

# GENERALIZED MULTI-SOURCE INFERENCE FOR TEXT CONDITIONED MUSIC DIFFUSION MODELS

Emilian Postolache<sup>1</sup>, Giorgio Mariani<sup>1</sup>, Luca Cosmo<sup>2</sup>, Emmanouil Benetos<sup>3</sup>, Emanuele Rodolà<sup>1</sup>

<sup>1</sup>Sapienza University of Rome

<sup>2</sup>Ca' Foscari University of Venice

<sup>3</sup>Queen Mary University of London

## ABSTRACT

Multi-Source Diffusion Models (MSDM) allow for compositional musical generation tasks: generating a set of coherent sources, creating accompaniments, and performing source separation. Despite their versatility, they require estimating the joint distribution over the sources, necessitating pre-separated musical data, which is rarely available, and fixing the number and type of sources at training time. This paper generalizes MSDM to arbitrary time-domain diffusion models conditioned on text embeddings. These models do not require separated data as they are trained on mixtures, can parameterize an arbitrary number of sources, and allow for rich semantic control. We propose an inference procedure enabling the coherent generation of sources and accompaniments. Additionally, we adapt the Dirac separator of MSDM to perform source separation. We experiment with diffusion models trained on Slakh2100 and MTG-Jamendo, showcasing competitive generation and separation results in a relaxed data setting.

**Index Terms**— Music Generation, Diffusion Models, Source Separation

## 1. INTRODUCTION

The task of musical generation has seen significant advancements recently, thanks to developments in generative models. The families of generative models showcasing state-of-the-art results are latent language models [1] and (score-based) diffusion models [2, 3, 4]. Latent language models map a continuous-domain (time or spectral) signal to a sequence of discrete tokens and estimate a density over such sequences autoregressively [5, 6] or via mask-modeling [7]. Diffusion models [8, 9], on the other hand, operate on continuous representations (time, spectral, or latent domains), capturing the gradient of the log-density perturbed by a noising process (Gaussian). Despite differences between these generative models, they typically share some mechanisms for conditioning on rich textual embeddings, obtained either using text-only encoders [10] or audio-text contrastive encoders [11, 12, 13]. Such a mechanism allows generating a musical track following a natural language prompt.

Generative models for music typically output only a final mixture. As such, generating the constituent sources is challenging. This implies that musical generative models are hard to employ in music production tasks, where the subsequent manipulation of sub-tracks, creation of accompaniments, and source separation is often required. Two existing approaches aim to address this issue. The first approach, called Multi-Source Diffusion Models (MSDM) [14], trains a diffusion model in time domain on (supervised) sets of coherent sources viewed as different channels without conditioning on textual information. Such a model allows for generating a set of coherent sources, creating accompaniments, and performing source

separation. Despite being a versatile compositional model for music, MSDM has three limitations: (i) It requires knowledge of separated coherent sources, which are hard to acquire. (ii) It architecturally assumes a fixed number of sources and their respective class type (e.g., Bass, Drums, Guitar, Piano). (iii) It is impossible to condition the sources on rich semantic information, as commonly done with text-conditioned music models. The second approach, based on supervised instruction prompting [15, 16], fine-tunes a latent diffusion model with instructions that allow adding, removing, and extracting sources present in a musical track. Although this approach addresses the issues (ii) and (iii) of MSDM, it does not solve the problem (i), necessitating pre-separated data. A strategy for scaling both models is training with data obtained by separating sources from mixtures using a pre-trained separator [17]. This approach, though, is not flexible because such separated data contains artifacts, and we are limited to the number and type of sources the separator can handle.

We develop a novel inference procedure for the task, called *Generalized Multi-Source Diffusion Inference (GMSDI)*, that can be used in combination with *any* text-conditioned (time-domain) diffusion model for music. Such a method: (i) Requires only mixture data for training, resulting in an unsupervised algorithm when paired with a contrastive encoder. (ii) Parameterizes an arbitrary number and type of sources. (iii) Allows for rich semantic control. To our knowledge, this is the first general algorithm for unsupervised compositional music generation. After developing the required background notions in Section 2, we develop the inference techniques in Section 3. We detail the experimental setup in Section 4 and show empirical results in Section 5. We conclude the paper in Section 6.

## 2. BACKGROUND

A musical track  $\mathbf{y}$  is a mixture of  $K$  instrumental and vocal sources  $\mathbf{x} = \{\mathbf{x}_k\}_{k \in [K]}$ . Therefore, we have  $\mathbf{y} = \sum_{k=1}^K \mathbf{x}_k$ , with  $K$  depending on the mixture. Fixing a source  $\mathbf{x}_k$ , we denote the complementary set with  $\mathbf{x}_{\bar{k}} = \{\mathbf{x}_l\}_{l \in [K]} - \{\mathbf{x}_k\}$ . While we typically do not have direct access to the audio constituents  $\{\mathbf{x}_k\}_{k \in [K]}$ , we are usually equipped with a text embedding  $\mathbf{z}$  which provides information about the sources. We can obtain  $\mathbf{z} = E_{\phi}^{\text{text}}(\mathbf{q})$  by encoding a text description  $\mathbf{q}$  with a text-only encoder  $E_{\phi}^{\text{text}}$ , or use a pre-trained contrastive audio-text encoder  $E_{\phi}^{\text{contr}}$  to extract embeddings both from the audio mixtures  $\mathbf{z} = E_{\phi}^{\text{contr}}(\mathbf{y})$  and from text descriptions  $\mathbf{z} = E_{\phi}^{\text{contr}}(\mathbf{q})$ .

### 2.1. Text-conditioned Score-based Diffusion Models

We work with continuous-time score-based [4] diffusion models. A text-conditioned score-based diffusion model  $S_{\theta}$  parameterizes the logarithm of the perturbed audio mixture density, conditioned on the

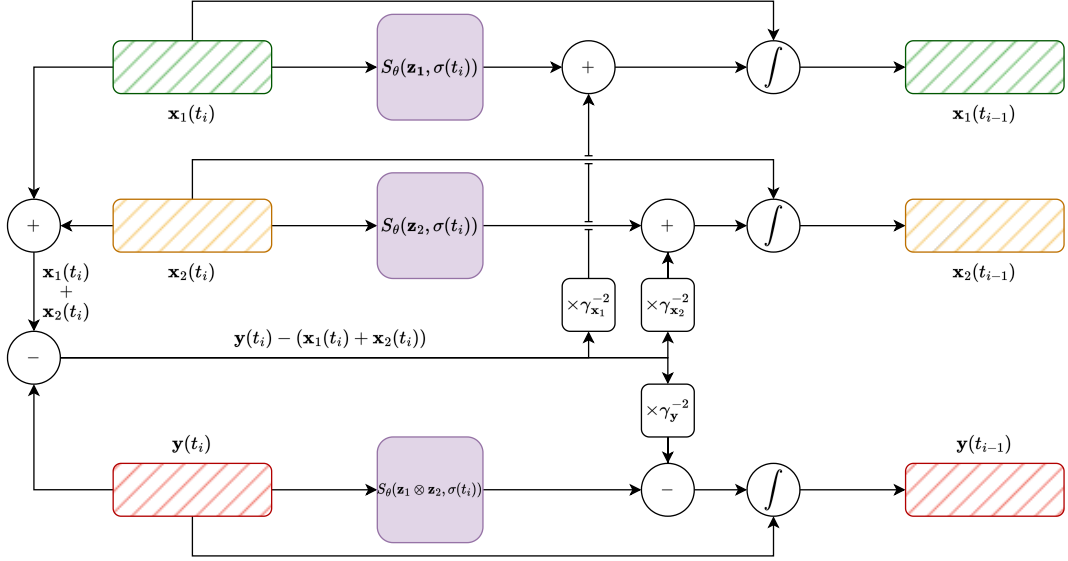


Fig. 1. Diagram for unconditional generation procedure with GMSDI, sampling two coherent sources.

textual embedding:

$$\nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t) | \mathbf{z}) \approx S_\theta(\mathbf{y}(t), \mathbf{z}, \sigma(t)), \quad (1)$$

where  $p(\mathbf{y}(t) | \mathbf{z}) = \int_{\mathbf{y}(0)} p(\mathbf{y}(t) | \mathbf{y}(0))p(\mathbf{y}(0) | \mathbf{z})$ , with

$$p(\mathbf{y}(t) | \mathbf{y}(0)) = \mathcal{N}(\mathbf{y}(t) | \mathbf{y}(0), \sigma^2(t)\mathbf{I}) \quad (2)$$

a Gaussian perturbation kernel depending on a noise schedule  $\{\sigma(t)\}_{t \in [0, T]}$ . We train  $S_\theta$  minimizing:

$$\mathbb{E}_{t \sim \mathcal{U}[0, T]} \mathbb{E}_{\mathbf{y}(0) \sim p(\mathbf{y}(0) | \mathbf{z})} \mathbb{E}_{\mathbf{y}(t) \sim p(\mathbf{y}(t) | \mathbf{y}(0))} [\mathcal{L}_{\text{SM}}],$$

where  $\mathcal{L}_{\text{SM}}$  is the denoising score-matching loss [2, 3]:

$$\mathcal{L}_{\text{SM}} = \|S_\theta(\mathbf{y}(t), \mathbf{z}, \sigma(t)) - \nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t) | \mathbf{y}(0))\|_2^2.$$

At inference time, we use classifier-free guidance [18], integrating

$$\begin{aligned} & S_\theta^*(\mathbf{y}(t), \mathbf{z}, \sigma(t)) \\ = & S_\theta(\mathbf{y}(t), \mathbf{z}, \sigma(t)) + w(S_\theta(\mathbf{y}(t), \mathbf{z}, \sigma(t)) - S_\theta(\mathbf{y}(t), \mathbf{z}^*, \sigma(t))), \end{aligned}$$

where  $\mathbf{z}^*$  is a fixed learned embedding modeling the unconditional  $\nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t))$ , and  $w \in \mathbb{R}$  is the embedding scale hyperparameter. We can use a *negative embedding* [19] instead of  $\mathbf{z}^*$  to better guide inference. With an abuse of notation, we will refer to  $S_\theta^*$  as  $S_\theta$ .

## 2.2. Multi-Source Diffusion Models

In [14], authors assume a fixed number  $K$  of coherent sources of known type  $\{\mathbf{x}_k\}_{k \in [K]}$  contained in the mixture  $\mathbf{y}$ . They train a *Multi-Source Diffusion Model (MSDM)*, an unconditional score-based diffusion model  $S_\theta^{\text{MSDM}}$  that captures the joint distribution of coherent sources:

$$\begin{aligned} & \nabla_{(\mathbf{x}_1(t), \dots, \mathbf{x}_K(t))} \log p(\mathbf{x}_1(t), \dots, \mathbf{x}_K(t)) \\ \approx & S_\theta^{\text{MSDM}}((\mathbf{x}_1(t), \dots, \mathbf{x}_K(t)), \sigma(t)). \end{aligned} \quad (3)$$

With the model, it is possible to perform music generation and source separation. *Total (unconditional) generation* integrates Eq. (3) directly, generating all coherent sources  $\{\mathbf{x}_k\}_{k \in [K]}$  composing a track. *Partial (conditional) generation* (i.e., accompaniment generation) fixes a known subset of sources  $\mathbf{x}_{\mathcal{I}} = \{\mathbf{x}_i\}_{i \in \mathcal{I}}$  ( $\mathcal{I} \subset [K]$ ) and generates the complementary subset  $\mathbf{x}_{\bar{\mathcal{I}}}$  ( $\bar{\mathcal{I}} = [K] - \mathcal{I}$ ) coherently. *Source separation* extracts all sources from an observable mixture  $\mathbf{y}(0)$ , integrating, for all  $k$ , the approximate posteriors  $\nabla_{\mathbf{x}_k(t)} \log p(\mathbf{x}_k(t) | \mathbf{y}(0))$ , modeled with Dirac delta likelihood functions. They propose a contextual separator using  $S_\theta^{\text{MSDM}}$  and a weakly supervised separator, using a model  $S_{\theta, k}$  for each source type. When constraining the last source, the weakly supervised separator samples from:

$$S_{\theta, k}(\mathbf{x}_k(t), \sigma(t)) - S_{\theta, K}(\mathbf{y}(0) - \sum_{k=1}^{K-1} \mathbf{x}_k(t), \sigma(t)). \quad (4)$$

## 3. GENERALIZED MULTI-SOURCE DIFFUSION INFERENCE

We train (or use) a text-conditioned diffusion model (Eq. (1))  $S_\theta(\mathbf{y}(t), \mathbf{z}, \sigma(t))$ , with pairs of audio mixtures  $\mathbf{y}(t)$  and associated text embeddings  $\mathbf{z}$ , containing information about the sources present in the mixture. We assume that each text embedding  $\mathbf{z}$  is of the form  $\mathbf{z}_1 \otimes \dots \otimes \mathbf{z}_K$  (more compactly  $\bigotimes_{k=1}^K \mathbf{z}_k$ ), where each  $\mathbf{z}_k$  describes a source  $\mathbf{x}_k$  present in  $\mathbf{y}$  and  $\otimes$  denotes an encoding of concatenated textual information (e.g.,  $\mathbf{z}_1 \otimes \dots \otimes \mathbf{z}_K = E_\phi^{\text{text}}(\mathbf{q}_1, \dots, \mathbf{q}_K)$ , with  $E_\phi^{\text{text}}(\mathbf{q}_k) = \mathbf{z}_k$ ). The idea is to leverage such text embeddings for parameterizing the individual source score functions:

$$\nabla_{\mathbf{x}_k(t)} \log p(\mathbf{x}_k(t) | \mathbf{z}_k) \approx S_\theta(\mathbf{x}_k(t), \mathbf{z}_k, \sigma(t)), \quad (5)$$

even if the model is trained only on mixtures. We devise a set of inference procedures for  $S_\theta$ , called *Generalized Multi-Source Diffusion Inference*, able to solve the tasks of  $S_\theta^{\text{MSDM}}$  in the relaxed data setting.

### 3.1. Total generation

In order to generate a coherent set of sources  $\{\mathbf{x}_k\}_{k \in [K]}$ , described by text embeddings  $\{\mathbf{z}_k\}_{k \in [K]}$ , we can sample from the conditionals  $p(\mathbf{x}_k(t) | \mathbf{x}_{\bar{k}}(t), \mathbf{y}(t), \mathbf{z}_1, \dots, \mathbf{z}_K, \mathbf{z}_1 \otimes \dots \otimes \mathbf{z}_K)$ :

$$\frac{p(\mathbf{x}(t), \mathbf{y}(t) | \mathbf{z}_1, \dots, \mathbf{z}_K, \mathbf{z}_1 \otimes \dots \otimes \mathbf{z}_K)}{p(\mathbf{x}_{\bar{k}}(t), \mathbf{y}(t) | \mathbf{z}_{\bar{k}}, \mathbf{z}_1 \otimes \dots \otimes \mathbf{z}_K)}. \quad (6)$$

First, we develop the numerator in Eq. (6) using the chain rule:

$$\begin{aligned} & p(\mathbf{x}(t), \mathbf{y}(t) | \mathbf{z}_1, \dots, \mathbf{z}_K, \mathbf{z}_1 \otimes \dots \otimes \mathbf{z}_K) \\ &= p(\mathbf{x}_k(t) | \mathbf{z}_k) p(\mathbf{y}(t), \mathbf{x}_{\bar{k}}(t) | \mathbf{x}_k(t), \mathbf{z}_{\bar{k}}, \mathbf{z}_1 \otimes \dots \otimes \mathbf{z}_K) \\ &= p(\mathbf{x}_k(t) | \mathbf{z}_k) p(\mathbf{y}(t) | \mathbf{x}(t)) p(\mathbf{x}_{\bar{k}}(t) | \mathbf{x}_k(t), \mathbf{z}_{\bar{k}}) \\ &\approx p(\mathbf{x}_k(t) | \mathbf{z}_k) p(\mathbf{y}(t) | \mathbf{x}(t)). \end{aligned} \quad (7)$$

We assume independence of the likelihood  $p(\mathbf{y}(t) | \mathbf{x}(t))$  from embeddings and approximate the last equality dropping the unknown term  $p(\mathbf{x}_{\bar{k}}(t) | \mathbf{x}(t), \mathbf{z}_{\bar{k}})$ . We substitute Eq. (7) in Eq. (6), take the gradient of the logarithm with respect to  $\mathbf{x}_k(t)$  and model the likelihood with isotropic Gaussians [20] depending on a variance  $\gamma_{\mathbf{x}_k}^2$ :

$$\begin{aligned} & \nabla_{\mathbf{x}_k(t)} \frac{\log p(\mathbf{x}_k(t) | \mathbf{z}_k) p(\mathbf{y}(t) | \mathbf{x}(t))}{\log p(\mathbf{x}_{\bar{k}}(t), \mathbf{y}(t) | \mathbf{z}_{\bar{k}}, \mathbf{z}_1 \otimes \dots \otimes \mathbf{z}_K)} \\ &= \nabla_{\mathbf{x}_k(t)} \log p(\mathbf{x}_k(t) | \mathbf{z}_k) + \nabla_{\mathbf{x}_k(t)} \log p(\mathbf{y}(t) | \mathbf{x}(t)) \\ &= \nabla_{\mathbf{x}_k(t)} \log p(\mathbf{x}_k(t) | \mathbf{z}_k) + \nabla_{\mathbf{x}_k(t)} \log \mathcal{N}(\mathbf{y}(t) | \sum_{l=1}^K \mathbf{x}_l(t), \gamma_{\mathbf{x}_k}^2 \mathbf{I}) \\ &= \nabla_{\mathbf{x}_k(t)} \log p(\mathbf{x}_k(t) | \mathbf{z}_k) + \frac{1}{\gamma_{\mathbf{x}_k}^2} (\mathbf{y}(t) - \sum_{l=1}^K \mathbf{x}_l(t)). \end{aligned} \quad (8)$$

Applying similar steps we obtain the score of the density on  $\mathbf{y}(t)$  conditioned on  $\mathbf{x}(t)$  (notice the opposite likelihood gradient):

$$\begin{aligned} & p(\mathbf{y}(t) | \mathbf{x}(t), \mathbf{z}_1, \dots, \mathbf{z}_K, \mathbf{z}_1 \otimes \dots \otimes \mathbf{z}_K) \\ &\approx \nabla_{\mathbf{y}(t)} \log p(\mathbf{y}(t) | \bigotimes_{l=1}^K \mathbf{z}_l) + \frac{1}{\gamma_{\mathbf{y}}^2} (\sum_{l=1}^K \mathbf{x}_l(t) - \mathbf{y}(t)). \end{aligned} \quad (9)$$

During inference, we sample from Eqs. (8) and (9) in *parallel*, replacing the gradients of the log-densities with score models (Eq. (5)):

$$\begin{cases} S_{\theta}(\mathbf{x}_k(t), \mathbf{z}_k, \sigma(t)) + \frac{1}{\gamma_{\mathbf{x}_k}^2} (\mathbf{y}(t) - \sum_{l=1}^K \mathbf{x}_l(t)) \\ S_{\theta}(\mathbf{y}(t), \bigotimes_{l=1}^K \mathbf{z}_l, \sigma(t)) + \frac{1}{\gamma_{\mathbf{y}}^2} (\sum_{l=1}^K \mathbf{x}_l(t) - \mathbf{y}(t)). \end{cases} \quad (10)$$

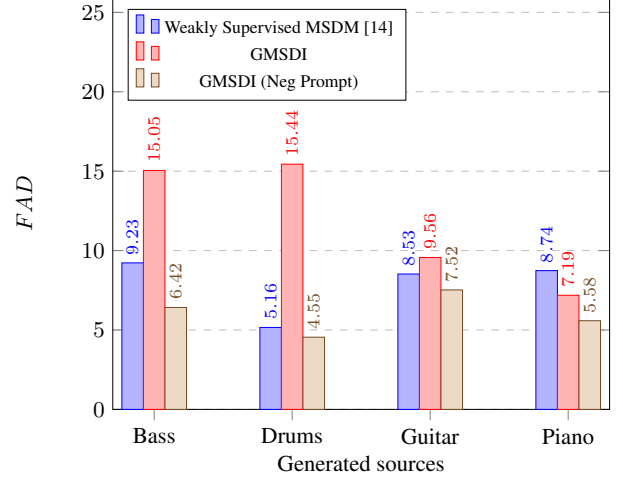
A diagram of the method is illustrated in Figure 1. Given a partition  $\{\mathcal{J}_m\}_{m \in [M]}$  of  $[K]$  containing  $M$  subsets (i.e.,  $\cup_{m \in [M]} \mathcal{J}_m = [K]$ ), we can perform inference more generally with:

$$\begin{cases} S_{\theta}(\sum_{j \in \mathcal{J}_m} \mathbf{x}_j(t), \bigotimes_{j \in \mathcal{J}_m} \mathbf{z}_j, \sigma(t)) + \frac{1}{\gamma_{\mathbf{x}_m}^2} (\mathbf{y}(t) - \sum_{l=1}^K \mathbf{x}_l(t)) \\ S_{\theta}(\mathbf{y}(t), \bigotimes_{l=1}^K \mathbf{z}_l, \sigma(t)) + \frac{1}{\gamma_{\mathbf{y}}^2} (\sum_{l=1}^K \mathbf{x}_l(t) - \mathbf{y}(t)). \end{cases} \quad (11)$$

### 3.2. Partial generation

We can generate accompaniments  $\mathbf{x}_{\mathcal{J}}$  for a given set of sources  $\mathbf{x}_{\mathcal{I}}$ , described by  $\{\mathbf{z}_i\}_{i \in \mathcal{I}}$ , by selecting a set of accompaniment text embeddings  $\{\mathbf{z}_j\}_{j \in \mathcal{J}}$ . We integrate Eqs. (10) for  $j \in \mathcal{J}$ :

$$\begin{cases} S_{\theta}(\mathbf{x}_j(t), \mathbf{z}_j(t), \sigma(t)) + \frac{1}{\gamma_{\mathbf{x}_j}^2} [\mathbf{y}(t) - (\alpha \sum_{i \in \mathcal{I}} \mathbf{x}_i(t) + \beta \sum_{l \in \mathcal{J}} \mathbf{x}_l(t))] \\ S_{\theta}(\mathbf{y}(t), \bigotimes_{l=1}^K \mathbf{z}_l, \sigma(t)) + \frac{1}{\gamma_{\mathbf{y}}^2} [(\alpha \sum_{i \in \mathcal{I}} \mathbf{x}_i(t) + \beta \sum_{l \in \mathcal{J}} \mathbf{x}_l(t)) - \mathbf{y}(t)], \end{cases} \quad (12)$$



**Fig. 2.** FAD (lower is better) between generated sources and Slakh100 test data (200 chunks,  $\sim 12$ s each). Neg Prompt indicates the presence of negative prompting.

with  $\mathbf{x}_i(t)$  ( $i \in \mathcal{I}$ ) sampled from the perturbation kernel in Eq. (2) conditioned on  $\mathbf{x}_i$  and  $\alpha, \beta \in \mathbb{R}$  scaling factors. Using Eq. (11), we can generate the accompaniment mixtures  $\sum_{j \in \mathcal{J}} \mathbf{x}_j$  directly.

### 3.3. Source separation

Source separation can be performed by adapting Eq. (4) to the text-conditioned model. Let an observable mixture  $\mathbf{y}(0)$  be composed by sources described by  $\{\mathbf{z}_k\}_{k \in [K]}$ . We can separate the sources by choosing a constrained source (w.l.o.g. the  $K$ -th) and sampling, for  $k \in [K-1]$ , with:

$$S_{\theta}(\mathbf{x}_k(t), \mathbf{z}_k, \sigma(t)) - S_{\theta}(\mathbf{y}(0) - \sum_{l=1}^{K-1} \mathbf{x}_l(t), \mathbf{z}_K, \sigma(t)). \quad (13)$$

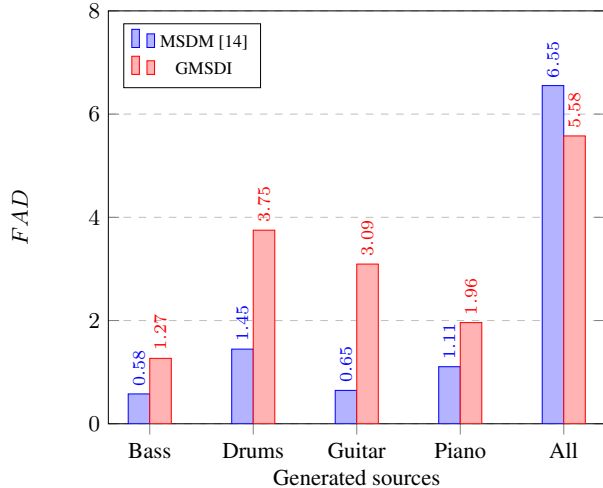
We call this method *GMSDI Separator*. We also define a *GMSDI Extractor*, where we extract the  $k$ -th source  $\mathbf{x}_k$  with:

$$S_{\theta}(\mathbf{x}_k(t), \mathbf{z}_k, \sigma(t)) - S_{\theta}(\mathbf{y}(0) - \mathbf{x}_k(t), \bigotimes_{l \neq k} \mathbf{z}_l, \sigma(t)), \quad (14)$$

constraining the mixture  $\sum_{l \neq k} \mathbf{x}_l(t)$ , complementary to  $\mathbf{x}_k(t)$ .

## 4. EXPERIMENTAL SETUP

To validate our theoretical claims, we train two time-domain MoSai-like [8] diffusion models. The first model is trained on Slakh2100 [21]. Slakh2100 is a dataset used in source separation, containing 2100 multi-source waveform music tracks obtained by synthesizing MIDI tracks with high-quality virtual instruments. We train the diffusion model on mixtures containing the stems Bass, Drums, Guitar, and Piano (the most abundant classes). To condition the diffusion model, we use the `t5-small` pre-trained T5 text-only encoder [10], which inputs the concatenation of the stem labels present in the mixture (e.g., “Bass, Drums” if the track contains Bass and Drums). Given that we know the labels describing the sources inside a mixture at training time, such an approach is weakly supervised. The window size is  $2^{18}$  at 22kHz ( $\sim 12$ s).



**Fig. 3.** FAD (lower is better) results on total and partial generation, with respect to Slakh2100 test mixtures (200 chunks,  $\sim 12$ s each).

The second model is trained on a more realistic dataset, namely MTG-Jamendo [22]. MTG-Jamendo is a music tagging dataset containing over 55000 musical mixtures and 195 tag categories. We train our diffusion model on the `raw_30s/audio-low` version of the dataset, using the first 98 shards for training and the last 2 for validation. The model window is of  $2^{19}$  samples ( $\sim 24$ s) at 22kHz. We condition the model with the pre-trained checkpoint `music_audioset_epoch_15_esc_90.14.pt`<sup>1</sup> of the LAION CLAP contrastive encoder [13]. At training time, we condition the diffusion model with embeddings  $E_{\phi}^{\text{contr}}(\mathbf{y})$  obtained from the training mixtures  $\mathbf{y}$  themselves, resulting in an unsupervised model. At inference time, we use ADPM2<sup>2</sup> [23] with  $\rho = 1$  for generation and AEuler<sup>2</sup> with  $s_{\text{chum}} = 20$  for separation.

## 5. EXPERIMENTAL RESULTS

First, we want to understand whether the model trained on Slakh2100 mixtures can parameterize single sources well. We sample, for each stem, 200 chunks of  $\sim 12$ s, conditioning with embeddings of single stem labels (e.g., “Bass”). Then, we compute the Fréchet Audio Distance (FAD) [24] with VGGish embeddings between such samples and 200 random Slakh2100 test chunks of the same source. In Figure 2, we compare our model against the weakly supervised version of MSDM [14], where a model learns the score function for each stem class (a setting requiring access to clean sources). We notice that single-stem prompting is insufficient for obtaining good FAD results, especially for Bass and Drums, causing silence to be generated. We find negative prompts (Section 2.1) essential for obtaining non-silent results using “Drums, Guitar, Piano” (Bass), “Bass” (Drums), “Bass, Drums” (Guitar), “Bass, Drums” (Piano). In all settings above, we use 150 sampling steps.

Following, we ask how well the model can perform coherent synthesis with GMSDI. In Figure 3, we compute the FAD between 200 random Slakh2100 test mixture chunks ( $\sim 12$ s each) and mixture chunks obtained by summing the model’s generated stems (unconditional) or the generated stems together with the conditioning

<sup>1</sup><https://github.com/LAION-AI/CLAP>

<sup>2</sup><https://github.com/crowsonkb/k-diffusion>

**Table 1.** Grid search over embedding scale  $w$  on 100 chunks ( $\sim 12$ s each) of Slakh2100 test set. Results in SI-SDR<sub>i</sub> (dB – higher is better). The source in parenthesis is the constrained source.

Model	$w = 3.0$	$w = 7.5$	$w = 15.0$	$w = 24.0$
GMSDI Extractor	7.66	<b>9.61</b>	6.00	-0.62
GMSDI Separator (Bass)	8.10	6.72	-1.09	-20.60
GMSDI Separator (Drums)	<b>9.44</b>	8.69	-1.48	-21.62
GMSDI Separator (Guitar)	5.82	4.37	-2.27	-17.49
GMSDI Separator (Piano)	7.60	6.41	-2.68	-16.90

**Table 2.** Quantitative results for source separation on the Slakh2100 test set. Results in SI-SDR<sub>i</sub> (dB – higher is better).

Model	Bass	Drums	Guitar	Piano	All
Demucs + Gibbs (512 steps) [27]	17.16	19.61	17.82	16.32	<b>17.73</b>
Weakly Supervised MSDM [14]	19.36	20.90	14.70	14.13	17.27
MSDM [14]	17.12	18.68	15.38	14.73	16.48
GMSDI Separator	9.76	15.57	9.13	9.57	11.01
GMSDI Extractor	11.00	10.55	9.52	10.13	10.30
Ensamble	11.00	15.57	9.52	10.13	<b>11.56</b>

tracks (conditional). On total generation (All), we set  $\gamma_{\mathbf{y}} = \infty$  and reach  $\sim 1$  lower FAD point, using 600 sampling steps. On partial generation, we sample using 300 steps, setting  $\gamma_{\mathbf{y}} \ll \infty$ , to inform the generated mixture about the conditioning sources. In this scenario, MSDM tends to generate silence. To enforce non-silent results with MSDM, we sample 100 examples for each conditioning chunk and select the sample with the highest  $L_2$  norm.

For source separation, we employ the SI-SDR improvement (SI-SDR<sub>i</sub>) [25] as an evaluation metric and follow the evaluation protocol of [14]. First, we perform a grid search (Table 1) to find a good embedding scale  $w$ . For the GMSDI Separator, we do not use negative prompting, while for the GMSDI Extractor, we only use negative prompts for Bass and Drums. We evaluate on the full Slakh2100 test set with  $w = 3$  and constrained Drums for GMSDI Separator and  $w = 7.5$  for GMSDI Extractor, showcasing results in Table 2. Training only with mixtures (plus associated labels), the ensemble of the two separators reaches 11.56 dB, being zero-shot, i.e., we do not target source separation during training [26].

We release qualitative examples for the Slakh2100 and MTG-Jamendo models on our demo page<sup>3</sup>.

## 6. CONCLUSIONS

We have proposed GMSDI, a compositional music generation method working with any time-domain text-guided diffusion model. The method obtains reasonable generation and separation metrics on Slakh2100, enabling unsupervised compositional music generation for the first time. In future work, we want to extend the technique to latent diffusion models and narrow the gap with supervised methods.

## 7. ACKNOWLEDGEMENTS

This work is supported by the ERC Grant no.802554 (SPECGEO) and PRIN 2020 project no.2020TA3K9N (LEGO.AI). L.C. is supported by the IRIDE grant from DAIS, Ca’ Foscari University of Venice. E.B. is supported by a RAEng/Leverhulme Trust Research Fellowship [grant no. LTRF2223-19-106].

<sup>3</sup><https://github.com/gladia-research-group/gmsdi>

## 8. REFERENCES

- [1] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu, “Neural discrete representation learning,” in *Advances in Neural Information Processing Systems*, 2017, vol. 30.
- [2] Yang Song and Stefano Ermon, “Generative modeling by estimating gradients of the data distribution,” *Advances in neural information processing systems*, vol. 32, 2019.
- [3] Jonathan Ho, Ajay Jain, and Pieter Abbeel, “Denoising diffusion probabilistic models,” *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [4] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole, “Score-based generative modeling through stochastic differential equations,” in *International Conference on Learning Representations*, 2020.
- [5] Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzett, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al., “Musiclm: Generating music from text,” *arXiv preprint arXiv:2301.11325*, 2023.
- [6] Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez, “Simple and controllable music generation,” *arXiv preprint arXiv:2306.05284*, 2023.
- [7] Hugo Flores Garcia, Prem Seetharaman, Rithesh Kumar, and Bryan Pardo, “Vampnet: Music generation via masked acoustic token modeling,” *arXiv preprint arXiv:2307.04686*, 2023.
- [8] Flavio Schneider, Zhijing Jin, and Bernhard Schölkopf, “Mousai: Text-to-music generation with long-context latent diffusion,” *arXiv preprint arXiv:2301.11757*, 2023.
- [9] Haohe Liu, Qiao Tian, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley, “Audioldm 2: Learning holistic audio generation with self-supervised pretraining,” *arXiv preprint arXiv:2308.05734*, 2023.
- [10] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 5485–5551, 2020.
- [11] Ilaria Manco, Emmanouil Benetos, Elio Quenton, and György Fazekas, “Learning music audio representations via weak language supervision,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 456–460.
- [12] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang, “Clap learning audio concepts from natural language supervision,” in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [13] Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov, “Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation,” in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.
- [14] Giorgio Mariani, Irene Tallini, Emilian Postolache, Michele Mancusi, Luca Cosmo, and Emanuele Rodolà, “Multi-source diffusion models for simultaneous music generation and separation,” *arXiv preprint arXiv:2302.02257*, 2023.
- [15] Yuancheng Wang, Zeqian Ju, Xu Tan, Lei He, Zhizheng Wu, Jiang Bian, and Sheng Zhao, “Audit: Audio editing by following instructions with latent diffusion models,” *arXiv preprint arXiv:2304.00830*, 2023.
- [16] Bing Han, Junyu Dai, Xuchen Song, Weituo Hao, Xinyan He, Dong Guo, Jitong Chen, Yuxuan Wang, and Yanmin Qian, “Instructme: An instruction guided music edit and remix framework with latent diffusion models,” *arXiv preprint arXiv:2308.14360*, 2023.
- [17] Chris Donahue, Antoine Caillon, Adam Roberts, Ethan Manilow, et al., “Singsong: Generating musical accompaniments from singing,” *arXiv preprint arXiv:2301.12662*, 2023.
- [18] Jonathan Ho and Tim Salimans, “Classifier-free diffusion guidance,” in *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- [19] Guillaume Sanchez, Honglu Fan, Alexander Spangher, Elad Levi, Pawan Sasanka Ammanamanchi, and Stella Biderman, “Stay on topic with classifier-free guidance,” *arXiv preprint arXiv:2306.17806*, 2023.
- [20] Vivek Jayaram and John Thickstun, “Source separation with deep generative priors,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 4724–4735.
- [21] Ethan Manilow, Gordon Wichern, Prem Seetharaman, and Jonathan Le Roux, “Cutting music source separation some slack: A dataset to study the impact of training data quality and quantity,” in *2019 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2019, pp. 45–49.
- [22] Dmitry Bogdanov, Minz Won, Philip Tovstogan, Alastair Porter, and Xavier Serra, “The mtg-jamendo dataset for automatic music tagging,” in *International Conference on Machine Learning*, 2019.
- [23] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu, “Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 5775–5787, 2022.
- [24] Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi, “Fréchet audio distance: A reference-free metric for evaluating music enhancement algorithms,” in *Interspeech*, 2019.
- [25] Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R. Hershey, “Sdr-half-baked or well done?,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 626–630.
- [26] Jordi Pons, Xiaoyu Liu, Santiago Pascual, and Joan Serra, “Gass: Generalizing audio source separation with large-scale data,” *arXiv preprint arXiv:2310.00140*, 2023.
- [27] Ethan Manilow, Curtis Hawthorne, Cheng-Zhi Anna Huang, Bryan Pardo, and Jesse Engel, “Improving source separation by explicitly modeling dependencies between sources,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 291–295.