

# MotionBricks: Scalable Real-Time Motions with Modular Latent Generative Model and Smart Primitives

TINGWU WANG<sup>†</sup>, NVIDIA, USA  
OLIVIER DIONNE<sup>†</sup>, NVIDIA, Canada  
MICHAEL DE RUYTER, NVIDIA, USA  
DAVID MINOR, NVIDIA, Canada  
DAVIS REMPE, NVIDIA, USA  
KAIFENG ZHAO, NVIDIA, Switzerland and ETH Zürich, Switzerland  
MATHIS PETROVICH, NVIDIA, Switzerland  
YE YUAN, NVIDIA, USA  
CHENRAN LI, NVIDIA, USA  
ZHENGYI LUO, NVIDIA, USA  
BRIAN ROBISON, NVIDIA, USA  
XAVIER BLACKWELL, NVIDIA, USA  
BERNARDO ANTONIAZZI, NVIDIA, USA  
XUE BIN PENG, NVIDIA, Canada and Simon Fraser University, Canada  
YUKE ZHU, NVIDIA, USA and The University of Texas at Austin, USA  
SIMON YUEN, NVIDIA, USA



Fig. 1. MotionBricks enables real-time motion control across animation and robotics. All motions are generated by our unified latent neural backbone using the smart primitive interface. Top: We showcase MotionBricks in a UE5 demo covering diverse locomotion styles, acrobatics, and object-scene interactions. Bottom: We deploy MotionBricks on the Unitree G1 robot for real-world whole-body control. The links to the dataset, code, and videos are available on the project webpage: <https://nvlabs.github.io/motionbricks>.

Despite transformative advances in generative motion synthesis, real-time interactive motion control remains dominated by traditional techniques.

<sup>†</sup>Joint first authors. Corresponding author: Tingwu Wang, [tingwuw@nvidia.com](mailto:tingwuw@nvidia.com).



This work is licensed under a Creative Commons Attribution 4.0 International License.

In this work, we identify two key challenges in bridging research and production: 1) *Real-time scalability*: Industry applications demand real-time generation of a vast repertoire of motion skills, while generative methods

© 2026 Copyright held by the owner/author(s).  
ACM 1557-7368/2026/7-ART  
<https://doi.org/10.1145/3811334>

exhibit significant degradation in quality and scalability under real-time computation constraints, and 2) *Integration*: Industry applications demand fine-grained multi-modal control involving velocity commands, style selection, and precise keyframes, a need largely unmet by existing text- or tag-driven models. Moreover, a systematic motion design interface for generative models remains absent. To overcome these limitations, we introduce MotionBricks: a large-scale, real-time generative framework with a two-fold solution. First, we propose a large-scale modular latent generative backbone tailored for robust real-time motion generation, effectively modeling a dataset of over 350,000 motion clips with a single model. Second, we introduce *smart primitives* that provide a unified, robust, and intuitive interface for authoring both navigation and object interaction. Notably, MotionBricks applies to new downstream tasks in a zero-shot manner, where no fine-tuning or task-specific tagging is required. Applications can be designed in a plug-and-play manner like assembling bricks without expert animation knowledge, enabling an accessible interface for applications in animation and robotics. Quantitatively, we show that MotionBricks produces state-of-the-art motion quality on open-source and proprietary datasets of various scales, while also achieving a real-time throughput of 15,000 FPS with 2ms latency. We demonstrate the flexibility and robustness of MotionBricks in a complete production-level animation demo, covering navigation and object-scene interaction across various styles with a unified model. To showcase our framework’s application beyond animation, we deploy MotionBricks on the Unitree G1 humanoid robot to demonstrate its flexibility and generalization for real-time robotic control.

CCS Concepts: • **Computing methodologies** → **Motion processing**; *Neural networks*; *Procedural animation*; • **Applied computing** → Computer games.

Additional Key Words and Phrases: character animation, motion synthesis, game engine, generative models, robotics

#### ACM Reference Format:

Tingwu Wang, Olivier Dionne, Michael De Ruyter, David Minor, Davis Rempe, Kaifeng Zhao, Mathis Petrovich, Ye Yuan, Chenran Li, Zhengyi Luo, Brian Robison, Xavier Blackwell, Bernardo Antoniazzi, Xue Bin Peng, Yuke Zhu, and Simon Yuen. 2026. MotionBricks: Scalable Real-Time Motions with Modular Latent Generative Model and Smart Primitives. *ACM Trans. Graph.* 45, 4 (July 2026), 22 pages. <https://doi.org/10.1145/3811334>

## 1 Introduction

Runtime interactive motion control has been one of the foundational pillars supporting the expanding array of animation and robotics applications. However, the typical process of designing runtime behavior systems is highly labor-intensive. Animation or behavior graphs, and their variants, have been widely adopted in the animation and robotics industries [Dynamics 2025; Kovar and Gleicher 2004; Min and Chai 2012; Safonova and Hodgins 2007; Unitree 2025; Unity 2025], largely due to their ability to construct complex, hierarchical state machines with granular quality control. In a conventional animation graph, pre-recorded motion clips are organized into discrete states, with transitions between them governed by game events and user input. Consequently, animation graphs suffer from severe scalability limitations during the process of asset and motion graph design. For example, a modern AAA animation state machine, such as that used in *Assassin’s Creed*, may require managing over 15,000 animations, 5,000 states, and nested graphs up to 12 levels deep [Holden 2018], making such systems nearly impossible to maintain or extend. This overwhelming complexity has effectively made high-quality runtime behavior systems the

exclusive domain of large, well-resourced studios, highlighting an urgent need for new, more accessible approaches for the broader community.

Recent advances in generative models have led to remarkable breakthroughs across a wide range of content creation tasks, including motion synthesis [Holden et al. 2017; Peng et al. 2018; Tevet et al. 2022]. Notably, two formulations have been extensively studied for interactive runtime motion control: (1) Models that directly receive control commands as input. This includes works spanning from non-parametric methods such as motion matching [Buttner 2015, 2019; Clavet 2016], to deterministic models [Holden et al. 2017; Lu et al. 2024; Starke et al. 2023, 2019, 2020, 2021; Zhang et al. 2018], to generative methods [Gou et al. 2025; Shi et al. 2024]. Control commands include, for example, tags, root trajectories, velocities, and object attributes. (2) Text-to-motion models [Alexanderson et al. 2023; Cohan et al. 2024; Guo et al. 2024; Pinyoanuntapong et al. 2024b; Tevet et al. 2022; Zhang et al. 2024a; Zhou et al. 2024] use text prompts as the main control modality. The former requires manual tagging or attribute preparation during training, and it is quite common to use a predefined “one-hot” vector to represent the control commands, which ultimately limits scalability as the types of control commands vary significantly across different motion categories. The latter, however, treats animation as a monolithic generation task from text prompts, which has yet to provide the fine-grained control needed for reliable industrial-level applications and often struggles to balance scalability, quality, and speed. Ultimately, runtime motion control is a complex task that bridges low-level motion synthesis and high-level behavior design, where a successful system needs to excel in both requirements.

To achieve fine-grained control and scalability in real-time, we present MotionBricks, a framework designed for production use that seamlessly combines a powerful low-level latent neural backbone with a flexible, universal high-level behavior system built upon our proposed *smart primitives*. For the low-level neural backbone, we adopt motion in-betweening as the foundational paradigm. We address real-time scalability challenges with a novel structured latent design that enables our backbone to significantly improve model capacity and spatial precision. Our model achieves state-of-the-art motion quality while maintaining throughput far exceeding real-time requirements. We further organize the model into a modular setup that supports progressive, coarse-to-fine motion generation, improving robustness and quality, while allowing for a transparent workflow for inspecting and refining intermediate results. Throughout the pipeline, we support flexible combinations of constraint specifications, enabling maximum adaptability for high-level smart primitive designs. For the high-level behavior system, we introduce two smart primitives: *smart locomotion* and *smart object*. Smart locomotion provides a robust mechanism for generating proxy keyframes across arbitrary navigation styles and velocity-heading commands, while safeguarding against command dead zones and producing natural movement details via neural root trajectory refinement. Smart objects offer an intuitive interface for proxy object interactions with keyframes, allowing precise and natural control over interaction with the neural backbone. This enables building highly complex scenes in a plug-and-play manner without the need for animation graphs. Both primitives expose a user-friendly, scalable interface

that unifies control into keyframe commands directed to the low-level backbone. Notably, all applications and demos shown in this paper are built using a fixed pre-trained neural backbone, without requiring further fine-tuning or additional tagging.

The core contributions of MotionBricks are summarized as follows: 1) We propose a modular generative neural backbone, featuring a novel multi-head tokenizer and a progressive coarse-to-fine generation pipeline. Our design achieves state-of-the-art motion quality on various internal and open-source datasets, while maintaining 2ms latency and 15,000 FPS throughput, far exceeding real-time requirements; 2) We introduce smart primitives, a flexible high-level behavior system built on top of our neural backbone. We implement smart locomotion and smart object primitives in Unreal Engine, and demonstrate their effectiveness through a complete production-level demo, showcasing a wide variety of navigation and interaction skills with a unified interface, all without fine-tuning or additional tagging; and 3) We deploy MotionBricks on the Unitree G1 humanoid robot, demonstrating that our framework bridges the gap between virtual character animation and physical robot control, enabling a unified approach to motion synthesis across both domains. The open-source project for MotionBricks can be found on our webpage.

## 2 Related Work

In this section, we provide an overview of relevant literature, covering both established traditional methods widely adopted in real-world applications and the latest advancements in data-driven neural network techniques.

*Traditional Methods.* Motion graphs, or animation behavior graphs, are the prevalent approach for runtime animation and certain robotics applicability in the industry [Arikan and Forsyth 2002; Dynamics 2025; Engine 2025; Kovar 2002; Lee et al. 2002; Unitree 2025]. Animation graphs are state machines consisting of motion states such as walking, running, idle, etc. Motion is generated by replaying clips within the state or blending between states during transitions [Unity 2025] based on user controls or game events. Academic research has further advanced the animation behavior graph, as seen in works such as Kovar and Gleicher [2004]; Min and Chai [2012]; Safonova and Hodgins [2007]. A particularly notable development is presented in Lee et al. [2010], which enables motion transitions at a fine-grained, frame-by-frame level rather than the traditional node-by-node approach. These academic innovations, Lee et al. [2010] in particular, ultimately inspired motion matching [Buttner 2015, 2019; Clavet 2016]. In motion matching, future frames are retrieved from a motion database to match current state and user controls. This approach eliminates the rigid, node-based constraints of traditional animation graphs, enabling significantly more flexible and nuanced motion synthesis. While efficiency improvements have been achieved through neural distillation and better search algorithms [Holden et al. 2020; Yi and Jee 2019], motion matching is typically still used as a subcomponent for navigation within broader animation graphs and does not overcome the fundamental scalability limitations. Ultimately, animation graphs become increasingly fragile and labor-intensive as the number of states and transitions grows, posing significant challenges for large-scale applications.

*Pre-Generative Methods.* Prior to the advent of powerful generative models [Achiam et al. 2023; Ho et al. 2020], researchers applied feedforward deterministic neural networks to motion synthesis [Fragkiadaki et al. 2015; Holden et al. 2016]. Phase-Functioned Neural Networks (PFNN) [Holden et al. 2017] is a seminal work in this direction, which predicts navigation motion on different terrains autoregressively. This approach inspired numerous models for runtime animation methods [Lu et al. 2024; Starke et al. 2023, 2024, 2019, 2020, 2021; Zhang et al. 2018]. However, deploying these models in real-world settings presents persistent challenges, such as foot sliding, object penetration and excessive smoothing. These models are also often specialized, tailored for a small set of tasks. As a result, their applicability in production remains restricted.

Another line of research formulates motion synthesis as an in-betweening or inpainting problem given the starting and ending keyframe constraints [Harvey et al. 2020; Oreshkin et al. 2024; Qin et al. 2022]. However, since these methods are typically tailored for a narrow set of skills, they inherit the same scalability challenges as the aforementioned autoregressive approaches. Moreover, as target keyframes are usually unavailable during runtime animation, these methods are predominantly utilized as offline tools for animation authoring.

*Generative Methods.* Recently, diffusion models [Ho and Ermon 2016; Song et al. 2020] have achieved remarkable success, particularly in image [Esser et al. 2024; Rombach et al. 2022] and video generation [Google 2025; OpenAI 2024], inspiring widespread adoption across various domains. Of particular relevance are motion diffusion models (MDMs) [Alexanderson et al. 2023; Cohan et al. 2024; Li et al. 2025; Tevet et al. 2025, 2022; Zhang et al. 2024a; Zhou et al. 2024]. In the pioneering work [Tevet et al. 2022], the authors train a motion diffusion model on the HumanML3D [Guo et al. 2022a] conditioned on text descriptions. This work is further extended in Shafir et al. [2023] to handle longer sequences and multiple interacting humans. Runtime diffusion guidance [Ho and Salimans 2022] is incorporated in Shafir et al. [2023], enabling the enforcement of spatial constraints during inference without retraining. In Li et al. [2024a, 2023], MDMs have also been extended to support object interactions. While MDMs offer strong scalability and motion quality, a significant challenge remains: inference speed. Standard diffusion models typically require several seconds, or even minutes, to generate motion clips, rendering them impractical for real-time applications. Although recent research explores auto-regressive diffusion models that amortize computation [Chen et al. 2024; Han et al. 2024; Li et al. 2024b; Shi et al. 2024] or use DDIM [Song et al. 2020] or techniques such as flow matching [Gou et al. 2025; Lipman et al. 2022], these approaches still face a significant trade-off between real-time performance and output quality, and barely achieve real-time speed, leaving them impractical for integration into already demanding animation or robotics pipelines. Another promising line of research centers on token-based generative models, which discretize continuous representations into discrete tokens [Esser et al. 2021; Razavi et al. 2019; Tian et al. 2024; Van Den Oord et al. 2017], and subsequently employ transformers for prediction. This approach draws inspiration from the advances seen in discrete language, image, and video generative models [Luo et al. 2024; Sun et al. 2024; Yu et al.

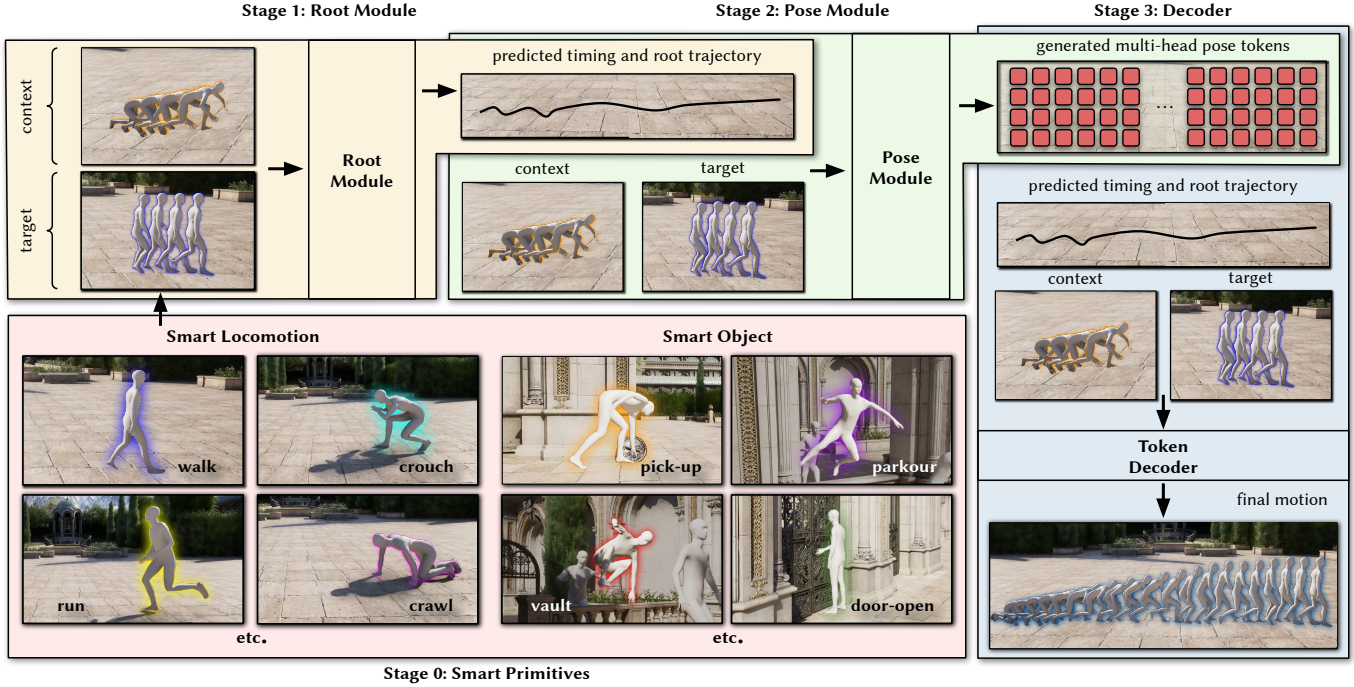


Fig. 2. MotionBricks’s inference pipeline consists of four stages. Given user commands or game events, smart primitives generate target keyframes. The root module first predicts timing and root trajectory, followed by the pose module that models the distribution of multi-head latent pose tokens. Finally, the decoder produces continuous motion conditioned on pose tokens, root trajectories, and keyframes.

2023a,b]. Works such as Guo et al. [2022b]; Jiang et al. [2023]; Zhang et al. [2024b] generate motion from text by treating motion tokens analogously to language tokens. It is common to use an iteratively masked token prediction strategy, in which tokens are predicted by progressively finalizing the most confident predictions [Guo et al. 2024; Pinyoanuntapong et al. 2024a,b]. Concurrently, SONIC [Luo et al. 2025] adopts a similar token-based latent generation approach for humanoid robot control, though with a different tokenizer and network architecture, and is limited to locomotion without support for object interaction. Overall, these methods face a similar trade-off between quality and speed as diffusion models, and have yet to achieve industry-grade interactivity, motion quality, or scalability.

### 3 Overview

We present MotionBricks, an intuitive framework for fine-grained, scalable, and robust runtime motion control. An overview of our framework is shown in Figure 2. The pipeline consists of four components: the smart primitive module, tokenizer, root module, and pose module.

We first train a conditional structured multi-head tokenizer that encodes  $T$  frames of motion into  $T/4$  discrete tokens using a temporal encoder and decoder model. Then we train a generative backbone comprising a root module and a pose module. The root module predicts timing and initial root trajectories from in-betweening constraints, while the pose module then models encoded pose token distributions conditioned on the in-betweening constraints as well as the root predictions. Crucially, training is agnostic to downstream

control modalities. No predefined control types or task-specific tags are required, as all smart primitives communicate through a unified keyframe interface. Once trained, the backbone requires no fine-tuning. Users configure smart primitives: smart locomotion for navigation motions, and smart object for scene object interactions. These smart primitives can be intuitively set up from motion libraries and interactive authoring tools, essentially creating a fully connected motion graph where transitions are generated with the neural backbone. At runtime, smart primitives generate target keyframes from user commands and game events. These keyframes, along with the context frames, are fed into the generative backbone to progressively generate timing, root trajectories, and pose tokens, which are finally decoded into continuous motion. The framework operates autoregressively, and replanning is triggered whenever the control signals are changed or the future motion buffer does not have enough frames left. The runtime inference loop is summarized in Algorithm 1. We use simplified notation for clarity; the formal definitions are provided in Sections 4 to 6.

*Flexible In-betweening Constraint Handling.* In MotionBricks, we adopt an in-betweening formulation: context keyframes represent the character’s recent states, while target keyframes are goals supplied by smart primitives. As shown in Figure 3,  $\mathcal{T}_1$ ,  $\mathcal{T}_2$ , and  $\mathcal{T}_3$  denote the sets of provided constraints for local root, global root, and pose, respectively, covering both context and target keyframes. We denote the keyframe constraints as  $\mathcal{T} = \{\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3\}$ . While context keyframes (indices 0, 1, 2, 3) are always available, target constraints

**Algorithm 1** MotionBricks Inference Loop

---

**Require:** Token decoder, root module  $\mathcal{F}$ , pose module  $\mathcal{P}$ , smart primitives  $\mathcal{SP}$

- 1: Initialize motion buffer  $\mathcal{B}$ , character state  $\mathcal{S}$
- 2: **while** running **do**
- 3:   Read game event or user commands  $\mathcal{C}$
- 4:   **if**  $\mathcal{C}$  changed **or**  $|\mathcal{B}|$  running low **then**
- 5:     % Stage 0: construct keyframe constraints  $\mathcal{T}$
- 6:      $\mathcal{T} \leftarrow \mathcal{SP}(\mathcal{S}, \mathcal{C})$
- 7:     % Stage 1: predict frame count  $T$ , root trajectory  $\{r\}$
- 8:      $T, \{r\} \leftarrow \mathcal{F}(\mathcal{T})$
- 9:     % Stage 2: predict pose tokens  $\{z_q\}$
- 10:      $\{z_q\} \leftarrow \mathcal{P}(\{r\}, \mathcal{T}, T)$
- 11:     % Stage 3: decode continuous motion  $\{r, p, q, v, c\}$
- 12:      $\{r, p, q, v, c\} \leftarrow \text{decoder}(\{z_q\}, \{r\}, \mathcal{T}, T)$
- 13:     update  $\mathcal{B}$  with generated motion
- 14:   **end if**
- 15:   Pop next frame from  $\mathcal{B}$ ; update  $\mathcal{S}$
- 16: **end while**

---

(indices 4, 5, 6, 7) vary by task. Different tasks require different levels of constraint density: navigation typically needs only 1 or 2 keyframes for style and root velocity to avoid over-constraining the motion, while object interaction requires denser specifications such as consecutive hand positions. To handle this variability, each module in MotionBricks accepts an arbitrary subset of constraints. When a constraint is not provided, we substitute it with a learnable mask embedding, allowing the model to gracefully handle partial specifications. This masking mechanism is applied consistently across the root module, pose module, and decoder. Throughout all modules, we support in-betweening segments ranging from 12 to 64 frames at 30 FPS.

#### 4 Structured Multi-headed Tokenizer

In MotionBricks, we propose a structured conditional multi-head tokenizer that serves as the cornerstone of our generative neural backbone. We map raw motion data into compact, abstract, structured latent representations and introduce a root-pose disentanglement strategy alongside a flexible conditioning mechanism tailored for diverse application needs. Together, these designs greatly improve the scalability, flexibility, and precision of the generative model under runtime constraints.

*State Representation.* We consider motion segments of length  $T$ , where  $T$  is randomly sampled from 12 to 64 frames in increments of 4. For each frame at time  $t$ , we define the motion state as a tuple  $(r_g, r_l, p, q, v, c)$ , comprising: (1) global root values  $r_g$ : projected global positions and the cosine and sine of the heading angle of the pelvis joint; (2) local root values  $r_l$ : projected positional velocity and angular velocity of the pelvis joint in global coordinates; (3) joint positions  $p$  and rotations  $q$  for all joints (excluding the projected root joint); (4) joint velocities  $v$  and contact labels  $c$ . Note that local and global root values are interconvertible; we retain both to

enable flexible constraint handling in the neural backbone. To properly process motions such as crawling and flipping, we represent joint rotations in global coordinates to avoid canonicalization with ill-defined root headings. Consequently, we do not apply heading canonicalization during training, but instead augment samples with random rotations to learn motion skills across all directions. Similarly, joint positions  $p$  are expressed as global coordinates relative to the root position, and velocities  $v$  are computed in the global frame, both without heading canonicalization, to maintain a consistent motion representation.

#### 4.1 Encoder

Drawing inspiration from foot-scaling and time-warping in traditional animation [Bereznyak 2024], we encode root and pose separately rather than jointly. Although root and pose are strongly correlated, pose intent can be largely disentangled from root intent: for instance, a person may walk at slightly varying speeds while maintaining the same gait pattern. Therefore, the encoder only encodes the local pose states without the root information as follows:

$$\{z_e^t\}_{t=1}^{T/4} = \text{enc} \left( \{p^t, q^t\}_{t=1}^T \right). \quad (1)$$

The latent embedding  $z_e^t$  is generated by progressively down-sampling the features at rates of 2 and 4 using either a 1D convolutional network [Kiranyaz et al. 2021] or a transformer with positional encoding [Vaswani et al. 2017] to ensure temporal consistency. For notational convenience, we overload  $t$  to index both tokens (left-hand side,  $T/4$  total because of the 4-frame downsampling) and frames (right-hand side,  $T$  total); we adopt this convention throughout the paper for cleaner presentation.

Rather than tokenizing or discretizing the feature space with one giant codebook or manually partitioning features by body part as in prior work [Aberman et al. 2020; Pinyoanuntapong et al. 2024b], we argue that such approaches are insufficient to capture the full diversity of human movement required for real-world applications. Instead, we rely on the network to learn this latent decomposition in a fully data-driven manner. The continuous latent embedding  $z_e^t$  is quantized into the discrete latent embedding  $z_q^t$  separately along the feature dimension using  $K$  discrete codebooks  $\{e_1, e_2, \dots, e_K\}$ :

$$z_q^t = \begin{Bmatrix} z_{q,1}^t \\ z_{q,2}^t \\ \vdots \\ z_{q,K}^t \end{Bmatrix} = \begin{Bmatrix} \arg \min_{e_1 \in \mathcal{E}_1} \|z_{e,1}^t - e_1\|_2^2 \\ \arg \min_{e_2 \in \mathcal{E}_2} \|z_{e,2}^t - e_2\|_2^2 \\ \vdots \\ \arg \min_{e_K \in \mathcal{E}_K} \|z_{e,K}^t - e_K\|_2^2 \end{Bmatrix}. \quad (2)$$

This strategy encourages the encoder to learn the subtle compositionality and disentanglement of human movement, without imposing manual bias on how motion should be partitioned. Additionally, multi-head tokenization yields a more robust latent manifold, enabling graceful degradation when individual tokens are mispredicted, as demonstrated later in our experiments.

#### 4.2 Decoder

Mirroring the encoder, the decoder progressively up-samples the discrete pose tokens  $z_q^t$  at rates of 2 and 4. Root trajectory conditions are incorporated at each up-sampling layer through skip

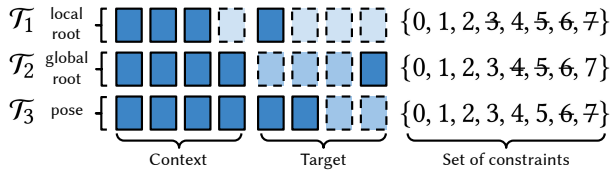


Fig. 3. In MotionBricks, the root module, pose module, and decoder accept a flexible number of constraints from smart primitives. We denote the sets of constraints for local root, global root, and pose as  $\mathcal{T}_1$ ,  $\mathcal{T}_2$ , and  $\mathcal{T}_3$ , respectively. Solid boxes indicate provided constraints; dashed boxes indicate optional or missing entries.

connections at their corresponding temporal positions, encouraging coherent full-body motion generation. To support flexible combinations of constraints from our smart primitives as shown in Figure 3, we allow keyframes produced by the primitives to be passed directly to the decoder. During training, we randomly sample between 0 and 10 keyframes and incorporate them at each layer in the same manner as root conditions. The decoder is formally defined as:

$$\{r_l^t, p^t, q^t, v^t, c^t\}_{t=1}^T = \text{dec} \left( \{z_q^t\}_{t=0}^{T/4}, \{\hat{r}_l^t\}_{t=0}^T, \{\check{p}\}, \{\check{q}\} \right), \quad (3)$$

where the local root trajectory  $\hat{r}_l^t$  is sampled from the dataset during training and predicted by the neural backbone at inference. We denote  $\{\check{p}\}$  and  $\{\check{q}\}$  as variable-length sets of keyframe constraints on joint positions and rotations, respectively. The local root trajectory  $\hat{r}_l^t$  and keyframe constraints  $\{\check{p}\}, \{\check{q}\}$  are provided to the decoder as input and additionally injected into hidden states at each up-sampling layer via skip connections. Since the pose tokens  $\{z_q^t\}_{t=0}^{T/4}$  have temporal dimension  $T/4$  while the local root trajectory  $\hat{r}_l^t$  has temporal dimension  $T$ , root features are temporally stacked by the downsampling factor 4 before concatenation with the pose token embeddings along the feature dimension. This stacking process is applied to the skip connection features of  $\hat{r}_l^t$  at each up-sampling layer, where the temporal downsampling factor progressively decreases from 4 to 2 to 1. Sparse keyframe constraints  $\{\check{p}\}, \{\check{q}\}$  are zero-padded to length  $T$  and processed in the same manner as root trajectory features. Additionally, a boolean availability mask determines whether keyframe embeddings from the skip connection or the decoder hidden states are selected at each layer. We refer readers to our open-source repository for implementation details. Notably, the decoder can also refine the predicted root trajectory, producing cleaner footsteps and smoother transitions.

### 4.3 Quantizer and Training

We train the tokenizer with a standard VQ-VAE quantizer loss, e.g.,  $L = \sum_{t=0}^{T/4} \sum_{k=1}^K \|sg(z_{q,k}^t) - e_k\|_2^2 + \beta \|z_{q,k}^t - sg(e_k)\|_2^2$ , where  $sg$  denotes the stop-gradient operator [Van Den Oord et al. 2017]. In practice, we use a running-mean codebook update, which yields more stable training [Razavi et al. 2019]. Alternatively, an FSQ tokenizer [Mentzer et al. 2023] can be used, which generates similar performance. As discussed in Section 3, during inference the decoder uses the first 4 frames as context keyframes, with 1 to 4 target keyframes provided by the smart primitives; together these form

the keyframe constraints  $\{\check{p}\}, \{\check{q}\}$  in Equation 3. However, during training, a random number (between 0 and 10) of keyframe constraints are sampled at arbitrary temporal positions to encourage the model to learn generalizable and robust features. The same variable-constraint sampling strategy is applied to the root and pose module training as well.

## 5 Neural Backbone

The neural backbone of MotionBricks is a progressive modular architecture that generates in-betweening motions driven by smart primitives in a coarse-to-fine manner. It comprises two components: the root module and the pose module. The root module predicts the number of frames between keyframes and produces an initial estimate of the root trajectory. The pose module then models the distribution of pose tokens conditioned on the root trajectory and keyframes. This modular design not only iteratively improves motion quality, but also provides a transparent interface for inspecting, modifying, and refining intermediate results.

### 5.1 Root Module

The root module follows a two-step design, taking the constraints shown in Figure 3 as input. The first step predicts the number of in-between frames; the second step estimates the initial root trajectory conditioned on the timing prediction and hidden states from the first step.

*Step 1: Hidden State Initialization and Timing Prediction.* We use a transformer encoder for step 1, whose input consists of three types of embeddings: (1) a sequence of 16 learnable frame-slot embeddings  $\{h_1^t\}_{t=1}^{16}$ , each serving as a query position for predicting the root trajectory at that temporal location, covering a maximum of 64 frames after  $4\times$  downsampling; (2) a global embedding encoding the keyframe constraints, denoted as  $f(\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3)$ ; and (3) a timing embedding  $g(T_1)$ , encoding the ground-truth frame count information  $T_1$  when available, or a learnable mask embedding otherwise. Both  $f$  and  $g$  are learnable 3-layer MLPs that project inputs to the same dimension as the frame-slot embeddings. The transformer outputs a distribution over in-between frame counts at 4-frame resolution, from which we sample or take the argmax during inference, yielding the predicted frame count  $T_2$ . When the ground-truth frame count is provided, the predicted distribution is masked out to only allow ground-truth frame count. We also support bit-wise masking over the distribution of valid frame counts, enabling flexible timing control at inference.

*Step 2: Conditional Root Trajectory Estimation.* Intermediate frame-slot embeddings from step 1 are passed to step 2, denoted as  $\{h_2^t\}_{t=1}^{16}$ . We mask out frame-slot embeddings beyond the predicted frame count  $T_2$ . A temporal transformer then processes this sequence to predict the initial root trajectory conditioned on keyframe and timing constraints.

Finally, the root module is formally defined as:

$$\begin{aligned} \{h_2\}, T_2 &= \mathcal{F}_1(\{h_1\}; g(T_1); f(\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3)), \\ \{r_g\} &= \mathcal{F}_2(\{h_2\}; g(T_2); f(\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3)). \end{aligned} \quad (4)$$

In Equation 4,  $\mathcal{F}_1$  and  $\mathcal{F}_2$  denote the transformer encoders for step 1 and step 2, respectively. In practice, we also augment the frame-slot embeddings  $h_1$  and  $h_2$  with positional encoding [Vaswani et al. 2017] for temporal information. The root module outputs global root values  $\{r_g\}_{t=1}^{T_2}$ , enabling direct control over root trajectories to satisfy global constraints. Due to the 4 times downsampling, each final embedding is decoded into 4 consecutive frames of root trajectory  $r_g$ . In subsequent pose and decoder modules, we convert these global root values to local root values for finer-grained motion generation.

## 5.2 Pose Module

The pose module is the largest transformer in our framework, modeling the distribution of pose tokens conditioned on root trajectories and keyframe constraints. Similar to the root module, we construct frame-slot embeddings, where each encodes: (1) local and global root values  $r_l, r_g$ , sampled from the dataset during training and predicted by the root module at inference; (2) local pose keyframe constraints  $p, q$ ; and (3) pose token embeddings, using a cosine scheduled mix of ground truth and masked-token embeddings during training, and all initialized to masked-token embeddings at inference. For local pose keyframe constraints, similar to the decoder (Section 4), we randomly sample 0 to 10 keyframes during training and extract them from  $\mathcal{T}_3$  at inference. Missing entries are replaced with learnable masked-pose embeddings. We adopt the masked token modeling strategy [Luo et al. 2024; Pinyoanuntapong et al. 2024b; Yu et al. 2023a,b] which iteratively predicts the pose tokens based on confidence. In practice, a single forward pass at inference typically produces high-quality motion. However, we still use cosine scheduling as an effective curriculum learning approach during training.

The pose module is formally defined as follows:

$$\begin{aligned} \{h_p\} &= \mathcal{P}(\{\phi(\{r_l\}; \{r_g\}; \{z_q\})\}), \\ p(z_{q,k}^t) &= \sigma\left(f_k\left(h_{p,k}^t\right)\right), k = 1, \dots, K \end{aligned} \quad (5)$$

where  $z_q$  is the latent token embedding vector in Equation 4.1 and  $\phi$  is the input embedding function that applies separate linear projections to the three input types and combines them via an MLP.  $\mathcal{P}$  denotes the transformer encoder, and  $h_{p,k}^t$  is the  $k$ -th partition of the final hidden state  $h_p^t$  along the feature dimension. Each partition is projected through a linear layer  $f_k$  followed by softmax  $\sigma$  to produce the token distribution for the  $k$ -th codebook head. The finalized pose tokens are then decoded along with the predicted root from the root module, as shown in Equation (3).

## 6 Smart Primitives

In this section, we introduce two smart primitives: (1) *smart locomotion* for navigation with arbitrary styles and velocity/heading commands, and (2) *smart object* for interacting with static scenes and dynamic objects. Together, these primitives form a unified and flexible framework powered by a shared neural backbone, providing a complete solution for building rich, interactive scenes.

A key advantage of MotionBricks is zero-shot generalization to diverse downstream tasks, enabled by our scalable generative

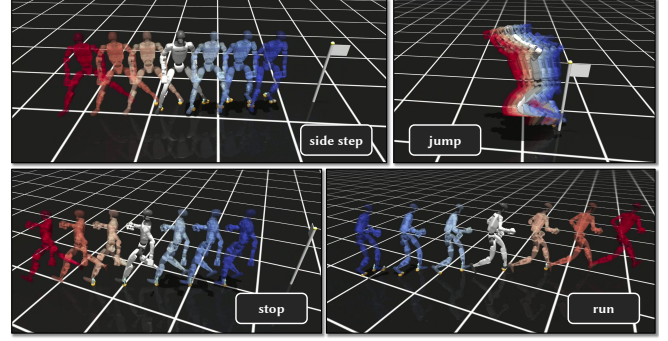


Fig. 4. Our root-disentangled decoder automatically warps motion to different root trajectories with the same pose tokens, enabling both robust and precise control. Root trajectories are interpolated linearly from red to blue characters. The neutral-colored character shows the original motion.

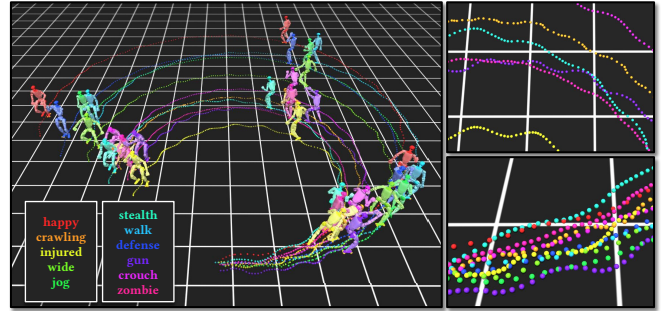


Fig. 5. Neural root refinement automatically adjusts trajectories based on locomotion style. Despite identical initial root states and navigation commands, each style (e.g., happy, crawling, stealth, crouch) produces a distinct root trajectory. Right column: zoomed views showing fine-grained trajectory details.

neural backbone. Unlike prior methods that require task-specific conditioning during training, such as velocity commands, one-hot task labels, or learned task descriptors, our model requires no such supervision. All task specifications are instead defined at runtime, enabling broad applicability without retraining.

### 6.1 Smart Locomotion

Controllable locomotion presents two key challenges: (1) generating natural motions from a wide range of velocity and heading commands, which is nontrivial when the “correct” velocity for a given navigation context is unknown or ambiguous, or when users specify infeasible targets (control dead zone); and (2) supporting diverse locomotion styles with smooth transitions. Smart locomotion addresses these with two mechanisms: *progressive root trajectory refinement*, which reliably and intelligently balances user input with neural correction, and *styled keyframe placement* with repetitive planning for natural style transitions. We introduce each component in temporal order of execution.

*Critically Damped Spring Model for Initial Root Trajectory Estimation.* As defined earlier, root states  $r_g$  comprise the projected global position and heading angle of the pelvis joint. Inspired by the

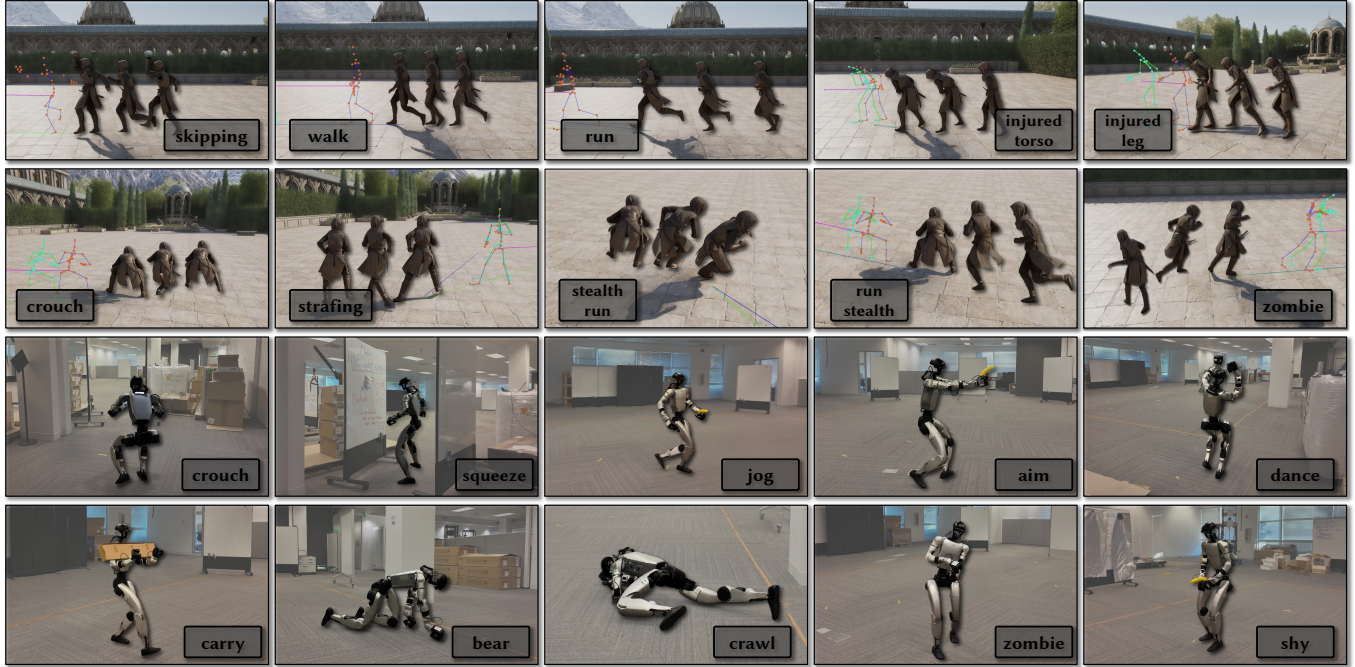


Fig. 6. Diverse locomotion styles generated by smart locomotion in UE5 (top two rows) and on the Unitree G1 robot (bottom two rows). Styles include walking, running, crouching, strafing, stealth, crawling, and expressive movements, all generated from a single neural backbone.

damped spring model commonly used in motion matching [Clavet 2016; Holden et al. 2020], we generate an initial guess of the root trajectory. We first compute a naive target  $r_{g,1}$  by linearly extrapolating the current state using the user’s commanded velocity and heading (from gamepad, keyboard, etc.) over a 1.0s horizon. We then apply a critically damped spring to smooth the trajectory  $r(t)$  from the current root  $r_0$  toward the target  $r_{g,1}$ :

$$r(t) = e^{-\gamma t} ((r_0 - r_{g,1}) + (v_0 + \gamma(r_0 - r_{g,1}))t) + r_{g,1}, \quad (6)$$

where  $\gamma$  is the damping coefficient, and  $r_0$  and  $v_0$  are the current position and velocity. Setting  $t = 1.0$  s yields the spring-smoothed target  $r_{g,2}$ .

*Style Control via Keyframe Selection.* To control locomotion style, we leverage our in-betweening formulation by placing keyframes, sampled from short reference clips or human-authored poses of the desired style, onto the spring-smoothed trajectory  $r_{g,2}$ .

Traditional locomotion controllers must carefully align keyframes with footstep phases and explicitly handle style transitions to avoid artifacts. In contrast, our method automatically generates high-quality, natural motions with smooth style transitions. We attribute this to two factors: (1) a powerful generative backbone that has learned diverse motion skills and can adapt keyframes to context, and (2) a replanning mechanism that re-plans after a fixed interval, preventing the output from being rigidly locked to arbitrarily placed keyframes. We demonstrate in both our UE5 demo and robotics deployment that even with a single keyframe and artificially inserted velocity/heading values, MotionBricks produces natural, coherent motions. We also support keyframe generation via 2D blendspaces,

showcasing the flexibility of our approach for integration with existing industrial animation tools. We discuss the differences between traditional 2D blendspaces and the keyframe blending approach in MotionBricks in Appendix D.

Table 1. Progressive root trajectory refinement in smart locomotion.

stage	symbol	description
naive	$r_{g,1}$	direct extrapolation
spring model	$r_{g,2}$	critically damped spring smoothing
root module	$r_{g,3}$	refinement with style and timing
decoder	$r_{g,4}$	refinement with full-body motion

*Neural Root Refinement.* While the spring-smoothed root  $r_{g,2}$  provides a reasonable initial estimate, it does not account for constraints imposed by context and target keyframes. We therefore feed  $r_{g,2}$  along with keyframe information into our root module to produce the refined root trajectory  $r_{g,3}$ . The root module generates trajectories with natural, realistic movement details (as opposed to the mechanistic spring-smoothed output), and leverages timing prediction to avoid artifacts such as invalid or abrupt velocities, effectively safeguarding the output from control dead zones. As shown in Figure 5, even with identical initial root states and navigation commands from keyboard, the final root trajectories are automatically adjusted to produce natural motions that respect style constraints. Finally,  $r_{g,3}$  is passed to the decoder, which produces the final motion while further refining root values in coordination with full-body details



Fig. 7. Smart objects enable diverse scene and object interactions from a flexible number of keyframes. Top rows: Examples including ledge climbing, vaulting, sitting, falling from heights, and object pickup. Bottom row: Keyframe authoring interface using standard translation and rotation gizmos in UE5.

such as footsteps and velocities, yielding the output  $r_{g,4}$ . Crucially, the decoder’s disentangled handling of root and pose enables robust adaptation to varied inputs. As shown in Figure 4, the final motions adapt gracefully to interpolated root trajectories, enabling flexible post-processing and root constraint adjustments in downstream applications.

As summarized in Table 1, the root trajectory undergoes progressive refinement from user input to final motion, striking the balance between responsiveness of user input and neural corrections. In Figure 6, we showcase diverse locomotion styles generated by smart locomotion in both UE5 and on the Unitree G1 robot.

## 6.2 Smart Object

Object interaction in games and simulations traditionally relies on template-based solutions: pre-canned animation clips or pre-defined robotic skill routines tied to specific object positions and orientations. This approach suffers from two major limitations: (1)

significant authoring burden, as different objects and configurations require dedicated animation clips or task-specific processing, and (2) rigid, inflexible playback that struggles with varied object placements or smooth transitions into and out of interactions.

Our smart object primitive addresses these limitations by leveraging the in-betweening capability of our neural backbone. Rather than relying on pre-canned animation clips, we define each interactive object via two lightweight components: *intent keyframes* that guide the character toward key poses (serving as either soft guidance or hard constraints depending on the interaction phase), and an *interaction binding* that anchors the interaction logic to the object.

*Smart Object Intent Keyframes.* Keyframes define the essential poses for an interaction behavior. For example, a climbing interaction may consist of multiple keyframe sets: (1) a contact pose where the character’s hands reach the wall edge, and (2) an exit pose where

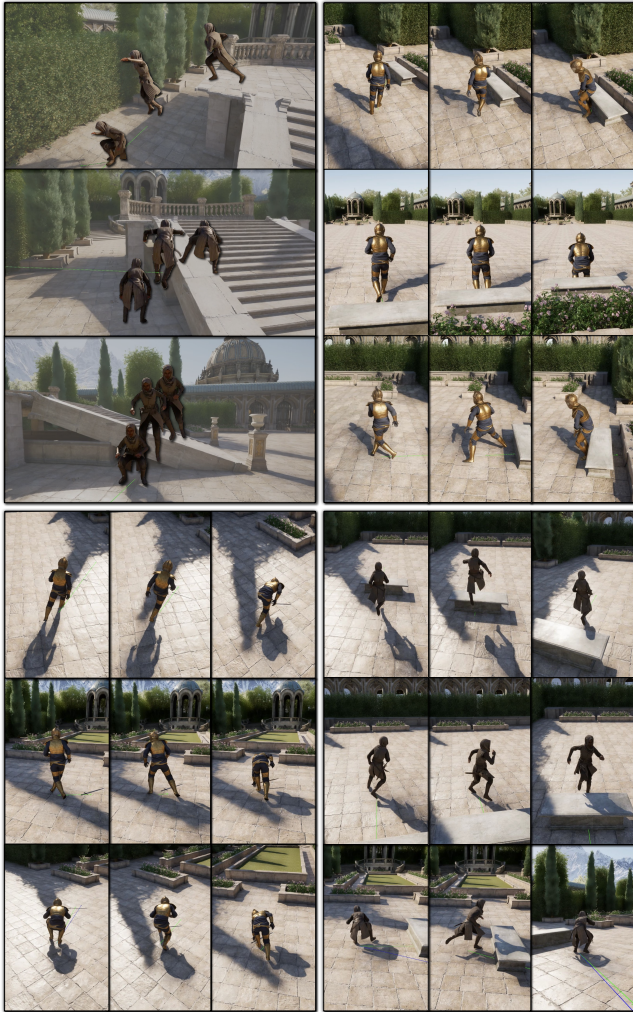


Fig. 8. Automatic motion variations from smart objects using the same keyframe setup. Clockwise from top-left: falling, sitting, vaulting, and sword pickup. Each skill shows three variants, with keyframes displayed chronologically from left to right.

the character stands on the ledge after climbing. As shown in Figure 7, we support a variable number of keyframe sets per behavior, with each set containing one or more keyframes. Users can obtain keyframes from a motion library or author them manually using standard translation, rotation, and IK gizmo tools.

To control how strictly keyframes are enforced, each keyframe set includes a drop-frame attribute  $\tau$ . Setting  $\tau = 0$  treats the keyframe as a *hard constraint*: the generated motion must fully reach the target pose before proceeding, which is essential for precise contact interactions like hand placement during climbing. Setting  $\tau > 0$  treats the keyframe as *soft guidance*: the model uses the keyframe to shape the motion’s intent, but may transition to the next phase up to  $\tau$  frames early without reaching the exact pose. This is useful for preparation keyframes (e.g., poses preparing the hands and feet for

Table 2. Overview of datasets used in our experiments.

Dataset	Hours	Train clips	Test clips	Joints
350k	700	315,162	35,018	27
70k	140	62,132	35,018	27
HumanML3D	28.6	23,206	2,578	22
LaFAN1-G1	4.6	2,362	262	34

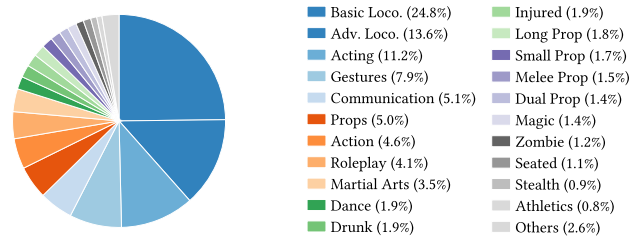


Fig. 9. Distribution of the 350k dataset across its 36 categories. “Basic.Loco.”, “Adv.Loco.”, and “Ext.Loco.” denote basic, advanced, and extreme locomotion, respectively. “Human Inter.” denotes human interaction. Categories comprising fewer than 0.5% of total clips are grouped under “Other”.

a vaulting motion) or exit poses, where strict enforcement would produce unnatural, over-constrained motion. We also note that a single smart object can define multiple keyframe sets for skill variety, and we support neurally generated keyframes on the fly. Additionally, each keyframe includes a boolean flag  $\omega$  indicating whether it should be rotated around the interaction pivot at runtime, as discussed below.

*Interaction Binding.* Each smart object defines an *interaction binding* that manages three aspects of interaction: detection, placement, and keyframe anchoring, as shown below. (1) Detection: We derive an interaction mesh from the object’s physical geometry, which serves as a trigger volume. At runtime, box collision tracing determines whether the character is within interaction range and can initiate the interaction. (2) Object sockets and placement: For portable objects, the binding defines sockets specifying where the object attaches to the character (e.g., hands) and placement logic determining where the object should be positioned once released, for instance, snapping to a valid surface detected via ray cast. (3) Keyframe anchoring: For scene and static object interactions, the mesh’s world transform acts as a pivot to anchor and rotate keyframes (for those keyframes with  $\omega = 1$ ) relative to the object. This enables a character to approach the same ledge from different angles, or interact with objects at varying positions, using identical keyframe definitions, a key advantage of our generative approach over template-based playback. Our implementation leverages UE5’s built-in modules (box traces, collision detection, socket systems), demonstrating that our method integrates seamlessly with existing game engine infrastructure. As shown in Figure 7, diverse scene and object interactions can be authored in a scalable way. Furthermore, Figure 8 demonstrates automatic motion variations. Falling motions adapt to different heights, and scene object interactions such as pickup,

vaulting, and sitting, adjust to varied angles, velocities, and styles without additional authoring.

## 7 Experiments

In this section, we evaluate MotionBricks across different datasets in both animation and robotics applications. Section 7.1 describes the datasets used in our experiments. Section 7.2 presents numerical comparisons with recent state-of-the-art methods. Section 7.3 provides ablation studies analyzing the contribution of different components. Section 7.4 discusses interactive animation, robotics applications, and engineering details.

### 7.1 Dataset

We evaluate our method on four datasets of varying scales. Our primary dataset is a proprietary motion capture collection containing approximately 700 hours of high-quality motion data with 350k motion clips, covering diverse actions including locomotion, combat, sports, and object interactions. Without further specification, we call it **350k** dataset. The 350k dataset covers 9,300 unique skills across 36 categories, captured from over 163 performers, as shown in Figure 9. The test set comprises 10% of the full dataset and consists of two 5% subsets: one split randomly, and the other split by motion skill category to maximize distributional difference from the training set. The dataset, along with download and usage instructions, is available on our project website. Additional statistics are provided in Appendix E.

For public benchmark comparison, we use **HumanML3D** [Guo et al. 2022a], an open-source dataset with 28.6 hours and 14,616 motion sequences. We removed broken clips and segment motions that are too long, so the final number of training clips is 23,206 and test clips is 2,578. Additionally, we evaluate on **LaFAN1-G1**, a variant of LaFAN1 [Harvey et al. 2020] retargeted to the Unitree G1 humanoid skeleton, which contains 4.6 hours with 2.4k clips and has become a common benchmark for robotics research. Since most LaFAN1 motion clips are quite long, we cut them into 6s clips. We also added four end-effector joints to the G1 skeleton, which has 30 joints (29 hinge joints and a free pelvis joint) so that in total we have 34 joints. We also train on a 140-hour subset of our proprietary dataset (referred to as **70k**) with 62k clips to study scaling behavior. For all four datasets, we use 30 FPS and retrain our method and baselines for fair comparison.

In total, as shown in Table 2, we have 4 datasets covering different scales of the datasets and skeletons of different number of joints and morphologies across animation and robotics domains. We use the same motion representation as described in Section 4.

### 7.2 Quantitative Evaluation

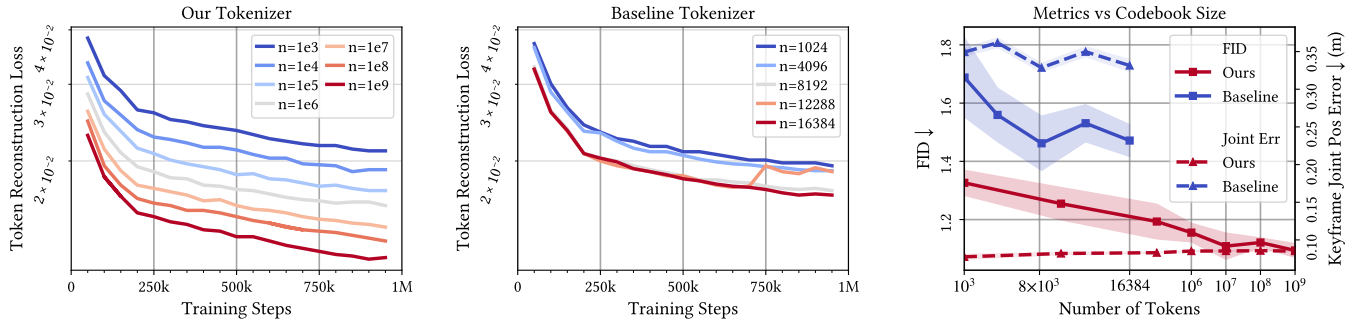
To demonstrate the effectiveness of our method, we compare with the following state-of-the-art methods: 1) Cond. In-betweening [Kim et al. 2022], 2) Delta-interpolator [Oreshkin et al. 2024], 3) Two-stage Trans. [Qin et al. 2022], 4) CondMDI [Cohan et al. 2024], 5) MMM [Pinyoanuntapong et al. 2024b], and 6) Closd-DiP [Tevet et al. 2025]. The first three are deterministic methods using transformers to predict in-between motions directly, with the two-stage transformer method employing a cascade of two transformers for

refinement. CondMDI and Closd-DiP are generative diffusion models. The former uses a U-Net architecture specifically designed for motion in-betweening, while the latter employs a closed-loop framework between physical simulation and a diffusion planner. In our experiments, we use only the diffusion planner of Closd-DiP for a fair comparison. MMM also uses a token-based approach but employs a single-head tokenizer without root-pose disentanglement or progressive design. We modified the open-source code repositories to support evaluation on our datasets, using the same motion representation and training settings as the original methods. All methods are formulated for the in-betweening task for our evaluation.

*Evaluation Metrics.* We randomly sample the context and target keyframes from the test set and generate motions between the two keyframes. We evaluate methods across seven categories of metrics: **(1) Speed:** We report frames per second (FPS) and latency in milliseconds to measure real-time performance. We take the latency from the original methods’ paper directly if the values are provided; otherwise, we use the modified open-source code to estimate the FPS and latency. **(2) Distribution:** We use Maximum Mean Discrepancy (MMD) [Jayasumana et al. 2024] and Fréchet Inception Distance (FID) [Heusel et al. 2017] to measure how well the generated motion distribution matches the ground truth distribution. To obtain the MMD and FID scores, we retrain a motion autoencoder following TMR [Petrovich et al. 2023], omitting the text branch, and compute both metrics in the learned latent space. Following [Jayasumana et al. 2024], we scale the distribution metrics for more human-readable values. We provide the groundtruth FID and MMD values for reference. **(3) Human Evaluation:** We conduct a user study with 40 participants from diverse backgrounds, including artists, technical specialists, engineers, and scientists from the animation and robotics domain. For each trial, participants view motions generated by 3 randomly selected methods on the same in-betweening setup and select a single winner. Win rate is computed as the percentage of trials where a method is selected as the best. Participants also rate each motion on a 1–5 Likert scale for overall quality. **(4) Diversity:** Following [Tevet et al. 2022], we run each generation 20 times and compute two diversity metrics: joint-level diversity, which measures the variability of joint positions across generated motions, and latent-level diversity, which similarly measures the mean pairwise distance in the learned latent space. **(5) Smoothness:** Joint Jitter and Root Jitter (both in  $m/s^2$ ) measure the acceleration magnitude of joint positions and root trajectory, respectively, as was proposed in [Gou et al. 2025]. It aligns well with the human perception of smoothness, and lower values indicate smoother motion. **(6) Physical Plausibility:** Foot Skate (m/frame) measures undesirable sliding of feet during ground contact, Penetration (m) measures ground penetration depth, and Foot Contact Accuracy (%) measures how often foot contact predictions align with ground truth labels. More accurate foot prediction will be helpful for postprocessing in downstream applications when necessary. **(7) Precision:** Context/Target Keyframe errors (m) measure joint position deviation from the input keyframes, Context/Target Root errors (m) measure root position deviation, and Reaching Success (%) measures whether the generated motion successfully reaches the target keyframe within a threshold. Note that instead of using mean

Table 3. Quantitative comparison on the 350k dataset. ↓ indicates lower is better, ↑ indicates higher is better. Best results are in **bold**, while underlined indicates the second-best or competitive runner-up results.

Method	Speed		Distribution		Human Eval		Diversity		Smoothness		Physical Plausibility			Precision				
	FPS ↑	Latency ↓	MMD ↓	FID ↓	Win ↑	Score ↑	Jnt ↑	Latent ↑	Jnt Jit. ↓	Root Jit. ↓	Foot Sk. ↓	Pene. ↓	Foot Acc. ↑	Ctx KF ↓	Tgt KF ↓	Ctx Root ↓	Tgt Root ↓	Reach ↑
Ground Truth	N/A	N/A	0.0533	0.022	N/A	N/A	N/A	N/A	3.47	1.96	0.0008	0.0009	N/A	N/A	N/A	N/A	N/A	N/A
Cond. Inbtwn. (2022)	<b>27,000</b>	<u>2.4ms</u>	0.1093	1.594	0.8%	1.83	0.00	0.00	16.88	12.65	0.018	0.034	51.8%	0.094	<u>0.078</u>	0.050	<u>0.051</u>	87.7%
Two-stage (2023)	<u>24,600</u>	2.6ms	0.1094	1.643	5.0%	2.18	0.00	0.00	38.39	33.72	0.038	0.037	43.7%	0.293	0.102	0.215	0.061	85.7%
Delta-interp. (2024)	18,800	3.4ms	0.1107	1.774	1.3%	1.92	0.00	0.00	24.92	20.15	0.028	0.032	50.6%	0.288	0.130	0.152	0.066	81.9%
CondMDI (2024)	1,930	33.2ms	<u>0.1080</u>	1.213	15.6%	3.13	2.03	16.42	16.19	14.62	0.012	<u>0.009</u>	62.7%	<b>0.045</b>	0.143	0.022	0.118	61.3%
CondMDI + CFG (2024)	1,050	60.5ms	<u>0.1082</u>	<u>1.201</u>			2.05	16.72	17.03	15.39	0.012	0.009	63.2%	0.049	0.121	0.024	0.097	65.9%
MMM (2024)	3,600	18.1ms	0.1176	1.544			<b>4.38</b>	<b>26.51</b>	5.40	3.11	0.005	0.011	86.5%	0.311	0.377	0.006	0.071	34.8%
MMM + 5 steps	2,700	24.5ms	0.1225	1.562	<u>19.9%</u>	<u>3.19</u>	3.21	<u>20.97</u>	5.36	3.10	<u>0.004</u>	0.011	<u>86.8%</u>	0.300	0.365	0.006	0.071	36.8%
MMM + 10 steps	1,400	46.2ms	0.1234	1.564			2.96	19.62	<u>5.35</u>	<u>3.09</u>	<u>0.004</u>	0.011	<u>86.8%</u>	0.300	0.364	<u>0.005</u>	0.070	36.6%
CloSD-DiP (2025)	4,200	15.3ms	0.1076	1.292			1.84	13.40	14.03	11.94	0.015	<u>0.010</u>	52.7%	0.071	0.129	0.042	0.091	75.7%
CloSD-DiP + CFG	2,800	23ms	0.1084	1.288	15.1%	2.90	1.70	12.45	18.33	15.98	0.022	0.012	50.5%	0.139	0.163	0.089	0.113	54.9%
CloSD-DiP + 5 steps	14,500	4.4ms	0.1081	1.338			1.51	10.86	17.56	15.54	0.017	0.012	47.4%	<u>0.069</u>	0.122	0.040	0.082	78.2%
CloSD-DiP + CFG + 5 steps	10,500	6.1ms	<u>0.1080</u>	1.319			1.40	10.08	20.85	18.60	0.023	0.013	47.7%	0.126	0.181	0.080	0.117	52.0%
<b>Ours</b>	15,000	<b>2ms</b>	<b>0.1056</b>	<b>1.054</b>	<b>86.5%</b>	<b>4.06</b>	<u>3.75</u>	<u>20.67</u>	<b>3.38</b>	<b>1.82</b>	<b>0.003</b>	<b>0.008</b>	<b>92.6%</b>	<u>0.052</u>	<b>0.076</b>	<b>0.001</b>	<b>0.023</b>	<b>99.6%</b>

Fig. 10. Scalability comparison between our multi-head tokenizer and a single-head baseline. Left and middle: Token reconstruction loss during training for varying codebook sizes for our method and the baseline. Our tokenizer continues to improve with larger codebooks (up to  $10^9$  tokens), while the baseline plateaus quickly. Right: Trade-off between FID (distribution quality) and keyframe joint position error as codebook size increases. Our method achieves better FID with larger codebooks while maintaining low keyframe error.

per joint position error (MJPE or MPJPE), we use the **max** joint position error to measure the precision of the generated motion. We argue that max per joint position error is more sensitive and aligned with the human perception. We define Reaching Success (%) as the percentage of the generated motion that successfully reaches the target keyframe’s root position within 5 cms and 15 degrees of the target keyframe’s heading angle.

*Results Summary.* As shown in Table 3 and Table 4, our method achieves consistently better performance than the baseline methods on all datasets. Key improvements on the 350k dataset include: better distribution matching including FID and MMD, lower jitter, better keyframe precision and near-perfect reaching success. Our algorithm also receives the highest score and win rate in human evaluation. We further provide visual in-betweening comparisons in Figure 12, where MotionBricks exhibits notably better physical plausibility, naturalness, and keyframe adherence than the baselines. More results are provided in Appendix F. For diversity metrics, MotionBricks outperforms baselines on larger datasets, while on smaller datasets such as HumanML3D and LaFAN1, baseline

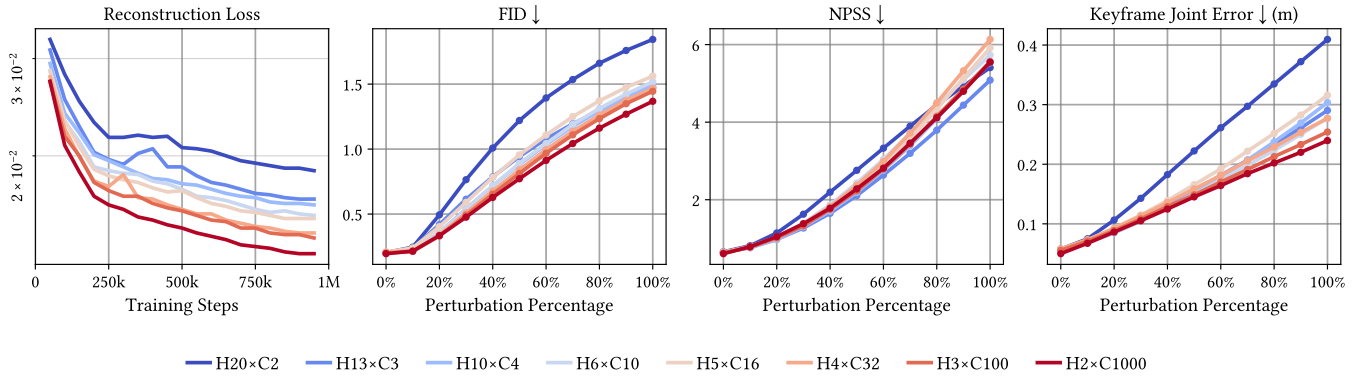
methods perform comparably or marginally better. Overall, there is no significant trade-off between diversity and precision or speed. We observe that non-generative methods achieve reasonable precision but have bad FID, MMD scores, and poor human evaluation scores, while diffusion-based methods improve distribution quality but struggle with satisfying spatial constraints. We also note that having more iterations in MMM gives minor improvements and classifier-free guidance for CondMDI and CloSDiP gives minor improvements to the metrics. Our method also achieves one of the best runtime efficiency among all methods, with only 2ms latency to generate the motion between two keyframes, achieving 15,000 FPS throughput.

### 7.3 Ablation Studies and Analysis

*Scalability of MotionBricks and baseline.* As shown in Figure 10, we compare the scalability of our structured multi-head tokenizer against a standard single-head VQ-VAE baseline. Our method (left) demonstrates consistent performance improvements as the number of tokens increases, while the baseline (middle) plateaus quickly

Table 4. Quantitative comparison on **LaFANI-G1**, **HumanML3D**, and **Bones-70k** datasets. ↓ indicates lower is better, ↑ indicates higher is better. Best results are in **bold**, second best are underlined.

Dataset	Method	Distribution		Diversity		Smoothness		Physical Plausibility			Precision				
		MMD ↓	FID ↓	Jnt ↑	Latent ↑	Jnt Jit. ↓	Root Jit. ↓	Foot Sk. ↓	Pene. ↓	Foot Acc. ↑	Ctx KF ↓	Tgt KF ↓	Ctx Root ↓	Tgt Root ↓	Reach ↑
LaFANI-G1	Cond. Inbtwn. (2022)	0.2897	1.693	0.00	0.00	40.73	38.17	0.030	0.019	83.8%	0.374	0.453	0.264	0.369	11.5%
	Two-stage (2022)	0.3083	1.985	0.00	0.00	85.06	81.31	0.036	0.021	83.4%	0.535	0.498	0.400	0.399	8.1%
	Delta-interp. (2024)	0.3060	1.846	0.00	0.00	45.27	43.19	0.028	0.020	86.1%	0.404	0.451	0.306	0.360	11.6%
	CondMDI (2024)	0.2860	<u>1.004</u>	0.75	4.08	17.89	16.93	0.014	<u>0.004</u>	82.8%	<u>0.057</u>	0.592	0.029	0.555	6.6%
	MMM (2024)	0.2920	1.189	<b>1.78</b>	<b>17.66</b>	<u>7.00</u>	<u>4.05</u>	<u>0.006</u>	<u>0.004</u>	<b>90.5%</b>	0.289	0.331	<u>0.008</u>	<u>0.103</u>	37.3%
	Closd-DiP (2025)	0.2839	1.165	<u>1.72</u>	<u>12.07</u>	13.96	12.41	0.013	0.005	85.5%	0.080	<b>0.178</b>	0.049	0.128	<u>50.7%</u>
	<b>Ours</b>	<b>0.2835</b>	<b>0.891</b>	<u>1.76</u>	11.42	<b>4.53</b>	<b>2.53</b>	<b>0.003</b>	<b>0.003</b>	<u>88.9%</u>	<b>0.053</b>	<u>0.179</u>	<b>0.006</b>	<b>0.100</b>	<b>61.1%</b>
HumanML3D	Cond. Inbtwn. (2022)	0.1107	1.497	0.00	0.00	10.07	8.04	0.019	0.016	62.7%	0.126	<b>0.090</b>	0.051	0.050	92.6%
	Two-stage (2022)	0.1131	1.709	0.00	0.00	19.78	16.51	0.030	0.017	61.6%	0.232	0.132	0.144	0.073	75.6%
	Delta-interp. (2024)	0.1133	1.585	0.00	0.00	10.86	8.74	0.031	0.015	60.1%	0.288	0.151	0.132	0.090	68.0%
	CondMDI (2024)	<b>0.1101</b>	<u>1.054</u>	1.49	<u>11.09</u>	5.82	4.87	0.009	<b>0.002</b>	77.4%	<b>0.037</b>	0.147	0.016	0.121	52.3%
	MMM (2024)	0.1238	1.068	<u>1.89</u>	<u>11.23</u>	<u>2.92</u>	<u>1.94</u>	<b>0.005</b>	<u>0.002</u>	<b>89.2%</b>	0.242	0.293	<u>0.003</u>	<u>0.037</u>	61.2%
	Closd-DiP (2025)	0.1104	1.192	1.28	9.09	6.12	5.06	0.012	0.003	69.9%	0.062	<u>0.105</u>	0.035	0.070	80.8%
	<b>Ours</b>	0.1114	<b>0.914</b>	<b>2.56</b>	<b>13.47</b>	<b>2.08</b>	<b>1.29</b>	<u>0.006</u>	<u>0.002</u>	<u>83.5%</u>	<u>0.050</u>	0.141	<b>0.001</b>	<b>0.022</b>	<b>98.9%</b>
Bones-70k	Cond. Inbtwn. (2022)	0.1073	1.559	0.00	0.00	18.18	13.91	0.019	0.031	49.4%	0.088	<u>0.091</u>	0.044	<u>0.062</u>	83.5%
	Two-stage (2022)	0.1081	1.604	0.00	0.00	39.32	34.11	0.039	0.037	44.3%	0.347	0.128	0.257	0.076	77.2%
	Delta-interp. (2024)	0.1183	2.024	0.00	0.00	39.32	32.79	0.048	0.042	44.8%	0.556	0.499	0.226	0.129	41.4%
	CondMDI (2024)	<u>0.1062</u>	<u>1.191</u>	2.11	<u>16.41</u>	16.62	15.10	0.012	<u>0.009</u>	60.7%	<b>0.049</b>	0.190	0.024	0.160	48.5%
	MMM (2024)	0.1227	1.494	<u>2.38</u>	<u>16.57</u>	<u>5.27</u>	<u>3.01</u>	<u>0.005</u>	0.011	85.6%	0.306	0.375	<u>0.006</u>	0.074	40.0%
	Closd-DiP (2025)	<u>0.1062</u>	1.305	1.86	13.13	14.14	12.10	0.015	0.011	50.8%	0.065	0.136	0.036	0.102	68.6%
	<b>Ours</b>	<b>0.1050</b>	<b>1.090</b>	<b>3.48</b>	<b>20.09</b>	<b>3.24</b>	<b>1.80</b>	<b>0.004</b>	<b>0.009</b>	<b>84.6%</b>	<u>0.054</u>	<b>0.087</b>	<b>0.001</b>	<b>0.023</b>	<b>99.6%</b>

Fig. 11. Ablation study on multi-head tokenization with fixed total codebook capacity ( $\sim 10^6$  tokens). “H” denotes the number of heads, “C” denotes the codebook size per head in the experiment name. Left: tokenizer’s reconstruction loss during training. Right three plots: FID, NPSS, and keyframe error under token perturbation, simulating lower-bound motion quality in real applications.

regardless of codebook size. In the right figure, we analyze the trade-off between distribution quality (FID) and keyframe precision as we scale from  $10^3$  to  $10^9$  tokens. Our method achieves better FID with more tokens but shows slightly increased keyframe errors at extreme scales, suggesting an optimal operating point around  $10^6$  to  $10^7$  tokens that balances both objectives. In contrast, the baseline exhibits no meaningful improvement in either FID or keyframe error with increased capacity, and maintains significantly worse absolute metrics across all scales. This demonstrates that our structured multi-head design is essential for effectively utilizing model capacities.

*Multi-head tokenization.* In Figure 11, we ablate different combinations of the number of heads and tokens per head while keeping the total number of tokens fixed at approximately  $10^6$ . We evaluate the tokenizer’s reconstruction loss, FID, the normalized power spectrum similarity score (NPSS) [Gopalakrishnan et al. 2019], and

keyframe errors under token perturbation. We introduce token perturbation to simulate prediction errors and assess how gracefully the model degrades, i.e., the lower-bound of the generated motion quality in potential real applications. We observe that increasing tokens per head raises reconstruction error, though the downstream metrics (FID, NPSS, keyframe errors) remain relatively stable. However, the perturbation analysis reveals a trade-off: too many tokens per head degrades NPSS (temporal coherence), while too few tokens per head hurts FID and keyframe precision. Based on these results, we recommend 128 to 256 tokens per head as the optimal configuration, combined with a total codebook capacity of approximately  $10^9$  tokens for the best overall performance.

*Scaling behavior with the size of the dataset.* In Figure 13, we investigate how model performance scales with dataset size by training on progressively larger subsets of our 350k dataset. This analysis

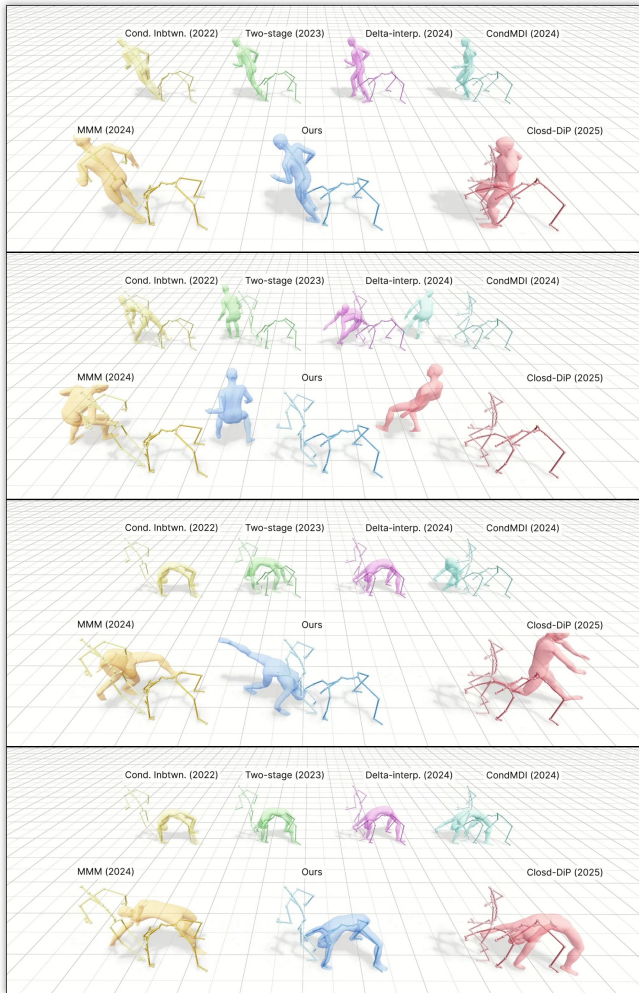


Fig. 12. An in-betweening test case transitioning from running left to a backbend pose. Time progresses from top to bottom; algorithm names are labeled in each subfigure. MotionBricks is shown as the blue character.

provides guidance for practitioners training on datasets of varying sizes. Both our method and baselines show improved tokenizer and generative model metrics with increasing data. However, the baseline struggles significantly on larger datasets due to its limited capacity and scalability, which becomes overwhelmed and exhibits degraded performance. Interestingly, FID scores worsen with larger datasets for both methods. We hypothesize that models trained on smaller datasets may overfit and produce less diverse outputs that happen to match the training distribution more closely, resulting in artificially lower FID scores. In contrast, keyframe precision consistently improves with more data, and our method maintains this improvement while the baseline’s keyframe errors increase as the dataset size grows. This demonstrates that our architecture is better suited to leveraging large-scale motion data while still working well with small datasets.

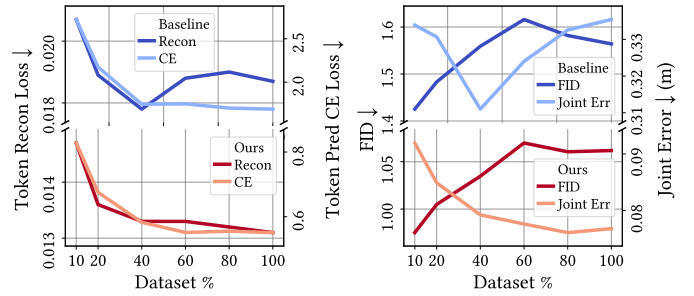


Fig. 13. Scaling behavior with dataset size from 10% to 100% of the 350k dataset. Left: Evaluation losses including the reconstruction loss for the tokenizer, and cross-entropy for token prediction. Right: FID and keyframe joint position error. Our method scales effectively with more data, while the baseline struggles at larger scales.

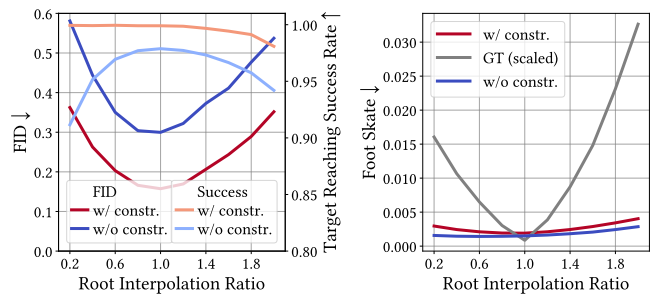


Fig. 14. Root trajectory interpolation analysis. Root interpolation ratio of 1.0 is the original trajectory;  $<1.0$  compresses,  $>1.0$  stretches. Left: FID and target reaching success rate remain stable across interpolation ratios, especially with keyframe constraints. Right: Foot skate stays low even under significant root manipulation, demonstrating the decoder’s robustness.

*Root interpolation designs.* Our decoder supports root trajectory interpolation, where we keep the predicted pose tokens fixed but override the root values with interpolated trajectories. This is a practical utility for runtime animation, enabling precise global keyframe matching without introducing foot sliding artifacts as mentioned earlier in Section 6.1. As shown in numerical results in Figure 14 and visual results in Figure 4, our design successfully hits target keyframes while maintaining FID quality and avoiding foot skating degradation when interpolating root trajectories. In Figure 14, we also compare decoding with and without pose keyframe constraints (as shown in Figure 3). Providing pose constraints improves both FID and target reaching success, demonstrating the value of our flexible constraint handling. This robustness to root manipulation stems from our modular architecture, where the pose decoder is trained to be agnostic to the source of root information.

*Additional Ablation Studies.* We provide additional ablation studies in the appendix: (1) the impact of varying the number of GPUs during training (Appendix A), (2) the choice of tokenizer design between FSQ and VQ-VAE (Appendix B), and (3) the effect of replanning frequency during deployment (Appendix C).

## 7.4 Interactive Animation, Robotics Applications and Engineering Details

As shown in Figure 1 and the demo video, we showcase our method in UE5 and deploy it on the Unitree G1 humanoid robot.

We train the tokenizer, root module, and pose module on 4 nodes with 8 GPUs each (32 GPUs total); training with 16 GPUs yields similar visual quality. Training runs for 2 million updates with a batch size of 256 per GPU, using Adam optimizer with learning rate  $5 \times 10^{-5}$  and a cosine schedule with 10k warmup steps decaying to  $2 \times 10^{-6}$ . The tokenizer uses a U-Net architecture with a downsampling rate of 4 (2 downsampling layers), where each layer consists of three residual 1D convolutional layers with 1024 channels and progressively larger kernel sizes (23.5M parameters total). The encoder and decoder have a symmetric architecture. We also experimented with transformer-based encoder and decoder designs, which achieve similar performance but are significantly slower during inference. FSQ [Mentzer et al. 2023] and VQ-VAE yield similar performance (see Appendix B for a detailed comparison). We also include standard foot sliding and velocity losses during tokenizer training to further improve motion quality, particularly important for retargeted datasets where artifacts like foot skating are more common. The pose module is a transformer with embedding dimension 1024, 16 attention heads, and 16 layers (150M parameters). The root module uses a similar architecture with fewer layers and 512 embedding dimension and 12 attention heads (50M parameters). Training takes approximately 7 days for the tokenizer, 3 days for the root module, and 7 days for the pose module on H100 GPUs, though Figure 13 suggests performance plateaus earlier. Training on L40S GPUs yields similar performance but is approximately  $2\times$  slower. During inference, we achieve 2ms latency and 15,000 FPS throughput on a desktop with an RTX 5090 GPU. For both the UE5 and G1 applications, we use the same model architecture and training settings.

For the UE5 application, we utilize UE5’s plugin system to implement both our neural backbone and smart primitives. Models are exported via ONNX and loaded in a native C++ plugin using TensorRT. Since the training skeleton differs from the demo characters, we use UE5’s built-in real-time retargeter [Games 2025] to retarget motions on the fly to the “Messenger” and “Guard” characters. The authoring time for smart navigation and smart object is similar, each taking less than 10 minutes for non-expert users. More details about the UE5 plugin are provided in Appendix G.

By combining our method with a physical tracking controller [Chen et al. 2025a; Luo et al. 2025; Wang et al. 2020; Yin et al. 2025; Zeng et al. 2025; Zhang et al. 2025], we drive real robots by supplying the generated kinematic motions as tracking targets. In our experiments, we use the tracking controller from Luo et al. [2025] to control the G1 humanoid robot. For G1 applications, we retarget our proprietary dataset to the G1 skeleton using a modified version of GMR [Araujo et al. 2025; Ze et al. 2025]. Automatic retargeting introduces some artifacts, but the results remain acceptable due to the scale of the dataset. For deployment, we use the Jetson Orin platform and Unitree’s official SDK, with models exported via ONNX and loaded using TensorRT similar to the UE5 application. Due to

Orin’s limited compute, latency increases to 5ms per inference. Re-planning is triggered at 10 Hz or whenever commands change, also similar to the strategy used in the UE5 demo. This lazy replanning strategy is more important for robotics deployment, as it avoids interference and competition of resources with the low-level control policy, which is sensitive to latency.

## 8 Discussion

In this work, we introduced MotionBricks, a real-time motion synthesis framework designed to bridge the gap between recent advances in generative models and the practical demands of animation and robotics production. By combining a scalable latent neural backbone with intuitive smart primitive interfaces, our approach enables diverse task-agnostic motion generation without fine-tuning or tagging. Through experiments on datasets of varying scales and demonstrations in both UE5 and on the Unitree G1 robot, we showed that a unified framework can address motion control across virtual characters and physical robots.

### 8.1 Limitations and Future Work

**Dataset Scale:** Our dataset of 350,000 motions, while substantial, remains small with limited coverage of rare motion types (e.g., only one clip for vaulting at 0.5m). Additionally, object interactions are not explicitly modeled. Continued dataset expansion [Team 2025] and leveraging video foundation models [Ni et al. 2025; Wen et al. 2025] offer promising paths forward.

**Visual Planning:** Real-world deployment lacks access to privileged simulation information such as ground-truth object poses and terrain geometry, disabling our smart object primitive. Robots must instead rely on noisy sensor input (RGB-D, SLAM). Developing kinematic planners driven by visual input [Chen et al. 2025b; He et al. 2025; Zhu et al. 2026] would benefit both robotics and animation domains.

**Physical Awareness:** Kinematic planners like MotionBricks may generate physically implausible motions (self-collisions) or motions exceeding hardware constraints (e.g., dangerous flips). Co-training the kinematic planner and tracking policy could address this, enabling the planner to produce feasible motions while the policy adapts to the planner’s output distribution.

**Retargeting:** Adapting motions across diverse morphologies remains challenging. Runtime retargeting (used in our UE5 demo) is fast but low-quality; offline retargeting (used for G1) is accurate but requires months of iteration. Recent work [Araujo et al. 2025; Kim et al. 2025; Yang et al. 2025; Ze et al. 2025] shows promise in bridging this gap.

### Acknowledgments

We thank Cyrus Hogg, John Malaska, Will Telford, Jon Shepard, Dmitry Korobchenko, Anna Minx, Edy Lim, Eugene Jeong, Sam Wu, Ehsan Hassani, Charles Zhou, Freya Li, Ling Li, Qiao Wang, and Lina Song for their support and guidance throughout the project. We thank Alberto Guerra, Boon Cotter, Byeong Gyu Park, Craig Christian, Gabriele Leone, Yenal Kal, Pierre Fleau, Miguel Guerrero, and Marc-Andre Carbonneau from the art team for their help in creating the Unreal Engine 5 demo and its assets. We also thank Hung

Yu Ling, Yue Zhao, Danila Krivenkov, Evgenii Tumanov, Morteza Ramezanali, Winston Chen, Yu-Shiang Wong, Jiefeng Li, Yeongho Seol, Jun Saito, Michael Buttner, Haotian Zhang, Yifeng Jiang, Chen Tessler, Sanja Fidler, Umar Iqbal, Jan Kautz, Yan Chang, and Jim Fan for their helpful discussions and feedback.

## References

- Kfir Aberman, Peizhuo Li, Dani Lischinski, Olga Sorkine-Hornung, Daniel Cohen-Or, and Baoquan Chen. 2020. Skeleton-aware networks for deep motion retargeting. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 62–1.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- Simon Alexanderson, Rajmund Nagy, Jonas Beskow, and Gustav Eje Henter. 2023. Listen, denoise, action! audio-driven motion synthesis with diffusion models. *ACM Transactions on Graphics (TOG)* 42, 4 (2023), 1–20.
- Joao Pedro Araujo, Yanjie Ze, Pei Xu, Jiajun Wu, and C Karen Liu. 2025. Retargeting matters: General motion retargeting for humanoid motion tracking. *arXiv preprint arXiv:2510.02252* (2025).
- Okan Arıkan and David A Forsyth. 2002. Interactive motion generation from examples. *ACM Transactions on Graphics (TOG)* 21, 3 (2002), 483–490.
- Alex Bereznyak. 2024. Shadow of HyperPose: New Animation System. In *ACM SIGGRAPH 2024 Talks*. 1–2.
- Bones Studio. 2026. BONES-SEED: Skeletal Everyday Embodied Dataset. <https://bones.studio/datasets>. Open-source motion dataset, approximately 140k motion clips.
- Michael Buttner. 2015. Motion Matching-The Road to Next-Gen Animation. In *Nucl. ai Conference*.
- Michael Buttner. 2019. Machine Learning for Motion Synthesis and Character Control. In *Interactive 3D Graphics and Games (I3D) 2019*. <https://www.youtube.com/watch?v=zuvmQxcCOM4>
- Rui Chen, Mingyi Shi, ShaoLi Huang, Ping Tan, Taku Komura, and Xuelin Chen. 2024. Taming diffusion probabilistic models for character control. In *ACM SIGGRAPH 2024 Conference Papers*. 1–10.
- Sirui Chen, Yufei Ye, Zi-ang Cao, Jennifer Lew, Pei Xu, and C Karen Liu. 2025b. Hand-eye autonomous delivery: Learning humanoid navigation, locomotion and reaching. *arXiv preprint arXiv:2508.03068* (2025).
- Zixuan Chen, Mazeyu Ji, Xuxin Cheng, Xuanbin Peng, Xue Bin Peng, and Xiaolong Wang. 2025a. GMT: General Motion Tracking for Humanoid Whole-Body Control. *arXiv preprint arXiv:2506.14770* (2025).
- Simon Clavet. 2016. Motion matching and the road to next-gen animation. *Proc. of GDC 2016* (2016).
- Setareh Cohan, Guy Tevet, Daniele Reda, Xue Bin Peng, and Michiel van de Panne. 2024. Flexible motion in-betweening with diffusion models. In *ACM SIGGRAPH 2024 Conference Papers*. 1–9.
- Boston Dynamics. 2025. Use Choreographer with Spot. <https://support.bostondynamics.com/s/article/Use-Choreographer-with-Spot-72036>.
- Unreal Engine. 2025. Graphing in Animation Blueprints. <https://dev.epicgames.com/documentation/en-us/unreal-engine/graphing-in-animation-blueprints-in-unreal-engine>.
- Epic Games. 2025. Blend Spaces in Unreal Engine. <https://dev.epicgames.com/documentation/en-us/unreal-engine/blend-spaces-in-unreal-engine>. Accessed: 2026-01-22.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. 2024. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 12873–12883.
- Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. 2015. Recurrent network models for human dynamics. In *Proceedings of the IEEE international conference on computer vision*. 4346–4354.
- Epic Games. 2025. IK Rig Animation Retargeting in Unreal Engine. <https://dev.epicgames.com/documentation/en-us/unreal-engine/ik-rig-animation-retargeting-in-unreal-engine>.
- Google. 2025. Veo 3. <https://deepmind.google/technologies/veo/>
- Anand Gopalakrishnan, Ankur Mali, Dan Kifer, Lee Giles, and Alexander G Ororbia. 2019. A neural temporal model for human motion prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12116–12125.
- Ruiyu Gou, Michiel van de Panne, and Daniel Holden. 2025. Control Operators for Interactive Character Animation. *ACM Transactions on Graphics (TOG)* 44, 6 (2025), 1–20.
- Chuan Guo, Yuxuan Mu, Muhammad Gohar Javed, Sen Wang, and Li Cheng. 2024. Momask: Generative masked modeling of 3d human motions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1900–1910.
- Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. 2022a. Generating diverse and natural 3d human motions from text. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5152–5161.
- Chuan Guo, Xinxin Zuo, Sen Wang, and Li Cheng. 2022b. Tm2t: Stochastic and tokenized modeling for the reciprocal generation of 3d human motions and texts. In *European Conference on Computer Vision*. Springer, 580–597.
- Bo Han, Hao Peng, Minjing Dong, Yi Ren, Yixuan Shen, and Chang Xu. 2024. Amd: Autoregressive motion diffusion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 2022–2030.
- Félix G Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. 2020. Robust motion in-betweening. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 60–1.
- Tairan He, Zi Wang, Haoru Xue, Qingwei Ben, Zhengyi Luo, Wenli Xiao, Ye Yuan, Xingye Da, Fernando Castañeda, Shankar Sastry, et al. 2025. VIRAL: Visual Sim-to-Real at Scale for Humanoid Loco-Manipulation. *arXiv preprint arXiv:2511.15200* (2025).
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* 30 (2017).
- Jonathan Ho and Stefano Ermon. 2016. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*. 4565–4573.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems* 33 (2020), 6840–6851.
- Jonathan Ho and Tim Salimans. 2022. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598* (2022).
- Daniel Holden. 2018. Character control with neural networks and machine learning. In *Proc. of GDC*, Vol. 1. 2.
- Daniel Holden, Oussama Kanoun, Maksym Perepichka, and Tiberiu Popa. 2020. Learned motion matching. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 53–1.
- Daniel Holden, Taku Komura, and Jun Saito. 2017. Phase-functioned neural networks for character control. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 42.
- Daniel Holden, Jun Saito, and Taku Komura. 2016. A deep learning framework for character motion synthesis and editing. *ACM Transactions on Graphics (ToG)* 35, 4 (2016), 1–11.
- Sadeep Jayasumana, Srikumar Ramalingam, Andreas Veit, Daniel Glasner, Ayan Chakrabarti, and Sanjiv Kumar. 2024. Rethinking fid: Towards a better evaluation metric for image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9307–9315.
- Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. 2023. Motiongpt: Human motion as a foreign language. *Advances in Neural Information Processing Systems* 36 (2023), 20067–20079.
- Chung Min Kim, Brent Yi, Hongsuk Choi, Yi Ma, Ken Goldberg, and Angjoo Kanazawa. 2025. PyRoki: A Modular Toolkit for Robot Kinematic Optimization. *arXiv preprint arXiv:2505.03728* (2025).
- Jihoon Kim, Taehyun Byun, Seungyoun Shin, Jungdam Won, and Sungjoon Choi. 2022. Conditional motion in-betweening. *Pattern Recognition* 132 (2022), 108894.
- Serkan Kiranyaz, Onur Avci, Osama Abdeljaber, Turker Ince, Moncef Gabbouj, and Daniel J Inman. 2021. 1D convolutional neural networks and applications: A survey. *Mechanical systems and signal processing* 151 (2021), 107398.
- L Kovar. 2002. Motion graphs. *ACM Trans. Graph.* 21, 3 (2002), 473–482.
- Lucas Kovar and Michael Gleicher. 2004. Automated extraction and parameterization of motions in large data sets. *ACM Transactions on Graphics (ToG)* 23, 3 (2004), 559–568.
- Jehee Lee, Jinxiang Chai, Paul SA Reitsma, Jessica K Hodgins, and Nancy S Pollard. 2002. Interactive control of avatars animated with human motion data. In *Proceedings of the 29th annual conference on Computer graphics and interactive techniques*. 491–500.
- Yongjoon Lee, Kevin Wampler, Gilbert Bernstein, Jovan Popović, and Zoran Popović. 2010. Motion Fields for Interactive Character Locomotion. In *ACM SIGGRAPH Asia 2010 Papers* (Seoul, South Korea) (SIGGRAPH ASIA '10). Association for Computing Machinery, New York, NY, USA, Article 138, 8 pages. doi:10.1145/1866158.1866160
- Jiefeng Li, Jinkun Cao, Haotian Zhang, Davis Rempe, Jan Kautz, Umar Iqbal, and Ye Yuan. 2025. GENMO: A GENeralist Model for Human Motion. *arXiv preprint arXiv:2505.01425* (2025).
- Jiaman Li, Alexander Clegg, Roozbeh Mottaghi, Jiajun Wu, Xavier Puig, and C Karen Liu. 2024a. Controllable human-object interaction synthesis. In *European Conference on Computer Vision*. Springer, 54–72.
- Jiaman Li, Jiajun Wu, and C Karen Liu. 2023. Object motion guided human motion synthesis. *ACM Transactions on Graphics (TOG)* 42, 6 (2023), 1–11.
- Tianyu Li, Calvin Qiao, Guanqiao Ren, KangKang Yin, and Sehoon Ha. 2024b. AAMDm: accelerated auto-regressive motion diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1813–1823.
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. 2022. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747* (2022).
- Jintao Lu, He Zhang, Yuting Ye, Takaaki Shiratori, Sebastian Starke, and Taku Komura. 2024. CHOICE: Coordinated human-object interaction in cluttered environments for pick-and-place actions. *arXiv preprint arXiv:2412.06702* (2024).

- Zhuoyan Luo, Fengyuan Shi, Yixiao Ge, Yujiu Yang, Limin Wang, and Ying Shan. 2024. Open-magvit2: An open-source project toward democratizing auto-regressive visual generation. *arXiv preprint arXiv:2409.04410* (2024).
- Zhengyi Luo, Ye Yuan, Tingwu Wang, Chenran Li, Sirui Chen, Fernando Castañeda, Zi-Ang Cao, Jiefeng Li, David Minor, Qingwei Ben, et al. 2025. Sonic: Supersizing motion tracking for natural humanoid whole-body control. *arXiv preprint arXiv:2511.07820* (2025).
- Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschanen. 2023. Finite scalar quantization: Vq-vae made simple. *arXiv preprint arXiv:2309.15505* (2023).
- Jiayuan Min and Jinxiang Chai. 2012. Motion graphs++ a compact generative model for semantic motion analysis and synthesis. *ACM Transactions on Graphics (TOG)* 31, 6 (2012), 1–12.
- James Ni, Zekai Wang, Wei Lin, Amir Bar, Yann LeCun, Trevor Darrell, Jitendra Malik, and Roel Herzig. 2025. From Generated Human Videos to Physically Plausible Robot Trajectories. *arXiv preprint arXiv:2512.05094* (2025).
- OpenAI. 2024. Sora: Creating Video From Text. <https://openai.com/sora>
- Boris N Oreshkin, Antonios Valkanas, Félix G Harvey, Louis-Simon Ménard, Florent Bouchet, and Mark J Coates. 2024. Motion In-Betweening via Deep Delta-Interpolator. *IEEE Transactions on Visualization and Computer Graphics* 30, 8 (2024), 5693–5704.
- Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. 2018. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–14.
- Mathis Petrovich, Michael J Black, and Gül Varol. 2023. Tmr: Text-to-motion retrieval using contrastive 3d human motion synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9488–9497.
- Ekkasit Pinyoanuntapong, Muhammad Usama Saleem, Korrawe Karunratanakul, Pu Wang, Hongfei Xue, Chen Chen, Chuan Guo, Junli Cao, Jian Ren, and Sergey Tulyakov. 2024a. Controlmm: Controllable masked motion generation. *arXiv preprint arXiv:2410.10780* (2024).
- Ekkasit Pinyoanuntapong, Pu Wang, Minwoo Lee, and Chen Chen. 2024b. Mmm: Generative masked motion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1546–1555.
- Jia Qin, Youyi Zheng, and Kun Zhou. 2022. Motion In-Betweening via Two-Stage Transformers. *ACM Trans. Graph.* 41, 6 (2022), 184–1.
- Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. 2019. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems* 32 (2019).
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- Alla Safonova and Jessica K Hodgins. 2007. Construction and optimal search of interpolated motion graphs. In *ACM SIGGRAPH 2007 papers*. 106–es.
- Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermano. 2023. Human motion diffusion as a generative prior. *arXiv preprint arXiv:2303.01418* (2023).
- Yi Shi, Jingbo Wang, Xuekun Jiang, Bingkun Lin, Bo Dai, and Xue Bin Peng. 2024. Interactive character control with auto-regressive motion diffusion models. *ACM Transactions on Graphics (TOG)* 43, 4 (2024), 1–14.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502* (2020).
- Paul Starke, Sebastian Starke, Taku Komura, and Frank Steinicke. 2023. Motion in-betweening with phase manifolds. *Proceedings of the ACM on Computer Graphics and Interactive Techniques* 6, 3 (2023), 1–17.
- Sebastian Starke, Paul Starke, Nicky He, Taku Komura, and Yuting Ye. 2024. Categorical codebook matching for embodied character controllers. *ACM Transactions on Graphics (TOG)* 43, 4 (2024), 1–14.
- Sebastian Starke, He Zhang, Taku Komura, and Jun Saito. 2019. Neural state machine for character-scene interactions. *ACM Transactions on Graphics (TOG)* 38, 6 (2019), 1–14.
- Sebastian Starke, Yiwei Zhao, Taku Komura, and Kazi Zaman. 2020. Local motion phases for learning multi-contact character movements. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 54–1.
- Sebastian Starke, Yiwei Zhao, Fabio Zinno, and Taku Komura. 2021. Neural animation layering for synthesizing martial arts movements. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–16.
- Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. 2024. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525* (2024).
- Tencent Hunyuan 3D Digital Human Team. 2025. HY-Motion 1.0: Scaling Flow Matching Models for Text-To-Motion Generation. *arXiv preprint arXiv:2512.23464* (2025).
- Guy Tevet, Sigal Raab, Setareh Cohan, Daniele Reda, Zhengyi Luo, Xue Bin Peng, Amit Haim Bermano, and Michiel van de Panne. 2025. CLoSD: Closing the Loop between Simulation and Diffusion for multi-task character control. In *The Thirteenth International Conference on Learning Representations*.
- Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano. 2022. Human motion diffusion model. *arXiv preprint arXiv:2209.14916* (2022).
- Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. 2024. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *Advances in neural information processing systems* 37 (2024), 84839–84865.
- Unitree. 2025. Unitree Boxing. <https://www.unitree.com/boxing>. Accessed: 2024-06-30.
- Unity. 2025. Unity User Manual (6.2). <https://docs.unity3d.com/Manual/class-BlendTree.html>.
- Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural discrete representation learning. *Advances in neural information processing systems* 30 (2017).
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- Tingwu Wang, Yunrong Guo, Maria Shugrina, and Sanja Fidler. 2020. UniCon: Universal Neural Controller For Physics-based Character Motion. *arXiv preprint arXiv:2011.15119* (2020).
- Boran Wen, Ye Lu, Keyan Wan, Sirui Wang, Jiahong Zhou, Junxuan Liang, Xinpeng Liu, Bang Xiao, Dingbang Huang, Ruiyang Liu, et al. 2025. Efficient and Scalable Monocular Human-Object Interaction Motion Reconstruction. *arXiv preprint arXiv:2512.00960* (2025).
- Lujie Yang, Xiaoyu Huang, Zhen Wu, Angjoo Kanazawa, Pieter Abbeel, Carmelo Sferazza, C Karen Liu, Rocky Duan, and Guanya Shi. 2025. Omniretarget: Interaction-preserving data generation for humanoid whole-body loco-manipulation and scene interaction. *arXiv preprint arXiv:2509.26633* (2025).
- Gwonjin Yi and Junghoon Jee. 2019. Search Space Reduction In Motion Matching by Trajectory Clustering. In *SIGGRAPH Asia 2019 Posters*. 1–2.
- Kangning Yin, Weishuai Zeng, Ke Fan, Minyue Dai, Zirui Wang, Qiang Zhang, Zheng Tian, Jingbo Wang, Jiangmiao Pang, and Weinan Zhang. 2025. UniTracker: Learning Universal Whole-Body Motion Tracker for Humanoid Robots. *arXiv:2507.07356 [cs.RO]* <https://arxiv.org/abs/2507.07356>
- Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama, Han Zhang, Huiwen Chang, Alexander G Hauptmann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, et al. 2023a. Magvit: Masked generative video transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10459–10469.
- Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. 2023b. Language Model Beats Diffusion—Tokenizer is Key to Visual Generation. *arXiv preprint arXiv:2310.05737* (2023).
- Yanjie Ze, João Pedro Araújo, Jiajun Wu, and C. Karen Liu. 2025. *GMR: General Motion Retargeting*. <https://github.com/YanjieZe/GMR> GitHub repository.
- Weishuai Zeng, Shunlin Lu, Kangning Yin, Xiaojie Niu, Minyue Dai, Jingbo Wang, and Jiangmiao Pang. 2025. Behavior Foundation Model for Humanoid Robots. *arXiv preprint arXiv:2509.13780* (2025).
- He Zhang, Sebastian Starke, Taku Komura, and Jun Saito. 2018. Mode-adaptive neural networks for quadruped motion control. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–11.
- Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. 2024a. Motiondiffuse: Text-driven human motion generation with diffusion model. *IEEE transactions on pattern analysis and machine intelligence* 46, 6 (2024), 4115–4128.
- Yaqi Zhang, Di Huang, Bin Liu, Shixiang Tang, Yan Lu, Lu Chen, Lei Bai, Qi Chu, Nenghai Yu, and Wanli Ouyang. 2024b. Motiongpt: Finetuned llms are general-purpose motion generators. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 7368–7376.
- Zhikai Zhang, Jun Guo, Chao Chen, Jilong Wang, Chenghui Lin, Yunrui Lian, Han Xue, Zhenrong Wang, Maoqi Liu, Huaping Liu, et al. 2025. Track Any Motions under Any Disturbances. *arXiv preprint arXiv:2509.13833* (2025).
- Wenyang Zhou, Zhiyang Dou, Zeyu Cao, Zhouyingcheng Liao, Jingbo Wang, Wenjia Wang, Yuan Liu, Taku Komura, Wenping Wang, and Lingjie Liu. 2024. Emdm: Efficient motion diffusion model for fast and high-quality motion generation. In *European Conference on Computer Vision*. Springer, 18–38.
- Shaoting Zhu, Ziwen Zhuang, Mengjie Zhao, Kun-Ying Lee, and Hang Zhao. 2026. Hiking in the Wild: A Scalable Perceptive Parkour Framework for Humanoids. *arXiv preprint arXiv:2601.07718* (2026).

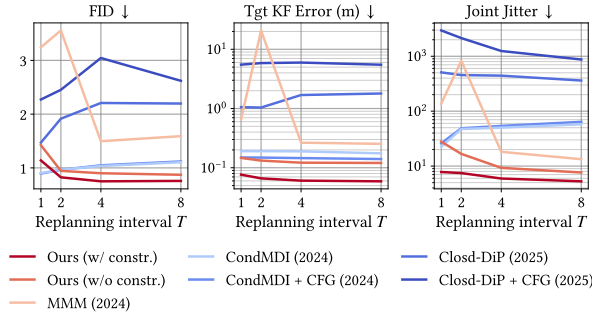


Fig. 15. Evaluation metrics under different replanning frequencies. Left: FID comparison between different replanning frequencies. Middle: Keyframe joint position error comparison between different replanning frequencies. Right: Joint jittering comparison between different replanning frequencies.

### A Scaling Behavior with Number of GPUs

Our main experiments utilize a large number of GPUs, which may not be accessible to individual researchers or smaller studios. To evaluate the impact of compute budget, we train with 1, 2, 4, 8, 16, 32, and 64 GPUs and report results in Figure 16. As shown in the left plot, throughput scales near-linearly with GPU count, directly reducing training time. In general, more GPUs lead to faster convergence and modestly better asymptotic performance, though occasional exceptions arise from random seed variations or distributional differences in the generative model training stage. The performance gap primarily manifests in tokenizer quality; however, these numerical differences translate to only minor perceptual changes in the generated motion. We conclude that training with fewer GPUs remains viable without significant degradation in visual quality.

### B Comparison between FSQ and VQ-VAE

Our tokenizer and pose module are agnostic to the choice of internal tokenization strategy. For instance, FSQ [Mentzer et al. 2023] has been adopted in recent work such as SONIC [Luo et al. 2025] and is also applicable to our framework. To draw a parallel between FSQ and VQ-VAE in our framework, we treat the latent dimensions in FSQ as the number of heads, and the levels per dimension as the codebook size per head. Similarly, for the pose module, we ask it to predict the corresponding token indices from the FSQ tokenizer.

As shown in Figure 17, FSQ and VQ-VAE yield similar performance in terms of FID, reconstruction loss, and cross-entropy loss given the same number of tokens and heads. However, VQ-VAE consistently achieves lower cross-entropy loss and better FID than FSQ, suggesting that VQ-VAE produces a latent space that is slightly easier for the generative model to model. We observe similar improvement in visual quality and robotics deployment. Based on these results, we use VQ-VAE as the default tokenizer for MotionBricks.

### C Replanning Frequency

In MotionBricks, replanning is governed by the drop-frame attribute  $\tau$  introduced in Section 6.2: a replan is triggered when fewer than  $\tau$  frames remain in the motion buffer. During deployment, replanning

also occurs whenever control commands change to allow for instant control reaction, as described in Section 7.4.

To quantify the impact of replanning frequency, we evaluate MotionBricks with replanning intervals ranging from every 1 frame to every 8 frames. As shown in Figure 15, discrete latent methods such as MotionBricks and MMM exhibit slightly worse FID, keyframe error, and joint jitter at higher replanning frequencies, whereas continuous diffusion methods show improved FID and joint jitter with more frequent replanning. We hypothesize that discrete latent representations struggle to capture subtle inter-step variations, and overly frequent replanning traps generation in its early phase, delaying the onset of the desired style and degrading fine details. We observe consistent trends in our visual demos and refer readers to the supplementary videos for further illustration. Based on these findings, we recommend a default replanning interval of 3 to 9 frames, combined with instant replanning upon command changes, as a practical engineering choice.

### D 2D Blendspace and Keyframe Blending in MotionBricks

A 2D blendspace is a standard technique in game engines such as Unreal Engine [Epic Games 2025] that blends between pre-recorded animation clips based on two continuous input axes, typically speed and direction. The neighboring clips are interpolated at runtime to produce smooth transitions for the given input parameters of speed and direction. While effective, 2D blendspaces rely heavily on heuristics and operate at the clip level, requiring careful authoring of sample coverage to avoid blending artifacts.

In contrast, MotionBricks performs keyframe-level blending: rather than interpolating between entire clips, we directly interpolate individual keyframes in joint space from reference poses and feed them to the neural backbone. Artifacts such as foot sliding or floating that would typically arise from naive interpolation are handled implicitly by MotionBricks. Direct keyframe interpolation in MotionBricks still produces natural, artifact-free motions for two reasons: our model is robust to imperfect keyframes, and the drop-frame replanning mechanism ensures that the blended keyframes are never rigidly enforced during generation. This frees us from the constraints of traditional 2D blendspaces, offering finer-grained control without the authoring overhead of populating a full blendspace grid. We refer readers to the supplementary videos for a detailed demonstration.

### E Dataset Statistics

Our full dataset contains 350k motion clips spanning 36 categories, 203 unique activities, 9,285 unique content types, and 163 unique performers. Categories, activities, and content types represent progressively finer levels of motion categorization, as illustrated in Figure 18.

We also release an open-source subset of approximately **140k motion clips**, publicly available as **BONES-SEED** [Bones Studio 2026] via Bones Studio, whose statistics are shown in Figure 19. This subset retains roughly 40% of motions, about half of the categories (17/36), 75% of activities (152/203), 41% of content types (3,829/9,285), and all 163 performers. We split this subset (143,792 motions) into three partitions: a training set (129,661 motions, 90.2%), a *Test-Unseen* set (7,166 motions, 5.0%), and a *Test-Heldout* set (6,965 motions,

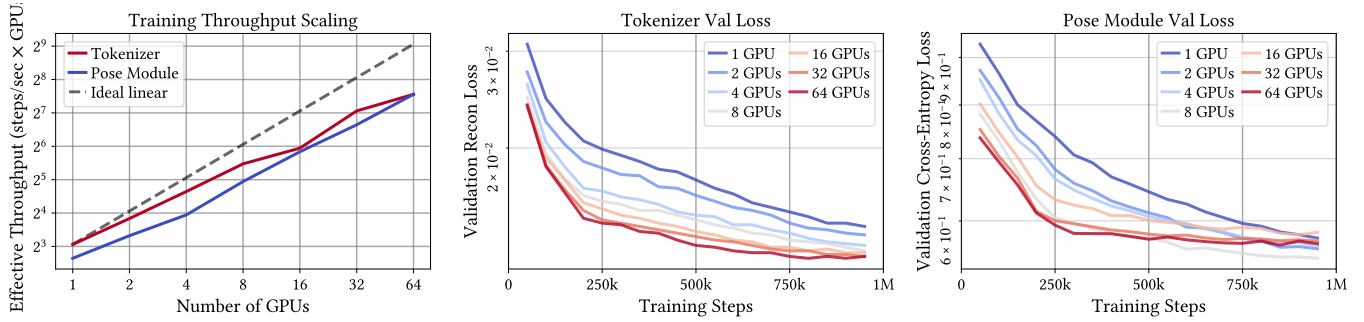


Fig. 16. Ablation study on GPU scaling during training. Left: ideal vs. achieved throughput scaling with the number of GPUs. Middle: tokenizer reconstruction loss on the validation set. Right: cross-entropy loss for token prediction on the validation set.

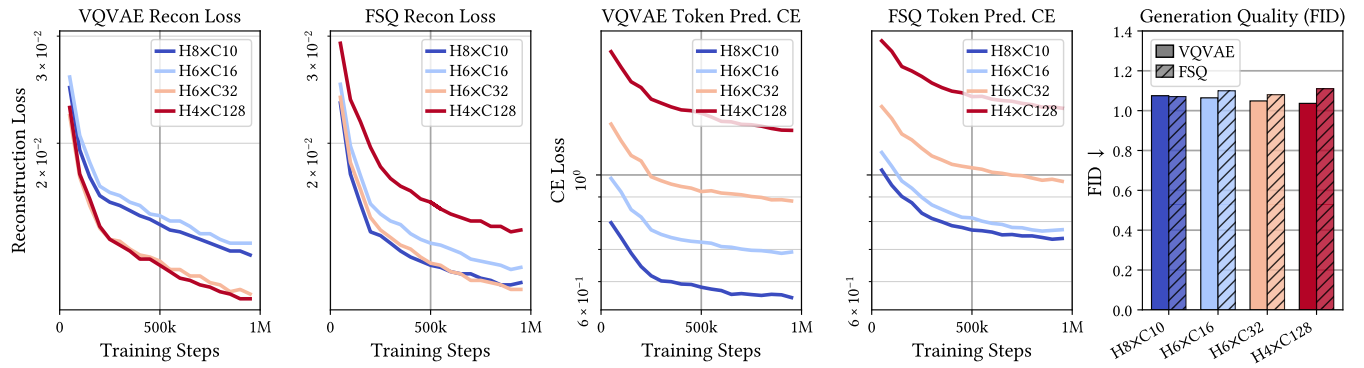


Fig. 17. Comparison between FSQ and VQ-VAE with matched token capacity. Left two plots: tokenizer reconstruction loss on the validation set during training. Middle two plots: cross-entropy loss for token prediction on the validation set. Right: FID comparison on the test set. “H” denotes the number of heads; “C” denotes the codebook size per head.

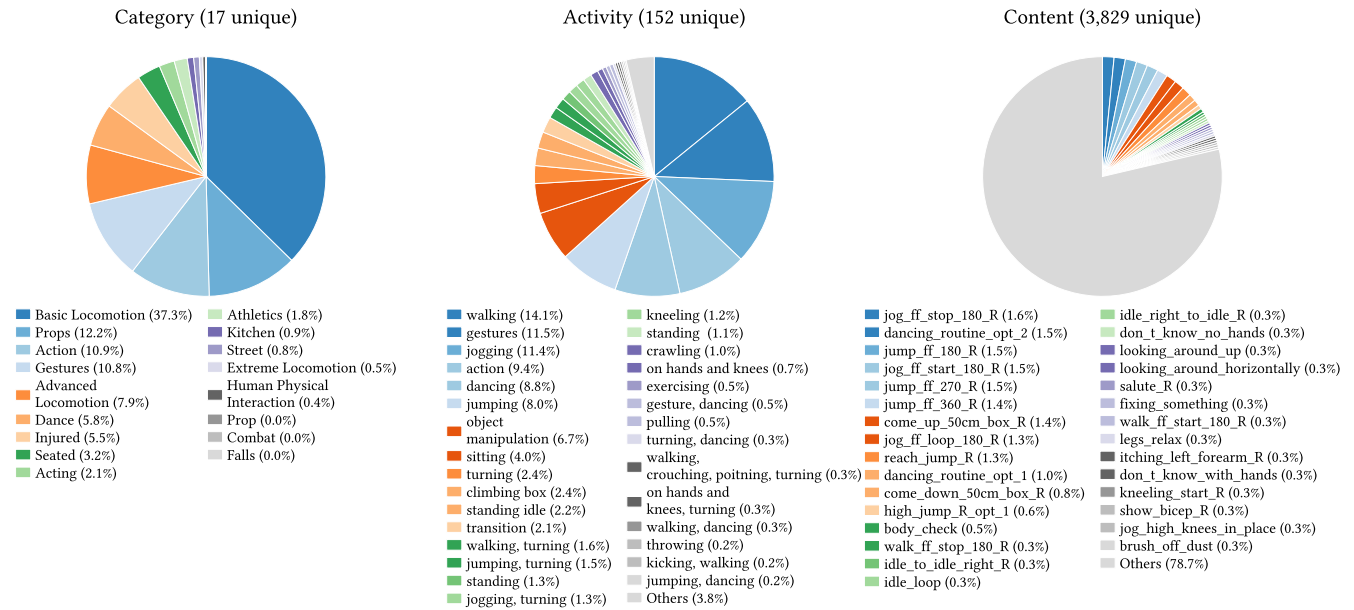


Fig. 18. Overall statistics and diversity of the 140k BONES-SEED open-source subset dataset by categories, activities and contents.

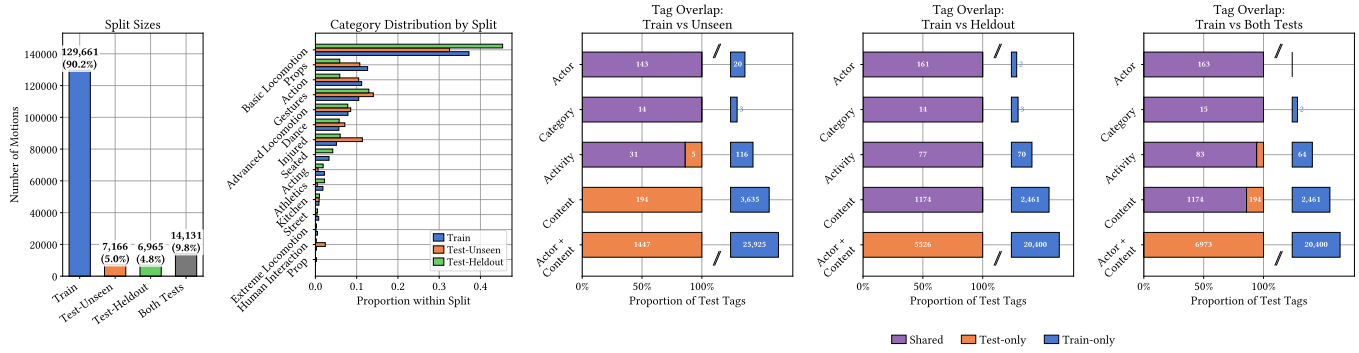


Fig. 19. Overview of the two test sets split from the 140k BONES-SEED open-source subset dataset.

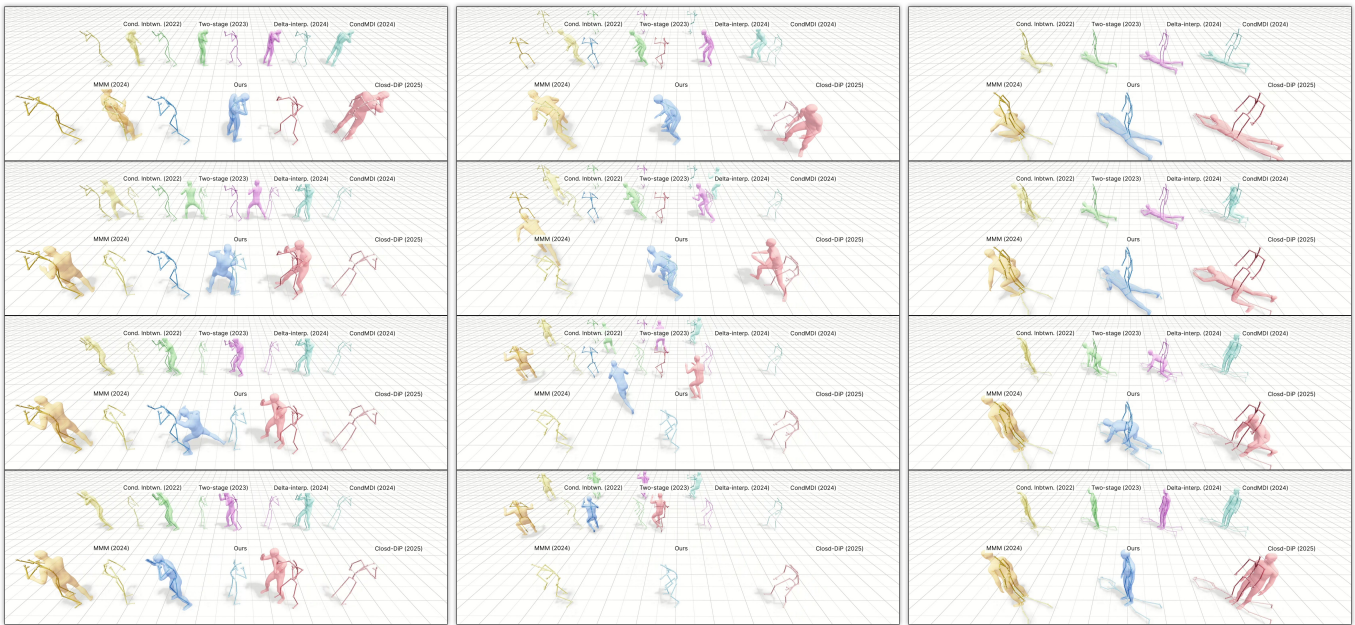


Fig. 20. Additional in-betweening visual comparisons across three test cases. Each column shows a different test case; time progresses from top to bottom. Algorithm names are labeled in each subfigure. MotionBricks is shown as the blue character.

4.8%). The two test sets are designed to evaluate different aspects of generalization. The Test-Unseen set has 0% content-type overlap with the training set: all 194 content types are novel, and 5 out of 36 activities are unseen during training. The Test-Heldout set shares 100% overlap at the category, activity, content, and actor levels, but contains novel actor-content combinations.

As shown in Figure 18, locomotion is by far the most represented category, with diverse variations in speed, direction, and style that enable robust and natural locomotion generation.

The least covered categories are object manipulation and terrain interaction. For manipulation, although many motions have been captured, the space of possible manipulation tasks is vast, and the dataset currently lacks finger and object motion data. Additionally, no samples include continuous terrain information, limiting the model’s ability to generate realistic motions on uneven surfaces.

While the dataset does contain motions with height variations, coverage remains sparse. In the most extreme case for example, vaulting over a 1m obstacle has only one or two captured clips, posing severe challenges for learning and generalization.

## F In-betweening Visual Comparisons

We provide additional isolated in-betweening visual comparisons in Figure 20. The skeletons denote keyframe constraints, and the characters denote snapshots of generations from each algorithm. From left to right, the three test cases are: (1) transitioning from yawning to a boxing pose, (2) transitioning from a stealth pose to running, and (3) transitioning from lying on the ground to a standing pose. MotionBricks consistently outperforms all baselines, which exhibit various artifacts such as implausible poses and poor

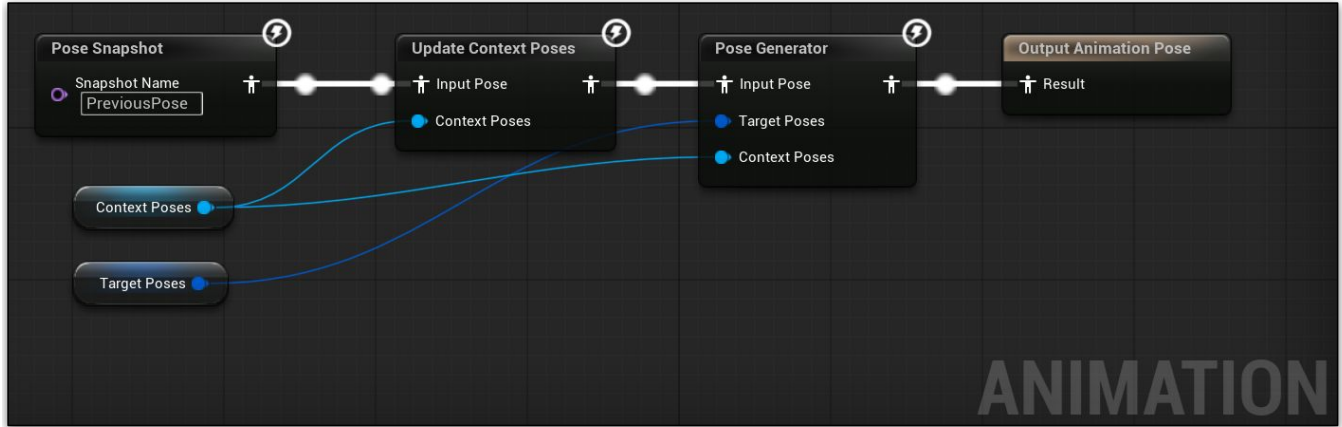


Fig. 21. The complete animation graph used in our UE5 demo. Context poses are updated each frame and passed alongside target keyframes to the pose generator node, which invokes the MotionBricks backbone only when new targets are received.

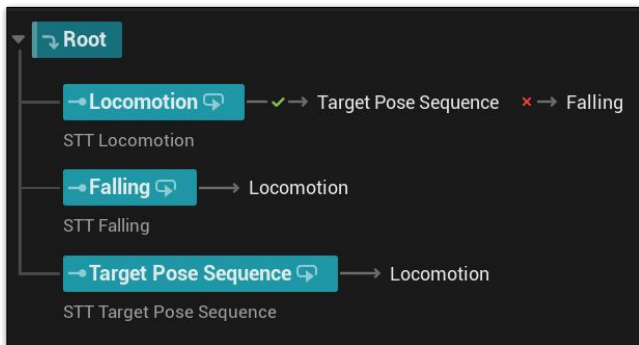


Fig. 22. The UE5 StateTree governing character behavior, managing transitions between locomotion, target pose sequence (object interaction), and falling states.

keyframe adherence. We refer readers to the demo videos for further illustration.

## G UE5 Demo Implementation

We describe the implementation of our UE5 demo, which showcases the full MotionBricks pipeline in a real-time interactive setting.

*Task Actors: A Unified Runtime Primitive.* Both smart locomotion and smart object are realized through a shared custom actor blueprint type that we call a *task actor*. Each task actor encapsulates the runtime logic, assets, and state needed to produce target keyframes for the neural backbone. At authoring time, a task actor samples and stores sparse target poses from user-provided clip assets, defining either the interaction keyframes for a scene object or the style poses for a parameterized locomotion blend-space. At runtime, the active task actor provides the target keyframes that MotionBricks consumes to generate the next motion sequence.

Because task actors leverage UE5 blueprints, each instance can implement its own gameplay logic independently. For example, a



Fig. 23. Runtime interaction detection via box traces in our UE5 demo. Green and red boxes indicate the detection and interaction trigger volumes, respectively. Examples shown include bench sitting, vaulting, object pickup, and platform jumping.

scene interaction task actor orients its keyframe sequence based on the character’s approach angle, using the keyframe anchoring mechanism described in Section 6.2. The key distinction between the two smart primitive types is ownership: the locomotion task actor resides on the character and is always active, while object interaction task actors are placed in the scene as standalone actors. While target keyframes could, in principle, be supplied by any source, including learned models or procedural methods, for our demo we provide them from a handful of selected clips.

*Interaction Detection and State Management.* During locomotion, the system performs box traces to detect nearby interactable objects, as shown in Figure 23. When a smart object is activated, it becomes the new keyframe provider, interrupting the current motion and triggering generation from the object’s target keyframes. This on-demand interruption relies on the deterministic replanning behavior

of MotionBricks: because the backbone produces consistent results regardless of when or how frequently replanning is triggered (see Appendix C), gameplay actions can seamlessly override the current motion without introducing artifacts. This property is essential to the smart primitive framework, as it enables responsive, interruptible character behavior driven entirely by game events.

As shown in Figure 22, the character’s high-level behavior is governed by a UE5 StateTree with three states: *locomotion*, where the locomotion task actor continuously supplies keyframes based on velocity, orientation, and replanning frequency; *target pose sequence*, where an object interaction task actor provides keyframes to execute the interaction; and *falling*, which dynamically spawns a task actor at runtime. For falls, the system measures the current height and runs a projectile simulation to determine the landing position, placing the appropriate landing keyframe based on the fall height. Upon completion, control returns to the locomotion state.

*Animation Graph and Inference Scheduling.* The animation graph for the entire demo is deliberately minimal, as shown in Figure 21. Each frame, the context pose buffer is updated with the previous frame’s pose, and both context and target poses are passed to a *pose generator* node. This node invokes the MotionBricks backbone

only when new target keyframes are received; otherwise, it samples directly from the previously generated motion buffer, so model inference does not occur every frame. Communication between gameplay and animation is bidirectional: gameplay can interrupt the current motion at any time by supplying new keyframes, while the pose generator notifies gameplay when the buffer is nearing its end and a new generation is needed.

## H Contributors

**Project Lead:** Tingwu Wang.

**Research:** Tingwu Wang, Olivier Dionne, Michael De Ruyter, David Minor, Davis Rempe, Kaifeng Zhao, Mathis Petrovich.

**Animation development:** Olivier Dionne, Tingwu Wang, Michael De Ruyter, David Minor, Brian Robison, Xavier Blackwell.

**Robotics development:** Tingwu Wang, Chenran Li, Zhengyi Luo, Ye Yuan.

**Game art, assets and level design:** Brian Robison, Xavier Blackwell, Bernardo Antoniazzi.

**Advising:** Simon Yuen, Xue Bin Peng, Yuke Zhu.

Received 22 January 2026; revised 15 April 2026; accepted 20 May 2026