exploration within

---

**Algorithm 1:** Action Prediction using MQ

---

**Network:** posterior network $f_\phi$, prior network $f_\theta$, quantizer $Q$, low-level controller $\pi_{\text{low}}$

**Definition:** stop gradient $sg(\cdot)$, latent $z$, margin $y$

**Input:** current state $s$, target reference state $\tilde{s}$

**Output:** action $a$, commitment loss $\mathcal{L}_{\text{commit}}$, margin-minimizing loss $\mathcal{L}_{\text{mm}}$

**Function** `ForwardActor`$(s, \tilde{s})$:

> $z \leftarrow f_\phi(s, \tilde{s})$    ▷ latent vector from posterior
> $z_p \leftarrow f_\theta(s)$    ▷ latent vector from prior
> $y \leftarrow z - sg(z_p)$    ▷ calculate margin
> $\bar{y}, \mathcal{L}_{\text{commit}} \leftarrow Q(y)$    ▷ calculate $\bar{y}$, $\mathcal{L}_{\text{commit}}$
> $\bar{z} \leftarrow \bar{y} + sg(z_p)$    ▷ reconstruct $z$
> $\mathcal{L}_{\text{mm}} \leftarrow \|\bar{z} - z_p\|^2$    ▷ calculate $\mathcal{L}_{\text{mm}}$
> $a \sim \pi_{\text{low}}(s, \bar{z})$    ▷ sample $a$

**return** $a$, $\mathcal{L}_{\textbf{commit}}$, $\mathcal{L}_{\textbf{mm}}$

---

the rich set of plausible motion space. Especially the discrete part of our hybrid representation leverages VQ-VAE [VDOV*17, ZZLH23], whose capacity, defined by the size of the codebook $\mathcal{B}$, is not fully utilized due to the low rate of code usage [MMAT23]. To address this, we employ hierarchical codebooks of Residual Vector Quantization (RVQ) [LKK*22, ZLO*21], which sequentially uses $N$ codebooks and retrieves discrete latent codes in an autoregressive manner. Our residual quantization process $Q$ for margin $y$ is represented as

$$\bar{y} = \sum_{n=1}^{N} e_n^*, \text{ where } e_n^* = \underset{e_n \sim \mathcal{B}_n}{\operatorname{argmin}} \|e_n - (y - \sum_{k=1}^{n-1} e_k^*)\|^2. \quad (5)$$

We find that even with the same total number of codes, a model with RVQ provides better imitation quality. In addition, we utilize three training techniques for RVQ used in the previous works [ZLO*21, AGL*22, YSZ*23, JCL*24]: *Exponential Moving Average (EMA) update* for optimizing the codes, *code reset* technique for eliminating unused codes, and the *quantizer dropout* to build priority on the multiple codebook system. Among the techniques, *quantizer dropout* creates a critical advantage for the task learning phase as further explained in Section 4.3.

### 4.3. Efficient Training for Policies

#### 4.3.1. Imitation Learning

Our training pipeline begins with the imitation learning phase where an agent learns latent motion prior via imitating the motion capture data. We find that conventional RL is less efficient at training complex latent models compared to simpler models, as shown in Figure 2. This is because latent models require additional terms, such as KL divergence or commitment losses, to regularize the latent space. To address this, we adopt *online distillation* [LCM*23] as a simple yet effective approach for training a low-level policy. In online distillation, a student policy is trained under the supervision from a pretrained expert policy. The online distillation enables the student policy to achieve highly accurate imitation, which would not be possible if trained from scratch with RL. Following the previous work [LCM*23], we first train an expert policy

using Proximal Policy Optimization (PPO) algorithm. As an expert policy, we use a simple encoder-decoder structure as illustrated in Figure 2. We explain further details about expert policy in the Appendix. Then we train an agent with supervised guidance from the pretrained expert imitation policy. We add the commitment loss and the margin-minimizing loss to the original training objective [LCM*23] as

$$\mathcal{L} = \mathcal{L}_{\text{action}} + \alpha\mathcal{L}_{\text{reg}} + \beta\mathcal{L}_{\text{commit}} + \gamma\mathcal{L}_{\text{mm}}, \quad (6)$$

where $\mathcal{L}_{\text{action}}$ refers to supervised guidance from expert action $\|a - a_{\text{expert}}\|^2$, while $\mathcal{L}_{\text{reg}}$ is a regularization term to minimize change between latent embeddings of neighboring frames. We additionally modify the regularization term as $\mathcal{L}_{\text{reg}} = \|\bar{y} - \bar{y}'\|^2 + \|z_p - z_p'\|^2$ where $\bar{y}'$ is the quantized margin and $z_p'$ is the prior latent vector cached from the previous time step.

Since our model employs RVQ, the latent representation better utilizes the limited number of codes compared to the conventional *discrete* model (Sec. 4.2). However, as a trade-off, RVQ exponentially increases training complexity during the task learning phase as our $\pi_{\text{high}}$ needs to explore the combinatorial action space composed of multiple codebooks. To mitigate this problem, we utilize *quantizer dropout* technique [ZLO*21, YSZ*23] while training the codebooks. Specifically, we randomly choose the maximum codebook index $M \in [1, N]$ where $N$ represents the total number of codebooks. We then only use $\{\mathcal{B}_1, \cdots, \mathcal{B}_M\}$ to reconstruct the target value, which is equivalent to Eq. (5), but with $N$ in the original equation substituted by $M$. As a result, our agent prioritizes codebooks with lower indices, thereby establishing a hierarchical structure within the overall set of codebooks.

#### 4.3.2. Task Learning

In the task learning phase, our agent efficiently navigates the latent space for motion prior and robustly adapts to each task scenario with a high-level policy $\pi_{\text{high}}$. The prioritization among codebooks allows $\pi_{\text{high}}$ to reduce the action space efficiently. $\pi_{\text{high}}$ can even perform only with the highest priority codebook, $\mathcal{B}_1$, in downstream tasks. This significantly improves over conventional discrete models, which must investigate all available codes to optimize the action space. Additionally, RVQ provides extra controllability during task learning. Empirically, we find that using more codebooks for $\pi_{\text{high}}$ allows the agent to exhibit more diverse motions, encouraging exploration during training. However, using too many codebooks can increase training complexity, so users can select the proper number to balance the results.

## 5. Experiments

We design several experiments to showcase the extreme versatility of our model with enhanced motion quality and training efficiency. We additionally provide the qualitative results with the supplementary video and further ablation studies in the Appendix.

### 5.1. Implementation Details

We train our agent in parallel environments on Isaac Gym simulator [MWG*21], with the internally implemented proportional

derivative (PD) controller. In each time step, the environment calculates state vector $s$, which can be written as

$$s = \{h, q_r, \dot{p}_r, \dot{q}_r, q, \dot{q}, p_{ee}\}, \tag{7}$$

where $h, q_r, \dot{p}_r, \dot{q}_r$ denote height, orientation, velocity, and angular velocity of the root in character coordinate, respectively. Also, $q, \dot{q}$ represents the position and velocity of joints, while $p_{ee}$ refers to the end-effector position in character coordinates. Also, we allow low-level policy $\pi_{low}$ to output action $a$ that acts as target joint positions $\hat{q}$ for PD control. Following the previous work [ZZLH23], we design $\pi_{low}$ to output the residual action relative to the current state of joint position $q$, which can be represented as $\hat{q} = q + a$. We train the models with LaFAN1 dataset [HYNP20] for imitation learning, which contains 3.5 hours of diverse motion capture data.
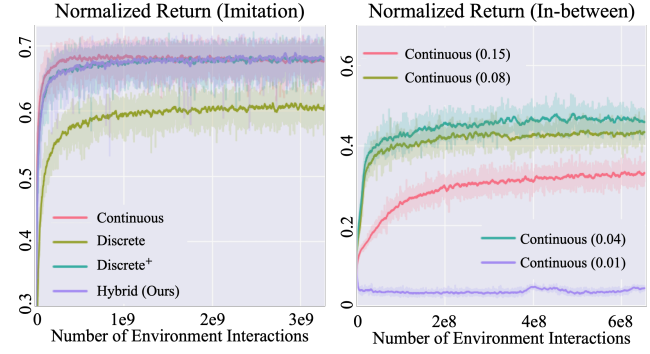
### 5.2. Baselines

We compare our *hybrid* latent space model with representative works that use latent space of motion prior. We adopt PULSE [LCK*23] and NCP [ZZLH23] from the released codes, and use them as a *continuous* and *discrete* models, respectively. We verify that our implementations perform best using most of the hyperparameters reported in the original papers. This indicates that our implementation faithfully reproduces the original versions of the previous works. Nonetheless, we carefully tune some hyperparameters for the baselines to achieve the best quality for imitation and task learning. Additionally, we find the *discrete* model can be further improved by incorporating Residual Vector Quantization (RVQ) layers. Therefore, we also implement the *discrete* latent model enhanced with RVQ, which we denote as *discrete*$^+$ model. We assign a codebook with the 8192 codes for the *discrete* model, while the RVQ models (both *discrete*$^+$ and *hybrid*) use the same total number of codes but in $N = 8$ codebooks with 1024 codes each. During the task learning, we manually tune the optimal hyperparameters for the exploration of $\pi_{high}$ in *continuous* model.
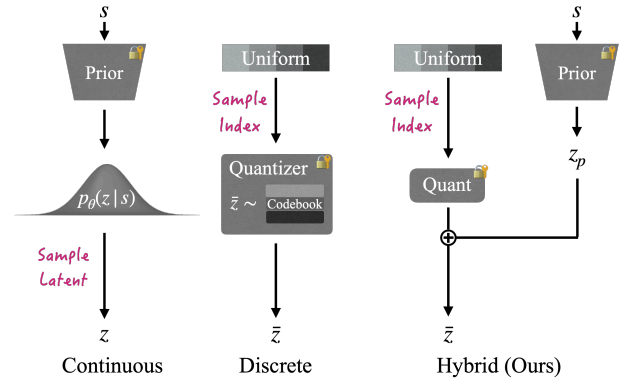
### 6. Discussions

### 6.1. Motion Imitation

We first examine the capacity of the various latent representations by testing the high-level policy to imitate the reference motions in the dataset. Figure 4 shows the quality of imitation in terms of the reward function. As suggested by previous works [LCM*23], continuous latent space can closely capture a diverse set of motions. In contrast, a finite set of discrete vectors often cannot comprehensively cover the dataset despite increased stability in various downstream tasks. Although increasing the number of codes can slightly improve imitation quality, we cannot achieve a comparable level of imitation with the *discrete* model. We find that using RVQ can overcome the prevalent limitation of *discrete* model, significantly improving the imitation quality. Despite using the same number of codes with *discrete* model, RVQ (*discrete*$^+$, *hybrid*) achieves the imitation quality on par with *continuous* latent space by efficiently facilitating the finite number of codes. Similar to *discrete* latent model, we find the imitation quality to be less accurate when we use the hybrid latent model with the standard VQ layer. We further provide a comparison between the hybrid latent model using RVQ and the standard VQ in the Appendix.
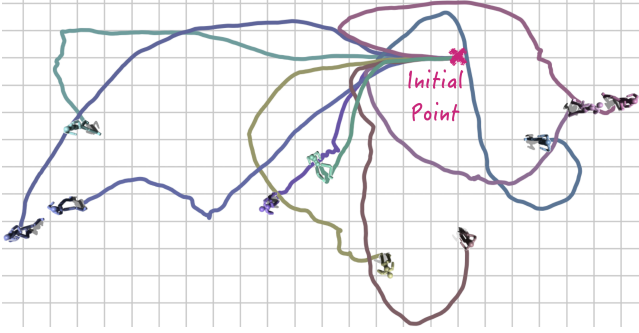


**Figure 4:** *Training curves during the imitation and the task learning phases, respectively. (Left) Discrete model with simple vector quantization cannot achieve accurate imitations. In contrast, discrete$^+$ model exhibits a similar training curve to ours. (Right) The value inside the parenthesis indicates experimented hyperparameters for the exploration rate. Continuous latent models are susceptible to the hyperparameters.*



**Figure 5:** *Prior sampling strategies for three different architectures. While continuous model samples are latent directly from the parameterized Gaussian distribution $p_\theta$, discrete and hybrid models uniformly sample the index from the pretrained codebook.*

### 6.2. Unconditional Motion Generation

We additionally evaluate the richness of the learned latent space by unconditional motion generation. In this setup, we randomly sample the latent vector $z$ from the prior distribution of each model, as depicted in Figure 5. The prior distribution of the *continuous* model is Gaussian distribution parameterized by its prior network. In contrast, the *discrete* and *hybrid* models leverage VQ-VAE [VDOV*17], which utilizes a uniform distribution as its prior distribution. We visualize the generated trajectories from 10 agents started with the same initial states in Figure 6. Surprisingly, our agents produce extremely smooth motions throughout the rollouts, even though the random sampling for the latent vector is completely independent at each timestep. We attribute this to our continuous latent vector $z_p$ from the prior network, which contributes to the temporal coherence. Due to the absence of such a continuous

**Figure 6:** *Unconditional motion generation from our model. We visualize trajectories from 30 seconds of 10 random rollouts with identical initial states.*

**Table 1:** *Evaluation on the unconditional motion generation. For discrete[+] and hybrid model with RVQ, we indicate number of the used codebook inside parentheses. We additionally mark the second best value with underlines.*

| Models | Distance ($d$) ↓ | Filtered (%) ↓ | Coverage (%) ↑ |
|---|---|---|---|
| *Continuous* | **0.798** (0.089) | 0.767 | 26.71 |
| *Discrete* | 2.048 (0.066) | 1.504 | 24.62 |
| *Discrete*[+] (1) | 1.516 (0.077) | 2.031 | 14.82 |
| *Discrete*[+] (2) | 1.708 (0.073) | 2.375 | 16.16 |
| *Discrete*[+] (3) | 1.856 (0.099) | 2.608 | 16.83 |
| *Hybrid* (1) | <u>0.820</u> (0.085) | <u>0.650</u> | 25.34 |
| *Hybrid* (2) | 0.883 (0.075) | **0.629** | <u>27.30</u> |
| *Hybrid* (3) | 0.935 (0.087) | 0.807 | **28.73** |

latent vector, we find that the *discrete* models largely suffer from high-frequency jittering.

We provide quantitative evaluations for the generated motions using the motion matching approach [PGH*22] in three metrics. The *distance d* measures the similarity between the generated and reference motions regarding the Euclidean distance between state transitions. For every pair of $(s, s')$, where $s$ and $s'$ are the current and the past proprioceptive states, we find the distance to the nearest state transitions $(s_m, s'_m)$ in the dataset $\mathcal{M}$

$$d(s, s') = \min_{(s_m, s'_m) \in \mathcal{M}} \|s - s_m\|^2 + \|s' - s'_m\|^2. \qquad (8)$$

Following the previous work [PGH*22], we normalize each dimension of the state vector with the mean and standard deviation calculated from the dataset. We also set a threshold value; if $d$ exceeds this threshold, we regard the queried transition as invalid and filter the calculated distance from the final statistics. We also report the *filtered rate* to show the ratio of valid motion in generations. Finally, the *coverage* counts the total number of visited transitions among datasets, representing the generated motion's diversity.

We measure 800,000 transitions, initializing each 30-second episode randomly from a state sampled from the reference. For models employing RVQ, we experiment with varying numbers of codebooks as described in Sec. 4.3.

The results in Table 1 indicate that our *hybrid* representation can generate plausible yet diverse motions from sampling the la-

tent variables. While it is expected that *continuous* models generate motions that best match the reference motion, our *hybrid* model can synthesize high-quality motion showing comparable trajectory distances. While RVQ also enhances the performance of *discrete* models, our proposed technique with MQ plays an essential role in generating plausible motions. Significantly, ours outperforms the *continuous model* regarding the validity and diversity of generated motions. Our model additionally offers user controllability for balancing the quality and the diversity of motion with RVQ layers, which the *continuous* model cannot provide. The user can increase the number of codebooks to enhance the diversity while sacrificing minimal imitation quality. This result also implies extra controllability for task learning. For example, if the high-level policy employs a large number of codebooks, the agent can achieve a high level of exploration. Conversely, using fewer codebooks allows the agent to behave more safely with smoother trajectories, resulting in more stable training. By balancing these properties, we can determine the optimal number of codebooks for each downstream task.

### 6.3. Downstream Tasks

Our hybrid model provides versatile latent space that an agent can efficiently navigate and robustly adapt to diverse downstream tasks. We deliberately demonstrate challenging scenarios where agents must deal with temporally or spatially sparse goal conditions. We further provide detailed explanations, including the reward functions, in the Appendix.
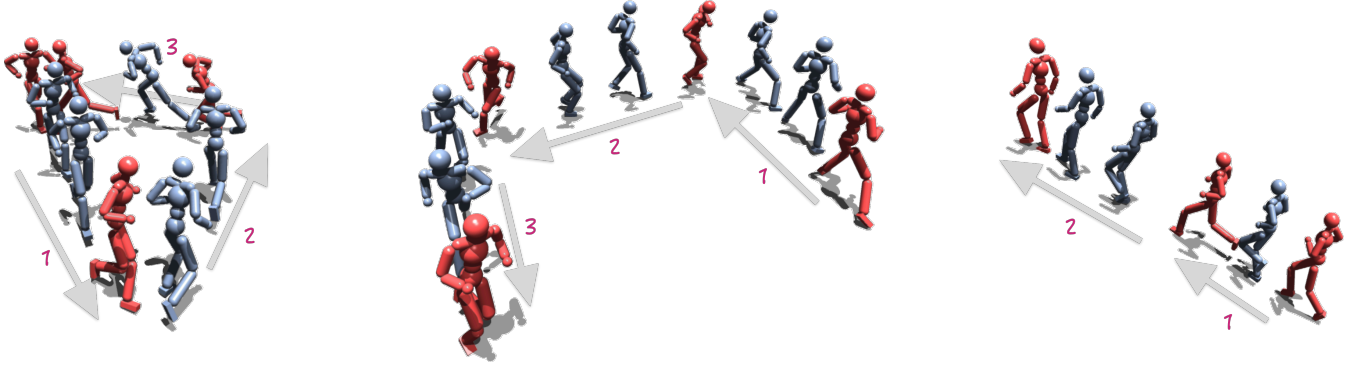
#### 6.3.1. Motion In-betweening

As our first downstream task, we present a motion in-betweening example in which agents must interpolate temporally sparse keyframes throughout the episodes. During the training, we randomly select zero-velocity poses as keyframes from the original motion sequence. Then, we provide the agent with a task goal vector $g$ based on the two keyframes for each local interval. We enhance the goal observation by augmenting the keyframe states with the time-to-arrival information for the next keyframe [HYNP20, GJW22]. We train our agent to mimic the ground-truth trajectories of the reference motion sequence using a full-body imitation reward [WGH20].

Our agent successfully interpolates keyframes even in novel scenarios that were not experienced during the training. To compare our model with the baselines, we quantitatively evaluate the generated motions with the conventional body tracking metrics [LCM*23]: mean per joint position error (MPJPE), MPJPE in global coordinates (G-MPJPE), root velocity error, and root acceleration error. The first two metrics indicate the motion fidelity compared to the ground truth trajectory, whereas the last two metrics represent the smoothness of motion. Table 2 summarizes the results of 2048 episodes from the test set. The result demonstrates that our model outperforms all other baselines. Notably, the performance of the *continuous* model reveals its failure to generate plausible motions, even with a dense full-body reward calculated from the reference trajectory.

Additionally, we empirically verify that the high-level policy of the *continuous* model is overly sensitive to hyperparameters during training. For example, as shown in Figure 4 (right), the agent

**Table 2:** *Evaluation on the motion in-betweening results. We randomly sample 5 keyframes from a 4 second-length motion sequence, i.e., the average time intervals between keyframes are 1 second. We report two sets of statistics; evaluated only with the keyframes ("Only Keyframes") and with full frames ("Full Frames"). The units for MPJPE and G-MPJPE are mm, while units for velocity and acceleartion error are mm/frame and mm/frame$^2$, respectively. We additionally mark the second best score with underlines.*

| Models | Only Keyframes | | | | Full Frames | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | MPJPE ↓ | G-MPJPE ↓ | Vel Error ↓ | Accel Error ↓ | MPJPE ↓ | G-MPJPE ↓ | Vel Error ↓ | Accel Error ↓ |
| *Continuous* | 157.3 (109.8) | 414.1 (540.3) | 48.44 (19.92) | 30.49 (10.95) | 181.8 (89.74) | 461.6 (502.9) | 43.79 (18.83) | 29.61 (10.87) |
| *Discrete* | 111.1 (88.59) | **304.9** (362.3) | 36.22 (18.24) | 16.02 (6.771) | 132.5 (79.59) | 348.7 (339.9) | 31.74 (16.57) | 15.69 (6.825) |
| *Discrete*$^+$ | <u>110.5</u> (88.11) | 329.5 (346.4) | <u>35.02</u> (18.71) | <u>13.08</u> (6.294) | <u>132.0</u> (77.24) | <u>369.8</u> (312.4) | <u>30.76</u> (16.70) | <u>13.26</u> (6.311) |
| *Hybrid* | **103.3** (87.98) | <u>316.7</u> (392.7) | **34.09 (19.97)** | **12.63 (6.323)** | **121.1 (79.35)** | **350.7 (356.6)** | **29.08 (17.21)** | **12.46 (6.176)** |



**Figure 7:** *Results of motion in-betweening example. Our simulated character (blue) generates plausible trajectories given zero-velocity keyframes (red). Arrows and numbers indicate the temporal flow of each episode.*

engages in excessive exploration when the noise level of the policy is set too high, while it focuses too much on exploitation when the noise level is set too low. In contrast, models with a discrete prior do not require such modification to balance exploration and exploitation. This implies that the discrete latent representation significantly enhances training efficiency for the task-learning phase. Among the models with discrete motion prior, our model generates minimal motion artifacts compared to *discrete* and *discrete*$^+$ models, given low values for the velocity and acceleration errors. We visualize the qualitative results from our model in Figure 7 and the supplementary video. Unlike the baselines, our agent consistently performs natural transitions between keyframes, even when the agent must rapidly shift the direction. Furthermore, our agent demonstrates natural motions in challenging scenarios, such as when it is required to walk backward.

### 6.3.2. Head-mounted Device Tracking

In addition, we explore a challenging task where an agent reconstructs full-body motion given spatially sparse tracking targets. In this setup, only the head trajectories are provided as a conditioning signal for the agents. The goal observation is provided as the state difference between the target and simulated head positions. The agent additionally observes future head trajectories for four additional time steps to resolve inherent ambiguity of the extremely sparse observation. During training, we calculate the reward based on how closely the agent's head state follows the target trajectory.
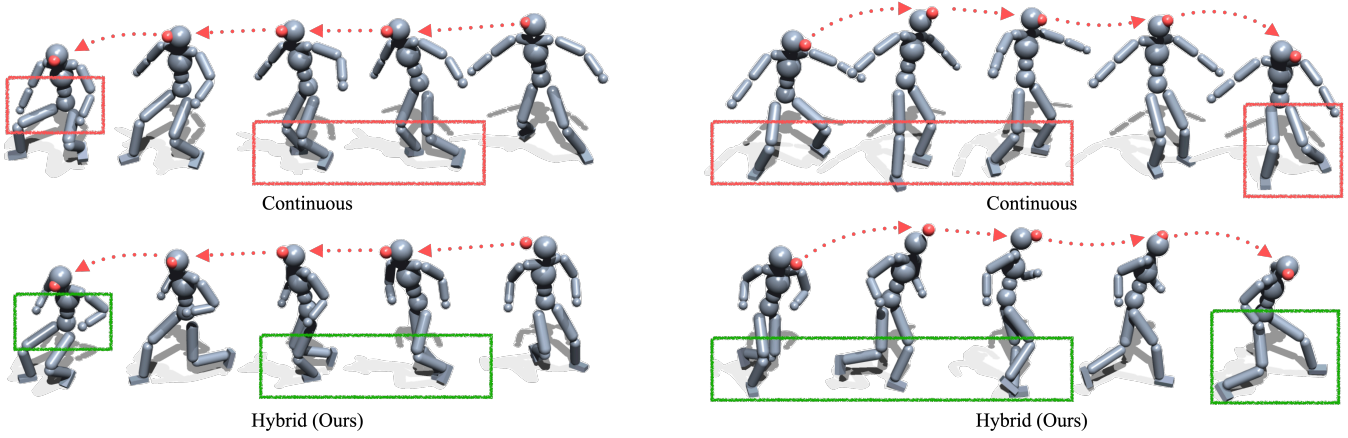
Similar to the temporally sparse setting, we find that our agent successfully discovers the optimal set of motions for head-mounted

**Table 3:** *Evaluation on the head-mounted device tracking. Similar to the data presented in Table 2, we utilize body tracking metrics to quantitatively evaluate the quality of motion.*
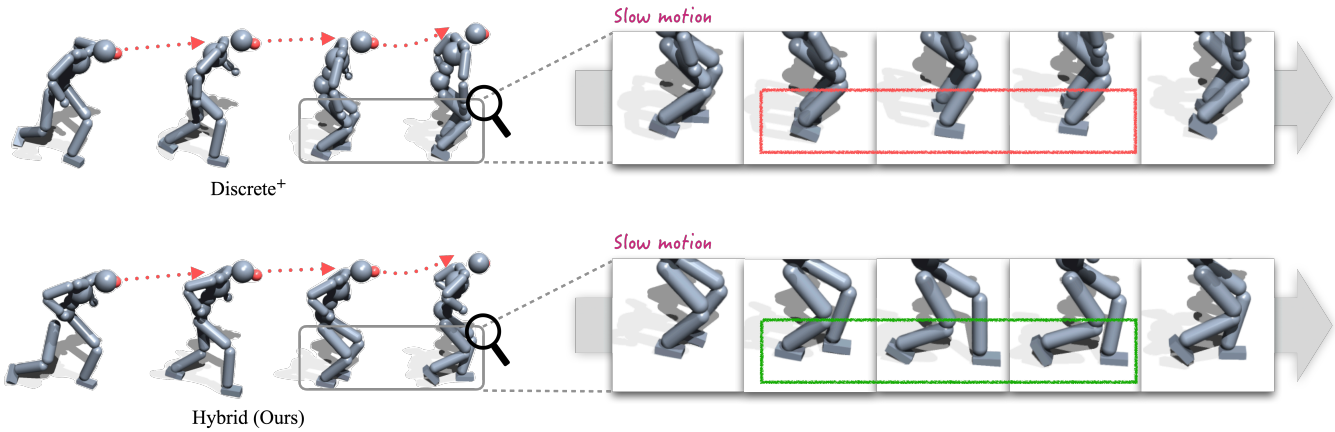
| Models | MPJPE ↓ | G-MPJPE ↓ | Vel Error ↓ | Accel Error ↓ |
| --- | --- | --- | --- | --- |
| *Continuous* | 168.7 (43.64) | 174.3 (108.5) | 29.10 (15.03) | 23.53 (11.22) |
| *Discrete* | <u>100.4</u> (47.91) | **110.9** (107.1) | 22.40 (10.55) | 15.34 (6.102) |
| *Discrete*$^+$ | 103.7 (50.39) | 116.1 (144.7) | 21.42 (10.26) | <u>14.19</u> (5.700) |
| *Hybrid* | **96.90** (57.88) | <u>112.9</u> (141.6) | **20.46** (11.57) | **13.48** (5.930) |

device tracking. Quantitative results in Table 3 and qualitative results in Figure 8 consistently demonstrate that our model effectively reuses the pretrained motion prior. The *hybrid* latent space enables the agent to track the head trackers more accurately than the *continuous* latent space, indicating enhanced training efficiency. Surprisingly, the gap between ours and the *discrete* models becomes more apparent in this scenario compared to the motion in-betweening. We attribute this to the robustness of our model under environments with sparse rewards. As shown in Figure 9, our model maintains a natural gait pattern, while *discrete*$^+$ model fails to do so. As clearly demonstrated with the supplementary video, *discrete* models often experience highly jittery movements as exploring the quantized latent space to fulfill the sparse goal. Such movement significantly degrades the visual quality of the generated motion. In contrast, our model does not suffer from these artifacts and consistently adheres to the learned motion priors.

**Figure 8:** *Qualitative comparison between continuous and hybrid models on head-mounted device tracking example. Our agent faithfully adheres to the learned motion prior while continuous model cannot.*



**Figure 9:** *Qualitative comparison between discrete$^+$ and hybrid models on head-mounted device tracking example. While discrete$^+$ model sometimes violates natural gaits, our model consistently exhibits realistic locomotion.*
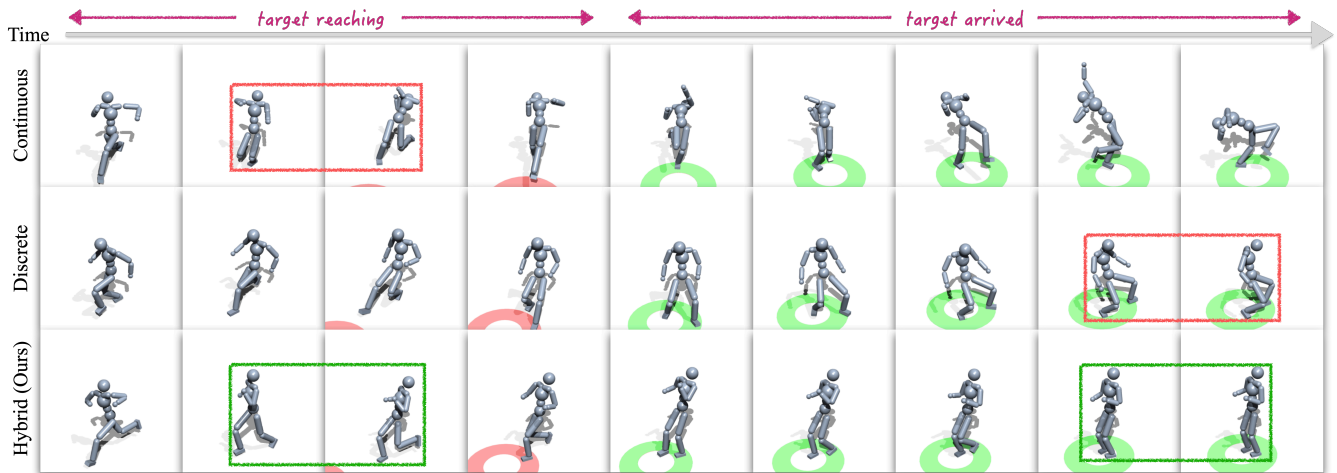
### 6.3.3. Point-goal Navigation

Lastly, we demonstrate a point-goal navigation example, which demonstrates a setting with spatio-temporally sparse goals. We provide the relative position to the goal in the ground plane as a goal observation. We employ *target location* reward that measures the relative distance and the heading velocity, similarly with the previous approach [PALVdP18]. Once the agent arrives at the goal, we give the maximum value for the reward to prevent meaningless motions in the goal location.
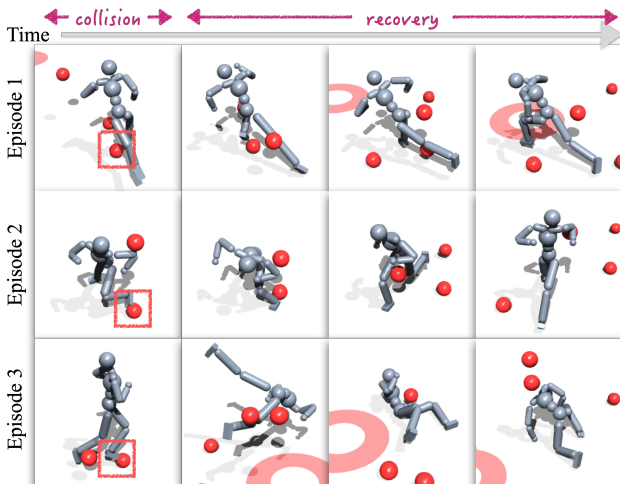
In this scenario, we can observe a similar tendency to the results from the other tasks. Note that this example does not include any ground-truth motions for additional guidance, and therefore the agent motions contain the pure exploration results to fulfill the task. The agents receive the maximum reward as long as they reach the goal, and there is no means to penalize unnatural motions. As demonstrated in Figure 10 and the accompanying video, only our *hybrid* model exhibit natural idle motions, whereas other models generate awkward movements. The result shows that our model can

faithfully generate natural motions with a versatile set of motion prior. Other works require additional reward designs to encourage desirable idle states, which involve non-trivial hand-crafted terms in the reward function.

We also test the robustness of our policy against physical perturbations. We introduce unexpected disruptions in the navigation environment and assess the agent's ability to fulfill the goal stably. To evaluate such ability, we test the policy trained on plain point-goal navigation environments in a zero-shot manner. For the external perturbations, spherical objects with a radius of 10 cm and a mass of 4.2 kg are thrown toward the agent every 0.3 seconds. The objects are thrown from a random location 1 meter away from the characters, with velocities randomly selected between 10 and 20 m/s. In Figure 11, we visualize the three different episodes of agents navigating with random perturbations. We find that our agent successfully maintains balance after collisions. Even when knocked down by severe impacts, as illustrated in the third row of

**Figure 10:** *Qualitative comparison between baselines (continuous, discrete) and hybrid model on the point-goal navigation example. Continuous model prioritizes speed over motion quality to maximize the task reward, which is higher when the goal is reached quickly. On the other hand, discrete model produces plausible motion when reaching the target location but struggles to maintain natural idle motion once it arrives.*



**Figure 11:** *Zero-shot adaptation to unexpected perturbations in a point-goal navigation task. Every 0.3 seconds, obstacles (red spheres) are thrown at the agent from random directions and velocities. Red boxes highlight the key moment of the collision that triggers the largest perturbation in an episode. Our agent consistently maintains natural movement despite significant disruptions.*

the figure, our agent recovers without deviating from the learned motion distribution.

These results demonstrate that our hybrid latent representation effectively handles test-time variations. We attribute this adaptability to the residual design, where the high-level policy modifies only a marginal portion of the latent vector produced by the prior network. This design enhances the agent's ability to respond to phys-

ical perturbations with greater stability. An ablation study from previous work [LCM*23] also supports that such residual modeling enhances performance at the task learning phase. Additional evidence is provided by the unconditional motion generation discussed in Sec. 6.2. As shown in Figure 5, our *hybrid* model uses a continuous latent vector as an offset to the discrete latent vector, which produces much more stable motions compared to *discrete* models.

## 7. Conclusion

In this paper, we present a novel architecture that utilizes integrated latent representation for the reusable motion prior. We suggest a hybrid representation that discretizes only the difference between the posterior and prior latent vectors. This approach fully utilizes the expressive yet temporally coherent continuous latent vector from the prior network during the task learning phase. In addition, we employ Residual Vector Quantization to improve the code usage and the training efficiency with the categorical latent prior. In the experiments, we demonstrate that our model is versatile in various scenarios. We show that *hybrid* latent space for motion prior enables efficient training in the task learning phase, allowing agents to exhibit temporally smooth motions consistently.

However, our models did not fully exploit the possible technical enhancement for recent variations of quantization or training pipelines. We believe state-of-the-art techniques such as finite scalar quantization [MMAT23] or group RVQ [YLH*23] can alleviate inherent limitations of vector quantization, such as overfitting problem [WF20]. Additionally, extending the hybrid latent model to alternative training pipelines, such as model-based approaches [YSCL22, YSZ*23], could enable a more comprehensive comparison with other state-of-the-art methods, providing valuable insights.

We suggest a few future directions to extend this work. Extending our hybrid latent model to the multi-subgoal scenarios could be one interesting research direction. For instance, future work could adapt the original PPO algorithm to incorporate a multi-task training scheme, as demonstrated in prior work [XSZK23]. Moreover, our model can be helpful in the motion composition task. By splitting the discrete latent space part-wise, one can achieve a dynamic combination of the part skills while maintaining the whole-body consistency with the continuous latent vector from the prior network.

## Acknowledgements

## Appendix

### A. Ablation Study

In this section, we provide an ablation study on our *hybrid* model.

### A.1. Effectiveness of Residual Vector Quantization

We first demonstrate the effectiveness of Residual Vector Quantization (RVQ) from the perspective of code usage and training efficiency. For comparison, we additionally train an ablated model where we use a conventional vector quantization layer instead of RVQ, namely *hybrid-vq*.
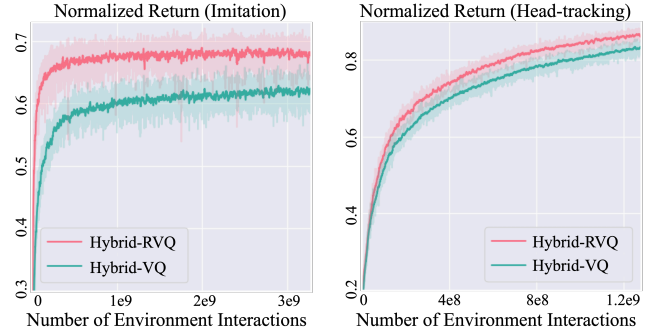
As shown in the training curve in Figure 12, *hybrid-vq* model fails to achieve a sufficient level of imitation. This implies that RVQ contributes to the enhanced code usage for the latent representation given a finite number of discrete codes.

Moreover, during task learning, we find that the *hybrid-vq* model exhibits lower performance in terms of normalized return compared to our model. Additionally, the *hybrid-vq* model requires more training time due to the computational burden of the argmin operation (Eq. (3)) with a larger discrete action space. In contrast, our model with RVQ handles this more efficiently.

### A.2. Number of Codebooks for Task Learning

In unconditional motion generation, we show that our *hybrid* model can balance the quality and the diversity of motion by controlling the number of codebooks used for prior sampling. We provide an additional evaluation to clearly demonstrate this balance in Table 4.

Along with this result, RVQ offers additional controllability through the number of active codebooks during the task learning



**Figure 12:** *Training curves for imitation learning (left) and task learning (head-mounted device tracking, right). We allocate 8192 codes for the hybrid model with a single VQ layer (Hybrid-VQ), while setting RVQ with 8 codebooks each containing 1024 codes for our main model (Hybrid-RVQ).*

**Table 4:** *Evaluation on the unconditional motion generation. The value inside the parenthesis indicates number of the used codebooks. We additionally mark the second best value with underlines.*

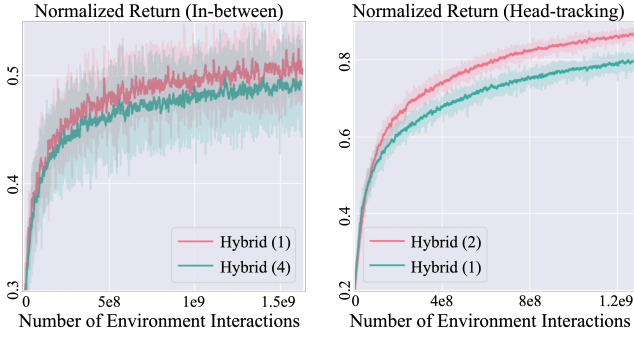| Models | Distance ($d$) ↓ | Filtered (%) ↓ | Coverage (%) ↑ |
|---|---|---|---|
| *Hybrid* (1) | **0.820** (0.085) | <u>0.650</u> | 25.34 |
| *Hybrid* (2) | <u>0.883</u> (0.075) | **0.629** | 27.30 |
| *Hybrid* (3) | 0.935 (0.087) | 0.807 | 28.73 |
| *Hybrid* (4) | 0.963 (0.090) | 0.738 | 29.29 |
| *Hybrid* (5) | 0.980 (0.092) | 0.774 | 29.94 |
| *Hybrid* (6) | 0.975 (0.081) | 0.735 | 30.30 |
| *Hybrid* (7) | 1.012 (0.083) | 0.772 | <u>30.83</u> |
| *Hybrid* (8) | 1.017 (0.089) | 0.749 | **30.99** |

phase. We show the result from different numbers of codebooks employed for task learning in Figure 13. As demonstrated in the results for motion in-betweening, using fewer codebooks leads to less diverse but more precise motion, allowing the high-level policy to be trained more stably. Conversely, in the tracking example, utilizing more codebooks can sometimes yield better results by enabling the high-level policy to explore the action space more actively during training. However, employing too many codebooks may excessively enlarge the action space, increasing the training complexity for task learning. Balancing this trade-off, we select the optimal number of active codebooks for each downstream task, achieving the best performance among the baselines.

### B. Imitation Learning Details

In this section, we provide an additional description of the imitation learning phase. We first elaborate on the further details about the imitation learning phase (Sec. B.1). Then, we provide additional information on the expert policy imported for training our imitation policy (Sec. B.2).

### B.1. Training

In Figure 2, our imitation policy includes a posterior network that outputs a latent vector $z$ given the current simulation state $s$ and the

**Figure 13:** *Training curves for motion in-betweening (left) and head-mounted device tracking (right). We indicate a number of active codebooks used for training in the parenthesis.*

target state $\tilde{s}$ sampled from the dataset. As represented in Eq. (7), the current simulation state $s$ primarily encodes the agent's proprioceptive state. In contrast, the target state $\tilde{s}$ represents the tracking objective. Instead of directly using the target state vector from the dataset, we represent $\tilde{s}$ as the state difference between the current state and the target state, following the approach in previous work [LCK*23, LCM*23]. Specifically, we first calculate the kinematic states for the $J$ rigid bodies of the agent, which include positions $p \in \mathbb{R}^{J \times 3}$, velocities $v \in \mathbb{R}^{J \times 3}$, orientations $\theta \in \mathbb{R}^{J \times 4}$, and angular velocities $\omega \in \mathbb{R}^{J \times 3}$. Simultaneously, we extract the same configurations from the target state of the reference motions, which we represent as $[\tilde{p}, \tilde{v}, \tilde{\theta}, \tilde{\omega}]$. Then, we calculate the state difference as follows:

$$\tilde{s} = [p - \tilde{p}, v - \tilde{v}, \theta \ominus \tilde{\theta}, \omega - \tilde{\omega}], \qquad (9)$$

where $\ominus$ denotes the operation to measure the difference between quaternions.

To improve sample efficiency, we also employ the early termination technique [PALVdP18], which is to terminate an episode when an agent violates a certain condition. For our imitation learning, we set an early termination condition that properly prevents large deviations from the ground-truth trajectories. Following a similar approach as in previous work [LCK*23, LCM*23], we set a termination condition based on the body positions. Specifically, we terminate an episode when the maximum distance from each body part of the agent to the corresponding body part in the target pose exceeds 0.5 meters, *i.e.,* ensuring that all body positions remain within 0.5 meters of the ground truth.

### B.2. Expert Policy

Following the previous work [LCM*23], we employ *online distillation* for training the imitation policy (Sec. 4.3). Unlike conventional Reinforcement Learning (RL), online distillation trains policy with supervised learning, resulting in much faster training compared to RL. To utilize online distillation, we need to train an expert policy in advance.

Basically, any method that achieves a reliable level of imitation can be adopted for training the expert policy. For example, one can employ an ensemble policy as an expert policy following previous works [LCK*23, LCM*23] that have achieved 100% imitation for extremely large datasets like AMASS [MGT*19]. However, we empirically find that our target dataset, LaFAN1 [HYNP20], does not require such a complex process for training the expert policy. Instead, we choose a simpler encoder-decoder structure (Section 3.1) to serve as the architecture for the expert policy. Due to the absence of the regularization on the latent space, this simple architecture cannot learn an effective latent space that is optimal for the high-level policy. However, we find that this type of architecture supports stable training of policy since such model does not require to optimize additional loss terms such as *Kullback-Leibler* loss in Eq. (1) or *commitment* loss in Eq. (3).

To train an expert policy, we leverage an imitation reward function from the previous work [WGH20]. Specifically, our imitation reward is multiplication of four terms, where each term calculates tracking error for root position $p_{\text{root}}$, joint position $q$, joint velocity $\alpha$, and the positions of end-effectors $p_{\text{ee}}$, respectively. The total imitation reward term $r_{\text{im}}$ can be written as

$$
\begin{aligned}
r_{\text{im}} &= r_{\text{root}} \cdot r_q \cdot r_\alpha \cdot r_{\text{ee}}, \\
r_{root} &= \exp(-10 \cdot \|p_{\text{root}} - \tilde{p}_{\text{root}}\|^2), \\
r_q &= \exp(-2 \cdot \|q - \tilde{q}\|^2), \\
r_\alpha &= \exp(-0.5 \cdot \|\alpha - \tilde{\alpha}\|^2), \\
r_{\text{ee}} &= \exp(-40 \cdot \|p_{\text{ee}} - \tilde{p}_{\text{ee}}\|^2),
\end{aligned}
\qquad (10)
$$

where symbols with the *tilde* indicate value from the reference state.

In addition to the imitation reward, we incorporate style reward from the adversarial training pipeline [PMA*21, LCK*23]. In this setting, a trainable discriminator $D_\phi$ determines whether the queried state transition $(s, s')$ comes from the simulated agent or the reference dataset. The discriminator outputs a value close to 1 if it predicts that the transition likely originates from the agent, and a value close to 0 for the opposite case. We express style reward $r_{\text{style}}$ and the total reward $r$ as

$$
\begin{aligned}
r_{\text{style}} &= -5 \cdot \log(1 - D_\phi(s, s')), \\
r &= 0.5 \cdot r_{\text{im}} + 0.5 \cdot r_{\text{style}}.
\end{aligned}
\qquad (11)
$$

### C. Task Learning Details

Expanding Section 6.3, we describe detailed settings for the downstream tasks. The explanation for each task includes mathematical expressions of task-specific goal observation, reward function, and the condition for early termination of the training.

### C.1. Motion In-betweening

In this example, an agent learns to generate interpolated motions given two keyframes and the queried time interval as goal observation $g$. In each episode, we randomly select 4 keyframes from a 4-second motion sequence sampled from the dataset.

**Goal observation** We employ the configuration of the goal observation vector used in the previous work [GJW22]. An agent observes two consecutive keyframes where one is the start keyframe $f_s$ and the other is the goal keyframe $f_g$ for each interval. Given an agent with $J$ rigid bodies, each keyframe $f$ contains body position $p_f \in \mathbb{R}^{J \times 3}$ and orientation $\theta_f \in \mathbb{R}^{J \times 4}$ in the global coordinate. At each timestep, the goal observation is calculated by measuring the difference between each keyframe and the simulated state, represented as

$$\bar{p}_f = p - p_f,$$
$$\bar{\theta}_f = \theta \ominus \theta_f, \tag{12}$$

where $p$ and $\theta$ denote the body position and orientation of the simulated agent. We repeat this process for the start keyframe $f_s$ and the goal keyframe $f_g$, which in result obtains a vector of $\hat{g} = [\bar{p}_{f_s}, \bar{\theta}_{f_s}, \bar{p}_{f_g}, \bar{\theta}_{f_g}]$. Followed by the previous work [HYNP20], we then augment the goal observation with a positional encoding of time-to-interval $e_{\text{time}}$, which we write as

$$g = \hat{g} + e_{\text{time}}. \tag{13}$$

**Reward** To encourage reliable motion, we guide the agent with the full-body imitation reward. We utilize the similar imitation reward term in Eq. (10). However, since the goal observation is sparse, it is impossible for the agent to strictly follow the ground-truth trajectories. Therefore, we modify the reward from a multiplication format to a summation format in order to encourage agents to behave flexibly. We represent the modified reward as

$$r = 0.1 \cdot r_{\text{root}} + 0.65 \cdot r_q + 0.1 \cdot r_\omega + 0.15 \cdot r_{\text{ee}}. \tag{14}$$

**Early Termination** To prevent an agent from interpolating in an unrealistic manner, we set an early termination condition that properly prevents large deviations from the ground-truth trajectories. Therefore, we employ the same condition used in the imitation learning phase.

## C.2. Head-mounted Device Tracking

In this scenario, we train an agent to reconstruct the full-body motion given the target head trajectories. For each episode, we sample 5 seconds of target head trajectories from the reference motion clips.

**Goal Observation** We can treat imitation learning as full-body tracking since it basically encourages agents to follow the target trajectories. In contrast to full-body tracking, this scenario assumes that only the target head trajectories are available for the agent. Therefore, we construct the goal observation vector $g$ by selecting the head parts from the target reference state $\tilde{s}$ in Eq. (9).

**Reward** We employ a similar reward structure used for motion in-betweening. However, we modify the total reward term to consider only the head part, allowing other parts to move freely without being penalized. Even though the reward term does not regulate the style of full-body motion, our model can still generate natural full-body motions.

**Early Termination** Similarly, with the goal observation and the reward, we only account for the head deviations in the early termination condition. We terminate an episode when the head of the simulated agent deviates from the ground truth head position more than 1m.

## C.3. Point-goal Navigation

In this scenario, the agent is trained to reach a target location that is randomly respawned within 10 meters from the agent's initial position. Additionally, we set the target location to be updated every 4 seconds. Since the maximum episode length is 8 seconds, the agent can experience at most one target change per training episode.

**Goal Observation** Since the target location is on the ground plane, we allow the agent to know its relative position toward the goal. The goal observation can be abbreviated as

$$g = p_{xy} - t_{xy}, \tag{15}$$

where $p_{xy}, t_{xy} \in \mathbb{R}^2$ indicate the locations of the agent and the target on the ground, respectively.

**Reward** We employ the reward function tailed for the *target location* task in the previous work [PALVdP18]. The reward consists of two parts: the first term $r_{\text{pos}}$ measuring positional error from the target location, and the second term $r_{\text{speed}}$ measuring the heading speed. Particularly, $r_{\text{speed}}$ encourages the agent to move toward the target location but does not provide an extra reward once the agent achieves the minimum target speed. Given the minimum target speed of $u = 1$ $(m/s)$, we can represent the total reward function as

$$r = 0.7 \cdot r_{\text{pos}} + 0.3 \cdot r_{\text{speed}},$$
$$r_{\text{pos}} = \exp[-0.5 \cdot \|p_{xy} - t_{xy}\|^2],$$
$$r_{\text{speed}} = \exp[-\max(0, u - \langle v_{xy}, \frac{p_{xy} - t_{xy}}{\|p_{xy} - t_{xy}\|}\rangle)^2], \tag{16}$$

where $v_{xy} \in \mathbb{R}^2$ denotes the current velocity of the agent on the ground.

**Early Termination** Unlike the previous tasks, this scenario does not include any reference trajectories sampled from the dataset. Therefore, we simply terminate the episode when the agent's head position falls below 0.5 meters, which potentially indicates a fall state.

## D. Hyperparameters

We provide hyperparameters for training. Note "→" indicates scheduling in Table 5-8. We also mark the parameters that are only applicable to certain model with the parenthesis, *e.g.*, $\mathcal{L}_{\text{KL}}$ for *continuous* model.

## References

[AGL*22] Ao T., Gao Q., Lou Y., Chen B., Liu L.: Rhythmic gesticulator: Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings. *ACM Transactions on Graphics (TOG) 41*, 6 (2022), 1–19. 5

**Table 5:** *Hyperparamters used for imitation learning.*

| Parameters | value |
|---|---|
| horizon length | 16 |
| number of environment | 8192 |
| batch size | 65536 |
| learning rate | 2e-4 |
| weight for $\mathcal{L}_{\text{expert}}$ | 10 |
| weight for $\mathcal{L}_{\text{KL}}$ (*continuous*) | $0.05 \to 0.5$ |
| weight for $\mathcal{L}_{\text{mm}}$ (*hybrid*) | $0.1 \to 1.0$ |
| weight for $\mathcal{L}_{\text{commit}}$ (*discrete*, *hybrid*) | 1.0 |
| weight for $\mathcal{L}_{\text{reg}}$ (*continuous*) | 0.005 |
| weight for $\mathcal{L}_{\text{reg}}$ (*discrete*, *hybrid*) | 0.05 |
| codebook size (*discrete*) | 8192 |
| codebook size (*discrete$^+$*, *hybrid*) | 1024 |
| Number of codebooks (*discrete$^+$*, *hybrid*) | 8 |

**Table 6:** *Hyperparamters used for motion in-betweening.*

| Parameters | value |
|---|---|
| horizon length | 16 |
| number of environment | 4096 |
| batch size | 32768 |
| learning rate | 5e-5 |
| $\gamma$ for PPO | 0.99 |
| $\tau$ for PPO | 0.95 |
| clip range for PPO | 0.1 |
| Number of active codebooks (*discrete$^+$*, *hybrid*) | 1 |

**Table 7:** *Hyperparamters used for head-mounted device tracking.*

| Parameters | value |
|---|---|
| horizon length | 16 |
| number of environment | 8192 |
| batch size | 65536 |
| learning rate | 5e-5 |
| $\gamma$ for PPO | 0.99 |
| $\tau$ for PPO | 0.95 |
| clip range for PPO | 0.1 |
| Number of active codebooks (*discrete$^+$*, *hybrid*) | 2 |

**Table 8:** *Hyperparamters used for point-goal navigation.*

| Parameters | value |
|---|---|
| horizon length | 16 |
| number of environment | 8192 |
| batch size | 65536 |
| learning rate | 5e-5 |
| $\gamma$ for PPO | 0.99 |
| $\tau$ for PPO | 0.95 |
| clip range for PPO | 0.1 |
| Number of active codebooks (*discrete$^+$*, *hybrid*) | 1 |

[BVV*15] BOWMAN S. R., VILNIS L., VINYALS O., DAI A. M., JOZEFOWICZ R., BENGIO S.: Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349* (2015). 3

[BWL*23] BAE J., WON J., LIM D., MIN C.-H., KIM Y. M.: Pmp: Learning to physically interact with environments using part-wise motion priors. In *ACM SIGGRAPH 2023 Conference Proceedings* (2023), pp. 1–10. 1, 2

[BXP*22] BHATNAGAR B. L., XIE X., PETROV I., SMINCHISESCU C., THEOBALT C., PONS-MOLL G.: Behave: Dataset and method for tracking human object interactions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (jun 2022), IEEE. 2

[CKD*15] CHUNG J., KASTNER K., DINH L., GOEL K., COURVILLE A. C., BENGIO Y.: A recurrent latent variable model for sequential data. *Advances in neural information processing systems 28* (2015). 3

[DD24] DILLER C., DAI A.: Cg-hoi: Contact-guided 3d human-object interaction generation. 2

[GJW22] GOPINATH D., JOO H., WON J.: Motion in-betweening for physically simulated characters. In *SIGGRAPH Asia 2022 Posters*. 2022, pp. 1–2. 7, 13

[GMJ*23] GUO C., MU Y., JAVED M. G., WANG S., CHENG L.: Momask: Generative masked modeling of 3d human motions. *arXiv preprint arXiv:2312.00063* (2023). 2

[GPB*18] GREGOR K., PAPAMAKARIOS G., BESSE F., BUESING L., WEBER T.: Temporal difference variational auto-encoder. *arXiv preprint arXiv:1806.03107* (2018). 3

[GvdOL*19] GÂRBACEA C., VAN DEN OORD A., LI Y., LIM F. S., LUEBS A., VINYALS O., WALTERS T. C.: Low bit-rate speech coding with vq-vae and a wavenet decoder. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (2019), IEEE, pp. 735–739. 3

[HGW*23] HASSAN M., GUO Y., WANG T., BLACK M., FIDLER S., PENG X. B.: Synthesizing physical character-scene interactions. In *ACM SIGGRAPH 2023 Conference Proceedings* (2023), pp. 1–9. 1, 2

[HMP*17] HIGGINS I., MATTHEY L., PAL A., BURGESS C. P., GLOROT X., BOTVINICK M. M., MOHAMED S., LERCHNER A.: beta-vae: Learning basic visual concepts with a constrained variational framework. *ICLR (Poster) 3* (2017). 3

[HYNP20] HARVEY F. G., YURICK M., NOWROUZEZAHRAI D., PAL C.: Robust motion in-betweening. 2, 6, 7, 12, 13

[IPOS13] IONESCU C., PAPAVA D., OLARU V., SMINCHISESCU C.: Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence 36*, 7 (2013), 1325–1339. 2

[JCL*24] JIANG B., CHEN X., LIU W., YU J., YU G., CHEN T.: Motiongpt: Human motion as a foreign language. *Advances in Neural Information Processing Systems 36* (2024). 5

[JZL*24] JIANG N., ZHANG Z., LI H., MA X., WANG Z., CHEN Y., LIU T., ZHU Y., HUANG S.: Scaling up dynamic human-scene interaction modeling. *arXiv preprint arXiv:2403.08629* (2024). 2

[KFPM21] KUMAR A., FU Z., PATHAK D., MALIK J.: Rma: Rapid motor adaptation for legged robots. *arXiv preprint arXiv:2107.04034* (2021). 1, 2

[KW13] KINGMA D. P., WELLING M.: Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013). 3

[LCK*23] LUO Z., CAO J., KITANI K., XU W., ET AL.: Perpetual humanoid control for real-time simulated avatars. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 10895–10904. 6, 12

[LCM*23] LUO Z., CAO J., MEREL J., WINKLER A., HUANG J., KITANI K., XU W.: Universal humanoid motion representations for physics-based control. *arXiv preprint arXiv:2310.04582* (2023). 2, 3, 5, 6, 7, 10, 12

[LKK*22] LEE D., KIM C., KIM S., CHO M., HAN W.-S.: Autoregressive image generation using residual quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 11523–11532. 2, 3, 5

[LTGN19] LUCAS J., TUCKER G., GROSSE R., NOROUZI M.: Understanding posterior collapse in generative latent variable models. 2

[LZCVDP20] LING H. Y., ZINNO F., CHENG G., VAN DE PANNE M.: Character controllers using motion vaes. *ACM Transactions on Graphics (TOG) 39*, 4 (2020), 40–1. 2

[LZL*24a] LIANG H., ZHANG W., LI W., YU J., XU L.: Intergen: Diffusion-based multi-human motion generation under complex interactions. *International Journal of Computer Vision* (2024), 1–21. 2

[LZL*24b] LIN J., ZENG A., LU S., CAI Y., ZHANG R., WANG H., ZHANG L.: Motion-x: A large-scale 3d expressive whole-body human motion dataset. *Advances in Neural Information Processing Systems 36* (2024). 2

[MGT*19] MAHMOOD N., GHORBANI N., TROJE N. F., PONS-MOLL G., BLACK M. J.: Amass: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF international conference on computer vision* (2019), pp. 5442–5451. 2, 12

[MHG*18] MEREL J., HASENCLEVER L., GALASHOV A., AHUJA A., PHAM V., WAYNE G., TEH Y. W., HEESS N.: Neural probabilistic motor primitives for humanoid control. *arXiv preprint arXiv:1811.11711* (2018). 1, 2

[MMAT23] MENTZER F., MINNEN D., AGUSTSSON E., TSCHANNEN M.: Finite scalar quantization: Vq-vae made simple. *arXiv preprint arXiv:2309.15505* (2023). 3, 5, 10

[MTA*20] MEREL J., TUNYASUVUNAKOOL S., AHUJA A., TASSA Y., HASENCLEVER L., PHAM V., EREZ T., WAYNE G., HEESS N.: Catch & carry: reusable neural controllers for vision-guided whole-body tasks. *ACM Transactions on Graphics (TOG) 39*, 4 (2020), 39–1. 1, 2

[MWG*21] MAKOVIYCHUK V., WAWRZYNIAK L., GUO Y., LU M., STOREY K., MACKLIN M., HOELLER D., RUDIN N., ALLSHIRE A., HANDA A., ET AL.: Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470* (2021). 5

[PALVdP18] PENG X. B., ABBEEL P., LEVINE S., VAN DE PANNE M.: Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions On Graphics (TOG) 37*, 4 (2018), 1–14. 2, 9, 12, 13

[PGH*22] PENG X. B., GUO Y., HALPER L., LEVINE S., FIDLER S.: Ase: Large-scale reusable adversarial skill embeddings for physically simulated characters. *ACM Transactions On Graphics (TOG) 41*, 4 (2022), 1–17. 2, 7

[PMA*21] PENG X. B., MA Z., ABBEEL P., LEVINE S., KANAZAWA A.: Amp: Adversarial motion priors for stylized physics-based character control. *ACM Transactions on Graphics (ToG) 40*, 4 (2021), 1–20. 2, 12

[PXW*23] PENG X., XIE Y., WU Z., JAMPANI V., SUN D., JIANG H.: Hoi-diff: Text-driven synthesis of 3d human-object interactions using diffusion models. *arXiv preprint arXiv:2312.06553* (2023). 2

[RVdOV19] RAZAVI A., VAN DEN OORD A., VINYALS O.: Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems 32* (2019). 3

[STKB23] SHAFIR Y., TEVET G., KAPON R., BERMANO A. H.: Human motion diffusion as a generative prior. *arXiv preprint arXiv:2303.01418* (2023). 2

[TKG*23] TESSLER C., KASTEN Y., GUO Y., MANNOR S., CHECHIK G., PENG X. B.: Calm: Conditional adversarial latent models for directable virtual characters. In *ACM SIGGRAPH 2023 Conference Proceedings* (2023), pp. 1–9. 2

[TRG*22] TEVET G., RAAB S., GORDON B., SHAFIR Y., COHEN-OR D., BERMANO A. H.: Human motion diffusion model. *arXiv preprint arXiv:2209.14916* (2022). 3

[VDOV*17] VAN DEN OORD A., VINYALS O., ET AL.: Neural discrete representation learning. *Advances in neural information processing systems 30* (2017). 3, 4, 5, 6

[WCL*22] WANG Z., CHEN Y., LIU T., ZHU Y., LIANG W., HUANG S.: Humanise: Language-conditioned human motion generation in 3d scenes. In *Advances in Neural Information Processing Systems (NeurIPS)* (2022). 2

[WF20] WU H., FLIERL M.: Vector quantization-based regularization for autoencoders. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2020), vol. 34, pp. 6380–6387. 10

[WGH20] WON J., GOPINATH D., HODGINS J.: A scalable approach to control diverse behaviors for physically simulated characters. *ACM Transactions on Graphics (TOG) 39*, 4 (2020), 33–1. 7, 12

[WGH21] WON J., GOPINATH D., HODGINS J.: Control strategies for physically simulated characters performing two-player competitive sports. *ACM Transactions on Graphics (TOG) 40*, 4 (2021), 1–11. 2

[WGH22] WON J., GOPINATH D., HODGINS J.: Physics-based character controllers using conditional vaes. *ACM Transactions on Graphics (TOG) 41*, 4 (2022), 1–12. 2, 3

[WLZ*23] WANG Y., LIN J., ZENG A., LUO Z., ZHANG J., ZHANG L.: Physhoi: Physics-based imitation of dynamic human-object interaction. *arXiv preprint arXiv:2312.04393* (2023). 2

[XK21] XU P., KARAMOUZAS I.: A gan-like approach for physics-based imitation learning and interactive character control. *Proceedings of the ACM on Computer Graphics and Interactive Techniques 4*, 3 (2021), 1–22. 2

[XSZK23] XU P., SHANG X., ZORDAN V., KARAMOUZAS I.: Composite motion learning with task control. *ACM Transactions on Graphics (TOG) 42*, 4 (2023), 1–16. 11

[XWWG24] XU S., WANG Z., WANG Y.-X., GUI L.-Y.: Interdreamer: Zero-shot text to 3d dynamic human-object interaction. *arXiv preprint arXiv:2403.19652* (2024). 2

[XXA*23] XU P., XIE K., ANDREWS S., KRY P. G., NEFF M., MCGUIRE M., KARAMOUZAS I., ZORDAN V.: Adaptnet: Policy adaptation for physics-based character control. *ACM Transactions on Graphics (TOG) 42*, 6 (2023), 1–17. 2

[XZY*24] XU L., ZHOU Y., YAN Y., JIN X., ZHU W., RAO F., YANG X., ZENG W.: Regennet: Towards human action-reaction synthesis. *arXiv preprint arXiv:2403.11882* (2024). 2

[YKK*23] YOUNES M., KIJAK E., KULPA R., MALINOWSKI S., MULTON F.: Maaip: Multi-agent adversarial interaction priors for imitation from fighting demonstrations for physics-based characters. *Proceedings of the ACM on Computer Graphics and Interactive Techniques 6*, 3 (2023), 1–20. 2

[YLG*23] YU L., LEZAMA J., GUNDAVARAPU N. B., VERSARI L., SOHN K., MINNEN D., CHENG Y., GUPTA A., GU X., HAUPTMANN A. G., ET AL.: Language model beats diffusion–tokenizer is key to visual generation. *arXiv preprint arXiv:2310.05737* (2023). 3

[YLH*23] YANG D., LIU S., HUANG R., TIAN J., WENG C., ZOU Y.: Hifi-codec: Group-residual vector quantization for high fidelity audio codec. *arXiv preprint arXiv:2305.02765* (2023). 3, 10

[YLK*21] YU J., LI X., KOH J. Y., ZHANG H., PANG R., QIN J., KU A., XU Y., BALDRIDGE J., WU Y.: Vector-quantized image modeling with improved vqgan. *arXiv preprint arXiv:2110.04627* (2021). 3

[YMG*23] YUAN Y., MAKOVIYCHUK V., GUO Y., FIDLER S., PENG X., FATAHALIAN K.: Learning physically simulated tennis skills from broadcast videos. *ACM Trans. Graph 42*, 4 (2023). 1, 2

[YSCL22] YAO H., SONG Z., CHEN B., LIU L.: Controlvae: Model-based learning of generative controllers for physics-based characters. *ACM Transactions on Graphics (TOG) 41*, 6 (2022), 1–16. 2, 3, 10

[YSZ*23] YAO H., SONG Z., ZHOU Y., AO T., CHEN B., LIU L.: Moconvq: Unified physics-based motion control via scalable discrete representations. *arXiv preprint arXiv:2310.10198* (2023). 2, 3, 5, 10

[YTB*24] YI H., THIES J., BLACK M. J., PENG X. B., REMPE D.: Generating human interaction motions in scenes with text control. *arXiv:2404.10685* (2024). 2

[YZAS21]   YAN W., ZHANG Y., ABBEEL P., SRINIVAS A.: Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157* (2021). 3

[ZLO*21]   ZEGHIDOUR N., LUEBS A., OMRAN A., SKOGLUND J., TAGLIASACCHI M.: Soundstream: An end-to-end neural audio codec. *IEEE/ACM Transactions on Audio, Speech, and Language Processing 30* (2021), 495–507. 2, 3, 5

[ZZLH23]   ZHU Q., ZHANG H., LAN M., HAN L.: Neural categorical priors for physics-based character control. *ACM Transactions on Graphics (TOG) 42*, 6 (2023), 1–16. 2, 3, 4, 5, 6

[ZZS*24]   ZHANG J., ZHANG J., SONG Z., SHI Z., ZHAO C., SHI Y., YU J., XU L., WANG J.: Hoi-m3: Capture multiple humans and objects interaction within contextual environment. *arXiv preprint arXiv:2404.00299* (2024). 2

[ZZW*23]   ZHAO K., ZHANG Y., WANG S., BEELER T., , TANG S.: Synthesizing diverse human motions in 3d indoor scenes. In *International conference on computer vision (ICCV)* (2023). 2