

PROMPTSEP: GENERATIVE AUDIO SEPARATION VIA MULTIMODAL PROMPTING

Yutong Wen^{1,2*}, Ke Chen¹, Prem Seetharaman¹, Oriol Nieto¹, Jiaqi Su¹, Rithesh Kumar¹
Minje Kim², Paris Smaragdis³, Zeyu Jin¹, Justin Salamon¹

¹Adobe Research, USA ²University of Illinois Urbana-Champaign, USA ³MIT, USA

ABSTRACT

Recent breakthroughs in language-queried audio source separation (LASS) have shown that generative models can achieve higher separation audio quality than traditional masking-based approaches. However, two key limitations restrict their practical use: (1) users often require operations beyond separation, such as sound removal; and (2) relying solely on text prompts can be unintuitive for specifying sound sources. In this paper, we propose PromptSep to extend LASS into a broader framework for general-purpose sound separation. PromptSep leverages a conditional diffusion model enhanced with elaborated data simulation to enable both audio extraction and sound removal. To move beyond text-only queries, we incorporate vocal imitation as an additional and more intuitive conditioning modality for our model, by incorporating Sketch2Sound as a data augmentation strategy. Both objective and subjective evaluations on multiple benchmarks demonstrate that PromptSep achieves state-of-the-art performance in sound removal and vocal-imitation-guided source separation, while maintaining competitive results on language-queried source separation.

Index Terms— Audio Source Separation, Diffusion Model.

1. INTRODUCTION

Audio source separation aims to isolate specific sounds from an audio mixture. Prior work has approached this task across various domains, where task-specific sound sources needed to be defined, including speech [1, 2], music [3, 4], and general sound events [5–7]. Target source extraction (TSE) is a variant where a specific source is specified by the users in the form of a conditioning signal for the model, such as class labels [8, 9], reference audio [10, 11], visual cues [12, 13], and other modalities [14, 15]. With recent breakthroughs in machine learning, language-queried audio source separation (LASS) has gradually emerged, where natural language serves as the conditioning input for target audio source separation [16–21].

Conventional source separation methods have largely been dominated by mask prediction models that minimize point-wise losses between masked audio and target audio signals [1, 22]. However, masking-based models often introduce distortions, artifacts, and leakages in separation results, as maintaining mask consistency on original audio becomes particularly challenging when multiple sounds overlap. Recently, generative models like diffusion and flow-matching are alternatives to the masking-based methods for separation tasks [23–27], achieving higher separation audio quality. And such methods are also adopted in the LASS settings [20, 21].

Despite promising progress, LASS systems still face two major limitations in real-world scenarios. First, most audio separation models only “extract” the target source from audio mixtures, while treating separation solely as an extraction operator is restrictive.

Users may also wish to remove specific sounds from an audio mixture (i.e., removal), rather than only isolating specific sounds (i.e., extraction). Supporting both extraction and removal within a single framework would better align with real-world needs. Masking-based approaches, however, struggle to generalize to such multi-operator use cases, as they often fail to deliver an accurate one-time mask for multiple sound targets for removal purpose [6, 10, 28]. In contrast, generative approaches offer an alternative. It is worth exploring whether generative models can synthesize high-quality audio outputs that support both extraction and removal operators by explicitly modeling the data distribution of high-fidelity audio samples.

Second, language may not be the most appropriate and effective query for audio sources. Some textual descriptions of sounds, such as “distortion”, “dark transition”, or “light flashing”, are too abstract or ambiguous to precisely specify the target sources. Additionally, a single sound can be described in many different ways, leading to high variability in prompt length, vocabulary, and phrasing. Multiple sources in a mixture may also satisfy a given description (e.g., “distortion” could refer to any harsh or intense sounds), causing unwanted sources to be extracted or removed. Most importantly, this limitation is inherent to language itself, as sounds are naturally perceived by human ears rather than captured by textual descriptions. To overcome this limitation, previous work [10] has explored audio prompts, example samples indicating the target source, to separate similar sounds in the mixture. But this approach is also unintuitive, as users must still provide appropriate reference audio samples.

To address the above limitations, we propose PromptSep, a latent diffusion model for open-vocabulary target audio source separation with multimodal cues. Specifically, we first extend the standard extraction-only separation process with **removal** operator. Second, we incorporate **vocal imitation** as another condition beyond text prompts, where users can mimic the target sound to guide the separation target. The two prompt modalities complement each other and provide more accurate condition. To achieve this, we explore how to augment vocal imitation data by leveraging existing sound effect generation models [29], datasets [30], and modules [27]. With PromptSep, users can perform either sound extraction or removal by leveraging either a textual description or a vocal imitation as a query. Vocal imitation alleviates the need to locate audio samples when the textual description is ambiguous. Our contributions are three-fold:

1. We extend extraction-only source separation with a conditional diffusion model and data simulation pipeline to support sound removal, offering more flexible separation operators.
2. We incorporate vocal imitation as an additional query modality via data augmentation and conditional modules, offering a more intuitive source control than text.
3. We provide a thorough evaluation and demonstrate PromptSep achieves superior performance on sound removal and vocal-imitation-guided extraction, and maintains competitive results on standard language-queried separation.

*Work partly done during an internship at Adobe Research.

2. METHOD

As shown in Figure 1, PromptSep takes an audio mixture as the primary input, along with two conditions: (1) a textual description of the separation target; and (2) a vocal imitation recording of the separation target. Both conditions can be used separately or jointly. Based on these inputs, the model outputs an audio track containing the sources specified by the conditions. In the following subsections, we introduce the construction of these conditions and the model architecture of PromptSep.

2.1. Separation Condition

Text Prompt The textual description of sounds varies in complexity. They may consist of just one keyword identifying the sound, or a long sentence containing multiple sound attributes. Consequently, text prompts can range from a single word to several sentences. To prepare the model for such a condition, we include text prompts of different lengths and descriptive styles for the same sounds during training, simulating the user inputs in real-world scenarios.

Moreover, we extend support for both the number of target sounds and separation operations. Previous LASS approaches typically rely on text prompts that specify a single target sound, and the models are also designed to handle only one target sound at a time. This constraint reduces practicality in real-world scenarios, where users want to extract multiple sounds simultaneously. In addition, users sometimes want to remove specific sounds rather than extract them. Indeed, such removal cases are particularly common, as it is often easier to describe the sounds to remove while preserving the remaining sounds. Previous approaches do not support this sort of interactions, to the best of our knowledge.

To address this gap, we first generalize the target of a single sound source to any subset of sources in the mixture. Then we introduce two text-based operators, **extraction** and **removal**, and use GPT-5 to generate 100 textual templates that paraphrase extraction or removal operations (50 for each). The specification of these templates can be viewed at our accompaniment webpage¹.

We conduct the data simulation pipeline to combine these templates with captions of the target sounds to generate text prompts of either extraction or removal operations. We combine them with corresponding input mixtures and output targets for the language-queried audio separation training.

Vocal Imitation As mentioned in Section 1, certain sounds (e.g., “distortion” and “buzzing”) can be too abstract to describe accurately using text alone, or too ambiguous for the model to only identify the single candidate. To enable more flexible and intuitive sound specification, we introduce an additional guidance modality for separation: vocal imitation. PromptSep can be conditioned on reference audio samples in which a user vocally imitates a target sound, thereby guiding the separation process. This approach provides a more natural and accessible way for users to describe the sounds.

While a few datasets, such as VimSketch [30], contain pairs of corresponding vocal imitation and sound effect samples, these pairs are not temporally aligned, and thus cannot serve as a good vocal imitation condition to identify the target sound in a mixture.

To address this, we leverage the sound effect generation model Sketch2Sound [29] for data augmentation. It can generate sound effect audio samples that match textual descriptions and are temporally aligned with vocal imitation prompts. Using around 12K real vocal imitation samples and their corresponding sound labels from the VimSketch dataset, we use Sketch2Sound to generate around 87K

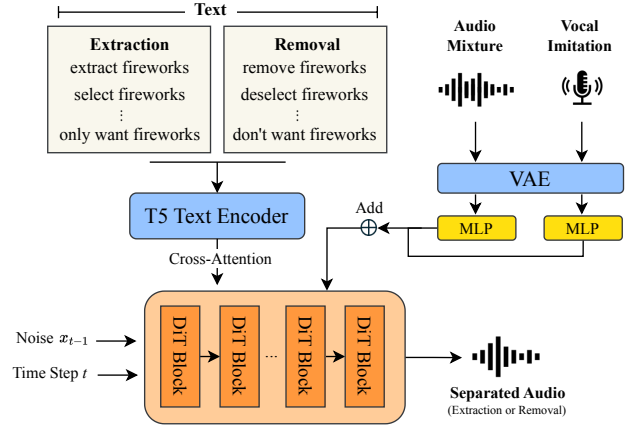


Fig. 1. The model architecture of PromptSep. Text and vocal imitation inputs can be used separately or combined.

temporally aligned sound effect samples as training data for PromptSep. To better simulate real-world conditions, we apply time-shift and pitch-shift augmentations as imitation variances. We also add ambient noise, from 4.36 hours of static noise data collections [31–33], to enhance the model generalization capability.

We note that the Sketch2Sound training data does not require paired vocal imitations and sound effects. Instead, it relies on RMS and pitch curves to guide sound generation, without using actual imitation samples. This leads to a possible exploration: whether curve-based features (RMS and pitch) or actual vocal-imitation raw waveform inputs provide more effective conditioning for audio separation. In Section 4, we address this through ablation studies.

2.2. Model Architecture

We implement PromptSep using a latent diffusion model via diffusion Transformer (DiT) with three input types (text, vocal imitation, and audio mixture), following the specifications of [29, 34, 35]. It contains three main components: (1) a pretrained variational autoencoder (VAE), following the architecture of Descript Audio Codec (DAC) [36], compresses 44.1 kHz mono audio samples into a sequence of continuous 128-dimensional embeddings at a temporal resolution of 40 Hz; (2) a pretrained FLAN-T5 encoder [37] encodes all text prompts; and (3) a DiT model is trained to generate new sequence of embeddings, which are decoded back to waveform via the VAE decoder to reconstruct target separated audio. For conditioning both the mixture and the vocal imitation samples, we adopt an in-place addition mechanism following [27, 29, 38] as they have the same latent dimensionality as the separated audio. Specifically, we apply a single MLP layer respectively to mixture and vocal imitation, and add the resulting embeddings to the noisy latent as input to DiT. The model has around 0.9B parameters.

As we allow the sum of multiple sounds as a target, a trivial solution to the learning objective is to replicate the input mixture to achieve a deceptively low loss. A light noise perturbation to the input can prevent this case, as replicating won’t achieve low loss values.

During training, we randomly drop condition signals for classifier-free guidance (CFG), with the drop rate 10% for text and mixture, but 90% for vocal imitation, as it gets easier to overfit through our experiments. We use a v -prediction framework for training and diffusion probabilistic models (DPM) solvers [39] for sampling. During inference, we set the CFG scale to 1.0.

¹<https://yutongwen.github.io/PromptSep/>

Models	SDRi \uparrow			L2 Mel \downarrow			F1 Decision Error \uparrow			CLAPScore \uparrow			CLAPScore _A \uparrow			FAD _{PANN} \downarrow		
	ACESC	FSD	ASFX	ACESC	FSD	ASFX	ACESC	FSD	ASFX	ACESC	FSD	ASFX	ACESC	FSD	ASFX	ACESC	FSD	ASFX
FlowSep [19]	-4.26	2.05	-2.75	3.06	13.80	4.93	0.88	0.45	0.55	0.24	0.11	0.14	0.74	0.45	0.63	23.70	50.39	38.68
SoloAudio [20]	2.42	14.75	5.15	8.35	2.26	4.73	0.74	0.80	0.53	0.26	0.30	0.19	0.62	0.77	0.61	21.17	5.79	8.82
PromptSep	1.74	10.89	5.65	5.04	7.60	4.23	0.82	0.60	0.60	0.24	0.22	0.18	0.62	0.58	0.66	12.00	19.75	3.19

Table 1. Results under the standard extraction setup, evaluated on AudioCaps + ESC50 (ACESC) [19], FreeSound (FSD) [20], and Adobe Audition Sound Effects (ASFX). Among three benchmarks, ASFX is the only evaluation set that is out-of-domain for all three models.

Models	SDRi \uparrow			L2 Mel \downarrow			F1 Decision Error \uparrow			CLAPScore _A \uparrow			FAD _{PANN} \downarrow		
	ACESC	FSD	ASFX	ACESC	FSD	ASFX	ACESC	FSD	ASFX	ACESC	FSD	ASFX	ACESC	FSD	ASFX
FlowSep [19]	-4.45	-12.44	-9.53	6.30	13.27	5.99	0.76	0.70	0.56	0.44	0.55	0.61	25.49	28.74	20.26
SoloAudio [20]	-1.08	-10.85	-5.50	12.40	37.84	10.87	0.59	0.29	0.36	0.30	0.20	0.45	27.20	87.79	18.54
PromptSep	1.17	-3.34	-3.20	6.40	9.13	4.86	0.80	0.87	0.75	0.54	0.71	0.72	16.27	15.99	3.81
FlowSep* [19]	-4.35	-13.14	-9.36	3.01	6.64	3.34	0.88	0.87	0.73	0.74	0.76	0.74	24.37	13.37	19.78
SoloAudio* [20]	2.26	-9.82	-3.77	8.60	35.31	8.70	0.74	0.44	0.54	0.62	0.32	0.59	23.57	78.45	14.20

Table 2. Results under the removal setup using negative text operators. We include a “upper-anchor” setup where models achieve the same removal effect by separating multiple target sounds as an upper bound of previous baselines (marked with * and in gray).

3. EXPERIMENTS

3.1. Training Datasets

Sound Event We train our model using an internal large collection of licensed sound effect datasets and publicly available, CC-licensed general audio corpora, in total consisting of approximately 1.1M audio samples. Each sound is accompanied with multiple versions of captions varying in length. As described in Section 2, we combine the captions with the text operator templates to form the final text input. Figure 1 shows example final text input.

Vocal Imitation We curate a new dataset, VimSketchGen, consisting of 87,171 pairs of aligned vocal imitations and sound effects. This dataset originally contains 12,453 vocal imitations from the VimSketch dataset, and each of them is paired with 7 corresponding sound effects generated using Sketch2Sound, with different median filter sizes $\in \{0, 3, 6, 9, 12, 15, 19\}$. All audio samples in VimSketchGen are 8-second stereo tracks sampled at 44.1 kHz.

Training Specification The input of PromptSep is a 10-second audio mixture by randomly combining 2 to 5 sound events from different categories. The mixing signal-to-noise ratio (SNR) is uniformly sampled between -3 and 10 dB. A random subset of sound events in the mixture is selected as the separation target, with the only exception that if the vocal imitation is chosen as condition, its corresponding sound event alone is used as the target. During training, the model is always conditioned on either text or vocal imitation, but not both simultaneously. While both conditions can be provided at inference time, their combined effect is not explored in this work and is left for future investigation.

3.2. Evaluation Datasets

AudioCaps [40] We follow the evaluation setup from FlowSep [19] to use the AudioCaps test set of 928 audio clips. We treat each audio clip as the target source and mix it with a noise clip randomly selected from the test set at an SNR randomly chosen between -15 dB and 15 dB. The first caption associated with the target audio is used as the query for separation.

ESC50 [41] It is also from the evaluation setup in FlowSep. We use the ESC50 evaluation set of 2000 audio clips. Similarly, each clip is mixed with another randomly selected clip at an SNR of 0 dB.

FSD-Mix [20] We use the test set of FSD-Mix, which contains 1,440 audio mixtures. Each mixture consists of 3 to 5 sound events, mixed at an SNR randomly selected between -10 and 10 dB.

Adobe Audition Sound Effects ² It serves as a completely out-of-domain benchmark, as both PromptSep and previous baselines are not trained on any training, validation, and test sets of it. We create 2000 mixtures and each contains 2 to 5 sound events randomly sampled from the test set, with an SNR between -10 and 10 dB.

VimSketchGen-Mix We use a subset of the VimSketchGen test split, containing 2000 sound event samples with their vocal imitation pairs. Each target sound is mixed with 1 to 3 interference sounds randomly selected from the AudioSet test set [42], using a SNR sampled between -3 and 10 dB. This evaluation set also contains the time-shift, pitch-shift, and noise interference augmentations.

3.3. Baselines

We compare our model against two generative language-queried audio separation models: FlowSep [19] and SoloAudio [20]. FlowSep performs separation in the latent space of a Mel-spectrogram VAE using flow matching, and reconstructs audio using a vocoder. It employs FLAN-T5 for text conditioning. SoloAudio, on the other hand, applies a modified DiT architecture in the VAE latent space and uses CLAP [43] for text embeddings.

3.4. Objective Metrics

We evaluate the performance of separation models using 6 objective metrics. Signal-to-distortion ratio improvement (SDRi) and L2 multi-resolution Mel-spectrogram distance [36] are the most conventional metrics to measure the signal-level differences between the output and the groundtruth. Following [19, 20, 44], we also include CLAPScore, CLAPScore_A, and Frechet Audio Distance (FAD) [45] with embeddings from PANNs [46], to assess the generation audio quality and the semantic correlation among the text prompt, separation audio, and groundtruth audio.

To better evaluate the separation decision error, we propose **F1 Decision Error** as a new metric to evaluate the ability of the model to identify the correct temporal regions of a target sound. Specifically, to obtain the F1 score of decision errors, we first compute the frame-wise RMS energy on both the separated audio and the groundtruth audio. These values are then binarized using a threshold (0.01) to obtain the activity sequences, determining the sound and unsound frames. Finally, we calculate the F1 score between the predicted activity sequence and the groundtruth activity sequence.

²<https://www.adobe.com/products/adobeauditiondcsfx>

Model	Extraction		Removal	
	REL \uparrow	OVL \uparrow	REL \uparrow	OVL \uparrow
Mixture	2.96 \pm 0.08	3.55 \pm 0.07	2.42 \pm 0.08	3.26 \pm 0.08
GT	3.94 \pm 0.07	4.17 \pm 0.06	3.27 \pm 0.08	4.06 \pm 0.06
FlowSep [19]	3.19 \pm 0.07	3.46 \pm 0.07	2.88 \pm 0.08	3.40 \pm 0.07
SoloAudio [20]	3.31 \pm 0.08	3.64 \pm 0.07	2.99 \pm 0.09	3.59 \pm 0.07
PromptSep	3.34 \pm 0.08	3.75 \pm 0.07	3.25 \pm 0.08	3.83 \pm 0.07

Table 3. Mean Opinion Scores (MOS) with standard error for text relevance (REL) and overall quality (OVL). Mixture and GroudTruth (GT) serve as the lower-anchor and upper-anchor results.

4. RESULTS

4.1. Language-queried Target Sound Extraction

Table 1 presents the results for the language-queried target sound extraction setup. Each model is provided with a text description of the target sounds and evaluated on its ability to extract the sounds from a mixture. We compare PromptSep against two baselines on four benchmarks. Due to the page limitation, we aggregate results from AudioCaps and ESC50 as ACESC using a weighted sum over their metrics, proportional to the number of samples in each dataset.

PromptSep achieves the best performance on nearly all metrics on ASFX, with the exception of CLAPScore, where SoloAudio surpasses us by 0.01. This highlights the strong generalization of our model, as ASFX is the only out-of-domain test set for all three models. FlowSep, trained on AudioSet, WavCaps, and VG-GSound, which share similar quality and sources with AudioCaps and ESC50, obtains the best L2 Mel distance, F1 decision error, and CLAPScore_A on ACESC. SoloAudio, trained primarily on FreeSound, achieves the best results on FSD. All of AudioCaps, ESC50, and FSD are out-of-domain for PromptSep, but it yields competitive performance: achieving or nearly matching the best scores in FAD, SDRi, F1, CLAPScore, and CLAPScore_A on ACESC, and in SDRi on FSD. These results further demonstrate the strong generalization ability of PromptSep.

Finally, we note that FlowSep performs poorly on SDRi. This is likely due to its separation process, which relies on vocoder to reconstruct the generated mel-spectrograms. While its output may preserve acoustic patterns and types of sound events, it can deviate substantially from the ground-truth waveform at the signal level.

4.2. Language-queried Target Sound Removal

Beyond the standard extraction setup, PromptSep also supports the sound removal. While our baselines were not explicitly trained with removal operations, they have been exposed to large-scale textual descriptions and may have implicitly learned some removal capability. Therefore, we also include their results for comparison. As presented in Table 2, PromptSep outperforms the baselines across all datasets and metrics, with the exception of the L2 Mel distance on ACESC. These highlight the strong performance of our model in language-queried target sound removal.

To further account for the limitations of baselines on sound removal, we design an alternative setup. Instead of directly removing the target sound, models are prompted to extract the remaining sounds by providing combined text descriptions of all non-target events, as an equivalent operation. Results under this configuration are shown in Table 2 (gray rows, with models marked by *). PromptSep continues to achieve most of the best scores, even compared against FlowSep* and SoloAudio*. This further demonstrates the strong performance of our model in sound removal.

Conditions	SDRi \uparrow	L2 Mel \downarrow	F1 Decision Error \uparrow	CLAPScore _A \uparrow	FAD \downarrow
Imitation	9.99	0.92	0.95	0.87	2.19
Pitch+RMS	7.17	3.30	0.84	0.71	6.66

Table 4. Results for the vocal imitation (Imitation) condition, along with the Pitch and RMS condition for ablation analysis.

4.3. Subjective Evaluation

We further conduct a subjective evaluation for both the extraction and removal setups by following the format of DCASE2024 Task 9³, using Relevance (REL) and Overall Sound Quality (OVL). The REL score measures how well the separated audio matches the given language query. The OVL score evaluates the perceived audio quality of the output, including factors such as clarity, naturalness, and artifacts. Both REL and OVL are rated by human annotators using a 5-point Likert scale.

We randomly select 100 samples from ASFX test set as it is the only out-of-domain benchmark for all three models. The same mixtures and text descriptions are used across both the extraction and removal setups to ensure consistency. The total number of participants is 100, with each of samples are rated by at least 4 participants. Results are shown in Table 3. PromptSep achieves the highest scores in both REL and OVL across both setups, demonstrating its strong performance in accurately separating target sounds and maintaining high audio quality.

4.4. Imitation-queried Target Sound Extraction

To the best of our knowledge, no existing system supports vocal imitations as a conditioning input for open-domain source separation. We evaluate our model on the VimSketchGen-Mix with no baseline. Results are presented in Table 4, where our model achieves an SDRi of 9.99 dB, an L2 multi-resolution Mel-spectrogram distance of 0.92, an F1 Decision Error of 0.95, a CLAPScore_A of 0.87, and a FAD score of 2.19. These results indicate strong separation performance and demonstrate that vocal imitation is an effective conditioning signal for source separation.

We also evaluate a variant of our model that uses only frame-wise pitch and RMS features extracted from the vocal imitation as the conditioning input, this setup is trained using the same median filter strategy as in [29], with the median filter size fixed to 8 during inference. While the pitch and RMS-based conditioning yields reasonably strong separation performance, it consistently underperforms the full vocal imitation condition across all evaluation metrics. We attribute this to the complexity of the mixtures, where overlapping sounds are common; in such cases, the raw vocal imitation provides richer information than the limited pitch and RMS features alone.

5. CONCLUSION

PromptSep offers a unified framework for sound extraction and removal that overcomes key limitations of existing LASS systems. Our approach supports both sound extraction and removal within a single model. Furthermore, the integration of vocal imitation as a query modality addresses the ambiguity and limitations of text prompts, offering a more intuitive interface for users. Through comprehensive evaluations, PromptSep demonstrates SOTA performance in sound removal and vocal-imitation-guided separation, while remaining competitive in standard LASS settings.

³<https://dcase.community/challenge2024/task-language-queried-audio-source-separation>

6. REFERENCES

- [1] Y. Luo et al., “Conv-TasNet: Surpassing Ideal Time–Frequency Magnitude Masking for Speech Separation,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, 2019.
- [2] D. Wang et al., “Supervised Speech Separation Based on Deep Learning: An Overview,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, 2018.
- [3] A. Défossez et al., “Music Source Separation in the Waveform Domain,” *CoRR*, vol. abs/1911.13254, 2019.
- [4] S. Rouard, F. Massa, et al., “Hybrid Transformers for Music Source Separation,” in *Proc. ICASSP*, 2023.
- [5] T. Ochiai et al., “Listen to What You Want: Neural Network-Based Universal Sound Selector,” *CoRR*, vol. abs/2006.05712, 2020.
- [6] Q. Kong* et al., “Universal Source Separation with Weakly Labelled Data,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, 2025.
- [7] I. Kavalero et al., “Universal Sound Separation,” in *2019 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2019, pp. 175–179.
- [8] M. Delcroix et al., “Few-Shot Learning of New Sound Classes for Target Sound Extraction,” *CoRR*, vol. abs/2106.07144, 2021.
- [9] K. Chen et al., “MDX-GAN: Enhancing Perceptual Quality in Multi-Class Source Separation via Adversarial Training,” in *Proc. ICASSP*, 2024.
- [10] K. Chen et al., “Zero-Shot Audio Source Separation through Query-Based Learning from Weakly-Labeled Data,” in *Proc. AAAI*, 2022.
- [11] P. Seetharaman et al., “Class-Conditional Embeddings for Music Source Separation,” in *Proc. ICASSP*, 2019.
- [12] C. Li et al., “Target Sound Extraction with Variable Cross-Modality Clues,” in *Proc. ICASSP*, 2023.
- [13] R. Gao et al., “Co-Separating Sounds of Visual Objects,” in *Proc. ICCV*, 2019.
- [14] P. Smaragdis, “User Guided Audio Selection from Complex Sound Mixtures,” in *Proc. of annual ACM symposium on User interface software and technology (UIST)*, 2009.
- [15] N. J. Bryan et al., “ISSE: An Interactive Source Separation Editor,” in *Proc. of the SIGCHI Conference on Human Factors in Computing Systems*, 2014.
- [16] H.-W. Dong et al., “CLIPSep: Learning Text-Queried Sound Separation with Noisy Unlabeled Videos,” 2023.
- [17] X. Liu et al., “Separate What You Describe: Language-Queried Audio Source Separation,” in *Proc. Interspeech*, 2022.
- [18] X. Liu et al., “Separate Anything You Describe,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, 2024.
- [19] Y. Yuan et al., “FlowSep: Language-Queried Sound Separation with Rectified Flow Matching,” in *Proc. ICASSP*, 2025.
- [20] H. Wang et al., “SoloAudio: Target Sound Extraction with Language-Oriented Audio Diffusion Transformer,” in *Proc. ICASSP*, 2025.
- [21] J. Hai et al., “DPM-TSE: A Diffusion Probabilistic Model for Target Sound Extraction,” in *Proc. ICASSP*, 2024.
- [22] T.-A. Hsieh et al., “TGIF: Talker Group-Informed Familiarization of Target Speaker Extraction,” *CoRR*, vol. abs/2507.14044, 2025.
- [23] G. Zhu et al., “A Review on Score-Based Generative Models for Audio Applications,” *CoRR*, vol. abs/2506.08457, 2025.
- [24] G. Zhu et al., “Music Source Separation with Generative Flow,” *IEEE Signal Process. Lett.*, 2022.
- [25] G. Mariani et al., “Multi-Source Diffusion Models for Simultaneous Music Generation and Separation,” in *Proc. ICLR*, 2024.
- [26] Y. C. Subakan et al., “Generative Adversarial Source Separation,” in *Proc. ICASSP*, 2018.
- [27] Y. Wen et al., “User-Guided Generative Source Separation,” in *Proc. ISMIR*, 2025.
- [28] Q. Kong et al., “Source Separation with Weakly Labelled Data: An Approach to Computational Auditory Scene Analysis,” in *Proc. ICASSP*, 2020.
- [29] H. F. García et al., “Sketch2Sound: Controllable Audio Generation via Time-Varying Signals and Sonic Imitations,” in *Proc. ICASSP*, 2025.
- [30] B. Kim et al., “VimSketch Dataset,” 2019.
- [31] K. Kinoshita et al., “The REVERB Challenge: A Common Evaluation Framework for Dereverberation and Recognition of Reverberant Speech,” in *Proc. WASPAA*. IEEE, 2013, pp. 1–4.
- [32] J. Eaton et al., “Estimation of Eoom Acoustic Parameters: The ACE Challenge,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 24, no. 10, pp. 1681–1693, 2016.
- [33] Z. Chen et al., “Structure from Silence: Learning Scene Structure from Ambient Sound,” in *5th Annual Conference on Robot Learning*, 2021.
- [34] Z. Evans et al., “Fast Timing-Conditioned Latent Audio Diffusion,” in *Proc. ICML*, 2024.
- [35] Z. Evans et al., “Long-Form Music Generation with Latent Diffusion,” *CoRR*, vol. abs/2404.10301, 2024.
- [36] R. Kumar et al., “High-Fidelity Audio Compression with Improved RVQGAN,” in *Proc. NeurIPS*, 2023.
- [37] C. Raffel et al., “Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer,” *CoRR*, vol. abs/1910.10683, 2019.
- [38] S.-L. Wu et al., “Music ControlNet: Multiple Time-Varying Controls for Music Generation,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, 2024.
- [39] C. Lu et al., “DPM-Solver++: Fast Solver for Guided Sampling of Diffusion Probabilistic Models,” *Machine Intelligence Research*, pp. 1–22, 2025.
- [40] C. D. Kim et al., “AudioCaps: Generating Captions for Audios in the Wild,” in *Proc. of Conference of the North American Chapter of the Association for Computational Linguistics*, 2019.
- [41] K. J. Piczak, “ESC: Dataset for Environmental Sound Classification,” in *Proc. ACM Multimed.*, 2015.
- [42] J. F. Gemmeke et al., “AudioSet: An Ontology and Human-Labeled Dataset for Audio Events,” in *Proc. ICASSP*, 2017.
- [43] Y. Wu* et al., “Large-Scale Contrastive Language-Audio Pretraining with Feature Fusion and Keyword-to-Caption Augmentation,” in *Proc. ICASSP*, 2023.
- [44] F. Xiao et al., “A Reference-Free Metric for Language-Queried Audio Source Separation Using Contrastive Language-Audio Pretraining,” *CoRR*, vol. abs/2407.04936, 2024.
- [45] K. Kilgour et al., “Fréchet Audio Distance: A Metric for Evaluating Music Enhancement Algorithms,” *CoRR*, vol. abs/1812.08466, 2018.
- [46] Q. Kong et al., “PANNs: Large-Scale Pretrained Audio Neural Networks for Audio Pattern Recognition,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, 2020.