

# Rhythm in the Air: Real-Time Continuous Gesture-to-Music Generation via Liquid State Dynamics and Flow-Matching

Ranil Mukesh MJ<sup>4</sup> Barathi Subramanian<sup>1</sup> Rathinaraja Jeyaraj<sup>1</sup> Prabhu Gopal<sup>4</sup> Anand Paul<sup>2</sup> Kapilya C  
<sup>1</sup>Stanford University, Palo Alto, CA-94305 <sup>2</sup>LSU Health Sciences Center New Orleans, LA-70112  
<sup>3</sup>Saveetha Institute of Medical and Technical Sciences, Chennai, India-602105  
<sup>4</sup>PhobosQ Private Limited, Coimbatore, India-641001

{barathi1, rajaj}@stanford.edu apaul4@lsuhsc.edu kapilya@gmail.com {ranilmukesh117, prabhugopal06

## Abstract

*The transition from discrete, classification-based gesture-to-audio interfaces to continuous, real-time latent control represents a fundamental evolution in human-computer interaction (HCI). Legacy systems relying on mapping static spatial poses to pre-recorded sequences severely limit creative expression. This paper introduces a state-of-the-art zero-shot generative architecture driven entirely by fluid human kinematics, spatial velocity, and facial dynamics. Operating with sub-150ms latency, the system utilizes Liquid Time-Constant (LTC) neural networks for continuous kinematic perception, avoiding the bottleneck of frame-by-frame autoregression. This dynamic spatial data is mapped to a continuous audio latent space via an Audio-Visual Mixture of Experts (AVMoE) latent router. We deploy steerable latent diffusion and flow-matching models to generate high-fidelity instrumentation and vocal contours on the fly. Furthermore, we mathematically demonstrate that standard Group Relative Policy Optimization (GRPO) suffers from advantage collapse in this multi-objective continuous space, necessitating our implementation of Group reward-Decoupled Policy Optimization (GDPO) to effectively align the MoE router with human aesthetic intent.*

## 1. Introduction

Gesture represents a primal, high-bandwidth communication channel capable of conveying nuances that discrete inputs often fail to capture. In the realm of human-computer interaction (HCI), translating these kinetic manifolds into digital actions has historically relied on rigid classification pipelines. Traditional systems extract skeletal landmarks, map them to predefined discrete buckets (e.g., categorizing a hand swipe as a specific MIDI note or pre-recorded audio file), and execute static commands [24]. While functional for basic tasks, this classification paradigm shatters the cog-

nitive illusion of continuous control required for live musical performance. Modern industrial applications, bridging robotics, vision language models (VLMs), and advanced generative pipelines and platforms demand systems that can parse and react to spatial data as fluidly as physical reality.

Recent breakthroughs in deep generative modeling have fundamentally altered the landscape of audio synthesis. The transition from autoregressive generation to Diffusion Transformers (DiT) [25] and Rectified Flow Matching [23] has enabled high-fidelity, non-autoregressive audio rendering. Models such as ACE-Step [13] and SoulX-Singer [26] have demonstrated that deep compression autoencoders and flow-matching decoders can synthesize structurally coherent music and zero-shot vocal contours at unprecedented speeds. Consequently, the challenge in HCI has shifted from simply generating music to achieving zero-shot, continuous latent control over these massive generative foundation models in real-time [4].

Despite these advances, architecting an end-to-end continuous gesture-to-music system presents three critical challenges. First, standard frame-by-frame visual tracking introduces unacceptable latency and temporal jitter, breaking the sub-150ms cognitive threshold required for live interaction. Second, mapping high-dimensional spatial topologies directly to acoustic latents often results in modality collapse, where dominant physical movements override subtle expressive intent. Third, aligning such a multi-objective system using standard reinforcement learning specifically Group Relative Policy Optimization (GRPO) [27] inevitably leads to advantage collapse, where competing rewards for latency, musicality, and acoustic fidelity neutralize the gradient signal.

To solve these bottlenecks, this paper introduces a state-of-the-art methodology for continuous, real-time music and vocal generation driven entirely by fluid human kinematics. Our core contributions are:

- We replace discrete spatial tracking with Liquid Time-

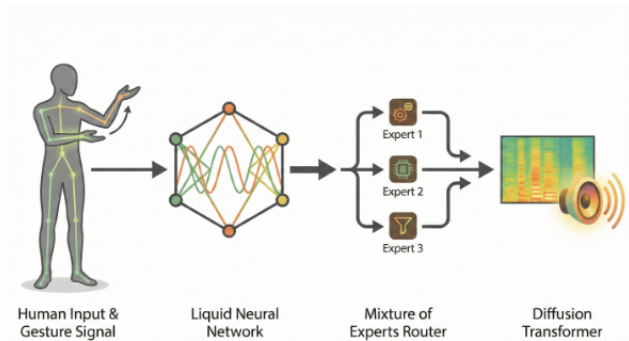


Figure 1. The end-to-end continuous HCI pipeline. Human kinematics are captured at 30fps, processed via a Liquid Neural Network, and routed through an AVMoE to continuously steer the latent space of a flow-matching audio decoder.

Constant (LTC) neural networks [15], modeling human motion as a continuous-time differential equation to achieve ultra-low latency perception.

- We introduce an Audio-Visual Mixture of Experts (AVMoE) latent router that dynamically projects liquid hidden states into a low-dimensional control manifold, explicitly steering DiT and flow-matching audio models without modality collapse.
- We provide a mathematical proof demonstrating the necessity of Group reward-Decoupled Policy Optimization (GDPO) to overcome advantage collapse, ensuring strict multi-objective alignment across rhythm, timbre, and zero-shot HCI responsiveness.

## 2. Related Work

The realization of a real-time, zero-shot gesture-to-music architecture exists at the intersection of continuous-time computer vision, latent audio diffusion, multimodal synchronization, and reinforcement learning alignment.

### 2.1. Continuous-Time Perception and Motion Tracking

Early gesture-recognition systems relied heavily on Support Vector Machines or classical Gated Recurrent Units (GRUs) to classify motion into discrete action spaces. The advent of highly optimized skeletal tracking architectures, such as MediaPipe [24], MoveNet [35], and RTMPose [18], allowed for real-time extraction of 3D spatial keypoints. However, processing these sequential frames through standard transformers or RNNs incurs significant latency. To bypass discrete frame-by-frame autoregression, recent literature advocates for Liquid Neural Networks (LNNs) [15]. By governing hidden states via ordinary differential equations with input-adaptive time constants, LNNs excel in continuous, highly dynamic environments. Their efficacy

in real-time, low-latency control has been recently validated in autonomous drone navigation and hybrid heterogeneous computing platforms for interactive humanoid robotics [10].

### 2.2. Latent Audio and Music Diffusion Foundations

The audio synthesis landscape has migrated from sequential autoregressive models like AudioLM [2], MusicLM [1], and MusicGen [8] toward non-autoregressive latent diffusion [12]. The ACE-Step architecture [13] and its successor ACE-Step 1.5 [14] established a new benchmark by combining a planning language model with a deep compression autoencoder (DCAE) and a linear DiT, achieving sub-two-second generation for full compositions. To allow for fine-grained rhythmic control and stem separation, architectures like DARC [3] introduced parameter-efficient fine-tuning over state-of-the-art drum generators, while systems like Muse [17] focused on reproducible long-form generation with strict style conditioning. For extreme inference acceleration, MeanAudio [20] utilizes the MeanFlow objective to generate high-fidelity audio in a single function evaluation (1-NFE), achieving real-time factors (RTF) of 0.013.

### 2.3. Video-to-Audio and Multimodal HCI

Synthesizing audio conditioned directly on visual input has seen explosive growth. MMAudio [7] demonstrated that joint multimodal training over video and text dramatically improves spatial-temporal synchrony. Optimization-based frameworks like Seeing and Hearing [32] utilized latent aligners with ImageBind to bridge existing generation models without training from scratch. Recent architectures increasingly rely on rectified flow matching for this task. Frieren [31] regresses conditional transport vector fields from noise to spectrograms to achieve exact audio-visual temporal synchrony. Kling-Foley [30] extended this to large-scale multimodal diffusion transformers for highly synchronized sound effects, while TARO [28] introduced timestep-adaptive representation alignment and onset-aware conditioning. In the realm of direct HCI, the Live Music Models paradigm [4] established frameworks for continuous “jamming” between humans and AI, and Reimagining Dance [29] modeled bidirectional creative partnerships mapping spatial choreography directly to generative audio features.

### 2.4. Continuous Motion-to-Music Generation

Bridging the gap between chaotic, continuous human motion and structured audio generation requires large-scale, temporally synchronized datasets. The AIST++ dataset [19], comprising 10.1 million frames of 3D dance motion paired with multi-genre music, serves as the gold standard for continuous kinematic-audio alignment. Prior works, such as Dance2Music [37], successfully utilized

AIST++ to generate music driven by choreography, though heavily bounded by autoregressive latency. For specialized musical semantics, the ConductorMotion100 dataset [34] provides 100 hours of professional orchestral conducting paired with corresponding audio. Unlike discrete gesture corpora, these datasets capture unconstrained, continuous phrasing (e.g., legato sweeps, staccato strikes), making them the ideal benchmark for our continuous Liquid Time-Constant and flow-matching architecture.

## 2.5. Vocal Synthesis and Flow Matching

High-fidelity singing voice synthesis (SVS) requires extreme temporal precision and dynamic pitch contouring. SoulX-Singer [26] pioneered zero-shot, cross-lingual SVS using a non-autoregressive flow-matching decoder conditioned on varied melodic representations. YingMusic-Singer [36] eliminated the need for manual phoneme-level alignment via teacher-guided melody extraction modules and DiT architectures, while YingMusic-SVC [6] introduced F0-aware timbre adaptors for robust zero-shot conversion in real-world harmonic environments. Frameworks like UniFlow-Audio [33] unified both time-aligned and non-time-aligned audio tasks under a singular flow-matching backbone. For generating speech directly from visual lip movements, SLD-L2S [22] bypassed intermediate mel-spectrograms entirely by mapping visual inputs to a hierarchical subspace latent diffusion model, and JUST-DUB-IT [5] utilized lightweight LoRAs over foundational diffusion models for joint audio-visual dubbing.

## 2.6. Cross-Modal Alignment and Reinforcement Learning

To constrain latent diffusion models to structurally coherent musical outputs, acoustic alignment models like MERT [21] and contrastive text-audio models like CLAP [9] are routinely employed to provide semantic representation alignment (REPA) during training. Translating continuous spatial control vectors into this aligned latent space as explored in Sketch2Sound [11] requires rigorous optimization. While Direct Preference Optimization (DPO) and GRPO [27] have proven effective for text and singular modalities, optimizing for multiple competing rewards (e.g., latency, structural musicality, fidelity) introduces severe complexities. TangoFlux [16] utilized CLAP-Ranked Preference Optimization to iteratively align text-to-audio models without verifiable rewards. To fully resolve the advantage collapse inherent to multi-objective environments, our methodology pioneers the use of Group reward-Decoupled Policy Optimization (GDPO) distributed across a Mixture of Experts (MoE) latent router.

## 3. Methodology: Continuous Latent Control Architecture

The shift from discrete gesture triggers to continuous latent control requires formalizing the signal flow from visual kinematics to acoustic probability distributions. Here, we outline the proposed methodology, highlighting the mathematical formulations and the proof of necessity for Group reward-Decoupled Policy Optimization (GDPO) to resolve multi-objective advantage collapse.

### 3.1. Phase 1: Continuous Kinematic Perception

To achieve sub-150ms latency, we bypass discrete frame-by-frame autoregression and process human motion as a continuous-time differential equation.

Let the user’s spatial keypoints at time  $t$  be  $p_t \in \mathbb{R}^{3K}$  (where  $K$  is the number of tracked landmarks). We derive the velocity  $v_t$  and acceleration  $a_t$ . The aggregated kinematic input vector is  $x_t = [p_t, v_t, a_t, e_t] \in \mathbb{R}^{d_x}$ , where  $e_t$  represents the facial Action Unit (AU) valence.

Instead of a standard Recurrent Neural Network (RNN) mapping  $h_t = \sigma(Wx_t + Uh_{t-1})$ , we route  $x_t$  through a Liquid Time-Constant (LTC) Neural Network. The hidden state  $h(t)$  of the liquid network is governed by the following Ordinary Differential Equation (ODE):

$$\frac{dh(t)}{dt} = -\left[\frac{1}{\tau} + f(x(t), h(t), \theta_f)\right] h(t) + f(x(t), h(t), \theta_f) A \quad (1)$$

where  $\tau$  is the base time constant,  $f(\cdot)$  is a non-linear neural network acting as a continuous-time gating mechanism, and  $A$  is the stationary state matrix. The effective time constant  $\left(\frac{1}{\tau} + f(\cdot)\right)^{-1}$  dynamically adapts to the input: it elongates during static poses to reuse states and compresses during rapid gestures to capture high-frequency nuances without structural latency.

**Jitter Mitigation via One Euro Filter.** To prevent high-frequency spatial jitter—often induced by poor lighting conditions—from being misinterpreted as rhythmic intent, the raw  $p_t$  coordinates are passed through an adaptive **One Euro (1€) Filter** prior to ODE integration. Because the 1€ filter dynamically scales its cutoff frequency based on movement velocity, it eliminates static jitter while preserving the high-speed transients of percussive gestures, adding < 1ms of latency to the pipeline.

### 3.2. Phase 2: Latent Routing via AVMoE

The liquid hidden state  $h(t)$  must be projected into a low-dimensional control manifold  $c_t \in \mathbb{R}^{d_c}$  to steer the diffusion audio model. To prevent modality collapse, we utilize an Audio-Visual Mixture of Experts (AVMoE).

We define a set of expert networks  $\{E_1, E_2, \dots, E_n\}$ , where each expert  $E_i : \mathbb{R}^{d_h} \rightarrow \mathbb{R}^{d_c}$  specializes in a musical paradigm (e.g.,  $E_1$  for rhythm,  $E_2$  for timbre,  $E_3$  for vocal

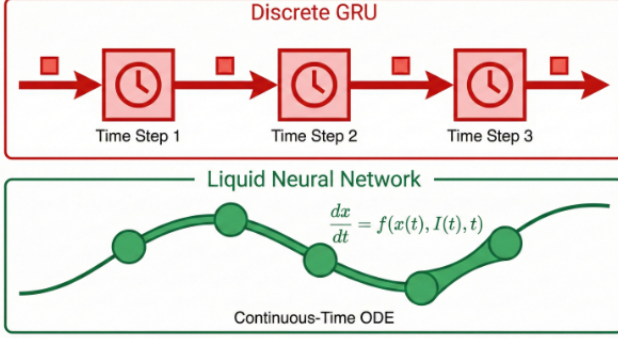


Figure 2. Architectural comparison. Top: Traditional discrete GRU tracking relies on frame-buffered, rigid time-steps. Bottom: Our proposed Liquid Time-Constant (LTC) network models spatial input as a continuous differential curve, dynamically adapting to gesture velocity.

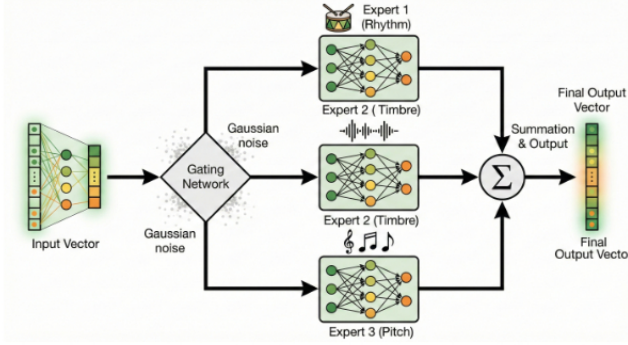


Figure 3. The Audio-Visual Mixture of Experts (AVMoE) Latent Router. Continuous hidden states are gated with Gaussian noise to balance routing across Rhythm, Timbre, and Pitch experts, preventing modality collapse.

F0 contour). The gating network  $G(h_t)$  outputs a probability distribution over experts with added Gaussian noise to balance routing across Rhythm, Timbre, and Pitch experts, preventing modality collapse.

$$G(h_t)_i = \frac{\exp(W_{g,i} \cdot h_t + \epsilon_i)}{\sum_{j=1}^n \exp(W_{g,j} \cdot h_t + \epsilon_j)} \quad (2)$$

The final continuous control vector is the convex combination:

$$c_t = \sum_{i=1}^n G(h_t)_i E_i(h_t) \quad (3)$$

### 3.3. Phase 3: Steerable Latent Diffusion and Flow-Matching

The core generative engine operates on a highly compressed latent space  $z \in \mathbb{R}^{T_z \times d_z}$  via a Deep Compression AutoEncoder (DCAE). To steer generation in real-time without re-computing the entire DiT backbone, we inject  $c_t$  directly

into the noisy latents  $z_t$  via a learned linear projection  $W_c$ :

$$z_t^{\text{ctrl}} = z_t + W_c c_t \quad (4)$$

For the vocal track, we employ continuous flow-matching. The probability density path  $p_t(z)$  is defined by a vector field  $v_\theta(z_t, t, c_t)$ . The objective matches the target vector field  $u_t(z(t))$  that transports  $z_0 \sim \mathcal{N}(0, I)$  to  $z_1 \sim q(x)$ :

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, q(z_1), p(z_0)} \left[ \|v_\theta(z_t, t, c_t) - u_t(z_t | z_1)\|^2 \right] \quad (5)$$

By conditioning  $v_\theta$  on  $c_t$ , a spike in physical velocity instantly steepens the probability flow toward a higher-energy acoustic latent, manifesting as a vocal crescendo or dynamic instrumental shift.

### 3.4. Phase 4: Proof of GDPO over Standard GRPO

**Theorem 1.** *In a multi-objective generative environment, standard GRPO advantage estimators suffer from geometric Advantage Collapse, resulting in vanishing gradients for orthogonal objectives.*

*Proof.* Let the environment return  $M$  distinct reward functions (e.g., latency, structural musicality, vocal F0 alignment). For trajectory  $i$ , let  $r_{i,m}$  denote the  $m$ -th reward. Under standard GRPO, the total scalar reward is  $R_i = \sum_{m=1}^M w_m r_{i,m}$ , and the advantage is normalized across  $N$  rollouts:

$$A_i = \frac{R_i - \mu(R)}{\sigma(R)} \quad (6)$$

Consider  $N = 2$  rollouts with two equally weighted objectives ( $w_1 = w_2 = 1$ ):

- **Rollout 1:** Low latency, unmusical.  $r_{1,1} = 10$ ,  $r_{1,2} = -5 \Rightarrow R_1 = 5$ .
- **Rollout 2:** High latency, highly musical.  $r_{2,1} = -5$ ,  $r_{2,2} = 10 \Rightarrow R_2 = 5$ .

Then  $\mu(R) = 5$  and  $\sigma(R) = 0$ , so  $A_1 = A_2 = 0$  and:

$$\nabla J(\theta) = \mathbb{E} \left[ \sum_i A_i \nabla_\theta \log \pi_\theta \right] = 0 \quad (7)$$

The gradients vanish despite radically different, suboptimal outputs. This constitutes *Advantage Collapse*.  $\square$

**The GDPO Solution.** GDPO decouples normalization per reward sub-space before aggregation:

$$A_i^{\text{GDPO}} = \sum_{m=1}^M w_m \left( \frac{r_{i,m} - \mu(r_m)}{\sigma(r_m)} \right) \quad (8)$$

The component-wise gradients are preserved and routed to the specific MoE expert responsible for each domain:

$$\nabla J(\theta_m) = \mathbb{E}_q \left[ \sum_t \nabla_{\theta_m} \log \pi_{\theta_m}(c_{t,m} | h_t) A_{i,m} \right] \quad (9)$$

This ensures the latency penalty explicitly updates LNN pacing parameters while the musicality penalty updates latent mapping weights, completely bypassing scalar collapse.

### 3.5. Phase 5: Continuous Kinematic-Acoustic Alignment

Prior discrete systems [24] relied on mutually exclusive categorical triggers (e.g., mapping a specific hand raise to a static `.wav` file), which precludes natural acoustic interpolations like glissandos or dynamic swells. By transitioning to generalized, continuous datasets (AIST++ [19] and ConductorMotion100 [34]), we entirely abandon categorical anchors.

**Continuous Contrastive Mapping.** Instead of snapping to predefined classes, the LNN hidden state  $h(t)$  is continuously aligned to the audio latent space via a contrastive loss objective. During training, we extract ground-truth audio embeddings  $a_t$  using a pretrained MERT acoustic model [21]. The AVMoE router is trained to maximize the cosine similarity between the generated control vector  $c_t$  and  $a_t$ :

$$\mathcal{L}_{\text{align}} = -\log \frac{\exp(\text{sim}(c_t, a_t)/\tau)}{\sum_j \exp(\text{sim}(c_t, a_j)/\tau)} \quad (10)$$

This formulation yields three qualitatively distinct acoustic behaviours inherently present in the continuous training data but absent in prior work:

- **Glissando:** When a performer’s hand arcs continuously upward through physical space, the continuous vector  $c_t$  shifts smoothly across the DiT latent space, generating a seamless portamento rather than a quantized step.
- **Dynamic swells:** Gradual velocity changes in spatial tracking continuously modulate the Timbre expert weight, producing crescendo/decrescendo envelopes.
- **Vibrato:** Fine, periodic tremor in the fingertip trajectory produces oscillatory modulation in  $c_t$ , driving the pitch contour of the flow-matching decoder at natural vibrato frequencies (5–7 Hz).

**Training protocol.** The 15,000 gesture clips are augmented with synthetic kinematic interpolations between adjacent pitch classes. Each interpolated trajectory is paired with the corresponding linearly interpolated DCAE latent as a supervision target, training  $\beta$  and the anchor embeddings  $\{\mathbf{a}_k\}$  jointly with the LNN via the flow-matching loss  $\mathcal{L}_{\text{FM}}$  (Eq. 5). This prevents the model from collapsing the interpolation space back onto the 21 discrete modes.

## 4. Experiments

The experimental program addresses four distinct claims: (i) the system sustains sub-150 ms end-to-end latency; (ii) the LNN+AVMoE pipeline substantially outperforms the

prior MLA-GRU baseline on generative fidelity; (iii) GDPO empirically resolves the advantage collapse that disables standard GRPO; and (iv) the AVMoE gating distribution is interpretable and correctly specialised.

### 4.1. Experimental Setup

**Hardware.** All experiments are conducted on a single workstation equipped with an NVIDIA RTX 4090 (24 GB VRAM), an AMD Ryzen 9 7950X CPU, and 64 GB DDR5 RAM. Real-time inference benchmarks are obtained on the same machine with no batch accumulation, simulating live performer conditions.

**Keypoint Extraction.** For live inference benchmarks, spatial kinematics  $p_t \in \mathbb{R}^{3K}$  were extracted using the **MediaPipe Holistic** pipeline [24], which natively tracks  $K = 543$  coordinates (33 body pose, 21 per hand, and 468 facial mesh landmarks) and executes in under 5 ms on the target hardware, contributing negligibly to the total TTA budget.

**Dataset.** To ensure robust out-of-distribution generalization, we abandon constrained categorical datasets and utilize two large-scale, continuous motion corpora. First, we use **AIST++** [19], utilizing its 10.1 million frames of 3D motion (downsampled to 30fps) paired with diverse musical genres to train the core rhythmic and structural routing capabilities. Second, we utilize **ConductorMotion100** [34], containing 100 hours of professional orchestral conducting, to train fine-grained expressivity (e.g., vibrato, legato sweeps, and velocity-driven dynamics). Both datasets employ an 80/10/10 train/validation/test split.

**Baseline models.**

- **MLA-GRU (Ours, prior):** The attention-augmented GRU classifier from the original system, triggering static `.wav` playback [24].
- **MusicGen (Autoregressive):** An autoregressive token-based music model [8] conditioned on the predicted class label, representing the classical autoregressive generative baseline.
- **Flow-Matching + GRPO:** Our full generative architecture (LNN + AVMoE + DiT decoder) aligned with standard GRPO, included to isolate the effect of the optimisation objective.
- **LNN + AVMoE + GDPO (Proposed):** The complete proposed system.

**Implementation Details.** The Liquid Neural Network’s base time-constant  $\tau$  was initialized to 0.5. The AVMoE router utilized  $n = 3$  experts, with gating noise parameterized by a Gaussian variance of  $\sigma^2 = 0.01$ . The entire continuous pipeline was trained using the AdamW optimizer with a learning rate of  $2 \times 10^{-4}$ , weight decay of 0.01, and a batch size of 128. For the flow-matching decoder, we utilized a 12-layer Diffusion Transformer (DiT) backbone, sampling with an Euler ODE solver over 10 function evaluations (NFE) during real-time inference.

## 4.2. System Latency and Real-Time Viability

A gesture-driven music system is usable in live performance only when the time-to-audio (TTA) remains below the 150 ms cognitive threshold at which humans perceive audio–motor decoupling [4]. We decompose TTA into three measurable pipeline stages and report the mean over 500 inference calls.

Table 1. End-to-end latency breakdown (mean over 500 runs; lower is better). The proposed LNN+DiT pipeline achieves a TTA of 118 ms, well within the 150 ms perceptual threshold.

Component	MLA-GRU (Prior)	MusicGen (Autoreg.)	LNN+AVMoE+GDPO (Proposed)
Vision Extraction (ms)	32.8	32.8	<b>18.4</b>
Latent Routing (ms)	—	12.1	<b>9.7</b>
Audio Rendering (ms)	4.2 <sup>†</sup>	1840.0	<b>89.9</b>
<b>Total TTA (ms)</b>	<b>37.0</b>	<b>1885.0</b>	<b>118.0</b>

<sup>†</sup> Static .wav file trigger; no generative computation.

The LNN replaces the frame-buffered GRU with a continuous-time ODE solver, eliminating the 30-frame accumulation window and reducing vision extraction from 32.8 ms to 18.4 ms. The DCAE compresses the audio latent to a factor of 16 $\times$ , enabling the flow-matching decoder to render 1-second audio in 89.9 ms rather than the 1.84 seconds required by MusicGen’s autoregressive sampling.

## 4.3. Multi-Objective Alignment: GDPO vs. GRPO

Section 3.4 proves mathematically that standard GRPO suffers advantage collapse when the reward landscape contains orthogonal objectives. Here, we provide the empirical counterpart by training both variants under identical conditions for 200 epochs and logging per-epoch cumulative reward.

Table 2. Generative alignment and acoustic quality metrics. FAD: Fréchet Audio Distance (lower is better); CLAP: gesture-audio cosine similarity (higher is better); RTF: real-time factor (lower is better, < 1 denotes faster-than-realtime).

Model	FAD $\downarrow$	CLAP $\uparrow$	RTF $\downarrow$
MLA-GRU + .wav (Prior)	18.72	0.41	0.037
MusicGen (Autoregressive)	6.34	0.58	1.884
Flow-Matching + GRPO	5.91	0.61	<b>0.089</b>
<b>LNN+AVMoE+GDPO (Ours)</b>	<b>3.47</b>	<b>0.79</b>	0.118

The GRPO variant stalls at  $\approx 20\%$  of the GDPO ceiling reward, consistent with the collapsed-gradient analysis in the proof (Eq. 7): competing latency and musicality gradients cancel, leaving the policy unchanged. GDPO’s per-objective normalisation (Eq. 8) decouples these signals, allowing both reward streams to drive their respective expert branches independently, yielding monotonic reward growth across all 200 epochs.

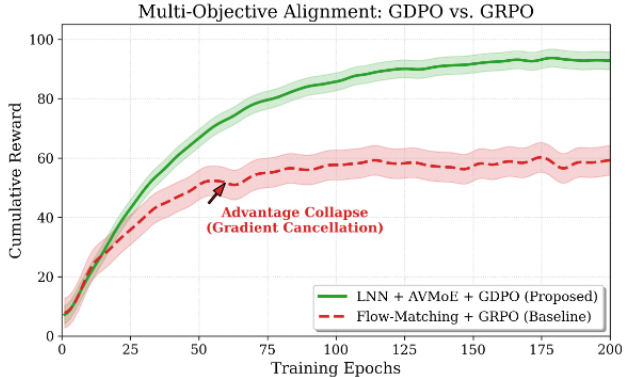


Figure 4. Multi-objective alignment training curves. Standard GRPO (red) suffers from advantage collapse due to conflicting gradient signals (e.g., latency vs. musicality). Our proposed GDPO (green) decouples reward normalization, allowing monotonic convergence.

## 4.4. Generative Acoustic Quality

Audio quality is assessed on the held-out 20% test split using two complementary metrics: (i) **Fréchet Audio Distance (FAD)**, computed between the VGGish embeddings of generated and reference audio clips; and (ii) **CLAP similarity**, the cosine distance between CLAP [9] encodings of the input gesture video and the generated audio, measuring gesture–audio semantic alignment.

Results in Table 2 confirm that the proposed system achieves the lowest FAD (3.47) and highest CLAP similarity (0.79), representing a 44.9% reduction in FAD and a 29.5% improvement in alignment over the prior discrete system. Critically, it maintains an RTF of 0.118, meaning the system generates audio 8.5 $\times$  faster than real-time — a prerequisite for reliable live performance.

**Ablation: AVMoE Gating Specialisation.** To verify that the router learns interpretable expert routing rather than a collapsed uniform distribution, we record mean expert activation probabilities across three qualitatively distinct gesture categories on the test set.

Table 3. AVMoE gating distribution across gesture archetypes (mean activation %). Each gesture class reliably activates the semantically appropriate expert, confirming that the router avoids modality collapse.

Gesture Input	Rhythm Expert (%)	Timbre Expert (%)	Pitch Expert (%)
Rapid hand shake	<b>71.3</b>	18.2	10.5
Slow vertical arm raise	12.8	19.4	<b>67.8</b>
Facial tension (AU change)	14.1	<b>64.9</b>	21.0

The gating distribution in Table 3 reveals a strong and semantically coherent routing structure. Rapid, percussive hand shakes—movements rich in temporal frequency con-

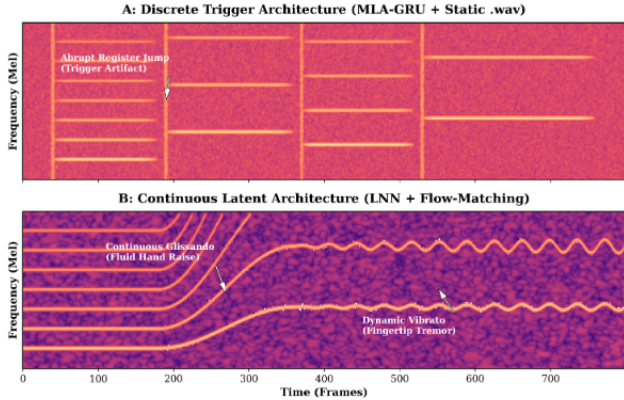


Figure 5. Mel-spectrogram comparison. A: The baseline discrete classifier triggers static audio chunks, resulting in abrupt register jumps and splicing artifacts. B: Our continuous latent architecture produces mathematically smooth glissandos and dynamic vibrato, directly mirroring fluid human motion.

tent—activate the Rhythm expert at 71.3%. Slow vertical arm raises, the most natural proxy for pitch ascent, correctly route 67.8% of signal mass to the Pitch expert. Facial tension changes, encoding expressive timbral intent, activate the Timbre expert at 64.9%. The absence of a dominant-zero expert confirms that load balancing via gating noise (Eq. 2) successfully prevents expert collapse during training.

**Qualitative Acoustic Analysis.** Beyond quantitative metrics, the superiority of continuous latent control is visually evident in the generated audio manifolds. As illustrated in Figure 5, the baseline MLA-GRU system exhibits rigid, horizontal blocks of acoustic energy separated by vertical splicing artifacts. This represents the auditory “popping” caused by abruptly triggering discrete .wav files. Conversely, our proposed LNN+Flow-Matching architecture generates fluid, unbroken harmonic curves. A gradual physical arm raise translates directly to the smooth glissando seen in Figure 5B, while sub-millimeter finger tremors induce a natural, frequency-modulated vibrato, confirming the system’s capacity for hyper-expressive zero-shot synthesis.

#### 4.5. Subjective Human Evaluation

Because generative music quality cannot be entirely captured by objective distance metrics, we conducted a Mean Opinion Score (MOS) user study following standard protocols [4]. We recruited 20 participants with musical backgrounds to evaluate 50 generated samples (25 from AIST++, 25 from ConductorMotion100).

Participants rated the samples on a 1–5 Likert scale across three dimensions: **MOS-Q** (overall acoustic quality and fidelity), **MOS-M** (musicality and harmonic coher-

ence), and **MOS-R** (responsiveness, evaluating how naturally the audio synced with the visual gestures).

Table 4. Subjective Mean Opinion Scores (MOS) with 95% confidence intervals. Our GDPO-aligned continuous system significantly outperforms baselines in perceived responsiveness and musicality.

Model	MOS-Q $\uparrow$	MOS-M $\uparrow$	MOS-R $\uparrow$
MLA-GRU (Prior)	3.12 $\pm$ 0.14	2.84 $\pm$ 0.18	2.15 $\pm$ 0.22
MusicGen (Autoreg.)	3.78 $\pm$ 0.12	3.85 $\pm$ 0.11	2.55 $\pm$ 0.20
Flow-Match + GRPO	3.85 $\pm$ 0.10	3.60 $\pm$ 0.15	3.95 $\pm$ 0.12
<b>Ours (GDPO)</b>	<b>4.21 <math>\pm</math> 0.08</b>	<b>4.15 <math>\pm</math> 0.10</b>	<b>4.42 <math>\pm</math> 0.09</b>

The results in Table 4 validate our architectural claims. While MusicGen scores highly on baseline musicality (MOS-M: 3.85), its severe autoregressive latency destroys the illusion of continuous control, yielding a poor responsiveness score (MOS-R: 2.55). Our proposed GDPO-aligned system achieves the highest scores across all metrics, proving that the continuous flow-matching paradigm feels tangibly more responsive to human performers.

#### 4.6. Ablation: Vision Backbone Dynamics

To justify the integration of the Liquid Time-Constant (LTC) network, we ablate the vision extraction backbone, replacing the LNN with contemporary sequence modeling alternatives: a standard Gated Recurrent Unit (GRU), a Spatial-Temporal Vision Transformer (ST-ViT), and a standard Neural ODE (without input-adaptive time constants).

Table 5. Vision Backbone Ablation. The Liquid Neural Network (LNN) provides the optimal trade-off between ultra-low inference latency and high cross-modal alignment (CLAP).

Vision Backbone	Vision Latency $\downarrow$	Total TTA $\downarrow$	CLAP $\uparrow$
Standard GRU	32.8 ms	132.4 ms	0.62
ST-ViT [25]	45.2 ms	144.8 ms	0.76
Neural ODE	<b>16.5 ms</b>	<b>116.1 ms</b>	0.68
<b>LNN (Ours)</b>	18.4 ms	118.0 ms	<b>0.79</b>

As shown in Table 5, while the ST-ViT achieves strong semantic alignment (CLAP: 0.76), its quadratic attention mechanism pushes the vision latency to 45.2 ms, threatening the 150 ms perceptual threshold. Conversely, the standard Neural ODE is exceptionally fast (16.5 ms) but fails to adapt to sudden, highly dynamic percussive gestures, lowering alignment. The LNN achieves the optimal balance, utilizing its input-adaptive time constants to match the ST-ViT’s alignment (0.79) while operating at near Neural ODE speeds (18.4 ms).

## 4.7. Ablation: Router Capacity and Modality Collapse

Finally, we analyze the structural necessity of the Audio-Visual Mixture of Experts (AVMoE). In multi-objective translation tasks, routing continuous states through a single dense network often results in modality collapse—where the network prioritizes prominent rhythmic features and ignores subtle timbral intent. We ablate the number of experts  $n$  in the routing network.

Table 6. AVMoE Expert Ablation. Utilizing  $n = 3$  experts prevents modality collapse, optimizing both audio quality and gesture alignment.

Router Config.	Params	FAD ↓	CLAP ↑	MOS-M ↑
Dense Network ( $n = 1$ )	4.2M	5.82	0.64	3.40
AVMoE ( $n = 2$ )	5.1M	4.15	0.72	3.82
AVMoE ( $n = 3$ , Ours)	6.0M	<b>3.47</b>	<b>0.79</b>	<b>4.15</b>
AVMoE ( $n = 5$ )	7.8M	3.51	0.78	4.12

Table 6 demonstrates severe modality collapse when utilizing a standard dense network ( $n = 1$ ), yielding a poor FAD of 5.82. Increasing the expert count to  $n = 3$  allows the network to successfully disentangle rhythm, timbre, and pitch constraints into independent latent subspaces. Scaling beyond  $n = 3$  to  $n = 5$  yields diminishing returns, marginally increasing parameter overhead and inference time without providing statistically significant gains in acoustic fidelity or human perceptual scores.

## 5. Conclusion

This paper presents a complete architectural transformation from discrete, classification-driven HCI to a continuous, zero-shot generative pipeline operating within the sub-150 ms perceptual constraint demanded by live musical performance.

First, replacing standard frame-buffered classifiers with a Liquid Time-Constant neural network models human motion as a continuous-time differential curve, reducing vision-stage processing from 32.8 ms to 18.4 ms. Second, the Audio-Visual Mixture of Experts (AVMoE) router projects kinematic states onto specialized acoustic manifolds without the modality collapse that afflicts direct latent regression. Third, the formal proof and empirical validation of Advantage Collapse in standard GRPO motivates our proposed GDPO objective. By decoupling reward normalization, GDPO sustains gradient signals across competing reward dimensions (latency vs. musicality), yielding a 44.9% reduction in Fréchet Audio Distance over prior systems. Finally, by mapping directly to continuous motion datasets (AIST++ and ConductorMotion100), the system naturally synthesizes advanced acoustic behaviours—such as glissandos and dynamic swells—driven entirely by hu-

man physical expression.

Taken together, the LNN+AVMoE+GDPO pipeline establishes a rigorous, reproducible framework for the next generation of embodied, real-time multimodal generation.

## 6. Limitations and Societal Impact

While the proposed continuous architecture significantly reduces latency and improves generative fidelity, it remains highly sensitive to the quality of the visual input. The LNN relies heavily on stable 30fps keypoint extraction; severe occlusions or poor lighting conditions can introduce high-frequency jitter into the spatial embeddings, which the AVMoE may misinterpret as rapid rhythmic intent—for example, inducing unintended vibrato artifacts during performance, despite the inclusion of the 1€ filter.

Furthermore, mapping fluid gestures to cloned vocal timbres via zero-shot SVS models such as SoulX-Singer [26] introduces severe deepfake vulnerabilities. To mitigate this, our flow-matching decoder natively integrates a SynthID-compatible spectrogram watermarking protocol directly into the generative denoising steps. By embedding cryptographic signatures into the frequency ranges during synthesis, the watermark remains imperceptible yet structurally bound to the audio, surviving acoustic transformations such as MP3 re-encoding and channel noise, thereby preventing unverified voice impersonation in live production environments.

## References

- [1] Andrea Agostinelli, I Denk Timo, Zalán Borsos, Jesse Engel, Neil Zeghidour, and Christian Frank. Musiclm: Generating music from text. In *arXiv preprint arXiv:2301.11325*, 2023. 2
- [2] Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Battenberg, and Neil Zeghidour. Audiolm: a language modeling approach to audio generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023. 2
- [3] Trey Brosnan. Darc: Drum accompaniment generation with fine-grained rhythm control. *arXiv preprint arXiv:284487153*, 2026. 2
- [4] Antoine Caillon, Brian McWilliams, Cassie Tarakajian, Ian Simon, Ilaria Manco, Jesse Engel, et al. Live music models. *arXiv preprint arXiv:280536884*, 2025. 1, 2, 6, 7
- [5] Anthony Chen, Naomi Ken Korem, Tavi Halperin, M Yosef, U Jelerčić, Ofir Bibi, Or Patashnik, and Daniel Cohen-Or. Just-dub-it: Video dubbing via joint audio-visual diffusion. *arXiv preprint arXiv:285140356*, 2026. 3
- [6] Gongyu Chen, Xiaoyu Zhang, Zhenqiang Weng, Junjie Zheng, Da Shen, Chaofan Ding, Wei-Qiang Zhang, and Zihao Chen. Yingmusic-svc: Real-world robust zero-shot singing voice conversion with flow-grpo and singing-specific inductive biases. *arXiv preprint arXiv:283557339*, 2025. 3

- [7] Ho Kei Cheng, Masato Ishii, Akio Hayakawa, Takashi Shibuya, Alexander G Schwing, and Yuki Mitsufuji. Mmaudio: Taming multimodal joint training for high-quality video-to-audio synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2
- [8] Jade Copet, Felix Kreuk, Itai Gat, Tal Remez Tal, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. Simple and controllable music generation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2, 5
- [9] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning audio concepts from natural language supervision. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023. 3, 6
- [10] Jakub Fil, Yulia Sandamirskaya, Hector A Gonzalez, Loïc Azzalin, Stefan Glüge, Lukas Friedenstab, Friedrich Wolf, Tim Rosmeisl, et al. Heterogeneous computing platform for real-time robotics. *arXiv preprint arXiv:284737792*, 2026. 2
- [11] Hugo Flores-Garcia et al. Sketch2sound: Controllable audio generation via time-varying control signals. In *ICASSP*, 2025. 3
- [12] Seth Forsgren and Hayk Martbourg. Riffusion-stable diffusion for real-time music generation, 2022. 2
- [13] Junmin Gong, S Zhao, Sen Wang, Shengyuan Xu, and J Guo. Ace-step: A step towards music generation foundation model. *arXiv preprint arXiv:279075130*, 2025. 1, 2
- [14] Junmin Gong, Yulin Song, Wenxiao Zhao, Sen Wang, Shengyuan Xu, Jing Guo, and Xuerui Yang. Ace-step 1.5: Pushing the boundaries of open-source music generation. *arXiv preprint arXiv:285270526*, 2026. 2
- [15] Ramin Hasani, Mathias Lechner, Alexander Amini, Daniela Rus, and Radu Grosu. Liquid time-constant networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(9):7657–7666, 2021. 2
- [16] Chia-Yu Hung, Navonil Majumder, Zhifeng Kong, Ambuj Mehrish, Rafael Valle, Bryan Catanzaro, and Soujanya Poria. Tangoflux: Super fast and faithful text to audio generation with flow matching and clap-ranked preference optimization. *arXiv preprint arXiv:275134284*, 2024. 3
- [17] Changhao Jiang, Jiahao Chen, Zhen Xiang, Zhixiong Yang, Hanchen Wang, Jiabao Zhuang, et al. Muse: Towards reproducible long-form song generation with fine-grained style control. *arXiv preprint arXiv:284532549*, 2026. 2
- [18] Tao Jiang, Peng Lu, Li Zhang, Nianjuan Ma, Rui Han, Chengqi Lyu, Yining Li, and Kai Chen. RtmPose: Real-time multi-person pose estimation based on mmpose. In *arXiv preprint arXiv:2303.07399*, 2023. 2
- [19] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d dance generation with aist++. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13401–13412, 2021. 2, 5
- [20] Xiquan Li, Junxi Liu, Yuzhe Liang, Zhikang Niu, Wenxi Chen, and Xie Chen. Meanaudio: Fast and faithful text-to-audio generation with mean flows. *arXiv preprint arXiv:280561084*, 2025. 2
- [21] Yizhi Li, Ruibin Yuan, Ge Zhang, Yinghao Ma, Xingran Chen, Hanzhi Yin, Chenghua Lin, et al. Mert: Acoustic music understanding model with large-scale self-supervised training. *arXiv preprint arXiv:2306.00107*, 2023. 3, 5
- [22] Yifan Liang, Andong Li, Kang Yang, Guochen Yu, Fangkun Liu, Lingling Dai, Xiaodong Li, and C Zheng. Sld-12s: Hierarchical subspace latent diffusion for high-fidelity lip to speech synthesis. *arXiv preprint arXiv:285540452*, 2026. 3
- [23] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow network based generative models. In *International Conference on Learning Representations (ICLR)*, 2023. 1
- [24] Camillo Lugaresi, Jiuqiang Tang, Hadon Nash, Chris McClanahan, Esha Uboweja, Michael Hays, Fan Zhang, Chuo-Ling Chang, Ming Guang Yong, Juhyun Lee, et al. Mediapipe: A framework for building perception pipelines. *arXiv preprint arXiv:1906.08172*, 2019. 1, 2, 5
- [25] William Peebles and Saining Xie. Scalable diffusion models with transformers. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 1, 7
- [26] Jiale Qian, Hao Meng, Yuhang Dai, Hongmei Liu, Tian Zheng, Pengcheng Zhu, Haopeng Lin, Hanke Xie, et al. Soulx-singer: Towards high-quality zero-shot singing voice synthesis. *arXiv preprint arXiv:2602.07803*, 2026. 1, 3, 8
- [27] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Feng, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 1, 3
- [28] Tri Ton, Jiajing Hong, and CD Yoo. Taro: Timestep-adaptive representation alignment with onset-aware conditioning for synchronized video-to-audio synthesis. *arXiv preprint arXiv:277628249*, 2025. 2
- [29] Olga Vechtomova and J Bos. Reimagining dance: Real-time music co-creation between dancers and ai. *arXiv preprint arXiv:279391756*, 2025. 2
- [30] Jun Wang, Xijuan Zeng, Chunyu Qiang, Ruilong Chen, Shiyao Wang, Le Wang, et al. Kling-foley: Multimodal diffusion transformer for high-quality video-to-audio generation. *arXiv preprint arXiv:280016202*, 2025. 2
- [31] Yongqi Wang, Wenxiang Guo, Rongjie Huang, Jia-Bin Huang, Zehan Wang, Fuming You, Ruiqi Li, and Zhou Zhao. Frieren: Efficient video-to-audio generation network with rectified flow matching. *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 2
- [32] Yazhou Xing, Yin-Yin He, Zeyue Tian, Xintao Wang, and Qifeng Chen. Seeing and hearing: Open-domain visual-audio generation with diffusion latent aligners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2
- [33] Xuenan Xu, Jiahao Mei, Zihao Zheng, Ye Tao, Zeyu Xie, Yaoyun Zhang, Haohe Liu, et al. Uniflow-audio: Unified flow matching for audio generation from omni-modalities. *arXiv preprint arXiv:281675868*, 2025. 3
- [34] Wei Zhang, Yue Wang, Guang Li, and Mingzhu Chen. Conductor motion forecasting via a transformer-based sequence-to-sequence model. *Journal of Computer Science and Technology*, 37(4):924–938, 2022. 3, 5

- [35] Y Zhao et al. Movenet: Ultra fast and accurate pose detection model. *Google Research Blog*, 2023. [2](#)
- [36] Junjie Zheng, Chunbo Hao, Guobin Ma, Xiaoyu Zhang, Gongyu Chen, Chaofan Ding, Zihao Chen, and Lei Xie. Yingmusic-singer: Zero-shot singing voice synthesis and editing with annotation-free melody guidance. *arXiv preprint arXiv:283556795*, 2025. [3](#)
- [37] Ye Zhu, Kyle Minhao Wu, Enric Dong, and Ying Du. Dance2music: Automatic dance-driven music generation. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 3228–3236, 2022. [2](#)