

Mousai: Efficient Text-to-Music Diffusion Models

Anonymous ACL submission

Abstract

Recent years have seen the rapid development of large generative models for text; however, much less research has explored the connection between text and another “language” of communication – *music*. Music, much like text, can convey emotions, stories, and ideas, and has its own unique structure and syntax. In our work, we bridge text and music via a text-to-music generation model that is highly efficient, expressive, and can handle long-term structure. Specifically, we develop a cascading latent diffusion approach that can generate multiple minutes of high-quality stereo music at 48kHz from textual descriptions. Moreover, our model features high efficiency, which enables real-time inference on a single consumer GPU with a reasonable speed. Through experiments and property analyses, we show our model’s competence over a variety of criteria compared with existing music generation models. Lastly, to promote the open-source culture, we provide a collection of open-source libraries with the hope of facilitating future work in the field.¹

1 Introduction

In recent years, natural language processing (NLP) has made significant strides in understanding and generating human language, due to the advancements in deep learning and large-scale pre-trained models (Radford et al., 2018; Devlin et al., 2019; Brown et al., 2020). While the majority of NLP research has focused on textual data, there exists another rich and expressive “language” of communication – *music*. Music, much like text, can convey emotions (Germer, 2011), stories (Chung, 2006), and ideas (Bicknell, 2002), and has its own unique structure and syntax (Swain, 1995).

In this paper, we further bridge the gap between text and music by leveraging the power of NLP

¹Our code and data are uploaded to the system, and will be released upon acceptance. Our anonymized music samples are available at <https://bit.ly/anonymous-mousai>.

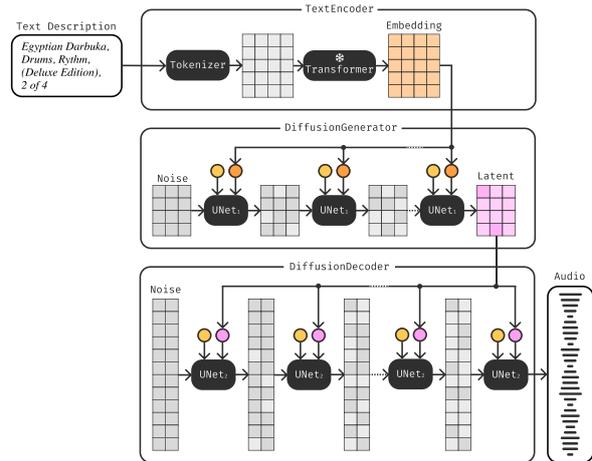


Figure 1: We propose a two-stage cascading diffusion method, where the first stage (the diffusion generator) compresses the music using a novel diffusion autoencoder, and the second stage (the diffusion decoder) generates music from the reduced representation conditioned on the encoding of a textual description.

techniques to generate music conditioned on textual input. Through our work, we not only aim to expand the scope of NLP applications, but also contribute to the interdisciplinary research at the intersection of language, music, and machine learning techniques.

However, like text, music generation has long been a challenging task, as it requires multiple aspects at different levels of abstraction (van den Oord et al., 2016; Dieleman et al., 2018). Existing audio generation models explore the use of recursive neural networks (Mehri et al., 2017), adversarial generative networks (Kumar et al., 2019; Kim et al., 2021; Engel et al., 2019; Morrison et al., 2022), autoencoders (Deng et al., 2021), and transformers (Yu et al., 2022). With the recent advancement in diffusion-based generative models in computer vision (Ramesh et al., 2022; Saharia et al., 2022), researchers in speech have also started to explore the use of diffusion models in tasks such as speech

059 synthesis (Kong et al., 2021; Lam et al., 2022; Leng
060 et al., 2022), although only a few these models can
061 apply well to the task of music generation.

062 Additionally, there are several long-standing chal-
063 lenges in the area of music generation: (1) music
064 generation at length, as most text-to-audio systems
065 (Forsgren and Martiros, 2022; Kreuk et al., 2022)
066 can only generate *a few seconds* of audio; (2) model
067 efficiency, as many need to run on GPUs for hours
068 to generate just one minute of audio (Dhariwal
069 et al., 2020; Kreuk et al., 2022); (3) lack of diver-
070 sity of the generated music, as many are limited by
071 their training methods taking in a single modality
072 (resulting in the ability to handle only single-genre
073 music, but *not diverse* genres) (Caillon and Esling,
074 2021; Pasini and Schlüter, 2022); and (4) easy con-
075 trollability by text prompts, as most are only con-
076 trolled by latent states (Caillon and Esling, 2021;
077 Pasini and Schlüter, 2022), the starting snippet of
078 the music (Borsos et al., 2022), or text but are lyrics
079 (Dhariwal et al., 2020) or descriptions of a daily
080 sound like dog barking (Kreuk et al., 2022).

081 A single model mastering all these aspects would
082 make a strong contribution to the music industry,
083 as it can enable the broader public to be part of
084 the creative process by allowing them to compose
085 music using an accessible text-based interface, as-
086 sist creators in finding inspiration, and provide an
087 unlimited supply of novel audio samples.

088 To address these challenges, we propose *Moûsai*,²
089 a novel text-conditional cascading diffusion model
090 illustrated in Figure 1. To ensure model efficiency,
091 our diffusion magnitude autoencoder can achieve
092 an audio signal compression rate of 64x. Together
093 with our design of a lightweight and specialized
094 1D U-Net architecture, our model enables a fast
095 inference speed on a single consumer GPU in min-
096 utes, and a training time of approximately one week
097 per stage on one A100 GPU, making it possible
098 to train and run the overall system using resources
099 available in most universities.

100 Remarkably, our diffusion-based model improves
101 significantly on previous models, as it can train on
102 a *variety* of music genres, generate *long-context*
103 music for several minutes with a *high quality of*
104 *48kHz* stereo music, runs real-time inference effi-
105 ciently within minutes, and can be easily controlled

²*Moûsai* is romanized ancient Greek for *Muses*, the sources
of artistic inspiration ([https://en.wikipedia.org/wiki/
Muses](https://en.wikipedia.org/wiki/Muses)), and also evokes a blend of *music* and *AI*.

106 by text. Our extensive evaluations on 11 criteria
107 also validate the quality of the generated music by
108 our model from multiple perspectives.

2 Related Work 109

Connecting Text and Music The connection be-
110 tween text and music lies in the intersection of NLP
111 and computational musicology. Previous work
112 looks into aspects such as the similarity of mu-
113 sic and linguistic structures (Papadimitriou and Ju-
114 rafsky, 2020), music and dialog (Berlingerio and
115 Bonin, 2018), and jointly modeling music and text
116 for emotion detection (Mihalcea and Strapparava,
117 2012). Apart from several work that generates mu-
118 sic from text (Dhariwal et al., 2020; Forsgren and
119 Martiros, 2022), we are the first to explore diffusion
120 models to interact text with music representations.
121

Generative Models Generative models aim to
122 learn a lower-dimension representation space, and
123 then reconstruct to the high-dimension space con-
124 ditioning on the given information (Rombach et al.,
125 2022; Yang et al., 2022; Kreuk et al., 2022; Ho
126 et al., 2022). Some effective methods earlier in-
127 clude auto-encoding (Hinton and Salakhutdinov,
128 2006; Kingma and Welling, 2014), or quantized
129 auto-encoding (van den Oord et al., 2017; Esser
130 et al., 2021; Lee et al., 2022). Recent proposals
131 focus on the quantized representation followed by
132 masked or autoregressive learning on tokens (Ville-
133 gas et al., 2022; Dhariwal et al., 2020; Kreuk et al.,
134 2022), and diffusion models (Ramesh et al., 2022;
135 Rombach et al., 2022; Saharia et al., 2022), which
136 leads to impressive performance. To the best of our
137 knowledge, we are the first to adapt the cascading
138 diffusion approach for audio generation.
139

Concurrent Work Upon the completion of our
140 work in Jan 2023, there came several powerful
141 generative music models, all led by large industry
142 labs (Agostinelli et al., 2023; Huang et al., 2023;
143 Copet et al., 2023). We do not include them in the
144 paper, as they count as concurrent work in the same
145 time or several months after our work, and also our
146 work is done in a university setting which cannot
147 compare with the performance of these large-scale
148 models supported by industry-level resources.
149

3 *Moûsai*: Efficient Long-Context Music Generation from Text 150

151 Our model *Moûsai* contains a two-stage training
152 process. In Stage 1, we use diffusion magnitude-
153

154 autoencoding (DMAE), which compresses the au-
 155 dio waveform 64x using a diffusion autoencoder.
 156 In Stage 2, we use a latent text-to-audio diffusion
 157 model, to generate a novel latent space by diffusion
 158 while conditioning on text embeddings obtained
 159 from a frozen transformer language model.

160 In the following, we first introduce the basic mod-
 161 ules of our models, and details of the two stages.

162 3.1 Modules

163 3.1.1 Latent Diffusion for Audio

164 ***v*-Objective Diffusion Process** We use the *v*-
 165 objective diffusion process as proposed by [Salimans and Ho \(2022\)](#). Suppose we have a sample
 166 \mathbf{x}_0 from a distribution $p(\mathbf{x}_0)$, some noise sched-
 167 ule $\sigma_t \in [0, 1]$, and some noisy data point $\mathbf{x}_{\sigma_t} =$
 168 $\alpha_{\sigma_t}\mathbf{x}_0 + \beta_{\sigma_t}\epsilon$. The *v*-objective diffusion tries to
 169 estimate a model $\hat{\mathbf{v}}_{\sigma_t} = f(\mathbf{x}_{\sigma_t}, \sigma_t)$ by minimizing
 170 the following objective:
 171

$$172 \mathbb{E}_{t \sim [0,1], \sigma_t, \mathbf{x}_{\sigma_t}} [\|f_{\theta}(\mathbf{x}_{\sigma_t}, \sigma_t) - \mathbf{v}_{\sigma_t}\|_2^2], \quad (1)$$

173 where $\mathbf{v}_{\sigma_t} = \frac{\partial \mathbf{x}_{\sigma_t}}{\partial \sigma_t} = \alpha_{\sigma_t}\epsilon - \beta_{\sigma_t}\mathbf{x}_0$, for which
 174 we define $\phi_t := \frac{\pi}{2}\sigma_t$, and obtain its trigonometric
 175 values $\alpha_{\sigma_t} := \cos(\phi_t)$, and $\beta_{\sigma_t} := \sin(\phi_t)$.

176 **DDIM Sampler for Denoising** The denoising step
 177 uses ODE samplers to turn noise into a new data
 178 point by estimating the rate of change. In this work,
 179 we adopt the DDIM sampler ([Song et al., 2021](#)),
 180 which we find to work well and have a reasonable
 181 tradeoff between the number of steps and audio
 182 quality. The DDIM sampler denoises the signal by
 183 repeated application of the following:

$$184 \hat{\mathbf{v}}_{\sigma_t} = f_{\theta}(\mathbf{x}_{\sigma_t}, \sigma_t) \quad (2)$$

$$185 \hat{\mathbf{x}}_0 = \alpha_{\sigma_t}\mathbf{x}_{\sigma_t} - \beta_{\sigma_t}\hat{\mathbf{v}}_{\sigma_t} \quad (3)$$

$$186 \hat{\epsilon}_{\sigma_t} = \beta_{\sigma_t}\mathbf{x}_{\sigma_t} + \alpha_{\sigma_t}\hat{\mathbf{v}}_{\sigma_t} \quad (4)$$

$$187 \hat{\mathbf{x}}_{\sigma_{t-1}} = \alpha_{\sigma_{t-1}}\hat{\mathbf{x}}_0 + \beta_{\sigma_{t-1}}\hat{\epsilon}_{\sigma_t}, \quad (5)$$

188 which estimates both the initial data point and the
 189 noise at the step σ_t , for some T -step noise schedule
 190 $\sigma_T, \dots, \sigma_0$ as a sequence evenly spaced between 1
 191 and 0.

192 **Diffusion Autoencoder for Audio Input** We pro-
 193 pose a new diffusion autoencoder that first encodes
 194 a magnitude spectrogram into a compressed rep-
 195 resentation, and later injects the latent into inter-
 196 mediate channels of the decoding modules. The
 197 standard method to do diffusion, such as the image

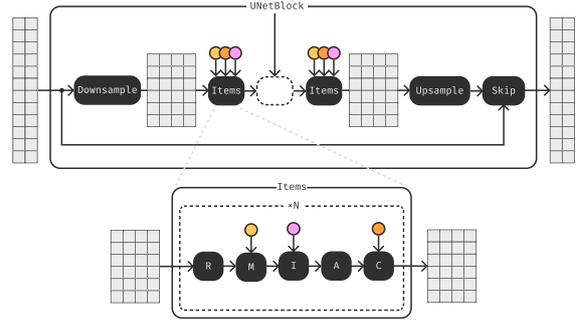


Figure 2: Our proposed 1D U-Net architecture. Each UNetBlock (top) consists of several U-Net items (bottom). In each U-Net item (bottom), we use a 1D convolutional ResNet (R), a modulation unit (M) to provide the diffusion noise level as a feature vector conditioning (●), an inject item (I) to inject external channels as conditioning (○), an attention item (A) to share time-wise information, and a cross-attention item (C) to condition on an external (text) embedding (●). Moreover, for the UNetBlocks, we can recursively nest them, which we indicate by the inner dashed region on the top.

198 diffusion model ([Rombach et al., 2022](#)), is to compress
 199 the input into a lower-dimensional representa-
 200 tion space and apply the diffusion process on the
 201 reduced latent space. We further compress and
 202 enhance the representation space by diffusion-based
 203 autoencoding ([Preechakul et al., 2022](#)), which is
 204 first introduced in computer vision, as a way to
 205 condition the diffusion process on a compressed latent
 206 vector of the input itself. Since diffusion serves as
 207 a more powerful generative decoder, and hence the
 208 input can be reduced to latent representations with
 209 higher compression ratios.

210 3.1.2 Efficient and Enriched 1D U-Net

211 Another crucial module in our model is the effi-
 212 cient 1D U-Net that we design. We identify that
 213 the vanilla U-Net architecture [Ronneberger et al. \(2015\)](#),
 214 originally introduced for medial image seg-
 215 mentation, has relatively limited efficiency and
 216 speed, as it uses an hourglass convolutional-only
 217 2D architecture with skip connections.

218 Hence, we propose a novel U-Net with only 1D
 219 convolutional kernels, which is more efficient than
 220 the original 2D architecture in terms of speed, and
 221 can be successfully used both on waveforms or on
 222 spectrograms if each frequency is considered as a
 223 different channel.

224 Moreover, we infuse our 1D U-Net with multi-
 225 ple new components, as illustrated in Figure 2: a
 226 ResNet residual 1D convolutional unit, a modula-

tion unit to alter the channels given features from the diffusion noise level, an inject item to concatenate external channels to the ones at the current depth, an attention item to share long-context structural information, and a cross-attention item to condition on the text embeddings. Note that inject items are applied only at a specific depth in the decoder in the first stage to condition on the latent representation of the music. Additionally, since attention and cross-attention items are for learning the structure and conditioning on text, we only use them for the second stage, text-conditioned music generation.

In summary, our novel 1D U-Net features more modern convolutional blocks, a variety of attention blocks, conditioning blocks, and improved skip connections, maintaining an efficient skeleton of the hourglass architecture.

3.2 Stage 1: Music Encoder by Diffusion Magnitude-Autoencoding (DMAE)

We design the first step of *Moûsai* to be learning a good music encoder to capture the latent representation space for music. Representation learning is crucial for generative models, as it can be drastically more efficient than handling the high-dimensional raw input data (Rombach et al., 2022; Yang et al., 2022; Kreuk et al., 2022; Ho et al., 2022; Villegas et al., 2022).

Overview To learn the representation space for music, we deploy a diffusion magnitude autoencoder (DMAE) shown in Figure 3. Specifically, we adopt our diffusion-based audio autoencoder, introduced in Section 3.1.1, to compress audio into a smaller latent space by 64x from the original waveform. To train the model, we first convert the waveform to a magnitude spectrogram, which is a better representation for audio models, and then we auto-encode it into a latent representation.

At the same time, we corrupt the original audio with a random amount of noise, and train our 1D U-Net (introduced in Section 3.1.2) to remove that noise. During the noise removal process, we condition the U-Net on the noise level and the compressed latent, which can have access to a reduced version of the non-noisy audio.

Model Architecture Our DMAE works as follows. Let \mathbf{w} be a waveform of shape $[c, t]$ for c channels and t timesteps, and $(\mathbf{m}_w, \mathbf{p}_w) = \text{stft}(\mathbf{w}; n =$

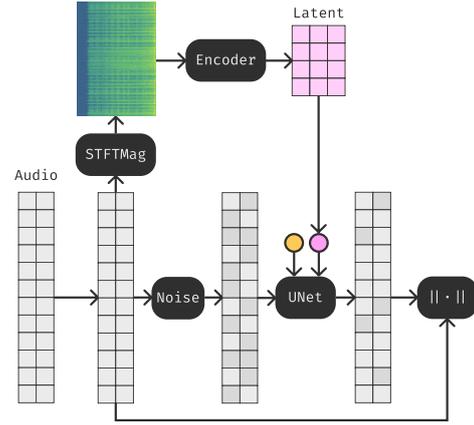


Figure 3: The training scheme of our diffusion magnitude autoencoder (DMAE). When denoising (bottom right), we condition the U-Net on the noise level (●) and compressed latent representation (●) from a reduced version of the non-noisy audio (the pink matrix).

1024, $h = 256$) be the magnitude and phase obtained from a short-time Fourier transform of the waveform with a window size of 1024 and hop-length of 256. Then the resulting spectrograms will have shape $[c \cdot n, \frac{t}{h}]$. We discard phase and encode the magnitude into a latent $\mathbf{z} = \mathcal{E}_{\theta_e}(\mathbf{m}_w)$ using a 1D convolutional encoder. The original waveform is then reconstructed by decoding the latent using a diffusion model $\hat{\mathbf{w}} = \mathcal{D}_{\theta_d}(\mathbf{z}, \epsilon, s)$, where \mathcal{D}_{θ_d} is the diffusion sampling process with starting noise ϵ and s is the number of decoding (sampling) steps. The decoder is trained with \mathbf{v} -objective diffusion while conditioning on the latent $f_{\theta_d}(\mathbf{w}_{\sigma_t}; \sigma_t, \mathbf{z})$, where f_{θ_d} is the proposed 1D U-Net, called repeatedly during decoding.

Since only the magnitude is used and phase is discarded, this diffusion autoencoder is simultaneously a compressing autoencoder and vocoder. By using the magnitude spectrograms, higher compression ratios can be obtained than autoencoding directly the waveform. We found that waveforms are less compressible and efficient to work with. Similarly, discarding phase is beneficial to obtain higher compression ratios for the same level of quality. The diffusion model can easily learn to generate a waveform with realistic phase even if conditioned only on the encoded magnitude.

In this way, the latent space for music can serve as the starting point for our text-to-music generator, which will be introduced next. To ensure this representation space fits the next stage, we apply a tanh function on the bottleneck, keeping the values in the range $[-1, 1]$. Note that we do not use

a more disentangled bottleneck, such as the one in VAEs (Kingma and Welling, 2014), as its additional regularization reduces the amount of allowed compressibility.

3.3 Stage 2: Text-to-Music Generator by Text-Conditioned Latent Diffusion (TCLD)

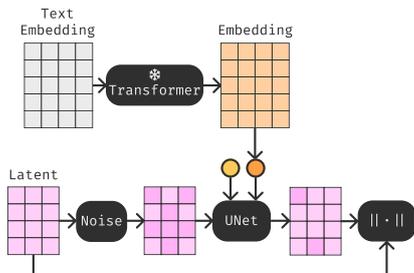


Figure 4: The training scheme of our text-conditioned latent diffusion (TCLD) generator. During the denoising process, we provide the U-Net a feature vector (●) and a text embedding (○).

Based on the learned music representation space, in this stage, we guide the music generation with text descriptions.

Overview We first encode the music source into the latent space using the DMAE encoder, and then we propose a text-conditioned latent diffusion (TCLD) which corrupt the latent space with a random amount of noise, and then train a series of U-Nets to remove the noise.

We illustrate the detailed process in Figure 4. Consistent with the previous stage, we use v -objective diffusion and the 1D U-Net architecture. To condition on the text embedding e , we use the U-Net configuration $f_{\theta_g}(z_{\sigma_t}; \sigma_t, e)$ to generate the compressed latent $z = \mathcal{E}_{\theta_e}(m_w)$. Then, the generator $\mathcal{G}_{\theta_g}(e, \epsilon, s)$ applies DDIM sampling and calls the U-Net s times to generate an approximate latent \hat{z} from the text embedding e and starting noise ϵ . The final generation stack during inference to obtain a waveform is

$$\hat{w} = \mathcal{D}_{\theta_d}(\mathcal{G}_{\theta_g}(e, \epsilon_g, s_g), \epsilon_d, s_d). \quad (6)$$

Text Conditioning To obtain the text embeddings, prior work on text-conditioning suggests either learning a joint data-text representation (Li et al., 2022; Elizalde et al., 2022; Ramesh et al., 2022), or using embeddings from pre-trained language model as direct conditioning (Saharia et al., 2022; Ho et al., 2022) of the latent model.

In our TCLD model, we follow the practice in Saharia et al. (2022) to use a pre-trained and frozen T5 language model (Raffel et al., 2020) to generate text embeddings from the given description. We use the classifier-free guidance (CFG) (Ho and Salimans, 2022) with a learned mask applied on batch elements with a probability of 0.1 to improve the strength of the text-embedding during inference.

Adapting the U-Net for Text Conditioning To enable the U-Net to condition on the text embedding e , we use the U-Net with the cross-attention blocks, which provide the conditioning text embedding, and multiple attention blocks, to ensure information sharing over the entire latent space, which is crucial to learn long-range audio structure. Given the compressed size of the latent space, we also increase the size of this inner U-Net to be larger than the first stage. And due to our efficiency design, it maintains a reasonable training and inference speed, even with large parameter counts.

4 Experimental Setup

4.1 Collection of the TEXT2MUSIC Dataset

To provide a fertile ground to train our text-to-music model on, we collect a new dataset, TEXT2MUSIC, which consists of 50K text-music pairs totaling 2,500 hours. We ensure a high quality of stereo music sampled at 48kHz and cover a wide variety of music spanning multiple genres, artists, instruments, and provenience. Many existing open-source music datasets, such as (Gillick et al., 2019; Hawthorne et al., 2019a), have limitations in terms of the specific musical instruments they encompass. While some datasets, like (Engel et al., 2017; Boulanger-Lewandowski et al., 2012), cover a broader array of instruments, they fall short in representing a wide variety of genres. This inadequacy underscores the need for a more comprehensive dataset that encompasses a rich tapestry of musical genres and diverse instrumentation.

As for the procedure to collect the music, we first check with the copyright regulations, which grants an exemption for using copyright infringing copies if the purpose is scientific research (Geiger et al., 2018; Delacroix, 2023), according to the EU regulation in Article 3 of the EU Directive on Copyright in the Digital Single Market (European Commission, 2016). Then, we follow Spotify’s top recommendations to collect seven very large playlists, each containing on average 7K pieces of music.

Genre	# Pieces	Percentage (%) in Dataset
Pop	5,498	27.29
Electronic	3,875	19.38
Rock	3,584	17.79
Metal	1,796	8.92
Hip Hop	818	4.06
Others	4492	22.56

Table 1: Our TEXT2MUSIC dataset covers a variety of music, including pop, electronic, rock, metal, hip pop, and others.

We iterate through every music sample in these playlists, for which we use the name of the music to search and download the music from YouTube, and we use the metadata to compose its corresponding text description, which contains the music title, author, album name, genre, and year of release.

In line with our spirit to open-source the model, we also open-source the data collection pipeline on GitHub,³ so future researchers can use it to facilitate new data collection.

We show the statistics about the diverse set of genres in our TEXT2MUSIC dataset in Table 1.

4.2 Implementation Details

Our diffusion autoencoder has 185M parameters, and text-conditional generator has 857M parameters, with more architecture details in Appendix A.3. We train the music autoencoder on random crops of length 2^{18} (~ 5.5 s at 48kHz), and the text-conditional diffusion generation model on fixed crops of length 2^{21} (~ 44 s at 48kHz) encoded in the 32-channels, 64x compressed latent representation. We use the AdamW optimizer (Loshchilov and Hutter, 2019) with a learning rate of 10^{-4} , β_1 of 0.95, β_2 of 0.999, ϵ of 10^{-6} , and weight decay of 10^{-3} . And we use an exponential moving average (EMA) with $\beta = 0.995$ and power of 0.7.

5 Evaluation

5.1 Assessment Criteria Overview

Evaluating music is a highly challenging task. We survey a large number of papers, and find that previous work adopts a variety of objective and subjective metrics,⁴ and the gist is that *no single metric is*

³Anonymous link. We will release it upon acceptance.

⁴The common metrics we surveyed include quality (Goel et al., 2022), fidelity (Goel et al., 2022; Hawthorne et al., 2019b; Hyun et al., 2022), musicality (Goel et al., 2022; Yu et al., 2022; Dhariwal et al., 2020), diversity (Goel et al., 2022; Dhariwal et al., 2020), and structure (Yu et al., 2022; Leng et al., 2022; Dhariwal et al., 2020).

perfect. After careful thinking, we design a comprehensive set of evaluation metrics covering three categories with a total of *11 metrics*, including both automatic and human evaluations. In the following, we will introduce the overall property analysis (Section 5.2), such as sample rate, efficiency and music type; text-music relevance (Section 5.3); music quality (Section 5.4); and long-term structure of the music (Section 5.5).

5.2 Property Analysis

Comparing the overall properties of various models in Table 2, we see a set of impressive properties of the *Moûsai* model: (1) We are among the very few that can control music generation easily by **text descriptions** of the type of music we want, as most other models do not take text as input (van den Oord et al., 2016; Caillon and Esling, 2021; Borsos et al., 2022), or take only lyrics or descriptions of daily sounds (e.g., “a dog barking”) (Kreuk et al., 2022; Dhariwal et al., 2020). The only other text-to-music model is the Riffusion model (Forsgren and Martiros, 2022), which only works with very short length of 5 seconds.

(2) Our model is also among the very few that enables **long-context** music generation for several minutes, among all others that can only generate seconds (van den Oord et al., 2016; Forsgren and Martiros, 2022; Kreuk et al., 2022; Pasini and Schlüter, 2022), except for Jukebox (Dhariwal et al., 2020) which generates songs given lyrics and takes very long to run inference.

(3) **Efficiency** is another highlight of our model, where we only need an inference time similar to the audio length on a consumer GPU, which is several minutes, while many other text-to-audio models take many GPU hours (Dhariwal et al., 2020; Kreuk et al., 2022). Our model is very friendly for research at university labs, as each model can be trained on a single A100 GPU in 1 week of training using a batch size of 32.

(4) Moreover, we also highlight the **diversity** of music we generate, as our model design enables multi-genre music training, instead of single-genre ones in previous models (Caillon and Esling, 2021; Pasini and Schlüter, 2022), and we can find rhythm, loops, riffs, and occasionally even entire choruses in our generated music.

Model	Sample Rate [↑]	Len. [↑]	Input (Text ✓)	Music (Diverse [↑])	Example	Infer. Time [↓]	Data
WaveNet (2016)	16kHz@1	Secs	None	Piano or speech	Piano	= Audio len.*	260
Jukebox (2020)	44.1kHz@1	Mins*	Lyrics, author, etc.	Song with the lyrics	Song	Hours	70K
RAVE (2021)	48kHz@2	Secs*	Latent	Single-genre Music	Strings	= Audio len.*	100
AudioLM (2022)	16kHz@1	Secs*	Beginning of the music	Piano or speech	Piano	Mins	40K
Musika (2022)	22.5kHz@2	Secs	Context vector	Single-genre Music	Piano	= Audio len.*	1K
Riffusion (2022)	44.1kHz@1	5s	Text (genre, author, etc.)	Music of any genre	Jazzy clarinet	Mins	–
AudioGen (2022)	16kHz@1	Secs*	Text (a phrase/sentence)	Daily sounds	Dog barks	Hours	4K
Moûsai (Ours)	48kHz@2	Mins*	Text (genre, author, etc.)	Music of any genre	African drums	= Audio len.	2.5K

Table 2: Comparison of our *Moûsai* model with previous music/audio generation models. We compare the followings aspects: (1) audio sample rate@the number of channels (**Sample Rate**[↑], where the higher the better), (2) context length of the generated music (**Len.**[↑], where the higher the more capable the model is to generate structural music; * indicates variable length, where we assume that autoregressive methods are variable by default, with an upper-bound imposed by attention), (3) input type (**Input**, where we feature using **Text** ✓ as the condition for the generation), (4) type of the generate music (**Music**, where the more **Diverse**[↑] genre, the better), (5) an example of the generated music type (**Example**), (6) inference time (**Infer. Time**[↓], where the shorter the better, and since the music length is seconds or minutes, the inference time equivalent to the audio length is the shortest, and we use * to show models that can run inference fast on CPU), and (7) total length of the music in the training data in hours (**Data**).

Model	CLAP (↑)	Inf. Time (s) (↓)	Inf. Mem. (G) (↓)
Riffusion	0.06	218.0	8.85
Moûsai	0.13	49.2	5.04

Table 3: Performance of our *Moûsai* and the Riffusion model in terms of the CLAP score, as well as the inference time (Inf. Time), and inference memory (Inf. Mem.) for a single 43-second music clip.

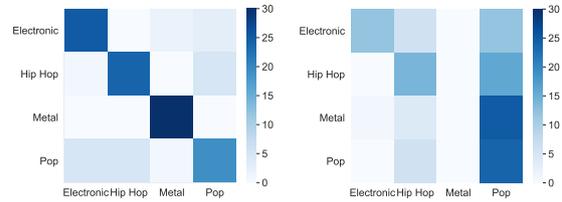
5.3 Evaluating the Text-Music Relevance

To assess how much the generated music corresponds to the given text prompt, we deploy both human and automatic evaluations.

Relevance & Distinctiveness by Human Evaluation We design a listener test where the annotators need to infer some coarse information of the text prompt behind a given piece of generated music. Since it is too challenging to infer the exact text prompt, we only ask annotators to infer the music genre indicated in the prompt.

To prepare the ground-truth prompts, we compose a list of 40 text prompts spanning across several common music genres: electronic, hip hop, metal, and pop. See Appendix C.1 for the entire list of prompts. Inspired by the two-alternative forced choice (2AFC) experiment design, we design a *four-alternative forced choice (4AFC)* paradigm, where the annotators need to categorize each music sample into exactly one of the four provided categories. See annotation details in Appendix C.1.

In Figure 5, we can see that our *Moûsai* model has the most mass on the diagonal (i.e., correctly identified), while the Riffusion model tends to generate generic samples that are mostly identified as pop



(a) Confusion matrix for the music pieces generated by *Moûsai*. (*y*-axis: true genre; *x*-axis: inferred genre.)

(b) Confusion matrix for the music pieces generated by the Riffusion model.

Figure 5: Results of the text-music relevance check, where the annotators are asked to infer the generated music by (a) our model and (b) the Riffusion model to their ground-truth genres: electronic, hip hop, metal, and pop. Perfect results are when the diagonal is dark.

for all ground-truth genres. This shows that the music generated by our model is both relevant to the test and distinct enough with the given genre against others.

Relevance by CLAP For automatic evaluation, we adopt the commonly used CLAP score (Wu et al., 2023) to quantify the alignment between the generated audio and the corresponding text. From Table 3, we can see that our model is two times better than Riffusion in terms of CLAP score, and also much faster in inference time.

5.4 Evaluating the Music Quality

We first introduce the four evaluation metrics, and then describe the evaluation results.

5.4.1 Metrics for Music Quality

To evaluate the quality of the generated music, we adopt four metrics: the automatic score by FAD, a

music Turing test, and human evaluation on musicality and audio clarity.

For automatic evaluation, we deploy the widely adopted **Fréchet Audio Distance (FAD)** (Kilgour et al., 2019) to assess the fidelity of the generated music distribution in comparison to the real music distribution (i.e., how **similar** the generated music is to the authentic music). To facilitate the computation of FAD, we employ the commonly used PANN model (Kong et al., 2020) as a means to effectively encode the music.

Then, we also set up three human evaluations, all on a scale of 1 (worst) to 5 (best). First, we let human annotators to assess the **authenticity/fidelity** of the generated music via a music Turing test, or fidelity (Goel et al., 2022; Hawthorne et al., 2019b; Hyun et al., 2022). Each time, we give the annotator two music samples, and ask them which one is real and which is generated. To provide a more fine-grained score, we also ask them how close the generated music they identified sounds like real music, on a scale of 1 (almost not similar at all) to 5 (highly similar). We keep their annotation score if they identify the generated music correctly, and otherwise we rate the music as 5, which means that the music perfectly passes the Turing test. Due to the space limit, we report the evaluation details in Appendix C.2.

The other two metrics we deploy are **musicality** and **audio clarity**. For musicality, we let human annotators rate the melodiousness and harmoniousness (Seitz, 2005) of the given music. And for audio clarity, or quality (Goel et al., 2022), we let them judge how close the quality is to a walkie-talkie (worst) or a high-quality studio sound system (best). The detailed setup of all our human evaluations are in Appendix C.2 and Appendix C.3.

5.4.2 Results

We show the evaluation results on all five metrics in Table 4. We can see that, on the automatic evaluation of FAD, our model has the best score, which is one magnitude smaller than previous models. Moreover, it also shows strong performance across the human evaluation metrics, outperforming the other two models on the music Turing test, harmoniousness, and sound clarity, as well as being comparable on the melodiousness metric.

Model	FAD (\downarrow)	Fidelity	Melody	Harmony	Clarity
Riffusion	0.0018	2.8	2.66	2.48	2.37
Musika	0.0020	3.04	3.21	3.04	2.88
<i>Moûsai</i>	0.00015	3.17	3.15	3.08	2.92

Table 4: Music quality scores for the three models.

5.5 Long-Term Structure of the Music

In music composition, the arrangement of a piece typically follows a gradual introduction, a main body with the core content, and a gradual conclusion, also called the sonata form (Webster, 2001). Accordingly, we look into whether our generated music also shows such long-term structure. Using the same text prompt, we can generate different segments/intervals of it by attaching the expression “1/2/3/4 out of 4” at the end of the text prompt, such as “Italian Hip Hop 2022, 3 of 4.” Specifically, we randomly generate 1000 music pieces, where the prompts have an even distribution of the four segment tags.

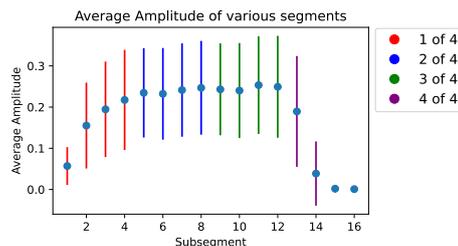


Figure 6: The average amplitude and variation across different segments of the generated music files.

We visualize the results in Figure 6, where we see the first segment shows a gradual increase in both the average amplitude and variance, followed by continuously high average amplitude and variance throughout Segments 2 and 3, and finally concluding with a gradual decline in the last segment.

6 Conclusion

In this work, we presented *Moûsai*, a novel text-to-music generation model using latent diffusion. We show that, in contrast to earlier approaches, our model can generate minutes of high-quality music in real-time on a consumer GPU, with good music quality and text-audio binding. In addition, we provide a collection of open-source libraries to facilitate future work in the field. We expect that the work will help pave the way towards higher-quality, longer-context text-to-music generation for future applications.

592 **Limitations and Future Work**

593 **Limited computational resources in an academic setting.**

594 We notice that there are some
595 concurrent work that are highly competitive
596 (Agostinelli et al., 2023), most of which are led
597 by industry labs. However, by the time we finished
598 this project in January 2023, our model is by far
599 the most performant given the limited resources in
600 the academic setting. The asymmetric distribution
601 of computational resources is making it almost im-
602 possible for academic labs to train state-of-the-art
603 generative models these days. We cannot compete
604 further with scaling up the model.

605 **Limited data.** Enhancing the scale of both data
606 and the model holds promising potential for yield-
607 ing significant improvements in quality. Following
608 (Dhariwal et al., 2020; Borsos et al., 2022), we sug-
609 gest training with 50k-100k hours instead of 2.5k.
610 Computer Vision studies like (Saharia et al., 2022)
611 show that utilizing larger pretrained language mod-
612 els for text embeddings plays an important role
613 in achieving superior quality outcomes. Drawing
614 upon this, we hypothesize that the application of
615 a larger pretrained language model to our second-
616 stage model can similarly contribute to enhanced
617 quality outcomes.

618 **Ethical Considerations**

619 Our work aims to bridge the gap between text and
620 music generation, enabling the creation of expres-
621 sive and high-quality music from textual descrip-
622 tions. While this research has the potential to ben-
623 efit various applications, such as music therapy,
624 entertainment, and education, we recognize that
625 it may also raise concerns in terms of copyright,
626 cultural appropriation, and the potential misuse of
627 generated content.

628 *Copyright and Intellectual Property:* Our model
629 may generate music that resembles existing copy-
630 righted works, which could lead to potential legal
631 disputes. First of all, for research-only use, it is
632 exempted from copyright infringement. For other
633 purposes, we suggest incorporating mechanisms
634 to detect and avoid generating music that closely
635 resembles copyrighted material.

636 *Economic Impact on Musicians and Composers:*
637 The widespread adoption of text-to-music genera-
638 tion models may have economic implications for
639 musicians and composers, potentially affecting

640 their livelihoods. We believe that our model should
641 be used as a tool to augment and inspire human
642 creativity, rather than replace it. We encourage col-
643 laboration between AI researchers, musicians, and
644 composers to explore new ways of integrating AI-
645 generated music into the creative process, ensuring
646 that the technology benefits all stakeholders.

647 In conclusion, we are committed to conducting
648 our research responsibly and ethically. We encour-
649 age the research community to engage in open dis-
650 cussions about the ethical implications of text-to-
651 music generation models and to develop guidelines
652 and best practices for their responsible use. By
653 addressing these concerns, we hope to contribute
654 to the development of AI technologies that benefit
655 society and promote creativity, while respecting the
656 rights and values of all stakeholders.

657 **References**

- 658 Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse
659 Engel, Mauro Verzetti, Antoine Caillon, Qingqing
660 Huang, Aren Jansen, Adam Roberts, Marco Tagliasac-
661 chi, Matt Sharifi, Neil Zeghidour, and Christian Frank.
662 2023. *Musilm: Generating music from text*. 662
- 663 Michele Berlingerio and Francesca Bonin. 2018. *To-
664 wards a music-language mapping*. In *Proceedings of
665 the Eleventh International Conference on Language Re-
666 sources and Evaluation (LREC 2018)*, Miyazaki, Japan.
667 European Language Resources Association (ELRA). 667
- 668 Jeanette Bicknell. 2002. Can music convey semantic
669 content? a kantian approach. *The Journal of Aesthetics
670 and Art Criticism*, 60(3):253–261. 670
- 671 Zalán Borsos, Raphaël Marinier, Damien Vincent,
672 Eugene Kharitonov, Olivier Pietquin, Matthew Shar-
673 ifi, Olivier Teboul, David Grangier, Marco Tagliasac-
674 chi, and Neil Zeghidour. 2022. *Audiolm: a lan-
675 guage modeling approach to audio generation*. *CoRR*,
676 abs/2209.03143. 676
- 677 Nicolas Boulanger-Lewandowski, Yoshua Bengio, and
678 Pascal Vincent. 2012. *Modeling temporal dependencies
679 in high-dimensional sequences: Application to poly-
680 phonic music generation and transcription*. 680
- 681 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
682 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
683 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
684 Askeel, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen
685 Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,
686 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris
687 Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
688 Gray, Benjamin Chess, Jack Clark, Christopher Berner,
689 Sam McCandlish, Alec Radford, Ilya Sutskever, and
690 Dario Amodei. 2020. *Language models are few-shot
691 learners*. In *Advances in Neural Information Processing*
692 691

692	<i>Systems</i> , volume 33, pages 1877–1901. Curran Associates, Inc.	748
693		749
694	Antoine Caillon and Philippe Esling. 2021. RAVE: A variational autoencoder for fast and high-quality neural audio synthesis . <i>CoRR</i> , abs/2111.05011.	750
695		751
696		752
697	Sheng-Kuan Chung. 2006. Digital storytelling in integrated arts education. <i>The International Journal of Arts Education</i> , 4(1):33–50.	753
698		754
699		755
700	Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. 2023. Simple and controllable music generation . <i>CoRR</i> , abs/2306.05284.	756
701		757
702		758
703		759
704	Sylvie Delacroix. 2023. Data rivers: carving out the public domain in the age of chat-gpt . Available at SSRN.	760
705		761
706	Kangle Deng, Aayush Bansal, and Deva Ramanan. 2021. Unsupervised audiovisual synthesis via exemplar autoencoders . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	762
707		763
708		764
709		765
710		766
711	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	767
712		768
713		769
714		770
715		771
716		772
717		773
718		774
719		775
720	Prafulla Dhariwal, Heewoo Jun, Christine Payne, Jong Wook Kim, Alec Radford, and Ilya Sutskever. 2020. Jukebox: A generative model for music . <i>CoRR</i> , abs/2005.00341.	776
721		777
722		778
723		779
724	Sander Dieleman, Aäron van den Oord, and Karen Simonyan. 2018. The challenge of realistic music generation: Modelling raw audio at scale . In <i>Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada</i> , pages 8000–8010.	780
725		781
726		782
727		783
728		784
729		785
730		786
731	Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. 2022. CLAP: learning audio concepts from natural language supervision . <i>CoRR</i> , abs/2206.04769.	787
732		788
733		789
734		790
735	Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Douglas Eck, Karen Simonyan, and Mohammad Norouzi. 2017. Neural audio synthesis of musical notes with wavenet autoencoders .	791
736		792
737		793
738		794
739	Jesse H. Engel, Kumar Krishna Agrawal, Shuo Chen, Ishaan Gulrajani, Chris Donahue, and Adam Roberts. 2019. Gansynth: Adversarial neural audio synthesis . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	795
740		796
741		797
742		798
743		799
744		800
745	Patrick Esser, Robin Rombach, and Björn Ommer. 2021. Taming transformers for high-resolution image synthesis . In <i>IEEE Conference on Computer Vision and Pat-</i>	801
746		802
747		803
	<i>tern Recognition, CVPR 2021, virtual, June 19-25, 2021</i> , pages 12873–12883. Computer Vision Foundation / IEEE.	
	European Commission. 2016. Proposal for a directive of the European parliament and of the council on copyright in the digital single market .	
	Seth* Forsgren and Hayk* Martiros. 2022. Riffusion - Stable diffusion for real-time music generation .	
	Christophe Geiger, Giancarlo Frosio, and Oleksandr Bulayenko. 2018. The exception for text and data mining (tdm) in the proposed directive on copyright in the digital single market-legal aspects. <i>Centre for International Intellectual Property Studies (CEIPI) Research Paper</i> , (2018-02).	
	Mark Germer. 2011. <i>Notes</i> , 67(4):760–765.	
	Daniel Gillick, Sayali Kulkarni, Larry Lansing, Alessandro Presta, Jason Baldrige, Eugene Ie, and Diego Garcia-Olano. 2019. Learning dense representations for entity retrieval. In <i>Computational Natural Language Learning (CoNLL)</i> .	
	Karan Goel, Albert Gu, Chris Donahue, and Christopher Ré. 2022. It’s raw! audio generation with state-space models . In <i>International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 7616–7633. PMLR.	
	Gal Greshler, Tamar Rott Shaham, and Tomer Michaeli. 2021. Catch-a-waveform: Learning to generate audio from a single short example . In <i>Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual</i> , pages 20916–20928.	
	Curtis Hawthorne, Andriy Stasyuk, Adam Roberts, Ian Simon, Cheng-Zhi Anna Huang, Sander Dieleman, Erich Elsen, Jesse Engel, and Douglas Eck. 2019a. Enabling factorized piano music modeling and generation with the MAESTRO dataset . In <i>International Conference on Learning Representations</i> .	
	Curtis Hawthorne, Andriy Stasyuk, Adam Roberts, Ian Simon, Cheng-Zhi Anna Huang, Sander Dieleman, Erich Elsen, Jesse H. Engel, and Douglas Eck. 2019b. Enabling factorized piano music modeling and generation with the MAESTRO dataset . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	
	Geoffrey E Hinton and Ruslan R Salakhutdinov. 2006. Reducing the dimensionality of data with neural networks. <i>science</i> , 313(5786):504–507.	
	Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey A. Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. 2022. Imagen video: High definition video generation with diffusion models . <i>CoRR</i> , abs/2210.02303.	

804	Jonathan Ho and Tim Salimans. 2022. Classifier-free diffusion guidance . <i>CoRR</i> , abs/2207.12598.	
805		
806	Qingqing Huang, Daniel S. Park, Tao Wang, Timo I. Denk, Andy Ly, Nanxin Chen, Zhengdong Zhang, Zhishuai Zhang, Jiahui Yu, Christian Havnø Frank, Jesse H. Engel, Quoc V. Le, William Chan, and Wei Han. 2023. Noise2music: Text-conditioned music generation with diffusion models . <i>CoRR</i> , abs/2302.03917.	
807		
808		
809		
810		
811		
812	Lee Hyun, Taehyun Kim, Hyolim Kang, Minjoo Ki, Hyeonchan Hwang, Kwanho Park, Sharang Han, and Seon Joo Kim. 2022. Commu: Dataset for combinatorial music generation . <i>CoRR</i> , abs/2211.09385.	
813		
814		
815		
816	Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. 2019. Fréchet audio distance: A metric for evaluating music enhancement algorithms .	
817		
818		
819	Minsu Kim, Joanna Hong, and Yong Man Ro. 2021. Lip to speech synthesis with visual context attentional GAN . In <i>Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual</i> , pages 2758–2770.	
820		
821		
822		
823		
824		
825	Diederik P. Kingma and Max Welling. 2014. Auto-encoding variational bayes . In <i>2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings</i> .	
826		
827		
828		
829		
830	Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D. Plumbley. 2020. Panns: Large-scale pretrained audio neural networks for audio pattern recognition .	
831		
832		
833		
834	Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. 2021. Diffwave: A versatile diffusion model for audio synthesis . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	
835		
836		
837		
838		
839	Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. 2022. Audiogen: Textually guided audio generation . <i>CoRR</i> , abs/2209.15352.	
840		
841		
842		
843	Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre de Brébisson, Yoshua Bengio, and Aaron C. Courville. 2019. Melgan: Generative adversarial networks for conditional waveform synthesis . In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada</i> , pages 14881–14892.	
844		
845		
846		
847		
848		
849		
850		
851		
852	Max W. Y. Lam, Jun Wang, Dan Su, and Dong Yu. 2022. BDDM: bilateral denoising diffusion models for fast and high-quality speech synthesis . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	
853		
854		
855		
856		
857		
858	Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. 2022. Autoregressive image generation using residual quantization . In <i>IEEE/CVF Con-</i>	
859		
860		
	<i>ference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 11513–11522. IEEE.	861
		862
		863
	Yichong Leng, Zehua Chen, Junliang Guo, Haohe Liu, Jiawei Chen, Xu Tan, Danilo P. Mandic, Lei He, Xiang-Yang Li, Tao Qin, Sheng Zhao, and Tie-Yan Liu. 2022. Binauralgrad: A two-stage conditional diffusion probabilistic model for binaural audio synthesis . <i>CoRR</i> , abs/2205.14807.	864
		865
		866
		867
		868
		869
	Manling Li, Ruochen Xu, Shuohang Wang, Luwei Zhou, Xudong Lin, Chenguang Zhu, Michael Zeng, Heng Ji, and Shih-Fu Chang. 2022. Clip-event: Connecting text and images with event structures . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 16399–16408. IEEE.	870
		871
		872
		873
		874
		875
		876
	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	877
		878
		879
		880
	Soroush Mehri, Kundan Kumar, Ishaan Gulrajani, Rithesh Kumar, Shubham Jain, Jose Sotelo, Aaron C. Courville, and Yoshua Bengio. 2017. Samplernn: An unconditional end-to-end neural audio generation model . In <i>5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings</i> . OpenReview.net.	881
		882
		883
		884
		885
		886
		887
	Rada Mihalcea and Carlo Strapparava. 2012. Lyrics, music, and emotions . In <i>Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning</i> , pages 590–599, Jeju Island, Korea. Association for Computational Linguistics.	888
		889
		890
		891
		892
		893
	Max Morrison, Rithesh Kumar, Kundan Kumar, Prem Seetharaman, Aaron C. Courville, and Yoshua Bengio. 2022. Chunked autoregressive GAN for conditional waveform synthesis . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	894
		895
		896
		897
		898
		899
	Isabel Papadimitriou and Dan Jurafsky. 2020. Learning Music Helps You Read: Using transfer to study linguistic structure in language models . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6829–6839, Online. Association for Computational Linguistics.	900
		901
		902
		903
		904
		905
	Marco Pasini and Jan Schlüter. 2022. Musika! fast infinite waveform music generation . <i>CoRR</i> , abs/2208.08706.	906
		907
		908
	Konpat Preechakul, Nattanat Chatthee, Suttisak Widadwongsa, and Supasorn Suwajanakorn. 2022. Diffusion autoencoders: Toward a meaningful and decodable representation . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 10609–10619. IEEE.	909
		910
		911
		912
		913
		914
		915
	Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understand-	916
		917

918	ing by generative pre-training. Technical report, OpenAI.	Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. 2017. Neural discrete representation learning . In <i>Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA</i> , pages 6306–6315.	973
919			974
920	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>J. Mach. Learn. Res.</i> , 21:140:1–140:67.		975
921			976
922			977
923			978
924			
925	Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with CLIP latents . <i>CoRR</i> , abs/2204.06125.	Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. 2022. Phenaki: Variable length video generation from open domain textual description . <i>CoRR</i> , abs/2210.02399.	979
926			980
927			981
928			982
			983
			984
929	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 10674–10685. IEEE.	James Webster. 2001. Sonata form. <i>The new Grove dictionary of music and musicians</i> , 23:687–698.	985
930			986
931			
932			987
933			988
934			989
			990
935	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation . In <i>Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015 - 18th International Conference Munich, Germany, October 5 - 9, 2015, Proceedings, Part III</i> , volume 9351 of <i>Lecture Notes in Computer Science</i> , pages 234–241. Springer.	Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. 2023. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation .	991
936			992
937			993
938			994
939			
940			995
941			996
			997
942	Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. 2022. Photo-realistic text-to-image diffusion models with deep language understanding . <i>CoRR</i> , abs/2205.11487.	Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. 2022. Diff-sound: Discrete diffusion model for text-to-sound generation . <i>CoRR</i> , abs/2207.09983.	992
943			993
944			994
945			
946			995
947			996
948			997
			998
949	Tim Salimans and Jonathan Ho. 2022. Progressive distillation for fast sampling of diffusion models . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	A More Data Details	999
950			
951			
952			
953			
954	Jay A Seitz. 2005. Dalcroze, the body, movement and musicality. <i>Psychology of music</i> , 33(4):419–435.	A.1 Data Collection Rationale	1000
955			
956	Jiaming Song, Chenlin Meng, and Stefano Ermon. 2021. Denoising diffusion implicit models . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	We have several desiderata when collecting the dataset: the data (1) must have text data paired with the music piece, and (2) must constitute a large size, which means that our data crawling procedure needs to be scalable, without tedious manual efforts to curate. Note that it is crucial to get a large-sized dataset in order to unleash the performance of audio generation diffusion models.	1001
957			1002
958			1003
959			1004
960			1005
			1006
			1007
			1008
961	Joseph P Swain. 1995. The concept of musical syntax. <i>The Musical Quarterly</i> , 79(2):281–308.	A.2 Training setup for the text-music pairs	1009
962			
963	A. M. TURING. 1950. I.—COMPUTING MACHINERY AND INTELLIGENCE. <i>Mind</i> , LIX(236):433–460.	For the textual description, we use metadata such as the title, author, album, genre, and year of release. Given that a song could span longer than 44s, we append a string indicating which chunk is currently being trained on, together with the total chunks the song is made of (e.g., <i>1 of 4</i>). This allows to select the region of interest during inference. Hence, an example prompt is like “ <i>Egyptian Darbuka, Drums, Rythm, (Deluxe Edition), 2 of 4.</i> ” To make the conditioning more robust, we shuffle the list of metadata and drop each element with a probability of 0.1. Furthermore, for 50% of the times we concatenate the list with spaces and the other 50% of the times we use commas to make	1010
964			1011
965			1012
966	Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. 2016. Wavenet: A generative model for raw audio . In <i>The 9th ISCA Speech Synthesis Workshop, Sunnyvale, CA, USA, 13-15 September 2016</i> , page 125. ISCA.		1013
967			1014
968			1015
969			1016
970			1017
971			1018
972			1019
			1020
			1021
			1022
			1023

the interface more robust during inference. Some example prompts in our dataset can be seen in Table 5.

Example Text Prompts in Our Dataset
Nr. 415 (Premium Edition), german hip hop, 2 of 7, 2012, XATAR, Konnekt
30 Años de Exitos, Mundanzas, 2 of 6, latin pop, Lupita D’Alessio, 2011
emo rap 2018 Runaway Lil Peep 4 of 5
Alone, Pt. II (Remixes) 2020 electro house Alone, Pt. II - Da Tweekaz Remix Alan Walker

Table 5: Example text prompts in our dataset.

A.3 Model Architecture and Parameters

Our diffusion autoencoder has 185M parameters, with 7 nested U-Net blocks of increasing channel count ([256, 512, 512, 512, 1024, 1024, 1024]), for which we downsample each time by 2, except for the first block ([1, 2, 2, 2, 2, 2, 2]). This makes the compression factor for our autoencoder to be 64x. Depending on the desired speed/quality tradeoff, more or less compression can be applied in this first stage. Following our single GPU constraint, we find that 64x compression factor is a good balance to make sure the second stage can work on a reduced representation. We discuss more about this tradeoff in Appendix E.4. The diffusion autoencoder only uses ResNet and modulation items with the repetitions [1, 2, 2, 2, 2, 2, 2]. We do not use attention, to allow decoding of variable and possibly very long latent representations. Channel injection only happens at depth 4, which matches the output of the magnitude encoder latent, after applying the tanh function.

Our text-conditional generator has 857M parameters (including the parameters of the frozen T5-base model) with 6 nested U-Net blocks of increasing channel counts ([128, 256, 512, 512, 1024, 1024]), and again downsampling each time by 2, except for the first block ([1, 2, 2, 2, 2, 2]). We use attention blocks at the depths [0, 0, 1, 1, 1, 1], skipping the first two blocks to allow for further downsampling before sharing information over the entire latent, instead use cross-attention blocks at all resolutions ([1, 1, 1, 1, 1, 1]). For both attention and cross-attention, we use 64 head features and 12 heads per layer. We repeat items with an increasing count towards the inner U-Net low-resolution and large-context blocks ([2, 2, 2, 4, 8, 8]), this allows good structural learning over minutes of audio.

B More experiments

B.1 Hardware Requirements

We use limited computational resources as available in a university lab. (3) **Efficiency** is another highlight of our model, where we only needs an inference time equivalent to the audio length on a consumer GPU, which is several minutes, while many other text-to-audio models take many GPU hours (Dhariwal et al., 2020; Kreuk et al., 2022). Our model is very friendly for research at university labs, as each of our models can be trained on a single A100 GPU in 1 week of training using a batch size of 32; this is equivalent to around 1M steps for both the diffusion autoencoder and latent generator. For inference, as an example, a novel audio source of ~ 43 s can be synthesized in less than 50s using a consumer GPU with a DDIM sampler and a high step count (100 generation steps and 100 decoding steps).

C More analysis

C.1 Annotation Details for the Genre Identification Test

Prompts We list all the text prompts composed for the four common music genres in Table 6.

Using these prompts, we generate music with both *Mou’sai* and the Riffusion model (Forsgren and Martiros, 2022), with a total of 80 pieces of music, two for each prompt.

To validate this quantitatively, we conducted a listener test with three perceivers (annotators) with diverse demographic backgrounds (both female and male, all with at least a Master’s degree of education). Each annotator listens to all 80 music samples we provide, and is instructed to categorize each sample into exactly one of the four provided genres.

Annotation We record how many times the perceiver correctly identifies the genre which the respective model was generating from. A large number (or score) means that the model often generated music that, according to the human perceiver, plausibly belonged to the correct category (when compared to the other three categories). To achieve a good score, the model needs to generate diverse and genre-specific music. We take the score as a quality score of the model when it comes to correctly performing text-conditional music generation.

In Figure 5, we display the confusion matrix of this genre identification test for both our model (left) and the Riffusion model (right). For our model, the annotators identify the right genres most of the time, whereas for the Riffusion model, the annotators often perceive the music as more generic, categorizing it as Pop.

C.2 Annotation Details for Turing Test

As for the details, we create 90 music samples, including 15 generated samples paired with 15 real music samples for each of the three models (Riffusion, Musika, and *Moûsai*). We recruit two annotators, pursuing Bachelor of Technology degree from the Indian Institute of Technology, Kharagpur, India. Additionally, the two annotators have pursued playing music as a hobby for the past 10 years.

We conducted a rigorous evaluation employing an experiment with a similar spirit to the Turing test (TURING, 1950) for natural language, but commonly called as the fidelity test in audio evaluation (Hyun et al., 2022) or speaker test (Greshler et al., 2021; Hawthorne et al., 2019b) in audio evaluation. Our methodology involved presenting a group of expert annotators with a total of 60 distinct folders, 15 corresponding to each of Mousai, Mousai (classical-only), Riffusion, and Musika models. Each folder containing two music files, one being the original and the other generated using a given model prompted with its corresponding metadata.

The annotators were provided with the task of determining the fidelity and providing a rating on a scale of 1 to 5, reflecting the perceived degree of authenticity of the generated audio. In cases where the annotators incorrectly identified the generated audio, the respective model was awarded 5 points. Conversely, if the annotators correctly identified the generated audio, the model’s rating was determined based on the score provided by the annotator. The annotators were compensated with 500 rupees (~6.5 dollars) for this 3 hour task (which is well above daily minimum wage in India).

Following are the exact instructions provided to the annotators

1. You will be presented with batches of two audio samples in subfolders of this folder named from 1 to 60. Each subfolder contains two audios named a.wav and b.wav.
2. Listen to each sample carefully.

3. It’s best to use headphones in a quiet environment if you can.
4. Some files may be loud, so it’s recommended to keep the volume moderate.
5. One of the audio samples in each pair is a real recording, while the other is a generated (synthetic) audio.
6. Listen to each pair of audio samples carefully.
7. Pay attention to the quality, characteristics, and nuances of each audio sample.
8. This folder contains a spreadsheet file called ‘Response_Task_2.xlsx’. Compare the samples to each other and provide a relative rating to the fake audio only out of 5, where 1 being the most fake and 5 being most real.

C.3 Annotation Details for Musicality

In order to ascertain the quality and artistic merit of the generated musical output, a rigorous human evaluation methodology was implemented. A total of 50 carefully curated folders, each containing three distinct audio files, were presented to human evaluators. These audio files were generated utilizing various models, all prompted by a specific prompt. We recruit two annotators, pursuing Bachelor of Technology degree from the Indian Institute of Technology, Kharagpur, India. Additionally, the two annotators have pursued playing music as a hobby for the past 10 years. The annotators were compensated with 500 rupees (~6.5 dollars) for this 3 hour task (which is well above daily minimum wage in India).

Following are the exact instructions provided to the annotators

1. Listen to the music and rate it based on three aspects: Quality, Melody, and Harmony.
2. It’s best to use headphones in a quiet environment if you can.
3. Some files may be loud, so it’s recommended to keep the volume moderate.
4. This folder contains folders subfolders through 1-50. Each subfolders contains three audio files named A.wav, B.wav, and C.wav . You need to listen to each of them and rate them (relative to each other) based on quality, melody, and harmony.
5. For Quality, consider how clear the audio sounds. Does it resemble a walkie-talkie (bad quality) or a high-quality studio sound system

(good quality)?

6. **Melodiousness** refers to the main pitch or note in the music. Pay attention to the rhythm and repetitiveness of the melody. A more rhythmic and repetitive melody is considered better, while the opposite is true for a less rhythmic melody.
7. **Harmoniousness** involves multiple notes played together to support the melody. Evaluate if these notes are in sync and enhance the effect of the melody. Higher scores should be given for good harmony and lower for poor harmony.
8. It is recommended view youtube videos: [this](#) or [this](#) short video explaining melody and harmony
9. This folder also contains a spreadsheet by the name "Response_Task_1.xlsx". Remember to provide ratings (out of 5) for each aspect of your evaluation in the file against appropriate folder number. Feel free to listen to each sample as many times before rating them.

D More Related Work

Audio Generation Audio generation is a challenging task. At the lowest level, we have digital waveforms that control air movement from speakers. Waveforms can be represented in different resolutions, or sample rates. Higher sample rates (e.g., 48kHz) allow for more temporal resolution and can represent higher frequencies, but at the same time it is computationally more demanding to generate. At higher levels of abstraction, we find qualitative properties such as texture (timbre) or pitch. Zooming out, we observe structure such as rhythm and melody that can span multiple seconds, or even structurally be composed into choruses that form minutes of interconnected patterns.

Audio can be represented with a single waveform (mono), two waveforms (stereo), or even more waveforms in the case of surround sound. Audio with two or more channels can give a sense of movement and spatialisation. From a modelling perspective, there are (1) unconditional models that generate novel samples from the training distribution without any additional information, or (2) conditional models that use a form of guidance, such as text, to control the generation. Models can be trained on a single modality (e.g., drums or piano) or on multiple modalities, which usually re-

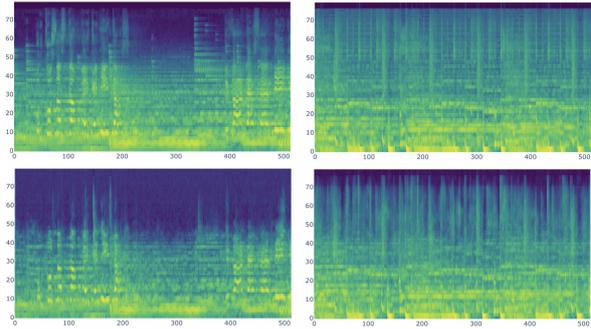


Figure 7: Mel spectrogram comparison between the true samples (top) and the auto-encoded samples (bottom); cf. text.

quire more parameters for an increased modelling capacity and decrease in speed.

E (Informal) Intuitions for Model Architecture and Training Setup

Sound types that our model is good at

Apart from the diversity and relevance, we also evaluate the sound quality of the music we generate. From the mel spectrograms we visualize in Figure 7, we can see that low-frequency sounds are handled rather well by our model. From the music samples we provide, it is apparent that our model performs well with drum-like sounds as frequently found in electronic, house, dubstep, techno, EDM, and metal music. This is likely a consequence of the lower amount of information required to represent low-frequency sounds.

E.1 Improving the Structure

We find that increasing the number of attention blocks (e.g., from a total of 4 – 8 to a total of 32+) in the latent diffusion model can improve the general structure of the songs, thanks to the long-context view. If the model is trained without attention blocks, the context provided by the U-Net is not large enough to learn any meaningful long-term structure.

E.2 Text-Audio Binding

We find that the text-audio binding works well with CFG higher than 3.0. Since the model is trained with metadata such as title, album, artist, genre, year, and chunk, the best keywords to control the generation appear to be frequent descriptive names, such as the genre of the music, or descriptions commonly found in titles, such as "remix", "(Deluxe

Edition)”, and possibly many more. A similar behavior has been observed and exploited in text-to-image models to generate better looking results.

E.3 Trade-Off between Speed and Quality

We find that 10 sampling steps in both stages can be enough to generate reasonable audio. We can achieve improved quality and reduced noise for high-frequency sounds by trading off the speed, i.e., increasing the number of sampling steps in the diffusion decoder, e.g., 50 – 100 steps). Increasing the number of sampling steps in the latent diffusion model (again in the order of 50 – 100 steps) will similarly improve the quality, likely due to the more detailed generated latents, and at the same time result in an overall better structured music. To make sure the results are comparable when varying the number of sampling steps, we use the same starting noise in both stages. In both cases, this suggests that using more advanced samplers could be helpful to improve on the speed-quality trade-off.

E.4 Trade-Off between Compression Ratio and Quality

We find that decreasing the compression ratio of the first stage (e.g., to 32x) can improve the quality of low-frequency sounds, but in turn will slow down the model, as the second stage has to work on higher dimensional data. As proposed later in Section 6, we hypothesize that using perceptually weighted loss functions instead of L2 loss during diffusion could help this trade-off, giving a more balanced importance to high frequency sounds even at high compression ratios.

Genre = Electronic

- Drops, Kanine Remix, Darkzy, Drops Remixes, bass house, (Deluxe) (Remix) 3 of 4
- Electronic, Dance, EDM (Deluxe) (Remix) 3 of 4
- Electro House (Remix), 2023, 3 of 4
- Electro Swing Remix 2030 (Deluxe Edition) 3 of 4
- Future Bass, EDM (Remix) 3 of 4, Remix
- EDM (Deluxe) (Remix) 3 of 4
- EDM, Vocal, Relax, Remix, 2023, 8D Audio
- Hardstyle, Drop, 8D, Remix, High Quality, 2 of 4
- Dubstep Insane Drop Remix (Deluxe Edition), 2 of 4
- Drop, French 79, BPM Artist, Vol. 4, Electronica, 2016

Genre = Hip Hop

- Real Hip Hop, 2012, Lil B, Gods Father, escape room, 3 of 4
- C’est toujours pour ceux qui savent, French Hip Hop, 2018 (Deluxe), 3 of 4
- Dejando Claro, Latin Hip Hop 2022 (Deluxe Edition) 3 of 4
- Latin Hip Hop 2022 (Deluxe Edition) 3 of 4
- Alternative Hip Hop Oh-My, 2016, (Deluxe), 3 of 4
- Es Geht Mir Gut, German Hip Hop, 2016, (Deluxe), 3 of 4
- Italian Hip Hop 2022 (Deluxe Edition) 3 of 4
- RUN, Alternative Hip Hop, 2016, (Deluxe), 3 of 4
- Hip Hop, Rap Battle, 2018 (High Quality) (Deluxe Edition) 3 of 4
- Hip Hop Tech, Bandlez, Hot Pursuit, brostep, 3 of 4

Genre = Metal

- Death Metal, 2012, 3 of 4
- Heavy Death Metal (Deluxe Edition), 3 of 4
- Black Alternative Metal, The Pick of Death (Deluxe), 2006, 3 of 4
- Kill For Metal, Iron Fire, To The Grave, melodic metal, 3 of 4
- Melodic Metal, Iron Dust (Deluxe), 2006, 3 of 4
- Possessed Death Metal Stones (Deluxe), 2006, 3 of 4
- Black Metal Venom, 2006, 3 of 4
- The Heavy Death Metal War (Deluxe), 2006, 3 of 4
- Heavy metal (Deluxe Edition), 3 of 4
- Viking Heavy Death Metal (Deluxe), 2006, 3 of 4

Genre = Pop

- (Everything I Do), I Do It For You, Bryan Adams, The Best Of Me, canadian pop, 3 of 4
 - Payphone, Maroon 5, Overexposed, Pop, 2021, 3 of 4
 - 24K Magic, Bruno Mars, 24K Magic, dance pop, 3 of 4
 - Who Is It, Michael Jackson, Dangerous, Pop (Deluxe), 3 of 4
 - Forget Me, Lewis Capaldi, Forget Me, Pop Pop, 2022, 3 of 4
 - Pop, Speak Now, Taylor Swift, 2014, (Deluxe), 3 of 4
 - Pop Pop, Maroon 5, Overexposed, 2016, 3 of 4
 - Pointless, Lewis Capaldi, Pointless, Pop, 2022, 3 of 4
 - Saved, Khalid, American Teen, Pop, 2022, 3 of 4
 - Deja vu, Fearless, Pop, 2020, (Deluxe), 3 of 4
-

Table 6: Text prompts composed for the four common music genres: electronic, hip hop, metal, and pop.