
MusicLM: Generating Music From Text

Andrea Agostinelli^{*1} Timo I. Denk^{*1}

Zalán Borsos¹ Jesse Engel¹ Mauro Verzetti¹ Antoine Caillon² Qingqing Huang¹ Aren Jansen¹
Adam Roberts¹ Marco Tagliasacchi¹ Matt Sharifi¹ Neil Zeghidour¹ Christian Frank¹

Abstract

We introduce MusicLM, a model for generating high-fidelity music from text descriptions such as “a calming violin melody backed by a distorted guitar riff”. MusicLM casts the process of conditional music generation as a hierarchical sequence-to-sequence modeling task, and it generates music at 24 kHz that remains consistent over several minutes. Our experiments show that MusicLM outperforms previous systems both in audio quality and adherence to the text descriptions. Moreover, we demonstrate that MusicLM can be conditioned on both text and a melody in that it can transform whistled and hummed melodies according to the style described in a text caption. To support future research, we publicly release MusicCaps, a dataset composed of 5.5k music-text pairs, with rich text descriptions provided by human experts. google-research.github.io/seanet/musiclm/examples

1. Introduction

Conditional neural audio generation covers a wide range of applications, ranging from text-to-speech (Zen et al., 2013; van den Oord et al., 2016) to lyrics-conditioned music generation (Dhariwal et al., 2020) and audio synthesis from MIDI sequences (Hawthorne et al., 2022b). Such tasks are facilitated by a certain level of temporal alignment between the conditioning signal and the corresponding audio output. In contrast, and inspired by progress in text-to-image generation (Ramesh et al., 2021; 2022; Saharia et al., 2022; Yu et al., 2022), recent work has explored generating audio from sequence-wide, high-level captions (Yang et al., 2022; Kreuk et al., 2022) such as “whistling with wind blowing”. While generating audio from such coarse captions represents a breakthrough, these models remain limited to simple acoustic scenes, consisting of few acoustic events over a

period of seconds. Hence, turning a single text caption into a rich audio sequence with long-term structure and many stems, such as a music clip, remains an open challenge.

AudioLM (Borsos et al., 2022) has recently been proposed as a framework for audio generation. Casting audio synthesis as a language modeling task in a discrete representation space, and leveraging a hierarchy of coarse-to-fine audio discrete units (or *tokens*), AudioLM achieves both high-fidelity and long-term coherence over dozens of seconds. Moreover, by making no assumptions about the content of the audio signal, AudioLM learns to generate realistic audio from audio-only corpora, be it speech or piano music, without any annotation. The ability to model diverse signals suggests that such a system could generate richer outputs if trained on the appropriate data.

Besides the inherent difficulty of synthesizing high-quality and coherent audio, another impeding factor is the scarcity of paired audio-text data. This is in stark contrast with the image domain, where the availability of massive datasets contributed significantly to the remarkable image generation quality that has recently been achieved (Ramesh et al., 2021; 2022; Saharia et al., 2022; Yu et al., 2022). Moreover, creating text descriptions of general audio is considerably harder than describing images. First, it is not straightforward to unambiguously capture with just a few words the salient characteristics of either acoustic scenes (e.g., the sounds heard in a train station or in a forest) or music (e.g., the melody, the rhythm, the timbre of vocals and the many instruments used in accompaniment). Second, audio is structured along a temporal dimension which makes sequence-wide captions a much weaker level of annotation than an image caption.

In this work, we introduce MusicLM, a model for generating high-fidelity music from text descriptions. MusicLM leverages AudioLM’s multi-stage autoregressive modeling as the generative component, while extending it to incorporate text conditioning. To address the main challenge of paired data scarcity, we rely on MuLan (Huang et al., 2022), a joint music-text model that is trained to project music and its corresponding text description to representations close to each other in an embedding space. This shared embedding space eliminates the need for captions at training time alto-

^{*}Equal contribution ¹Google Research ²IRCAM - Sorbonne Université (work done while interning at Google). Correspondence to: Christian Frank <chfrank@google.com>.

gether, and allows training on massive audio-only corpora. That is, we use the MuLan embeddings computed from the audio as conditioning during training, while we use MuLan embeddings computed from the text input during inference.

When trained on a large dataset of unlabeled music, MusicLM learns to generate long and coherent music at 24 kHz, for text descriptions of significant complexity, such as “*enchanting jazz song with a memorable saxophone solo and a solo singer*” or “*Berlin 90s techno with a low bass and strong kick*”. To address the lack of evaluation data for this task, we introduce MusicCaps, a new high-quality music caption dataset with 5.5k examples prepared by expert musicians, which we publicly release to support future research.

Our experiments show through quantitative metrics and human evaluations that MusicLM outperforms previous systems such as Mubert (Mubert-Inc, 2022) and Riffusion (Forsgren & Martiros, 2022), both in terms of quality and adherence to the caption. Furthermore, since describing some aspects of music with words can be difficult or even impossible, we show how our method supports conditioning signals beyond text. Concretely, we extend MusicLM to accept an additional melody in the form of audio (e.g., whistling, humming) as conditioning to generate a music clip that follows the desired melody, rendered in the style described by the text prompt.

We acknowledge the risks associated with music generation, in particular, the potential misappropriation of creative content. In accordance with responsible model development practices, we conduct a thorough study of memorization by adapting and extending the methodology of Carlini et al. (2022) used for text-based large language models. Our findings show that when feeding MuLan embeddings to MusicLM, the sequences of generated tokens significantly differ from the corresponding sequences in the training set.

The key contributions of this work are the following:

1. We introduce MusicLM, a generative model that produces high-quality music at 24 kHz which is consistent over several minutes while being faithful to a text conditioning signal.
2. We extend our method to other conditioning signals, such as a melody that is then synthesized according to the text prompt. Furthermore, we demonstrate long and coherent music generation of up to 5-minute long clips.
3. We release the first evaluation dataset collected specifically for the task of text-to-music generation: MusicCaps is a hand-curated, high-quality dataset of 5.5k music-text pairs prepared by musicians.

2. Background and Related Work

The state-of-the-art in generative modeling for various domains is largely dominated either by Transformer-based autoregressive models (Vaswani et al., 2017) or U-Net-based diffusion models (Ho et al., 2020). In this section, we review the related work with an emphasis on autoregressive generative models operating on discrete tokens, which share similarities with MusicLM.

2.1. Quantization

Modeling sequences of discrete tokens autoregressively has proven to be a powerful approach in natural language processing (Brown et al., 2020; Cohen et al., 2022) and image or video generation (Esser et al., 2021; Ramesh et al., 2021; Yu et al., 2022; Villegas et al., 2022). Quantization is a key component to the success of autoregressive models for continuous signals, including images, videos, and audio. The goal of quantization is to provide a compact, discrete representation, which at the same time allows for high-fidelity reconstruction. VQ-VAEs (Van Den Oord et al., 2017) demonstrated impressive reconstruction quality at low bitrates in various domains and serve as the underlying quantizer for many approaches.

SoundStream (Zeghidour et al., 2022) is a universal neural audio codec capable of compressing general audio at low bitrates, while maintaining a high reconstruction quality. To achieve this, SoundStream uses residual vector quantization (RVQ), allowing scalability to higher bitrate and quality, without a significant computational cost. More specifically, RVQ is a hierarchical quantization scheme composing a series of vector quantizers, where the target signal is reconstructed as the sum of quantizer outputs. Due to the composition of quantizers, RVQ avoids the exponential blowup in the codebook size as the target bitrate increases. Moreover, the fact that each quantizer is fitted to the residual of coarser quantizers introduces a hierarchical structure to the quantizers, where coarser levels are more important for high-fidelity reconstruction. This property is desirable for generation, since the past context can be defined by only attending to the coarse tokens. Recently, SoundStream was extended by EnCodec (Défossez et al., 2022) to higher bitrates and stereophonic audio. In this work, we rely on SoundStream as our audio tokenizer, since it can reconstruct 24 kHz music at 6 kbps with high fidelity.

2.2. Generative Models for Audio

Despite the challenge of generating high-quality audio with long-term consistency, a series of approaches have recently tackled the problem with some success. Jukebox (Dhariwal et al., 2020), for example, proposes a hierarchy of VQ-VAEs at various time resolutions to achieve high temporal

coherence, but the generated music displays noticeable artifacts. PerceiverAR (Hawthorne et al., 2022a), on the other hand, proposes to model a sequence of SoundStream tokens autoregressively, achieving high-quality audio, but compromising the long-term temporal coherence.

Inspired by these approaches, AudioLM (Borsos et al., 2022) addresses the trade-off between coherence and high-quality synthesis by relying on a hierarchical tokenization and generation scheme. Concretely, the approach distinguishes between two token types: (1) *semantic* tokens that allow the modeling of long-term structure, extracted from models pretrained on audio data with the objective of masked language modeling; (2) *acoustic* tokens, provided by a neural audio codec, for capturing fine acoustic details. This allows AudioLM to generate coherent and high-quality speech as well as piano music continuations without relying on transcripts or symbolic music representations.

MusicLM builds on top of AudioLM with three important additional contributions: (1) we condition the generation process on a descriptive text, (2) we show that the conditioning can be extended to other signals such as melody, and (3) we model a large variety of long music sequences beyond piano music (from *drum'n'bass* over *jazz* to *classical music*).

2.3. Conditioned Audio Generation

Generating audio from a text description (such as “*whistling with laughter in the background*”) has recently been tackled by several works. DiffSound (Yang et al., 2022) uses CLIP (Radford et al., 2021) as the text encoder and applies a diffusion model to predict the quantized mel spectrogram features of the target audio based on the text embeddings. AudioGen (Kreuk et al., 2022) uses a T5 (Raffel et al., 2020) encoder for embedding the text, and an autoregressive Transformer decoder for predicting target audio codes produced by EnCodec (Défossez et al., 2022). Both approaches rely on a modest amount of paired training data such as AudioSet (Gemmeke et al., 2017) and AudioCaps (Kim et al., 2019) (totalling less than 5k hours after filtering).

Closer to MusicLM, there are also works focusing on music generation conditioned on text. In Mubert (Mubert-Inc, 2022), the text prompt is embedded by a Transformer, music tags which are close to the encoded prompt are selected and used to query the song generation API. Based on the selected tags, Mubert generates a combination of sounds, which in turn were generated by musicians and sound designers. This is in contrast to Riffusion (Forsgren & Martiros, 2022), which fine-tunes a Stable Diffusion model (Rombach et al., 2022a) on mel spectrograms of music pieces from a paired music-text dataset. We use both Mubert and Riffusion as baselines for our work, showing that we improve the audio generation quality and adherence to the text description.

Symbolic representations of music (e.g., MIDI) can also be used to drive the generative process as a form of strong conditioning, as demonstrated by Huang et al. (2019); Hawthorne et al. (2019); Engel et al. (2020). MusicLM enables a more natural and intuitive way of providing a conditioning signal, for example through a hummed melody, which can also be combined with a text description.

2.4. Text-Conditioned Image Generation

Precursor to text-conditioned audio synthesis are the text-conditioned image generation models, which made significant progress in quality due to architectural improvements and the availability of massive, high-quality paired training data. Prominent Transformer-based autoregressive approaches include Ramesh et al. (2021); Yu et al. (2022), while Nichol et al. (2022); Rombach et al. (2022b); Saharia et al. (2022) present diffusion-based models. The text-to-image approaches have been extended to generating videos from a text prompt (Wu et al., 2022a; Hong et al., 2022; Villegas et al., 2022; Ho et al., 2022).

The closest to our approach among these works is DALL-E 2 (Ramesh et al., 2022). In particular, similarly to the way DALL-E 2 relies on CLIP (Radford et al., 2021) for text encoding, we also use a joint music-text embedding model for the same purpose. In contrast to DALL-E 2, which uses a diffusion model as a decoder, our decoder is based on AudioLM. Furthermore, we also omit the prior model mapping text embeddings to music embeddings, such that the AudioLM-based decoder can be trained on an audio-only dataset and the music embedding is simply replaced during inference by the text embedding.

2.5. Joint Embedding Models for Music and Text

MuLan (Huang et al., 2022) is a music-text joint embedding model consisting of two embedding towers, one for each modality. The towers map the two modalities to a shared embedding space of 128 dimensions using contrastive learning, with a setup similar to (Radford et al., 2021; Wu et al., 2022b). The text embedding network is a BERT (Devlin et al., 2019) pre-trained on a large corpus of text-only data, while we use the ResNet-50 variant of the audio tower.

MuLan is trained on pairs of music clips and their corresponding text annotations. Importantly, MuLan imposes only weak requirements on its training data quality, learning cross-modal correspondences even when the music-text pairs are only weakly associated. The ability to link music to unconstrained natural language descriptions makes it applicable for retrieval or zero-shot music tagging. In this work, we rely on the pretrained and frozen model of Huang et al. (2022).

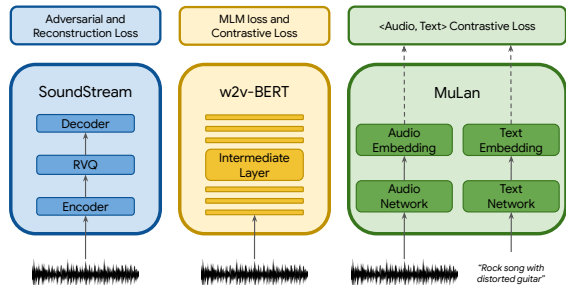


Figure 1. Independent pretraining of the models providing the audio and text representations for MusicLM: SoundStream (Zeghidour et al., 2022), w2v-BERT (Chung et al., 2021), and MuLan (Huang et al., 2022).

3. Method

In this section, we describe MusicLM and its components. Section 3.1 describes the models that provide the audio representations. Then, we show in Section 3.2 how we use these representations for text-conditioned music generation.

3.1. Representation and Tokenization of Audio and Text

We use three models for extracting audio representations that will serve for conditional autoregressive music generation, which are illustrated in Figure 1. In particular, by following the approach of AudioLM, we use the self-supervised audio representations of SoundStream (Zeghidour et al., 2022), as acoustic tokens to enable high-fidelity synthesis, and w2v-BERT (Chung et al., 2021), as semantic tokens to facilitate long-term coherent generation. For representing the conditioning, we rely on the MuLan music embedding during training and the MuLan text embedding at inference time. All three of these models are pretrained independently and then frozen, such that they provide the discrete audio and text representations for the sequence-to-sequence modeling.

SoundStream. We use a SoundStream model for 24 kHz monophonic audio with a striding factor of 480, resulting in 50 Hz embeddings. The quantization of these embeddings is learned during training by an RVQ with 12 quantizers, each with a vocabulary size of 1024. This results in a bitrate of 6 kbps, where one second of audio is represented by 600 tokens. We refer to these as *acoustic tokens*, denoted by A .

w2v-BERT. Similarly to AudioLM, we use an intermediate layer of the masked-language-modeling (MLM) module of a w2v-BERT model with 600M parameters. After pretraining and freezing the model, we extract embeddings from the 7th layer and quantize them using the centroids of a learned k-means over the embeddings. We use 1024 clusters and a sampling rate of 25 Hz, resulting in 25 *semantic tokens* for every second of audio, denoted by S .

MuLan. To train MusicLM, we extract the representation of the target audio sequence from the audio-embedding network of MuLan. Note that this representation is continuous and could be directly used as a conditioning signal in Transformer-based autoregressive models. However, we opt for quantizing the MuLan embeddings in such a way that both the audio and the conditioning signal have a homogeneous representation based on discrete tokens, aiding further research into autoregressively modeling the conditioning signal as well.

Since MuLan operates on 10-second audio inputs and we need to process longer audio sequences, we calculate the audio embeddings on 10-second windows with 1-second stride and average the resulting embeddings. We then discretize the resulting embedding by applying an RVQ with 12 vector quantizers, each with a vocabulary size of 1024. This process yields 12 MuLan audio tokens M_A for an audio sequence. During inference, we use as conditioning the MuLan text embedding extracted from the text prompt, and quantize it with the same RVQ as the one used for the audio embeddings, to obtain 12 tokens M_T .

Conditioning on M_A during training has two main advantages. First, it allows us to easily scale our training data, since we are not limited by the need of text captions. Second, by exploiting a model like MuLan, trained using a contrastive loss, we increase the robustness to noisy text descriptions.

3.2. Hierarchical Modeling of Audio Representations

We combine the discrete audio representations presented above with AudioLM to achieve text-conditioned music generation. For this, we propose a hierarchical sequence-to-sequence modeling task, where each stage is modeled autoregressively by a separate decoder-only Transformer. The proposed approach is illustrated in Figure 2.

The first stage is the *semantic modeling* stage, which learns the mapping from the MuLan audio tokens to the semantic tokens S , by modeling the distribution $p(S_t|S_{<t}, M_A)$, where t is the position in the sequence corresponding to a the time step. The second stage is the *acoustic modeling* stage, where the acoustic tokens A_q are predicted conditioned on both the MuLan audio tokens and the semantic tokens, modeling the distribution $p(A_t|A_{<t}, S, M_A)$.

Notably, to avoid long token sequences, AudioLM proposed to further split the acoustic modeling stage into a coarse and fine modeling stage. We rely on the same approach, where the coarse stage models the first four levels from the output of the SoundStream RVQ, and the fine stage models the remaining eight — we refer to Borsos et al. (2022) for details.

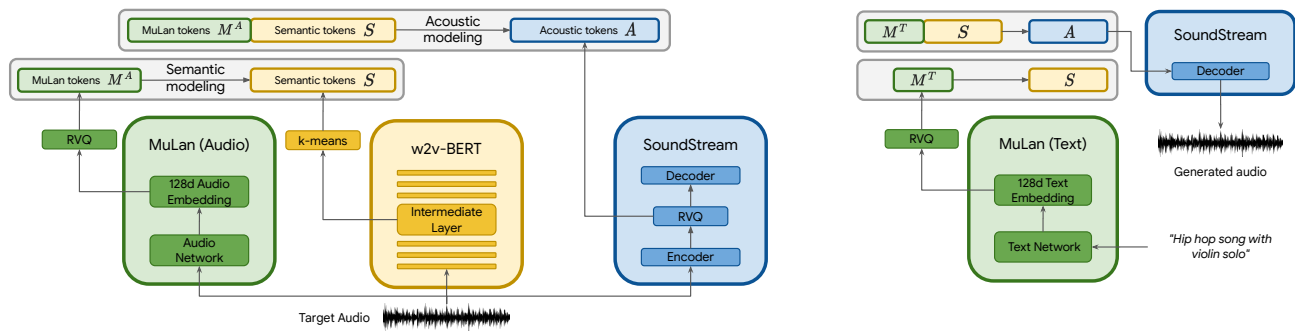


Figure 2. Left: During training we extract the MuLan audio tokens, semantic tokens, and acoustic tokens from the *audio-only* training set. In the semantic modeling stage, we predict semantic tokens using MuLan audio tokens as conditioning. In the subsequent acoustic modeling stage, we predict acoustic tokens, given both MuLan audio tokens and semantic tokens. Each stage is modeled as a sequence-to-sequence task using decoder-only Transformers. Right: During inference, we use MuLan text tokens computed from the text prompt as conditioning signal and convert the generated audio tokens to waveforms using the SoundStream decoder.

4. Experimental Setup

4.1. Models

We use decoder-only Transformers for modeling the semantic stage and the acoustic stages of AudioLM. The models share the same architecture, composed of 24 layers, 16 attention heads, an embedding dimension of 1024, feed-forward layers of dimensionality 4096, dropout of 0.1, and relative positional embeddings (Raffel et al., 2020), resulting in 430M parameters per stage.

4.2. Training and Inference

By relying on pretrained and frozen MuLan, we need audio-only data for training the other components of MusicLM. We train SoundStream and w2v-BERT on the Free Music Archive (FMA) dataset (Defferrard et al., 2017), whereas the tokenizers and the autoregressive models for the semantic and acoustic modeling stages are trained on a dataset containing five million audio clips, amounting to 280k hours of music at 24 kHz. Each of the stages is trained with multiple passes over the training data. We use 30 and 10-second random crops of the target audio for the semantic stage and the acoustic stage, respectively. The AudioLM fine acoustic modeling stage is trained on 3-second crops.

During inference, we make use of the joint embedding space between audio and text learned by MuLan, that is, we substitute M_A with M_T . We then follow the stages described above and obtain A given M_T . We use temperature sampling for the autoregressive sampling in all stages, with temperature of 1.0 for the semantic modeling stage, 0.95 and 0.4 for the coarse and fine acoustic modeling stages respectively. These temperature values were chosen based on subjective inspection to provide a good trade-off between diversity and temporal consistency of the generated music.

4.3. Evaluation Dataset

To evaluate MusicLM, we prepare MusicCaps, a high-quality music caption dataset, which we make publicly available.¹ This dataset includes 5.5k music clips from AudioSet (Gemmeke et al., 2017), each paired with corresponding text descriptions in English, written by ten professional musicians. For each 10-second music clip, MusicCaps provides: (1) a free-text *caption* consisting of four sentences on average, describing the music and (2) a list of music *aspects*, describing genre, mood, tempo, singer voices, instrumentation, dissonances, rhythm, etc. On average, the dataset includes eleven aspects per clip. See Appendix A for a few caption and aspect list examples.

MusicCaps complements AudioCaps (Kim et al., 2019), as they both contain audio clips from AudioSet with corresponding textual descriptions. However, while AudioCaps contains non-music content, MusicCaps focuses exclusively on music and includes highly detailed expert-provided annotations. The examples are extracted from both the train and eval split of AudioSet, covering a diverse distribution of genres, as detailed in Appendix A. MusicCaps also provides a *genre-balanced* split of the data with 1k examples.

4.4. Metrics

We compute different metrics to evaluate MusicLM, capturing two important aspects of music generation: the audio quality and the adherence to the text description.

Fréchet Audio Distance (FAD). The Fréchet Audio Distance (Kilgour et al., 2019) is a reference-free audio quality metric, which correlates well with human perception. Models producing samples with a low FAD score are expected

¹kaggle.com/datasets/googleai/musiccaps

to generate plausible audio. However, the generated samples might not necessarily adhere to the text description provided as conditioning.

We report the FAD based on two audio embedding models, both of which are publicly available: (1) Trill² (Shor et al., 2020), which is trained on speech data, and (2) VGGish³, (Hershey et al., 2017) which is trained on the YouTube-8M audio event dataset (Abu-El-Haija et al., 2016). Because of the difference in training data, we expect the models to measure different aspects of the audio quality (speech and non-speech, respectively).

KL Divergence (KLD). There is a many-to-many relationship between text descriptions and music clips compatible with them. It is therefore not possible to directly compare the generated music with the reference at the level of the audio waveform. To assess the adherence to the input text description, we adopt a proxy method similar to the one proposed in Yang et al. (2022); Kreuk et al. (2022). Specifically, we use a LEAF (Zeghidour et al., 2021) classifier trained for multi-label classification on AudioSet, to compute class predictions for both the generated and the reference music and measure the KL divergence between probability distributions of class predictions. When the KL-divergence is low, the generated music is expected to have similar acoustic characteristics as the reference music, according to the classifier.

MuLan Cycle Consistency (MCC). As a joint music-text embedding model, MuLan can be used to quantify the similarity between music-text pairs. We compute the MuLan embeddings from the text descriptions in MusicCaps as well as the generated music based on them, and define the MCC metric as the average cosine similarity between these embeddings.

Qualitative evaluation. Ultimately, we rely on subjective tests to evaluate the adherence of generated samples to the text description. We set up an A-vs-B human rating task, in which raters are presented with the text description and two samples of music generated by two different models, or one model and the reference music. There are five possible answers: strong or weak preference for A or B, and no preference. The raters are instructed not to take the music quality into account when making their decision, because this aspect of the evaluation is already covered by the FAD metric.

We consider the output of n different models, in addition to the reference music, thus a total of $n + 1$ conditions and $n(n + 1)/2$ pairs. To aggregate the results of the pairwise tests and rank conditions, we count the number of “wins”,

that is, how often a condition is strongly or weakly preferred. The samples are selected from the genre-balanced 1k subset of our evaluation data.

Training data memorization. Large language models have the capacity to memorize patterns seen in the training data (Carlini et al., 2020). We adapt the methodology used in Carlini et al. (2022) to study the extent to which MusicLM might memorize music segments. We focus on the first stage, responsible for semantic modeling. We select N examples at random from the training set. For each example, we feed to the model a prompt which includes the MuLan audio tokens M_A followed by a sequence of the first T semantic tokens S , with $T \in \{0, \dots, 250\}$, corresponding to up to 10 seconds. We use greedy decoding to generate a continuation of 125 semantic tokens (5 seconds) and we compare the generated tokens to the target tokens in the dataset. We measure exact matches as the fraction of examples for which generated and target tokens are identical over the whole sampled segment.

In addition, we propose a methodology to detect approximate matches, based on the observation that sequences of seemingly different tokens might lead to acoustically similar audio segments. Namely, we compute the histogram of semantic token counts over the corresponding vocabulary $\{0, \dots, 1023\}$ from both the generated and target tokens, and define a matching cost measure between histograms as follows. First, we compute the distance matrix between pairs of semantic tokens, which is populated by the Euclidean distances between the corresponding k-means centroids used to quantize w2v-BERT to semantic tokens (see Section 3.1). Then, we solve an optimal transport problem to find the matching cost between a pair of histograms using the Sinkhorn algorithm (Cuturi, 2013), considering only the sub-matrix corresponding to non-zero token counts in the two histograms. To calibrate the threshold used to determine whether two sequences might be approximate matches, we construct negative pairs by permuting the examples with target tokens and measure the empirical distribution of matching costs for such negative pairs. We set the match threshold τ to 0.85, which leads to less than 0.01% false positive approximate matches.

5. Results

We evaluate MusicLM by comparing it with two recent baselines for music generation from descriptive text, namely Mubert (Mubert-Inc, 2022) and Riffusion (Forsgren & Martiros, 2022). In particular, we generate audio by querying the Mubert API,⁴ and by running inference on the Riffusion model.⁵ We perform our evaluations on MusicCaps, the evaluation dataset we publicly release together with this paper.

²tfhub.dev/google/nonsemantic-speech-benchmark/trill/3

³tfhub.dev/google/vggish/1

⁴github.com/MubertAI (accessed in Dec 2022 and Jan 2023)

⁵github.com/riffusion/riffusion-app (accessed on Dec 27, 2022)

Table 1. Evaluation of generated samples using captions from the MusicCaps dataset. Models are compared in terms of audio quality, by means of Fréchet Audio Distance (FAD), and faithfulness to the text description, using Kullback–Leibler Divergence (KLD) and MuLan Cycle Consistency (MCC), and counts of wins in pairwise human listening tests (Wins).

MODEL	FAD _{TRILL} ↓	FAD _{VGG} ↓	KLD ↓	MCC ↑	WINS ↑
RIFFUSION	0.76	13.4	1.19	0.34	158
MUBERT	0.45	9.6	1.58	0.32	97
MUSICLM	0.44	4.0	1.01	0.51	312
MUSICCAPS	-	-	-	-	472

Comparison to baselines. Table 1 reports the main quantitative and qualitative results of this paper. In terms of audio quality, as captured by the FAD metrics, on FAD_{VGG} MusicLM achieves better scores than Mubert and Riffusion. On FAD_{TRILL}, MusicLM scores similarly to Mubert (0.44 vs. 0.45) and better than Riffusion (0.76). We note that, according to these metrics, MusicLM is capable of generating high-quality music comparable to Mubert, which relies on pre-recorded sounds prepared by musicians and sound designers. In terms of faithfulness to the input text description, as captured by KLD and MCC, MusicLM achieves the best scores, suggesting that it is able to capture more information from the text descriptions compared to the baselines.

We further supplement our evaluation of text faithfulness with a human listening test. Participants are presented with two 10-second clips and a text caption, and asked which clip is best described by the text of the caption on a 5-point Likert scale. We collect 1200 ratings, with each source involved in 600 pair-wise comparisons. Table 1 reports the total number of “wins”, that is, counting how often the human raters preferred a model in a side-by-side comparison. MusicLM is clearly preferred over both baselines, while there is still a measurable gap to the ground truth reference music. Full details of the listening study can be found in Appendix B.

Listening to examples in which the ground truth was preferred over MusicLM reveals the following patterns: (1) captions are extremely detailed, referring to more than five instruments or describing non musical aspects such as “wind, people talking”; (2) captions describe temporal ordering of the audio being played; (3) negations are used, which are not well captured by MuLan.

Overall, we conclude that: (1) our approach is able to capture fine-grained information from the rich free-text captions of MusicCaps; (2) the KLD and MCC metrics provide a quantitative measure of the faithfulness to the text description, which is in accordance with the human rating study.

Importance of semantic tokens. To understand the usefulness of decoupling semantic modeling from acoustic mod-

eling, we train a Transformer model which directly predicts coarse acoustic tokens from MuLan tokens, by modeling $p(A_t|A_{<t}, M_A)$. We observe that while the FAD metrics are comparable (0.42 FAD_{TRILL} and 4.0 FAD_{VGG}), KLD and MCC scores worsen when removing the semantic modeling stage. In particular the KLD score increases from 1.01 to 1.05, and the MCC score decreases from 0.51 to 0.49, indicating that semantic tokens facilitate the adherence to the text description. We also confirm this qualitatively by listening to the samples. In addition, we observe degradation in long term structure.

Information represented by audio tokens. We conduct additional experiments to study the information captured by the semantic and the acoustic tokens. In the first study, we fix the MuLan text tokens as well as the semantic tokens, running the acoustic modeling stage multiple times to generate several samples. In this case, by listening to the generated music, it is possible to observe that the samples are diverse, yet they tend to share the same genre, rhythmical properties (e.g., drums), and part of the main melody. They differ in terms of specific acoustic properties (e.g., level of reverb, distortion) and, in some cases, different instruments with a similar pitch range can be synthesized in different examples. In the second study, we fix only the MuLan text tokens and generate both the semantic and acoustic tokens. In this case, we observe a much higher level of diversity in terms of melodies and rhythmic properties, still coherent with the text description. We provide samples from this study in the accompanying material.

Memorization analysis. Figure 3 reports both exact and approximate matches when the length of the semantic token prompt is varied between 0 and 10 seconds. We observe that the fraction of exact matches always remains very small (< 0.2%), even when using a 10 second prompt to generate a continuation of 5 seconds. Figure 3 also includes results for approximate matches, using $\tau = 0.85$. We can see a higher number of matches detected with this methodology, also when using only MuLan tokens as input (prompt length $T = 0$) and the fraction of matching examples increases as the length of the prompt increases. We inspect these matches more closely and observe that those with the lowest matching score correspond to sequences characterized by a low level of token diversity. Namely, the average empirical entropy of a sample of 125 semantic tokens is 4.6 bits, while it drops to 1.0 bits when considering sequences detected as approximate matches with matching score less than 0.5. We include a sample of approximate matches obtained with $T = 0$ in the accompanying material. Note that acoustic modeling carried out by the second stage introduces further diversity in the generated samples, also when the semantic tokens match exactly.

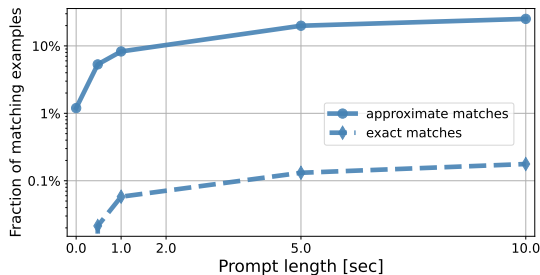


Figure 3. Memorization results for the semantic modeling stage. We compare the semantic tokens generated for 5 seconds of audio to corresponding tokens in the training set, considering exact and approximate matches.

6. Extensions

Melody conditioning. We extend MusicLM in such a way that it can generate music based on both a text description and a melody, which is provided in the form of humming, singing, whistling, or playing an instrument. This requires extending the conditioning signal in a way that captures the target melody. To this end, we create a synthetic dataset composed of audio pairs with matching melodies but different acoustics. To create such pairs, we use different versions of the same music clip, such as covers, instrumentals, or vocals. Additionally, we acquire data pairs of people humming and singing. We then train a joint embedding model such that when two audio clips contain the same melody, the corresponding embeddings are close to each other. For implementation details we refer to Appendix C.

To extract the melody conditioning for MusicLM, we quantize the melody embeddings with RVQ, and concatenate the resulting token sequences with the MuLan audio tokens M_A . During inference, we compute melody tokens from the input audio clip and concatenate them with the MuLan text tokens M_T . Based on this conditioning, MusicLM can successfully generate music which follows the melody contained in the input audio clip, while adhering to the text description.

Long generation and story mode. In MusicLM, generation is autoregressive in the temporal dimension which makes it possible to generate sequences longer than those used during training. In practice, the semantic modeling stage is trained on sequences of 30 seconds. To generate longer sequences, we advance with a stride of 15 seconds, using 15 seconds as prefix to generate an additional 15 seconds, always conditioning on the same text description. With this approach we can generate long audio sequences which are coherent over several minutes.

With a small modification, we can generate long audio sequences while changing the text description over time. Borrowing from Villegas et al. (2022) in the context of video generation, we refer to this approach as *story mode*. Con-

cretely, we compute M_T from multiple text descriptions and change the conditioning signal every 15 seconds. The model generates smooth transitions which are tempo consistent and semantically plausible, while changing music context according to the text description.

7. Conclusions

We introduce MusicLM, a text-conditioned generative model that produces high-quality music at 24 kHz, consistent over several minutes, while being faithful to the text conditioning signal. We demonstrate that our method outperforms baselines on MusicCaps, a hand-curated, high-quality dataset of 5.5k music-text pairs prepared by musicians.

Some limitations of our method are inherited from MuLan, in that our model misunderstands negations and does not adhere to precise temporal ordering described in the text. Moreover, further improvements of our quantitative evaluations are needed. Specifically, since MCC also relies on MuLan, the MCC scores are favorable to our method.

Future work may focus on lyrics generation, along with improvement of text conditioning and vocal quality. Another aspect is the modeling of high-level song structure like introduction, verse, and chorus. Modeling the music at a higher sample rate is an additional goal.

8. Broader Impact

MusicLM generates high-quality music based on a text description, and thus it further extends the set of tools that assist humans with creative music tasks. However, there are several risks associated with our model and the use-case it tackles. The generated samples will reflect the biases present in the training data, raising the question about appropriateness for music generation for cultures underrepresented in the training data, while at the same time also raising concerns about cultural appropriation.

We acknowledge the risk of potential misappropriation of creative content associated to the use-case. In accordance with responsible model development practices, we conducted a thorough study of memorization, adapting and extending a methodology used in the context of text-based LLMs, focusing on the semantic modeling stage. We found that only a tiny fraction of examples was memorized exactly, while for 1% of the examples we could identify an approximate match. We strongly emphasize the need for more future work in tackling these risks associated to music generation — we have no plans to release models at this point.

References

- Abu-El-Haija, S., Kothari, N., Lee, J., Natsev, P., Toderici, G., Varadarajan, B., and Vijayanarasimhan, S. Youtube8m: A large-scale video classification benchmark. *arXiv:1609.08675*, 2016.
- Borsos, Z., Marinier, R., Vincent, D., Kharitonov, E., Pietquin, O., Sharifi, M., Teboul, O., Grangier, D., Tagliasacchi, M., and Zeghidour, N. Audioldm: a language modeling approach to audio generation. *arXiv:2209.03143*, 2022.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., Oprea, A., and Raffel, C. Extracting training data from large language models, 2020. URL <https://arxiv.org/abs/2012.07805>.
- Carlini, N., Ippolito, D., Jagielski, M., Lee, K., Tramer, F., and Zhang, C. Quantifying memorization across neural language models, 2022. URL <https://arxiv.org/abs/2202.07646>.
- Chung, Y., Zhang, Y., Han, W., Chiu, C., Qin, J., Pang, R., and Wu, Y. W2v-bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training. *arXiv:2108.06209*, 2021.
- Cohen, A. D., Roberts, A., Molina, A., Butryna, A., Jin, A., Kulshreshtha, A., Hutchinson, B., Zevenbergen, B., Aguera-Arcas, B. H., ching Chang, C., Cui, C., Du, C., Adiwardana, D. D. F., Chen, D., Lepikhin, D. D., Chi, E. H., Hoffman-John, E., Cheng, H.-T., Lee, H., Krivokon, I., Qin, J., Hall, J., Fenton, J., Soraker, J., Meier-Hellstern, K., Olson, K., Aroyo, L. M., Bosma, M. P., Pickett, M. J., Menegali, M. A., Croak, M., Díaz, M., Lamm, M., Krikun, M., Morris, M. R., Shazeer, N., Le, Q. V., Bernstein, R., Rajakumar, R., Kurzweil, R., Thoppilan, R., Zheng, S., Bos, T., Duke, T., Doshi, T., Zhao, V. Y., Prabhakaran, V., Rusch, W., Li, Y., Huang, Y., Zhou, Y., Xu, Y., and Chen, Z. Lamda: Language models for dialog applications. *arXiv:2201.08239*, 2022.
- Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2013.
- Defferrard, M., Benzi, K., Vandergheynst, P., and Bresson, X. FMA: A dataset for music analysis. In *International Society for Music Information Retrieval Conference (ISMIR)*, 2017.
- Devlin, J., Chang, M., Lee, K., and Toutanova, K. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.
- Dhariwal, P., Jun, H., Payne, C., Kim, J. W., Radford, A., and Sutskever, I. Jukebox: A generative model for music. *arXiv:2005.00341*, 2020.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- Défossez, A., Copet, J., Synnaeve, G., and Adi, Y. High fidelity neural audio compression. *arXiv:2210.13438*, 2022.
- Engel, J. H., Hantrakul, L., Gu, C., and Roberts, A. DDSF: differentiable digital signal processing. In *International Conference on Learning Representations (ICLR)*, 2020.
- Esser, P., Rombach, R., and Ommer, B. Taming transformers for high-resolution image synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- Forsgren, S. and Martiros, H. Riffusion - Stable diffusion for real-time music generation, 2022. URL <https://riffusion.com/about>.
- Gemmeke, J. F., Ellis, D. P., Freedman, D., Jansen, A., Lawrence, W., Moore, R. C., Plakal, M., and Ritter, M. Audio set: An ontology and human-labeled dataset for audio events. In *IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2017.
- Hawthorne, C., Stasyuk, A., Roberts, A., Simon, I., Huang, C. A., Dieleman, S., Elsen, E., Engel, J. H., and Eck, D. Enabling factorized piano music modeling and generation with the MAESTRO dataset. In *International Conference on Learning Representations (ICLR)*, 2019.
- Hawthorne, C., Jaegle, A., Cangea, C., Borgeaud, S., Nash, C., Malinowski, M., Dieleman, S., Vinyals, O., Botvinick, M. M., Simon, I., Sheahan, H., Zeghidour, N., Alayrac, J., Carreira, J., and Engel, J. H. General-purpose, long-context autoregressive modeling with perceiver AR. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvári, C., Niu, G., and Sabato, S. (eds.), *International Conference on Machine Learning (ICML)*, 2022a.

- Hawthorne, C., Simon, I., Roberts, A., Zeghidour, N., Gardner, J., Manilow, E., and Engel, J. H. Multi-instrument music synthesis with spectrogram diffusion. *arXiv:2206.05408*, 2022b.
- Hershey, S., Chaudhuri, S., Ellis, D. P. W., Gemmeke, J. F., Jansen, A., Moore, C., Plakal, M., Platt, D., Saurous, R. A., Seybold, B., Slaney, M., Weiss, R., and Wilson, K. Cnn architectures for large-scale audio classification. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Ho, J., Salimans, T., Gritsenko, A., Chan, W., Norouzi, M., and Fleet, D. J. Video diffusion models. *arXiv:2204.03458*, 2022.
- Hong, W., Ding, M., Zheng, W., Liu, X., and Tang, J. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv:2205.15868*, 2022.
- Huang, C. A., Vaswani, A., Uszkoreit, J., Simon, I., Hawthorne, C., Shazeer, N., Dai, A. M., Hoffman, M. D., Dinculescu, M., and Eck, D. Music transformer: Generating music with long-term structure. In *International Conference on Learning Representations (ICLR)*, 2019.
- Huang, Q., Jansen, A., Lee, J., Ganti, R., Li, J. Y., and Ellis, D. P. W. Mulan: A joint embedding of music audio and natural language. In *International Society for Music Information Retrieval Conference (ISMIR)*, 2022.
- Kilgour, K., Zuluaga, M., Roblek, D., and Sharifi, M. Fréchet audio distance: A reference-free metric for evaluating music enhancement algorithms. In *INTERSPEECH*, 2019.
- Kim, C. D., Kim, B., Lee, H., and Kim, G. Audiocaps: Generating captions for audios in the wild. In *NAACL-HLT*, 2019.
- Kreuk, F., Synnaeve, G., Polyak, A., Singer, U., Défossez, A., Copet, J., Parikh, D., Taigman, Y., and Adi, Y. Audiogen: Textually guided audio generation, 2022.
- Mubert-Inc. Mubert. <https://mubert.com/>, <https://github.com/MubertAI/Mubert-Text-to-Music>, 2022.
- Nichol, A. Q., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., and Chen, M. GLIDE: towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning (ICML)*, 2022.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P. J., et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research (JMLR)*, 2020.
- Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., and Sutskever, I. Zero-shot text-to-image generation. In Meila, M. and Zhang, T. (eds.), *International Conference on Machine Learning (ICML)*, 2021.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. Hierarchical text-conditional image generation with clip latents. *arXiv:2204.06125*, 2022.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, June 2022a.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022b.
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E. L., Ghasemipour, S. K. S., Ayan, B. K., Mahdavi, S. S., Lopes, R. G., Salimans, T., Ho, J., Fleet, D. J., and Norouzi, M. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv:2205.11487*, 2022.
- Schroff, F., Kalenichenko, D., and Philbin, J. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2015.
- Shor, J., Jansen, A., Maor, R., Lang, O., Tual, O., de Chaumont Quiry, F., Tagliasacchi, M., Shavitt, I., Emanuel, D., and Haviv, Y. Towards Learning a Universal Non-Semantic Representation of Speech. In *INTERSPEECH*, 2020.
- van den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A. W., and Kavukcuoglu, K. Wavenet: A generative model for raw audio. In *ISCA*, 2016.
- Van Den Oord, A., Vinyals, O., et al. Neural discrete representation learning. *Advances in neural information processing systems (NeurIPS)*, 2017.

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems (NeurIPS)*, 2017.
- Villegas, R., Babaeizadeh, M., Kindermans, P.-J., Moraldo, H., Zhang, H., Saffar, M. T., Castro, S., Kunze, J., and Erhan, D. Phenaki: Variable length video generation from open domain textual description. *arXiv:2210.02399*, 2022.
- Wu, C., Liang, J., Ji, L., Yang, F., Fang, Y., Jiang, D., and Duan, N. Nüwa: Visual synthesis pre-training for neural visual world creation. In *European Conference on Computer Vision (ECCV)*, 2022a.
- Wu, H., Seetharaman, P., Kumar, K., and Bello, J. P. Wav2clip: Learning robust audio representations from CLIP. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022b.
- Yang, D., Yu, J., Wang, H., Wang, W., Weng, C., Zou, Y., and Yu, D. Diffsound: Discrete diffusion model for text-to-sound generation. *arXiv:2207.09983*, 2022.
- Yu, J., Xu, Y., Koh, J. Y., Luong, T., Baid, G., Wang, Z., Vasudevan, V., Ku, A., Yang, Y., Ayan, B. K., Hutchinson, B., Han, W., Parekh, Z., Li, X., Zhang, H., Baldridge, J., and Wu, Y. Scaling autoregressive models for content-rich text-to-image generation, 2022.
- Zeghidour, N., Teboul, O., de Chaumont Quitry, F., and Tagliasacchi, M. LEAF: A learnable frontend for audio classification. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=jM76BCb6F9m>.
- Zeghidour, N., Luebs, A., Omran, A., Skoglund, J., and Tagliasacchi, M. Soundstream: An end-to-end neural audio codec. *IEEE ACM Trans. Audio Speech Lang. Process.*, 30, 2022.
- Zen, H., Senior, A., and Schuster, M. Statistical parametric speech synthesis using deep neural networks. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2013.

A. MusicCaps Dataset

Together with this paper, we release MusicCaps, a high-quality music caption dataset.⁶ This dataset includes music clips from AudioSet (Gemmeke et al., 2017), paired with corresponding text descriptions in English. It contains a total of 5,521 examples, out of which 2,858 are from the AudioSet eval and 2,663 from the AudioSet train split. We further tag 1,000 examples as a balanced subset of our dataset, which is balanced with respect to the genres of the music contained. All examples in the balanced subset are from the AudioSet eval split.

Examples of free text captions:

- *“This folk song features a male voice singing the main melody in an emotional mood. This is accompanied by an accordion playing fills in the background. A violin plays a droning melody. There is no percussion in this song. This song can be played at a Central Asian classical concert.”*
- *“This is a live recording of a keyboardist playing a twelve bar blues progression on an electric keyboard. The player adds embellishments between chord changes and the piece sounds groovy, bluesy and soulful.”*
- *“A synth is playing an arpeggio pluck with a lot of reverb rising and falling in velocity. Another synth sound is playing pads and a sub bassline. This song is full of synth sounds creating a soothing and adventurous atmosphere. This song may be playing at a festival during two songs for a buildup.”*
- *“A low sounding male voice is rapping over a fast paced drums playing a reggaeton beat along with a bass. Something like a guitar is playing the melody along. This recording is of poor audio-quality. In the background a laughter can be noticed. This song may be playing in a bar.”*
- *“The electronic music features a section that repeats roughly every two seconds. It consists of a beat that’s made of a kick drum and claps. A buzzing synth sets the pulsation of the music by playing once every two beats. The whole music sounds like a loop being played over and over. Towards the end of the excerpt a crescendo-like buzzing sound can be heard, increasing the tension.”*

Examples of aspect lists:

- *“pop, tinny wide hi hats, mellow piano melody, high pitched female vocal melody, sustained pulsating synth lead, soft female vocal, punchy kick, sustained synth bass, claps, emotional, sad, passionate”*
- *“amateur recording, finger snipping, male mid range voice singing, reverb”*
- *“backing track, jazzy, digital drums, piano, e-bass, trumpet, acoustic guitar, digital keyboard song, medium tempo”*
- *“rubab instrument, repetitive melody on different octaves, no other instruments, plucked string instrument, no voice, instrumental, fast tempo”*
- *“instrumental, white noise, female vocalisation, three unrelated tracks, electric guitar harmony, bass guitar, keyboard harmony, female lead vocalisation, keyboard harmony, slick drumming, boomy bass drops, male voice backup vocalisation”*

⁶kaggle.com/datasets/googleai/musiccaps

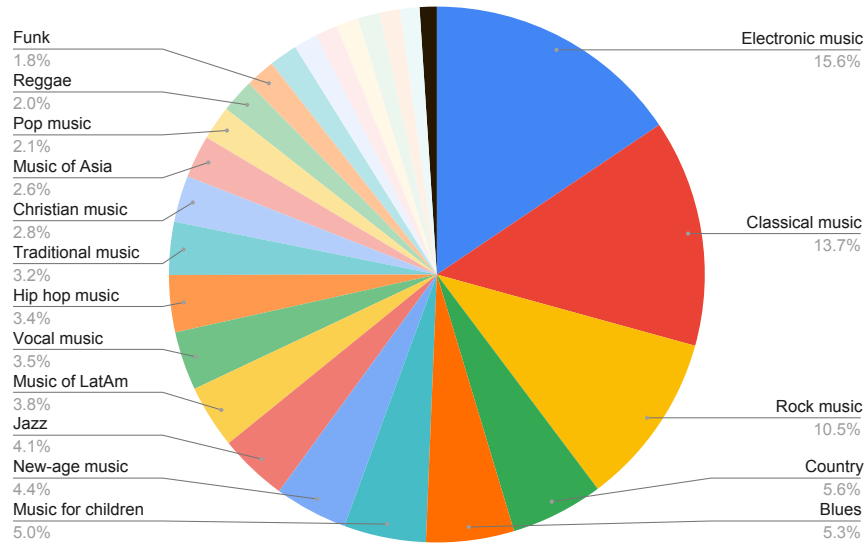


Figure 4. Genre distribution of all 5.5k examples of MusicCaps, according to an AudioSet classifier.

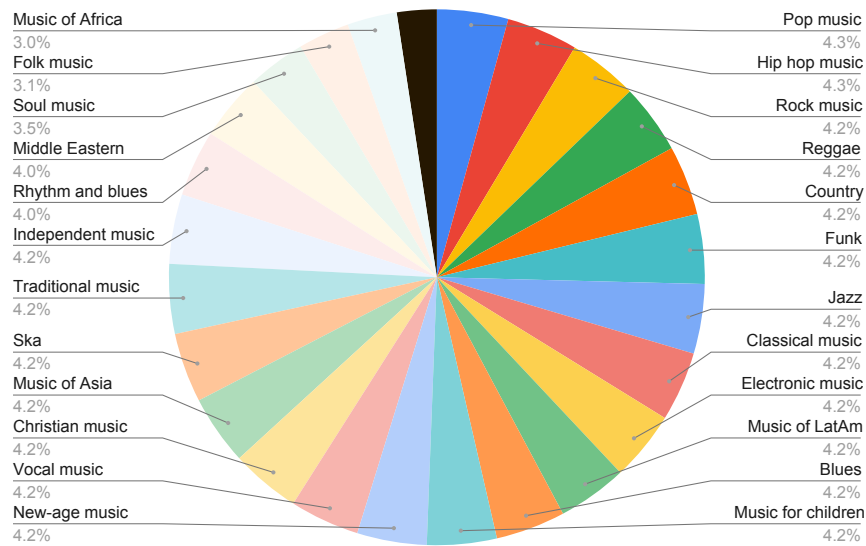


Figure 5. Genre distribution of a balanced 1k example subset of MusicCaps, according to an AudioSet classifier.

B. Qualitative Evaluation

Participants in the listening test were presented with two 10-second clips and a text caption, and asked which clip is best described the text of the caption on a 5-point Likert scale. They were also instructed to ignore audio quality and focus just on how well the text matches the music (similar to MuLan score). Figure 6 shows the user interface presented to raters.

We collected 1200 ratings, with each source involved in 600 pair-wise comparisons. Figures 7 and 8 show the granular results of pairwise comparisons between the models. According to a post-hoc analysis using the Wilcoxon signed-rank test with Bonferroni correction (with $p < 0.01/15$), the orderings shown in Figure 8 from raters are all statistically significant.

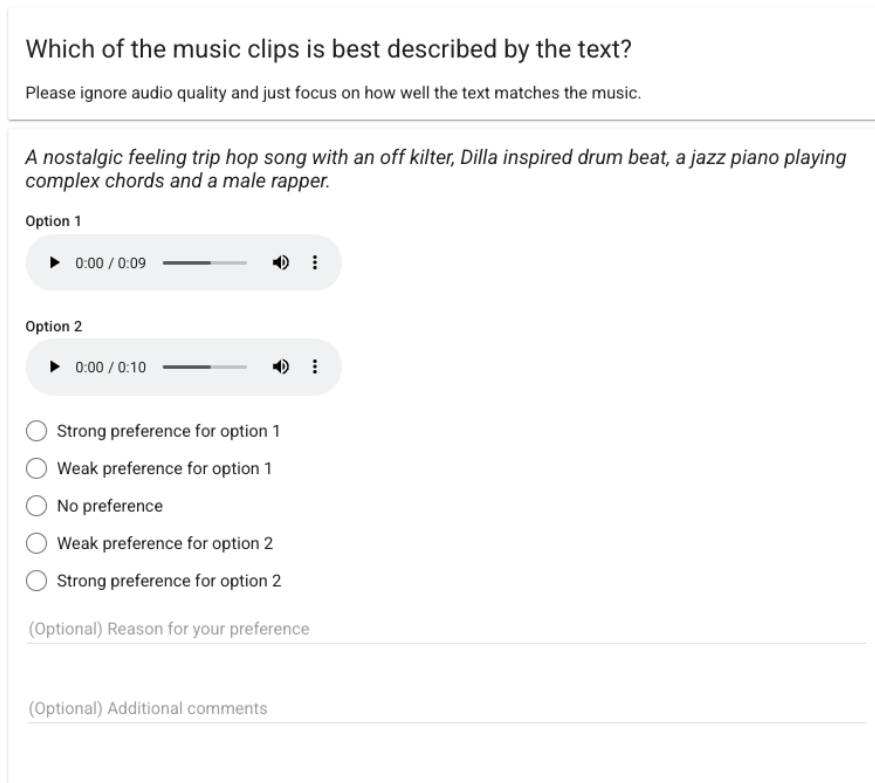


Figure 6. User interface for the human listener study.

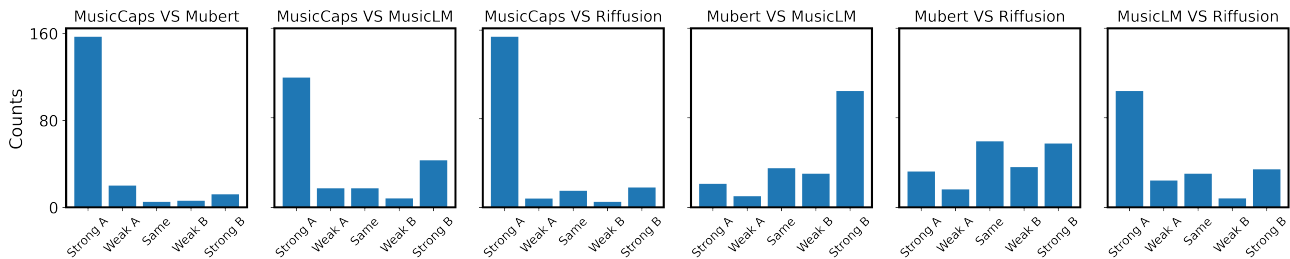


Figure 7. Pairwise comparisons from the human listener study. Each pair is compared on a 5-point Likert scale. Raters had a decisive model preference in all cases except Mubert vs. Riffusion.

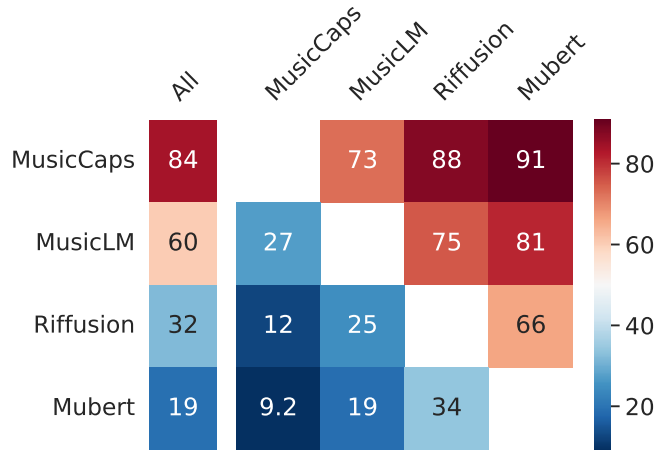


Figure 8. Win percentage from the human listener study. Each row indicates the % of times listeners found the music to better match the caption from that system to those from any other system (first column, $N = 1200$) and each system individually (other columns, $N = 600$). The ground truth data (MusicCaps) clearly is the best match to the captions, but followed closely by MusicLM, which even beats the ground truth in 27% of comparisons.

C. Melody Conditioning

We provide here implementation details of the model used for conditioning the music generation on melody. The model is based on a small ViT (Dosovitskiy et al., 2021) composed of 12 layers, 6 attention heads, embedding dimension of 512 and feed-forward layer of dimension 1024. The input to the model are the temporal frames of the mel spectrogram of the audio. We use semi-hard triplet loss (Schroff et al., 2015) to train the melody embedding model to generate 192 dimensional embeddings for each 4 seconds of audio. The model learns to generate embeddings which are representative of a melody while being invariant to acoustic properties related to the instruments being played. This is particularly advantageous, since this representation is complementary to the representation learned by the MuLan embeddings. Hence, our melody embeddings and the MuLan can be jointly and complementarily used for conditioning the music generation process. During training, we consider input audio with a duration of 10 seconds. We extract three melody embeddings, with a hop length of 3 seconds, discretize each of them to tokens with residual vector quantization (RVQ) and concatenate the resulting token sequences with the MuLan audio tokens M_A . We use an RVQ composed of 24 quantizers, each with a vocabulary size of 512.