

GestureEncoder: a Joint Embedding Model for Movement-Based Musical Performance

JASON BRENT SMITH and BRYAN PARDO, Northwestern University, USA

Artificial Intelligence has been used to great effect to support musicians in live performance through gesture recognition, enabling control of the output of interactive music systems using movement. However, musicians working with gesture recognition must often manually map each new gesture to specific musical parameters, a lengthy and constrained process that is usually infeasible for live improvisation. Additionally, the inner workings of AI systems are hidden from users, making it difficult for them to understand if the system’s mappings between gestures and music align with their own. This contrasts with pairs of human improvisers who can build a shared understanding of mappings from motion to music in real time. This paper presents *GestureEncoder*, a model that maps gestures to musical outcomes in a way grounded in human understanding: a user can define gestures alongside a verbal description, and the system generates a joint embedding space of positional gesture data and user-defined physical and musical descriptions. *GestureEncoder* enables a user to quickly and organically create a set of motion-based musical instructions for human-understandable human-machine improvisation, using a two-dimensional visualization to represent how gestures are mapped to physical and musical descriptions.

Additional Key Words and Phrases: Artificial Intelligence, Music, Explainable AI

ACM Reference Format:

Jason Brent Smith and Bryan Pardo. 2018. GestureEncoder: a Joint Embedding Model for Movement-Based Musical Performance. *J. ACM* 37, 4, Article 111 (August 2018), 5 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

In a live improvised musical performance, musicians who perform together collaborate to define a shared musical language built from gestures and their meanings. These gestures can take simple yet subjective forms, such as a band leader lifting and lowering their head to subtly signal to the rest of the musicians that a new section of a song is about to begin. These languages can be given more formal or regimented definitions: for example, *Conduction* is a trademarked (protocol/language) by Lawrence D. “Butch” Morris that relies on the live enactment of motion gestures by an improviser’s orchestra leader, which are then refined and iterated on to enable the group to act as a unit [21].

Generative Artificial Intelligence (AI) is a rapidly growing area of research that has been incorporated into many music production workflows, with standard inputs to these systems often being text that describes whole music passages to generate [2, 15]. AI has also been used to map gestures to musical parameters during live musical performances. For example, Google’s MediaPipe [24] platform can be integrated into live camera applications for music performance, using classifiers for hand signs (e.g., open hand) and motions (e.g., circular pattern) to create low-dimensional representations of gestures that can be linked to audio parameters [18–20]. Although these motion-controlled systems are more capable

Authors’ Contact Information: Jason Brent Smith, jason.smith1@northwestern.edu; Bryan Pardo, pardo@northwestern.edu, Northwestern University, Evanston, Illinois, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Manuscript submitted to ACM

of real-time performances, their “black-box” nature has limited performers’ and audiences’ ability to understand how the system *interprets* gestures. This work explores how human-AI musical improvisation can be made more understandable, while maintaining the ability for an AI to react to gestures in real time, by enabling the AI to form connections between a user’s performance gestures and their intended musical meaning.

We present *Gesture Encoder*, a model for a human-understandable connection between physical motion and musical intent for improvisation on a level that mirrors Morris’ *Conduction*. This model encodes low-dimensional representations of gesture data alongside user-defined labels for that gesture’s meaning (both on a literal, physical level as well as the gesture’s intended musical output) into a shared embedding space [6]. Using this shared embedding space, a user’s motion can be converted into natural-language prompts for generative audio models, thereby providing an explainable, user-controlled, and real-time capable connection between the user and the AI’s shared understanding of motion and sound. Our goal is to use shared embeddings so that a generative model traditionally controlled by text can be controlled by gestures, while preserving the text instructions to promote understanding of the gesture-to-music mapping.

Embedding spaces are organizations in which topics can be arranged according to distances based on similarity [6], and various inputs have been used to perform conditional audio generation based on embedding spaces. Shared embedding spaces have been used to connect music and language for the generation of audio and symbolic music based on text input [8, 11, 13] as well as controllable latent representations of audio [3, 4]. MusicLM uses three models for audio and text representation and allows descriptive text to control generated audio [1]. VampNet generates variations of a given input audio segment by modeling masked acoustic tokens [7]. Diffusion models such as Stable Audio Open Small [14] and Magenta Realtime [22] enable real-time text-controlled audio generation. Using a shared gesture-sound latent space, *GestureEncoder* generates human-understandable prompts for language models, taking as input a user’s motions and creating a “vocabulary” that represents the intention of those motions.

2 *Gesture Encoder*

Gesture Encoder is a joint embedding model developed in PyTorch that combines three modalities:

- **Gesture Data:** Sequences of positional data captured through the Google MediaPipe hand-tracking model [24].
- **Gesture Description:** Verbal descriptions of physical motion (e.g., “spiral”, “point up”).
- **Music Description:** An intended musical outcome of a gesture (e.g., “speed up”, “fade”).

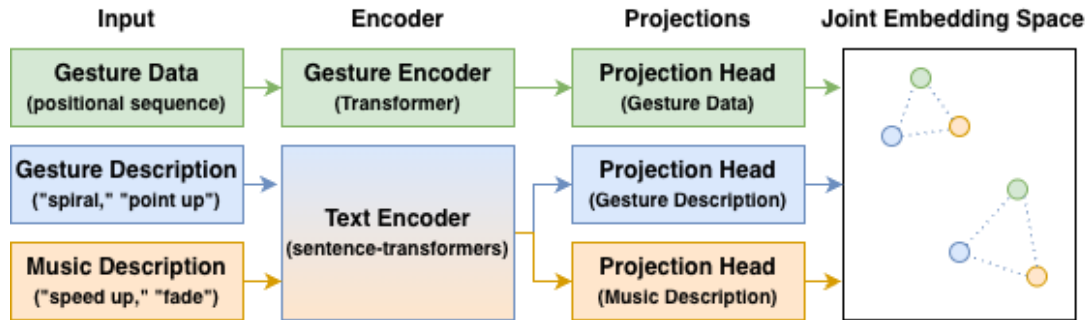


Fig. 1. *GestureEncoder* architecture. Gesture data sequences are encoded using a Transformer model, along with a language model for the gesture and music description labels. Three projection heads are combined to form a shared embedding space.

The user can record gestures and label them with motion descriptors and intended musical outcomes, thereby creating a data set of triplets: gesture, text, and music. Gestures are encoded using PyTorch `nn.TransformerEncoderLayer` and `nn.Transformer` layers with 256 dimensions. To encode the gesture description and music description labels, we use frozen, pre-trained Sentence Transformers¹ models [16] based on BERT architectures [10]. All three representations use a projection head with `nn.ReLU` and linear activation functions to form a 128-dimensional embedding. During training, *GestureEncoder* uses contrastive loss [9] by combining a PyTorch implementation² of InfoNCELoss [17] across the three domains through a weighted sum of InfoNCELoss (L) for each combination of gesture data (g), gesture description (d), and music description (m) predictions.

$$L = L_{g,t} + \alpha L_{g,m} + \beta L_{t,m} \quad (1)$$

3 Evaluation

We use a very limited data set of five gestures (Table 1) to demonstrate a feasible set of instructions that could be presented to *GestureEncoder* and labeled in one minute in a hypothetical performance situation. The average training time for these five gestures is 75.95 seconds over 300 epochs on an Apple M2 MacBook Pro 2022, which represents a small portion of the up to three-hour sessions of teaching “vocabulary” to musicians by Butch Morris when recording performances using *Conduction* [5].

Gesture Data	Gesture Description	Music Description
10 sec. video	Wave Hand	Speed Up
8 sec. video	Clap Hands	Stop
10 sec. video	Thumbs Up	Louder
10 sec. video	Point Left	Trombone
10 sec. video	Point Right	Trumpet

Table 1. A sample dataset of five user-defined, *Conduction*-style gestures.

Figure 2 shows a two-dimensional uniform manifold approximation and projection (UMAP) [12] visualization of the trained *GestureEncoder*, using cosine similarity. Each “triangle” represents a triplet of positional data, Gesture Description, and Music Description, with similar gestures grouped together. Each item is closer to the other members of its “triangle” than to any other item; as a result, on this small dataset, the retrieval accuracy and mean reciprocal rank (MRR) are both 100% for each combination of the three modalities. Gestures with an inherent similarity, such as how “point right” and “point left” are closer to each other than they are to the “thumbs up” gesture.

Geometric alignment (the average cosine similarity between pairs) and uniformity (the log of the average Gaussian kernel between pairs) are two metrics used to assess the quality of a representation learned with contrastive loss [23]. The alignment of each pair *GestureEncoder*’s modalities is 0.883 (gesture data, gesture description), 0.883 (gesture data, music description), and 0.926 (gesture description, music description), with each pair having a uniformity of -4.6. These results indicate greater alignment between the two text representations, as expected, since both were trained with the same encoder. Equal scores between the gesture and the two text modalities and the overall negative uniformity indicate an even spread of the “triplets” across the shared embedding space.

¹<https://pypi.org/project/sentence-transformers/>

²<https://github.com/RElbers/info-nce-pytorch>

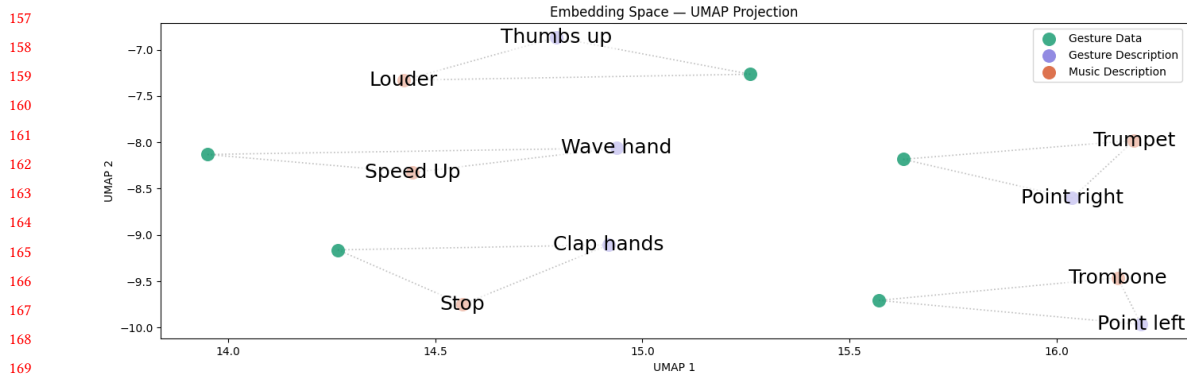


Fig. 2. The shared embedding space of *GestureEncoder* represented in two UMAP dimensions. For each gestural information (green), an equivalent motion descriptor (blue) and musical meaning (orange) are connected by a dotted line.

4 Conclusion

GestureEncoder is a prototype joint-embedding model that supports human-AI musical improvisation by enabling a user to define a "vocabulary" of physical movements alongside their intended musical outputs, allowing them to use gestures to control a language model in live performance settings. A user can define a set of performance gestures using a few video examples, along with labels indicating the gesture's semantic meaning and the intended musical outcome. As a result, *GestureEncoder* can recognize labeled performance gestures without explicit classification and compare their meanings in a semantic space that represents similarities between gestures and musical outcomes. In future work, we will expand on this by interpolating between known musical meanings to generate language-model-friendly musical instructions for previously unseen, undefined gestures. We also aim to incorporate *GestureEncoder* into live musical applications, leveraging its ability to map motion to human-understandable musical descriptions to serve as live prompting input for a generative audio model, while visualizing the shared embedding space.

Acknowledgments

This work was supported by NSF Award Number 2300633. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

References

- [1] Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. 2023. Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325* (2023).
- [2] Misagh Azimi and Mo H Zareei. 2025. Live Improvisation with Fine-Tuned Generative AI: A Musical Metacreation Approach. In *Proceedings of the International Conference on New Interfaces for Musical Expression*. 389–393.
- [3] Antoine Caillon and Philippe Esling. 2021. RAVE: A variational autoencoder for fast and high-quality neural audio synthesis. *arXiv preprint arXiv:2111.05011* (2021).
- [4] Franco Caspe, Jordie Shier, Mark Sandler, Charalampos Saitis, and Andrew McPherson. 2025. Designing Neural Synthesizers for Low Latency Interaction. *arXiv preprint arXiv:2503.11562* (2025).
- [5] Farai Chideya. 2008. Butch Morris on the Art of 'Conduction'. *NPR News & Notes*. <https://www.npr.org/transcripts/19145728> Accessed April 30, 2026.
- [6] Adji B Dieng, Francisco JR Ruiz, and David M Blei. 2020. Topic modeling in embedding spaces. *Transactions of the Association for Computational Linguistics* 8 (2020), 439–453.

- [7] H Flores_Garcia, P Seetharaman, R Kumar, and B Pardo. 2023. VampNet: Music Generation via Masked Acoustic Token Modeling. 24th International Society for Music Information Retrieval Conference.
- [8] Qingqing Huang, Aren Jansen, Joonseok Lee, Ravi Ganti, Judith Yue Li, and Daniel PW Ellis. 2022. Mulan: A joint embedding of music audio and natural language. *arXiv preprint arXiv:2208.12415* (2022).
- [9] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. Supervised contrastive learning. *Advances in neural information processing systems* 33 (2020), 18661–18673.
- [10] Mikhail V Koroteev. 2021. BERT: a review of applications in natural language processing and understanding. *arXiv preprint arXiv:2103.11943* (2021).
- [11] Max WY Lam, Qiao Tian, Tang Li, Zongyu Yin, Siyuan Feng, Ming Tu, Yuliang Ji, Rui Xia, Mingbo Ma, Xuchen Song, et al. 2024. Efficient neural music generation. *Advances in Neural Information Processing Systems* 36 (2024).
- [12] Leland McInnes, John Healy, and James Melville. 2018. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).
- [13] Gautam Mittal, Jesse Engel, Curtis Hawthorne, and Ian Simon. 2021. Symbolic music generation with diffusion models. *arXiv preprint arXiv:2103.16091* (2021).
- [14] Zachary Novack, Zach Evans, Zack Zukowski, Josiah Taylor, CJ Carr, Julian Parker, Adnan Al-Sinan, Gian Marco Iodice, Julian McAuley, Taylor Berg-Kirkpatrick, et al. 2025. Fast Text-to-Audio Generation with Adversarial Post-Training. *arXiv preprint arXiv:2505.08175* (2025).
- [15] Y Yogi Tegar Nugroho and P Paulus Metta Dwi Manggala. 2024. The use of AI in creating music compositions: A case study on Suno application. In *7th Celt International Conference (CIC 2024)*. Atlantis Press, 177–189.
- [16] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics. <https://arxiv.org/abs/1908.10084>
- [17] Evgenia Rusak, Patrik Reizinger, Attila Juhos, Oliver Bringmann, Roland S Zimmermann, and Wieland Brendel. 2024. Infonce: Identifying the gap between theory and practice. *arXiv preprint arXiv:2407.00143* (2024).
- [18] Jason Smith and Jason Freeman. 2021. Effects of Deep Neural Networks on the Perceived Creative Autonomy of a Generative Musical System. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, Vol. 17. 91–98.
- [19] Jason Brent Smith and Jason Freeman. 2023. Effects of Visual Explanation on Perceived Creative Autonomy in an AI-Based Generative Music System. In *Companion Proceedings of the 28th International Conference on Intelligent User Interfaces (Sydney, NSW, Australia) (IUI '23 Companion)*. Association for Computing Machinery, New York, NY, USA, 25–28.
- [20] Jason Brent Smith and Jason Freeman. 2025. Adaptation and Perceived Creative Autonomy in Gesture-Controlled Interactive Music. In *Proceedings of the 25th international conference on New Interfaces for Musical Expression*.
- [21] Thomas Taylor Stanley. 2009. Butch Morris and the Art of Conduction. (2009).
- [22] Lyria Team, Antoine Caillon, Brian McWilliams, Cassie Tarakajian, Ian Simon, Ilaria Manco, Jesse Engel, Noah Constant, Yunpeng Li, Timo I Denk, et al. 2025. Live music models. *arXiv preprint arXiv:2508.04651* (2025).
- [23] Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International conference on machine learning*. PMLR, 9929–9939.
- [24] Fan Zhang, Valentin Bazarevsky, Andrey Vakunov, Andrei Tkachenka, George Sung, Chuo-Ling Chang, and Matthias Grundmann. 2020. Mediapipe hands: On-device real-time hand tracking. *arXiv preprint arXiv:2006.10214* (2020).