# Guiding Audio Editing with Audio Language Model

**Zitong Lan**    **Yiduo Hao**    **Mingmin Zhao**
University of Pennsylvania

## Abstract

Audio editing plays a central role in VR/AR immersion, virtual conferencing, sound design, and other interactive media. However, recent generative audio editing models depend on template-like instruction formats and are restricted to mono-channel audio. These models fail to deal with declarative audio editing, where the user declares what the desired outcome should be, while leaving the details of editing operations to the system. We introduce SmartDJ, a novel framework for stereo audio editing that combines the reasoning capability of audio language models with the generative power of latent diffusion. Given a high-level instruction, SmartDJ decomposes it into a sequence of atomic edit operations, such as adding, removing, or spatially relocating events. These operations are then executed by a diffusion model trained to manipulate stereo audio. To support this, we design a data synthesis pipeline that produces paired examples of high-level instructions, atomic edit operations, and audios before and after each edit operation. Experiments demonstrate that SmartDJ achieves superior perceptual quality, spatial realism, and semantic alignment compared to prior audio editing methods. Demos are available at project page.

## 1 Introduction

Imagine you recorded a forest on a rainy day and wanted to transform it into the soundscape of a sunny forest. Achieving this transformation requires many edits: removing mismatched elements like rainfall and adding new effects such as bright leaf rustling. Traditionally, audio editing is a *procedural process*: the user specifies how to achieve the goal, step by step, by removing rain, layering new samples, adjusting gain, and so on. Yet in practice, users would prefer just to issue a single, high-level instruction, e.g., *"make this sound like a sunny day"*, and rely on an intelligent audio editor to decide what edits are needed. We call this *declarative editing*: the user declares *what* the desired outcome should be, while leaving the *how*, i.e., the sequence of operations, to the system. Such a declarative paradigm would unlock a wide range of applications in VR/AR, gaming, cinematic post-production, and beyond, where designing and modifying audio scenes is central to immersive experiences.

Recent advances in deep generative models have opened exciting possibilities for text-to-audio generation and language-driven editing Evans et al. (2024a;b); Hai et al. (2024); Heydari et al. (2024); Huang et al. (2023a); Liu et al. (2023a; 2024a); Jia et al. (2024); Liang et al. (2024); Wang et al. (2023); Xu et al. (2024). However, existing audio editors remain limited in two key aspects. First, they rely on templated instructions such as "add the sound of birds" or "remove the sound of rain", which restricts their ability to handle high-level or abstract user instructions. Second, they operate only on monaural (single-channel) audio, discarding spatial features such as interaural time and level differences, which are primary cues for spatial hearing. Without these cues, even semantically correct edits sound flat and fail to deliver immersive experiences.

In contrast, declarative editing requires an editor to bridge from high-level goals to detailed operations automatically. When a user issues instructions such as *"place me in a concert hall"* or *"have this audio in a library"*, the system must reason about which sounds to remove, which to preserve, how to adjust loudness, when to introduce new events, and how to shift spatial position. Current audio editors based on diffusion models cannot achieve this, as they lack the reasoning capacity to interpret these high-level instructions. On the other hand, language models alone are also insufficient: while they can parse the user's request, they lack grounding in the audio itself, making it impossible to decide which sound elements should be suppressed, emphasized, or retained. This gap highlights the need for a framework that jointly reasons over language and audio, combining natural-language instruction understanding with audio-aware analysis and editing.
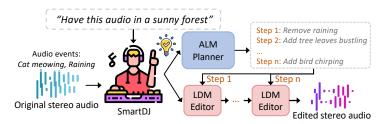
Figure 1: SmartDJ has an Audio Language Model (ALM) that acts as an edit planner to decompose the instructions into atomic steps, guiding the Latent Diffusion Model (LDM) editor to produce high-quality edited audio.

In this work, we present SmartDJ, the first framework for declarative audio editing by introducing Audio Language Models (ALM) into the audio editing loop. This approach is motivated by the recent success of MLLMs in multimodal grounding and reasoning Chu et al. (2023); Ghosh et al. (2025a); Kong et al. (2024a); Liu et al. (2023b); Li et al. (2023a); Zheng et al. (2025). As illustrate in Fig. 1, our key idea is to let the ALM act as a *planner*: it perceives the original audio while interpreting the user's high-level instruction and decomposes it into *a sequence of atomic edit operations*, such as adding or removing a sound, shifting direction, or adjusting volume, etc. These atomic edit operations are then executed sequentially by a conditional Latent Diffusion Model (LDM) as an audio *editor*, realizing the user's high-level goal. This separation of planning and editing transforms audio editing from a procedural task into a declarative one. Moreover, because the intermediate representation is expressed in natural language, users can easily inspect, refine, or override the planned edits interactively.

Training a system to perform declarative editing introduces a fundamental challenge: it requires paired examples of *high-level instructions*, their corresponding *atomic edit sequences*, and the *before-and-after edited audio*. Such data is rarely available, as natural soundscapes are complex and difficult to edit at scale. To address this, we construct a scalable data generation pipeline that provides *controllable audio scenes*. Each scene is assembled from independently parameterized sound events with attributes such as direction and loudness. Within this pipeline, an off-the-shelf LLM acts as the *designer*, producing diverse high-level editing instructions and the corresponding atomic operations. Audio signal processing then serves as the *composer*, rendering each operation by adjusting or mixing sound events to produce the corresponding audio after every edit. Together, this designer–composer pipeline mirrors the planner–editor structure of our model. It produces a large-scale corpus of realistic editing operations, providing the supervision needed to train and evaluate our model.

Experimental results show that SmartDJ delivers superior editing quality and better alignment with high-level user instructions from both subjective metrics and human evaluations. Our LDM also outperforms existing baselines for audio editing. Ablation studies show that our ALM can effectively reason about and decompose high-level user instructions into a sequence of editing actions.

In summary, our main contributions are as follows:

- We introduce SmartDJ, the first stereo audio editor capable of interpreting high-level user instructions with an ALM and executing them as precise atomic edit operations through an LDM.
- We introduce the first scalable pipeline for generating editable stereo audio scenes, combining high-level instructions with controllable events to enable reasoning-based audio editors.
- We conduct extensive experiments and user studies with different baseline methods and demonstrate that SmartDJ achieves the highest editing quality for both objective and subjective metrics.

## 2   RELATED WORK

**Audio generation and editing.** With the current advances in deep generative models, lots of methods have achieved high-quality audio generation from text and multi-modal conditions Chen et al. (2024); Evans et al. (2024a); Hai et al. (2024); Huang et al. (2023a); Liu et al. (2023a; 2024a); Wang et al. (2023). Recently, spatial audio generation has attracted more attention Evans et al. (2024a); Heydari et al. (2024); Liu et al. (2025); Sun et al. (2024), providing more realistic and immersive listening experiences. Parallel to these generation efforts, text-guided audio editing also emerged as a powerful tool for modifying existing audio recordings. Audit Wang et al. (2023) introduced an end-to-end diffusion model conditioned on both the input audio and simple, structured text commands, but its reliance on fixed editing templates limits flexibility to interpret high-level user prompts. WavCraft

Liang et al. (2024) leverages GPT-api to parse user instructions, yet it expects fully specified prompts (e.g., "extract baby crying from the audio" or "apply a low-pass filter to the wave crashing sound and add it back"). Some recent work Jia et al. (2024); Manor & Michaeli (2024); Xu et al. (2024) adapts image-editing techniques (e.g., DDPM inversion, Null-text inversions, attention map manipulation) to monaural audio. They require precise token-level swapping or deletion to complete the edit, struggle with high-level instructions. Besides, they offer no support for stereo audio editing that are important in many applications Lan et al. (2024); Xu et al. (2025); Liang et al. (2025). To summarize, none of the existing work can interpret high-level user input to complete the audio editing. Besides, existing frameworks remain confined to monaural outputs and are ill-suited for immersive spatial scenarios.

**Multimodal Large Language Model.** Large language models (LLMs) are remarkable in natural language processing tasks. Provided with multimodal inputs, such as images and audio, multimodal LLMs (MLLMs) et al. (2022); Liu et al. (2023b); Li et al. (2023a); Team (2025); Kong et al. (2024b); Ghosh et al. (2025b) demonstrate exceptional performance across a wide range of downstream visual-language and audio-language tasks. In the vision domain, LLaVA Liu et al. (2023b) enables LLMs to achieve general-purpose visual and language understanding by fine-tuning on a multimodal instruction-following dataset. In the audio domain, LTU Gong et al. (2023) and Audio Flamingo Kong et al. (2024b); Ghosh et al. (2025b) enhance LLMs with the ability to process non-speech sounds and non-verbal speech. With strong capabilities of MLLMs, researchers introduced them into the field of content generation Wu et al. (2024c;b); Koh et al. (2023); Fu et al. (2024); Huang et al. (2023b), grounding Lai et al. (2024); Hao et al. (2024); Cheng et al. (2024), world modeling Wu et al. (2024a); Ge et al. (2024); Mai et al. (2024), and embodied AI Driess et al. (2023); Li et al. (2023b); Mu et al. (2023). In image generation and editing, various works Koh et al. (2023); Fu et al. (2024); Huang et al. (2023b) use VLM to guide diffusion models. However, in the audio domain, existing methods have not exploited the reasoning capabilities of audio language models.

**Diffusion-based Image Editing.** In contrast to text-to-image generation, image editing focuses on altering specific elements or attributes within an image while preserving the contents of the remaining image. Diffusion models have been widely used in image editing tasks Couairon et al. (2022); Hertz et al. (2022); Hui et al. (2024) by altering the inversion process, which produces a latent representation that can reconstruct the image through the generative process. SDEdit Meng et al. (2022) first adds noise to the source image, and then subsequently denoises the image through the SDE to produce the target image. P2P Hertz et al. (2022) adjusts the cross-attention features according to the difference between the source and target captions to generate the target images. Based on this, IP2P Brooks et al. (2023) finetuned a diffusion model on edit image triplets to enable image editing with simple natural language instructions. Following works Geng et al. (2023); Hui et al. (2024) further scale up the dataset to support more capable and generalized models. Furthermore, Fu et al. (2024); Huang et al. (2023b) use vision-language models to guide diffusion models for image editing tasks.

## 3 METHOD

In this section, we first define the task and notations. We then introduce our proposed framework SmartDJ that combines an Audio Language Model for interpreting high-level editing instructions with a diffusion model for executing sequential audio edits. Finally, we describe a scalable data generation pipeline powered by LLMs to support supervised training and evaluation.

### 3.1 PROBLEM DEFINITION AND NOTATIONS

Let $a_0$ denote the original audio waveform, which contains multiple audio events (e.g., *cat meowing*, *rainfall*), as shown in Fig. 1. We define a *high-level editing instruction* $\mathcal{P}$ as a natural language description of a desired transformation of the overall audio scene. Such instructions are typically declarative: they specify what the desired outcome should be (e.g., "*make it sound like a quiet morning in a sunny forest*" or "*transform this into an indoor library setting*"), but do not explicitly prescribe the individual operations. Since $\mathcal{P}$ is a high-level instruction, it must be decomposed into a detailed sequence of atomic editing steps $\mathcal{S} = \{s_1, s_2, \ldots, s_n\}$, where each step $s_i$ either modifies an existing audio event in $a_0$ or introduces a new event needed to fulfill $\mathcal{P}$. We denote $a_i (i = 1, 2, ..., n)$ as the intermediate audio after applying step $s_i$, where $a_0$ is the original input and $a_n$ is the final edited audio.

Specifically, each step $s_i$ either modifies an existing audio event or introduces a new event required to satisfy $\mathcal{P}$. The atomic editing operations considered in this work are:

- `Add:` Mix a new sound event into the scene (e.g., inject bird chirps).
- `Remove:` Delete an existing sound event (e.g., remove car engine noise).
- `Extract:` Isolate a particular sound event from the original audio while removing others.
- `Turn volume up/down:` Adjust the volume of a specific event.
- `Change direction:` Modify the spatial location of an event.

The goal of the editing is to produce a target edited audio clip $a_n$ by applying the sequence of edits $\mathcal{S}$ to $a_0$. Importantly, this editing formulation must preclude shortcut solutions that ignore the original input (e.g., re-generating an entirely new clip from scratch). Instead, the edited audio $a_n$ must preserve all unedited content from $a_0$ while achieving the requested audio scene transformation. Please refer to Appendix A.2 for more details.

## 3.2 SMARTDJ FRAMEWORK

SmartDJ consists of an Audio Language Model (ALM) and a Latent Diffusion Model (LDM). We leverage the ALM as a planner to interpret high-level instructions and generate a sequence of atomic editing steps. The LDM editor then executes these steps sequentially to transform the original audio. As illustrated in Fig. 2, the ALM takes the original audio $a_0$ and the high-level editing instruction $\mathcal{P}$ as input, and outputs a sequence of atomic editing steps $\mathcal{S} = \{s_1, s_2, ..., s_n\}$. These editing steps are then executed sequentially by the LDM, producing intermediate results $a_1, a_2, \ldots, a_n$, where $a_n$ is the final edited audio. The overall process is formulated as:



Figure 2: SmartDJ framework overview

$$\{s_1, s_2, ..., s_n\} = \mathrm{ALM}(a_0; \mathcal{P}) \tag{1}$$

$$a_i = \mathrm{LDM}(a_{i-1}; s_i), \, i = 1, 2, ..., n \tag{2}$$

## 3.3 AUDIO LANGUAGE MODEL FOR ATOMIC EDITING STEPS GENERATION

The Audio Language Model (ALM) takes as input the original audio clip and the high-level editing instruction and generates a sequence of atomic editing steps. As shown at the top of Fig. 3, we first encode $a_0$ using a pretrained audio encoder (i.e., CLAP Wu et al. (2023)) to obtain an audio embedding $z_a$, which is injected into the ALM via adapter layers. In parallel, the instruction $\mathcal{P}$ is tokenized and encoded as a sequence of embeddings $(p_1, p_2, \ldots, p_k)$, which serve as the textual context for the ALM.

Our ALM is trained in an auto-regressive manner to generate the token sequence corresponding to the atomic editing steps $\mathcal{S}$, by minimizing the following objective:

$$\mathcal{L}_{\mathrm{ALM}} = -\sum_{t=1}^{l} \log P_\theta(r'_t = r_t \mid z_a, r_{1:t-1}, p_{1:k}), \tag{3}$$

where $r$ and $r'$ are the ground truth and predicted text tokens for the atomic editing steps $\mathcal{S}$, $l$ is the length of the tokens, and $\theta$ denotes the model parameters. To enable efficient fine-tuning, we freeze the parameters of the CLAP audio encoder, apply Low-Rank Adaptation (LoRA) Hu et al. (2022) to a small subset of the LLM layersGhosh et al. (2025b), and fully fine-tune the adapter layers.

**Separate training.** We train the ALM and LDM as independent modules rather than end-to-end. This enables human-in-the-loop editing, where users can easily intervene at the level of generated natural-language-based atomic steps before LDM editor inference. Besides, it makes training and computation more efficient and promotes modularity. This design allows for different ALMs or LDMs to be swapped in or out with minimal re-training effort, making it both practical and extensible for diverse editing scenarios.

## 3.4 SEQUENTIAL STEREO AUDIO EDITING WITH LATENT DIFFUSION MODEL

The Latent Diffusion Model (LDM) in our framework performs audio editing conditioned on the atomic editing steps $\mathcal{S}$. To support this, we adopt a latent diffusion architecture Hai et al. (2024); Rombach et al. (2022) and extend it to enable editing of stereo audio with spatial effects.
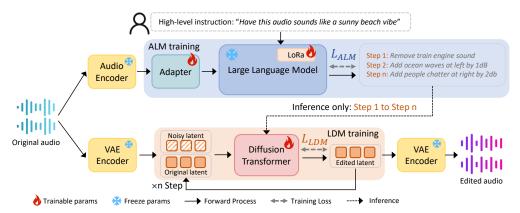
Figure 3: **SmartDJ framework**. Our method incorporates an ALM as an edit planner that understands both the original audio and the high-level instructions to produce atomic edit steps. These atomic steps are then fed into an LDM editor to edit the audio sequentially. The ALM and LDM modules are trained separately.

**Stereo Audio VAE.** Given a stereo audio signal $a \in \mathbb{R}^{2 \times L}$, where $L$ is the number of time-domain samples at two audio channels (left and right). The audio Variational Autoencoder (VAE) encodes $a$ into a latent representation $\hat{a} \in \mathbb{R}^{C \times L'}$, where $C$ and $L'$ denote the number of latent channels and the temporal length of the latent sequence. Similar to DAC Kumar et al. (2023) and Stable Audio Open Evans et al. (2024b), our audio VAE is based on a 1D-CNN autoencoder with a continuous VAE bottleneck and snake activation functions. The resulting latent $\hat{a}$ has a dimension of $C = 128$ and length of $L' = L/480$, resulting a compression ratio of $7.5\times$.

**Latent Diffusion Model:** Our diffusion model conditions on both the text description $s_i$ at the $i$-th editing step and the latent representation of the audio from the previous step, $\hat{a}_{i-1}$, to generate the updated latent $\hat{a}_i$. We use the FLAN-T5 Chung et al. (2024) text encoder $E_{\text{text}}()$ to convert $s_i$ to text embeddings. At each editing step, we initialize a randomly noised latent $\hat{a}'_i \in \mathbb{R}^{C \times L'}$, which is concatenated with $\hat{a}_{i-1}$ to form the input $[\hat{a}_{i-1}; \hat{a}'_i] \in \mathbb{R}^{2C \times L'}$ to the diffusion model. The model is conditioned on $E_{\text{text}}(s_i)$ via cross-attention layers, and the diffusion timestep $t$ is incorporated through a modified AdaLN module Hai et al. (2024) to reduce model parameters. We implement a Diffusion Transformer (DiT) that learns to denoise the latent $\hat{a}'_i$ across multiple timesteps by predicting the added noise. Let $\epsilon$ denote the true added Gaussian noise, and let $\epsilon_\theta(\cdot)$ be the predicted noise output by the model. The training objective is to minimize the following denoising loss:

$$\mathcal{L}_{\text{LDM}} = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), t, s_i, \hat{a}_{i-1}, \hat{a}'_i} \| \epsilon - \epsilon_\theta(t, E_{\text{text}}(s_i), [\hat{a}_{i-1}; \hat{a}'_i]) \|_2 . \qquad (4)$$

During inference, we use DDIM sampling Song et al. (2021) with classifier-free guidance (CFG), which has proven effective for text-guided generation and editing Ho & Salimans (2022). CFG steers the denoising process by interpolating between conditional and unconditional model predictions:

$$\tilde{\epsilon}_\theta = \omega \cdot \epsilon_\theta(t, E_{\text{text}}(s_i), [\hat{a}_{i-1}; \hat{a}'_i]) + (1 - \omega) \cdot \epsilon_\theta(t, \varnothing, [\hat{a}_{i-1}; \hat{a}'_i]), \qquad (5)$$

where $\omega$ is the guidance scale and $\varnothing$ denotes the text embedding of an empty string.

## 3.5 High-level Instruction Audio Editing Dataset Curation

Since no public dataset features audio editing conditioned on high-level instructions, we develop a scalable data generation pipeline, as illustrated in Fig. 4a. For each data point, we first randomly sample $K$ single-event audio clips from public datasets, each labeled with tags such as {"*car engine*", "*bell ring*", "*goat bleat*", ...}. We feed these labels into GPT-4o and prompt it to act as a sound designer: it designs a high-level editing instruction $\mathcal{P}$ that transforms the original mix into a new audio scene (e.g., *"Make this sound like a countryside morning"* or *"Make it sound like a busy train station on a sunny afternoon"*). It then decomposes $\mathcal{P}$ into a sequence of atomic edits $\mathcal{S} = \{s_1, s_2, ..., s_n\}$, including both `add` operations and modifications to existing events (e.g., `remove`, `turn up/down`, `change sound direction`).

Once the sound designer has provided the sequential edit operations, an audio signal processing based composer takes over. To generate edit audio pairs, the composer first synthesizes the initial audio $a_0$ by superposing the $K$ audio clips. Spatial effects are rendered with direction-dependent phase and amplitude on two channels Lan et al. (2024). For each atomic edit $s_i$, we proceed as follows:
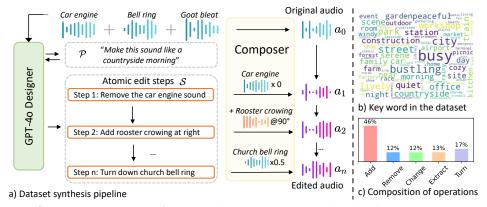
Figure 4: **Scalable data synthesis.** a) pipeline: we sample audio clips from databases with text labels and compose them into original audio $a_0$; These text labels are then fed into GPT-4o, which is prompted to design a high-level instruction $\mathcal{P}$ and generate corresponding atomic steps $\mathcal{S}$. We compose the target audio $a_1, a_2, ..., a_n$ following the atomic steps sequentially with rule-based composer. b) Key words in the high-level instruction. c) The proportion of each single-step edit operation in the dataset.

- If $s_i$ modifies an existing event, we update its volume level or sound direction.
- If $s_i$ is an *add* operation, we retrieve a new clip from the database with a matching label.

Since each sampled audio event is *independently editable*, an edit step $s_i$ that modifies an existing event can be converted into event-level parameter adjustments, without altering any other sources in the mixture $a_{i-1}$. For example, to simulate *"remove the car engine sound"* (Fig. 4), we set the car engine's volume to zero to generate $a_1$. To simulate *"add rooster crowing at right"*, we retrieve a rooster clip, apply the specified spatial effect, and superpose it on $a_1$ to obtain $a_2$. To simulate *"turn down the dog bark"*, we reduce the volume of the corresponding clip. This allows us to generate a complete editing trajectory $a_1, a_2, ..., a_n$ by progressively updating event-level parameters and re-composing the audio scene. Our resulting dataset contains high-level editing instructions, atomic edit sequences, and the audio for the full editing trajectory. More details available in Appx. A.3.

## 4 EXPERIMENT

### 4.1 SETUP

**Dataset**. We use a combination of datasets including AudioCaps Kim et al. (2019), VGGSound Chen et al. (2020), FSD50k Fonseca et al. (2021), ESC50 Piczak (2015), and WavCaps Mei et al. (2024). We adopt a series of dataset cleaning pipelines following previous work Hai et al. (2024); Sun et al. (2024); Wang et al. (2023) by filtering events with noisy data labels or low clap scores. Each audio is trimmed or padded into 10 seconds with a sampling rate of 24 K. We sample 2-5 audio events and use GPT-4o to create 50k training pairs and 1k evaluation pairs of high-level audio editing data to train audio language model and evaluate the whole editing pipeline. We present the keyword in the high-level instructions in Fig. 4b. We also expand the size of single-step editing data pairs $(s_t, a_{t-1}, a_t)$ to 0.5M, where each step is a single atomic operation to train our LDM audio editor. Fig. 4c shows the proportion of each operation. Please find more details in the Appx. A.1.

**Metrics**. To evaluate edit quality and diversity, we use common metrics in audio generation and editing Liu et al. (2023a); Wang et al. (2023), including Fréchet Distance (FD), Kullback-Leibler divergence (KL), Fréchet Audio Distance (FAD), Inception Score (IS), and Log-Spectral Distance (LSD). We use CLAP score to measure the semantic similarities between the edited audio and the text prompt. For spatial audio, we calculate GCC MSE (GCC) based on Generalized Cross-Correlation with Phase Transform (GCC-PHAT), and use StereoCRW Chen et al. (2022) to produce stereo audio features to evaluate CRW MSE (CRW) and Fréchet Stereo Audio Distance (FSAD).

**Baselines**. To evaluate the high-level instruction based audio editing task, we first train an end-to-end version of Audit Wang et al. (2023) that directly predicts the final-step edited audio $a_n$ conditioned on the original audio $a_o$ and high-level instruction $\mathcal{P}$ in one step without ALM. We extend the mono-channel Audit to our binaural setting, where we stack the left and right channels of the mel-spectrograms as the model inputs. We also evaluate various zero-shot and training-required

Figure 5: Examples of ALM's output detailed steps. Our ALM module identifies events in the original audio clips and reasons on the given high-level instruction to produce aligned editing steps.

editing methods based on the ALM's outputs to perform multi-step sequential editing. The zero-shot methods include SDEdit Meng et al. (2022), DDIM Inversion Mokady et al. (2022), ZETA Manor & Michaeli (2024), and AudioEditor Jia et al. (2024). For a fair comparison, we replace these methods' generation backbone with the current SOTA methods. In AudioEditor, we replace Affusion with BEWO Sun et al. (2024) to support binaural editing. In SDEdit, DDIM Inversion, and ZETA, we use Stable-Audio-Open Evans et al. (2024b) as the backbone. For sequential editing with a trainable editor, we also use an Audit baseline trained on single-step audio editing. In addition, we evaluate the same set of baseline methods on single-step audio editing tasks. More details available in Appx. B.1.

**Implementation details.** SmartDJ ALM is initialized from Audio Flamingo 2 Ghosh et al. (2025a) with 3B parameters. During training, we freeze the AF-CLAP module and fine-tune the adapter layers and LLM with a LoRA module for 20 epochs with a batch size of 24. For LDM, it uses velocity prediction with Zero-SNR, and CFG rescaling technique Lin et al. (2024) to adjust the magnitude of the predicted velocity and avoid over-exposure. It is trained on single-step editing data with a batch size of 256 of 500k iterations. 10% text is replaced with empty strings to enable unconditional modeling. The learning rates for the ALM and LDM training are 1e-5 and 5e-5 with AdamW. All experiments are conducted with four NVIDIA L40S GPUs. More details available in Appx. B.2.

## 4.2 RESULTS

**High-level instruction audio editing task.** We first show inference examples from our ALM module in Fig. 5. The ALM-generated atomic editing steps accurately align both the original audio contents and the high-level editing instructions. For example, it correctly removes audio event *engine roar* when transferring audio scene into a *busy family home* vibe. It also removes *people whistle* and adds *pages turning* to enhance the immersion of being in an *old library*. More examples in Appx. C.1.

| Framework | Method | Training | Speed | FD ↓ | FAD ↓ | KL↓ | LSD↓ | IS ↑ | CLAP↑ |
|-----------|--------|----------|-------|------|-------|-----|------|------|-------|
| w/o ALM | Audit | ✓ | 2.07s | 39.6 | 10.07 | 3.13 | 1.96 | 3.29 | 0.125 |
| w/ ALM | SDEdit | ✗ | 301s (74.6s) | <u>27.3</u> | <u>3.73</u> | 3.26 | 2.25 | 6.66 | 0.188 |
| | DDIM | ✗ | 331s (82.1s) | 34.3 | 9.49 | 4.07 | 2.23 | 3.97 | 0.076 |
| | ZETA | ✗ | 356s (88.2s) | 28.8 | 3.75 | 2.93 | 2.24 | 7.72 | <u>0.224</u> |
| | AE | ✗ | 406s (101s) | 27.6 | 5.02 | 3.22 | 2.11 | **8.91** | 0.211 |
| | Audit | ✓ | 11.6s (2.07s) | 29.4 | 5.71 | **2.81** | <u>1.51</u> | 3.95 | 0.197 |
| | SmartDJ (Ours) | ✓ | 13.1s (2.40s) | **14.7** | **1.53** | <u>2.85</u> | **1.42** | <u>8.36</u> | **0.238** |

Table 1: Quantitative results of the whole pipeline from high-level instructions to audio editing. Speed in () is the time for a single-step edit. AE denotes AudioEditor; DDIM denotes DDIM Inversion.

We present the results of the high-level instruction audio editing task in Tab. 1. The first row is the end-to-end Audit baseline, which iss directly trained on high-level instructions and the final target audio, showing the worst performance overall. Since no prior method can interpret the high-level instructions, we use the same set of ALM-generated atomic steps to guide all audio editing models in the multi-step evaluation, including the baselines and our method. We compare the edited audios from each method with 1k reference audios. Our method achieves the lowest metric in FD, FAD, LSD, and comparably low in KL, indicating the smallest discrepancy from the reference audios. It also delivers a high IS metric, showing strong audio quality and diversity. In addition, the highest CLAP score demonstrates the best semantic alignment between the edited audio and instruction.

**Inference time.** We also provide inference speed analysis in Tab. 1. The inference speed reported outside the brackets denote the total time to complete the entire high-level instruction edit, while the values inside indicate the per-round inference time of LDM. On average, SmartDJ's ALM requires 4.8s to generate one set of atomic editing instructions. With multi-round reasoning, SmartDJ
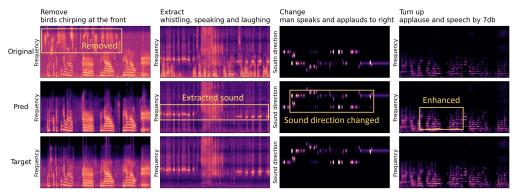
Figure 6: Examples on `Remove`, `Extract`, `Change`, `Turn up` operations. SmartDJ edited audio are closely aligned with the ground truth. Yellow boxes highlight edited sound components.

| Method | Add | | | | | | | | Remove/Extract | | | | | | |
|--------|------|-------|------|------|------|-------|-------|--------|------|-------|------|------|------|-------|--------|
| | FD ↓ | FAD ↓ | KL ↓ | LSD ↓ | IS ↑ | GCC ↓ | CRW ↓ | FSAD ↓ | FD↓ | FAD↓ | KL↓ | LSD↓ | GCC↓ | CRW↓ | FSAD↓ |
| SDEdit | 33.0 | 3.92 | 2.59 | 2.01 | 4.71 | 143.3 | 131.4 | 0.42 | 46.6 | 4.64 | 2.38 | 1.89 | <u>159.5</u> | 157.9 | 0.45 |
| DDIM | 36.9 | 6.28 | 2.64 | 2.02 | 4.57 | <u>131.2</u> | <u>116.1</u> | <u>0.11</u> | 53.9 | 5.69 | 2.75 | 1.85 | 159.6 | <u>138.5</u> | <u>0.29</u> |
| ZETA | 37.6 | <u>3.64</u> | 2.46 | 1.71 | 5.04 | 143.3 | 127.6 | 0.52 | <u>40.6</u> | 3.78 | <u>1.92</u> | 1.99 | 161.7 | 155.9 | 0.43 |
| AE | <u>30.5</u> | 3.66 | 2.13 | 1.84 | <u>5.51</u> | 133.5 | 145.8 | 0.55 | 47.1 | <u>3.65</u> | 2.43 | 2.01 | 164.0 | 165.6 | 0.58 |
| Audit | 36.5 | 4.49 | <u>1.95</u> | <u>1.42</u> | 4.37 | 145.7 | 136.3 | 0.31 | 54.5 | 7.56 | 2.12 | **1.65** | 189.2 | 204.5 | 1.07 |
| SmartDJ | **22.7** | **1.82** | **1.39** | **1.35** | **5.96** | **76.5** | **41.3** | **0.03** | **26.0** | **2.57** | **1.03** | <u>1.71</u> | **15.9** | **5.5** | **0.02** |

(a) Audio editing operation `add` and `remove/extract`

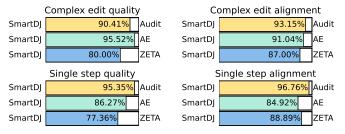| Method | Turn Up/Down | | | | | | | Change | | | | | | |
|--------|------|-------|------|------|------|------|-------|------|-------|------|------|------|------|-------|
| | FD↓ | FAD↓ | KL↓ | LSD↓ | GCC↓ | CRW↓ | FSAD↓ | FD↓ | FAD↓ | KL↓ | LSD↓ | GCC↓ | CRW↓ | FSAD↓ |
| Audit | 47.1 | 5.6 | 1.51 | 1.04 | 136.8 | 139.0 | 0.89 | 42.5 | 4.9 | 1.34 | 1.02 | 170.3 | 163.5 | 0.99 |
| SmartDJ | **11.8** | **1.0** | **0.27** | **1.01** | **23.3** | **2.86** | **0.01** | **13.0** | **0.88** | **0.33** | **1.00** | **59.6** | **36.6** | **0.02** |

(b) Audio editing operation `turn up/down`, `change sound direction`

Table 2: Quantitative results on all individual audio editing operations.

completes a full edit in 13.1s, which is significantly faster than training-free baselines. Our approach is slower compared with end-to-end Audit (2.07s) or Audit with ALM (11.6s), but this trade-off yields substantially better editing quality and alignment with target instructions.

**Single-step audio editing.** We present the results of single-step audio editing operations. Tab. 2a shows results on `add`. Our method consistently outperforms baseline methods in the edited audio, including better similarities to the ideal target audio (lowest FD, FAD and KL), and higher quality and diversity shown by the highest IS. Furthermore, stronger spatial metrics also indicate that the stereo audio characteristics are preserved better by SmartDJ.

Tab. 2a also shows the performance on `remove` and `extract` tasks. The results clearly indicate that our method delivers the best performance in aligning with the ground truth edited audio across both operations. We show the results on `turn up/down` and `change sound direction` tasks in Tab. 2b. Our method again shows stronger performance over Audit, which demonstrates SmartDJ has better fine-grained manipulation in audio event properties. This is because Audit's VAE operates on mel-spectrograms and discards phase information, which is critical for spatial cues. SmartDJ employs a diffusion transformer, which provides stronger long-range temporal modeling and richer cross-attention conditioning. These architectural upgrades allow edits to be both semantically precise and spatially coherent. Some examples of spectrogram visualization are shown in Fig. 6. More qualitative comparisons can be found in Appx. C.2.

**Human evaluations.** We evaluate the subjective preference via a user study. We provide users with data pairs consisting of the original audio, the editing prompt, SmartDJ edited result, and edited results from a random competing baseline. We ask the user to select the one that has higher audio quality, has better alignment with the text or spatial instructions, and aligns with the original audio. We conducted extensive evaluations with 19 participants and 20 (10 high-level audio editing, 10 single-step editing) data pairs per participant. We separate the evaluation for high-level instruction audio editing and single-step editing in Fig. 7. In both tasks, SmartDJ is much preferred over all competing methods. Our method delivers the best audio quality and the best alignment with both the high-level instruction and the single-step instruction from user perception. More details in Appx. C.3.

**Complex edit quality**
SmartDJ 90.41% Audit
SmartDJ 95.52% AE
SmartDJ 80.00% ZETA

**Complex edit alignment**
SmartDJ 93.15% Audit
SmartDJ 91.04% AE
SmartDJ 87.00% ZETA

**Single step quality**
SmartDJ 95.35% Audit
SmartDJ 86.27% AE
SmartDJ 77.36% ZETA

**Single step alignment**
SmartDJ 96.76% Audit
SmartDJ 84.92% AE
SmartDJ 88.89% ZETA

Figure 7: User study results. In both audio quality and text/audio alignment, SmartDJ is consistently preferred over baselines in the high-level instruction editing task and single-step tasks.

Figure 8: Similarity with original audio after multiple editing rounds. (Log Spec. Distance vs. Number of edit rounds; legend: AudioEditor, ZETA, Audit, SmartDJ)

| Study Objectives | Variation | FD↓ | FAD↓ | KL↓ | LSD↓ | IS↑ | CLAP↑ |
|---|---|---|---|---|---|---|---|
| ALM module | w/o ALM | 23.6 | 3.14 | 2.91 | 1.84 | 4.63 | 0.137 |
| | w/ ALM | **14.7** | **1.53** | **2.85** | **1.42** | **8.36** | **0.238** |
| Editing order | Add → Modify → Remove | 14.9 | 1.59 | 2.88 | 1.48 | 8.09 | 0.234 |
| | Random order | **14.7** | 1.56 | 2.88 | 1.47 | 8.16 | 0.237 |
| | Remove → Modify → Add | **14.7** | **1.53** | **2.85** | **1.42** | **8.36** | **0.238** |
| Extract operation | AudioSep Liu et al. (2024b) | 27.1 | 2.86 | 0.88 | **1.63** | N/A | N/A |
| | SmartDJ | **25.7** | **2.55** | **0.78** | 1.71 | N/A | N/A |

Table 3: Model ablations.

## 4.3 ABLATION STUDIES

**Audio quality over multi-round editing.** Since the high-level instruction audio editing task involves a sequence of editing, unchanged content must remain intact after multiple steps. To evaluate this, we design a "round-trip" edit experiment: we perform the operations "add the sound of $\mathcal{A}$" and "remove the sound of $\mathcal{A}$" on an audio clip for five rounds, where $\mathcal{A}$ is a pseudo audio label. For an ideal audio editor, this sequence of round-trip operation should exactly reconstruct the original audio. We measure the LSD between each round's edited output with the original audio clip (Fig. 8). SmartDJ consistently achieves the lowest LSD, indicating the smallest drift from the original content. This demonstrates that our method preserves the unedited audio content under repeated editing actions.

**Effectiveness of ALM.** We conduct an ablation by removing the ALM module and training a variant of the LDM end-to-end. As shown in Tab. 3, SmartDJ performs significantly worse without ALM, showing the importance of ALM's intermediate reasoning capabilities. ALM enables the model to produce semantically coherent edits aligned with the high-level instruction and the original audio.

**Editing order.** We adopt a simple ordering strategy of remove → modify (volume/direction) → add, which intuitively avoids removing newly inserted content. To test the impact of ordering, we also experiment with two alternatives: randomized and reversed order. As shown in Tab. 3, both alternatives result in only marginal degradation compared to our default ordering. This indicates that the editing order has minimal influence on the final results, suggesting that ALM-generated instructions rarely contain conflicting operations.

**Compare with sound separation model.** We compare SmartDJ with a recent sound separation model AudioSep Liu et al. (2024b) on the Extract operation. Ours achieves comparable or slightly better performance (Tab. 3). This suggests that while our model is designed for general-purpose audio editing, it is competitive on target sound separation tasks.

## 5 DISCUSSION

**Conclusion.** We presented SmartDJ, the first framework for high-level instruction guided stereo audio editing that utilizes the reasoning capability of audio language models and strong editing capabilities of the latent diffusion model. Our approach produces atomic editing steps and executes them sequentially to achieve perceptually realistic stereo audio transformations. Extensive evaluations on both subjective audio metrics and human perceptual studies demonstrate that SmartDJ outperforms prior methods, and preserves spatial fidelity in complex scenes.

**Limitations.** Supporting a new task-specific editing operation on the LDM requires retraining the diffusion model. However, these edits can usually be achieved through a combination of our proposed atomic steps. Besides, a future direction is to implement an end-to-end joint training strategy of ALM and LDM that combines reasoning with audio editing.

## 6  ETHICS STATEMENT.

Our work focuses on audio editing for research and creative applications such as immersive media, conferencing, and sound design. The dataset is generated from publicly available sound event libraries and synthetic mixing, without personal or sensitive recordings. We encourage responsible use aligned with academic and creative purposes.

## 7  REPRODUCIBILITY STATEMENT

We provide detailed descriptions of the model architectures, training objectives, and data generation pipeline in the main paper and appendix. Hyperparameters, training configurations, and dataset construction details are included to ensure reproducibility. Code, pretrained models, and the synthesized dataset will be released upon acceptance to facilitate replication of our results.

## REFERENCES

Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions, 2023. URL https://arxiv.org/abs/2211.09800.

Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 721–725. IEEE, 2020.

Ziyang Chen, David F Fouhey, and Andrew Owens. Sound localization by self-supervised time delay estimation. In *European Conference on Computer Vision*, pp. 489–508. Springer, 2022.

Ziyang Chen, Prem Seetharaman, Bryan Russell, Oriol Nieto, David Bourgin, Andrew Owens, and Justin Salamon. Video-guided foley sound generation with multimodal controls. *arXiv preprint arXiv:2411.17698*, 2024.

Tianheng Cheng, Lin Song, Yixiao Ge, Wenyu Liu, Xinggang Wang, and Ying Shan. Yolo-world: Real-time open-vocabulary object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16901–16911, June 2024.

Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.

Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance, 2022. URL https://arxiv.org/abs/2210.11427.

Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model, 2023. URL https://arxiv.org/abs/2303.03378.

Jean-Baptiste Alayrac et al. Flamingo: a visual language model for few-shot learning, 2022. URL https://arxiv.org/abs/2204.14198.

Zach Evans, CJ Carr, Josiah Taylor, Scott H Hawley, and Jordi Pons. Fast timing-conditioned latent audio diffusion. In *Forty-first International Conference on Machine Learning*, 2024a.

Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Stable audio open. *arXiv preprint arXiv:2407.14358*, 2024b.

Eduardo Fonseca, Xavier Favory, Jordi Pons, Frederic Font, and Xavier Serra. Fsd50k: an open dataset of human-labeled sound events. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:829–852, 2021.

Tsu-Jui Fu, Wenze Hu, Xianzhi Du, William Yang Wang, Yinfei Yang, and Zhe Gan. Guiding instruction-based image editing via multimodal large language models, 2024. URL `https://arxiv.org/abs/2309.17102`.

Zhiqi Ge, Hongzhe Huang, Mingze Zhou, Juncheng Li, Guoming Wang, Siliang Tang, and Yueting Zhuang. WorldGPT: Empowering LLM as multimodal world model. In *ACM Multimedia 2024*, 2024. URL `https://openreview.net/forum?id=G1tsqarGAw`.

Zigang Geng, Binxin Yang, Tiankai Hang, Chen Li, Shuyang Gu, Ting Zhang, Jianmin Bao, Zheng Zhang, Han Hu, Dong Chen, and Baining Guo. Instructdiffusion: A generalist modeling interface for vision tasks, 2023. URL `https://arxiv.org/abs/2309.03895`.

Sreyan Ghosh, Zhifeng Kong, Sonal Kumar, S Sakshi, Jaehyeon Kim, Wei Ping, Rafael Valle, Dinesh Manocha, and Bryan Catanzaro. Audio flamingo 2: An audio-language model with long-audio understanding and expert reasoning abilities. *arXiv preprint arXiv:2503.03983*, 2025a.

Sreyan Ghosh, Zhifeng Kong, Sonal Kumar, S Sakshi, Jaehyeon Kim, Wei Ping, Rafael Valle, Dinesh Manocha, and Bryan Catanzaro. Audio flamingo 2: An audio-language model with long-audio understanding and expert reasoning abilities, 2025b. URL `https://arxiv.org/abs/2503.03983`.

Yuan Gong, Hongyin Luo, Alexander H Liu, Leonid Karlinsky, and James Glass. Listen, think, and understand. *arXiv preprint arXiv:2305.10790*, 2023.

Jiarui Hai, Yong Xu, Hao Zhang, Chenxing Li, Helin Wang, Mounya Elhilali, and Dong Yu. Ezaudio: Enhancing text-to-audio generation with efficient diffusion transformer. *arXiv preprint arXiv:2409.10819*, 2024.

Yiduo Hao, Sohrab Madani, Junfeng Guan, Mohammed Alloulah, Saurabh Gupta, and Haitham Hassanieh. Bootstrapping autonomous driving radars with self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15012–15023, June 2024.

Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control, 2022.

Mojtaba Heydari, Mehrez Souden, Bruno Conejo, and Joshua Atkins. Immersediffusion: A generative spatial audio latent diffusion model. *arXiv preprint arXiv:2410.14945*, 2024.

Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.

Jiawei Huang, Yi Ren, Rongjie Huang, Dongchao Yang, Zhenhui Ye, Chen Zhang, Jinglin Liu, Xiang Yin, Zejun Ma, and Zhou Zhao. Make-an-audio 2: Temporal-enhanced text-to-audio generation. *arXiv preprint arXiv:2305.18474*, 2023a.

Yuzhou Huang, Liangbin Xie, Xintao Wang, Ziyang Yuan, Xiaodong Cun, Yixiao Ge, Jiantao Zhou, Chao Dong, Rui Huang, Ruimao Zhang, and Ying Shan. Smartedit: Exploring complex instruction-based image editing with multimodal large language models, 2023b. URL `https://arxiv.org/abs/2312.06739`.

Mude Hui, Siwei Yang, Bingchen Zhao, Yichun Shi, Heng Wang, Peng Wang, Yuyin Zhou, and Cihang Xie. Hq-edit: A high-quality dataset for instruction-based image editing. *arXiv preprint arXiv:2404.09990*, 2024.

Yuhang Jia, Yang Chen, Jinghua Zhao, Shiwan Zhao, Wenjia Zeng, Yong Chen, and Yong Qin. Audioeditor: A training-free diffusion-based audio editing framework. *arXiv preprint arXiv:2409.12466*, 2024.

Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In *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)*, pp. 119–132, 2019.

Jing Yu Koh, Daniel Fried, and Ruslan Salakhutdinov. Generating images with multimodal language models, 2023. URL https://arxiv.org/abs/2305.17216.

Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio flamingo: A novel audio language model with few-shot learning and dialogue abilities. *arXiv preprint arXiv:2402.01831*, 2024a.

Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio flamingo: A novel audio language model with few-shot learning and dialogue abilities, 2024b. URL https://arxiv.org/abs/2402.01831.

Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar. High-fidelity audio compression with improved rvqgan. *Advances in Neural Information Processing Systems*, 36:27980–27993, 2023.

Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9579–9589, June 2024.

Zitong Lan, Chenhao Zheng, Zhiwei Zheng, and Mingmin Zhao. Acoustic volume rendering for neural impulse response fields. *Advances in Neural Information Processing Systems*, 37:44600–44623, 2024.

Chunyuan Li, Zhe Gan, Zhengyuan Yang, Jianwei Yang, Linjie Li, Lijuan Wang, and Jianfeng Gao. Multimodal foundation models: From specialists to general-purpose assistants, 2023a. URL https://arxiv.org/abs/2309.10020.

Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation, 2023b. URL https://arxiv.org/abs/2312.16217.

Jinhua Liang, Huan Zhang, Haohe Liu, Yin Cao, Qiuqiang Kong, Xubo Liu, Wenwu Wang, Mark D Plumbley, Huy Phan, and Emmanouil Benetos. Wavcraft: Audio editing and generation with large language models. *arXiv preprint arXiv:2403.09527*, 2024.

Susan Liang, Dejan Markovic, Israel D Gebru, Steven Krenn, Todd Keebler, Jacob Sandakly, Frank Yu, Samuel Hassel, Chenliang Xu, and Alexander Richard. Binauralflow: A causal and streamable approach for high-quality binaural speech synthesis with flow matching models. *arXiv preprint arXiv:2505.22865*, 2025.

Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. Common diffusion noise schedules and sample steps are flawed. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 5404–5411, 2024.

Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503*, 2023a.

Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley. Audioldm 2: Learning holistic audio generation with self-supervised pretraining. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024a.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b. URL https://arxiv.org/abs/2304.08485.

Huadai Liu, Tianyi Luo, Qikai Jiang, Kaicheng Luo, Peiwen Sun, Jialei Wan, Rongjie Huang, Qian Chen, Wen Wang, Xiangtai Li, et al. Omniaudio: Generating spatial audio from 360-degree video. *arXiv preprint arXiv:2504.14906*, 2025.

Xubo Liu, Qiuqiang Kong, Yan Zhao, Haohe Liu, Yi Yuan, Yuzhuo Liu, Rui Xia, Yuxuan Wang, Mark D Plumbley, and Wenwu Wang. Separate anything you describe. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024b.

Xinji Mai, Zeng Tao, Junxiong Lin, Haoran Wang, Yang Chang, Yanlan Kang, Yan Wang, and Wenqiang Zhang. From efficient multimodal models to world models: A survey, 2024. URL https://arxiv.org/abs/2407.00118.

Hila Manor and Tomer Michaeli. Zero-shot unsupervised and text-based audio editing using ddpm inversion. *arXiv preprint arXiv:2402.10009*, 2024.

Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D Plumbley, Yuexian Zou, and Wenwu Wang. Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.

Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations, 2022. URL https://arxiv.org/abs/2108.01073.

Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models, 2022. URL https://arxiv.org/abs/2211.09794.

Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and Ping Luo. EmbodiedGPT: Vision-language pre-training via embodied chain of thought. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=IL5zJqfxAa.

Karol J Piczak. Esc: Dataset for environmental sound classification. In *Proceedings of the 23rd ACM international conference on Multimedia*, pp. 1015–1018, 2015.

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.

Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=St1giarCHLP.

Peiwen Sun, Sitong Cheng, Xiangtai Li, Zhen Ye, Huadai Liu, Honggang Zhang, Wei Xue, and Yike Guo. Both ears wide open: Towards language-driven spatial audio generation. *arXiv preprint arXiv:2410.10676*, 2024.

Qwen Team. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

Yuancheng Wang, Zeqian Ju, Xu Tan, Lei He, Zhizheng Wu, Jiang Bian, et al. Audit: Audio editing by following instructions with latent diffusion models. *Advances in Neural Information Processing Systems*, 36:71340–71357, 2023.

Jialong Wu, Shaofeng Yin, Ningya Feng, Xu He, Dong Li, Jianye HAO, and Mingsheng Long. ivideoGPT: Interactive videoGPTs are scalable world models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a. URL https://openreview.net/forum?id=4TENzBftZR.

Jiannan Wu, Muyan Zhong, Sen Xing, Zeqiang Lai, Zhaoyang Liu, Zhe Chen, Wenhai Wang, Xizhou Zhu, Lewei Lu, Tong Lu, Ping Luo, Yu Qiao, and Jifeng Dai. VisionLLM v2: An end-to-end generalist multimodal large language model for hundreds of vision-language tasks. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b. URL https://openreview.net/forum?id=nvYDPF4LJK.

Jianzong Wu, Xiangtai Li, Chenyang Si, Shangchen Zhou, Jingkang Yang, Jiangning Zhang, Yining Li, Kai Chen, Yunhai Tong, Ziwei Liu, and Chen Change Loy. Towards language-driven video inpainting via multimodal large language models, 2024c. URL `https://arxiv.org/abs/2401.10226`.

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)*, pp. 1–5. IEEE, 2023.

Manjie Xu, Chenxing Li, Dan Su, Wei Liang, Dong Yu, et al. Prompt-guided precise audio editing with diffusion models. *arXiv preprint arXiv:2406.04350*, 2024.

Shuyang Xu, Zhiyang Dou, Mingyi Shi, Liang Pan, Leo Ho, Jingbo Wang, Yuan Liu, Cheng Lin, Yuexin Ma, Wenping Wang, et al. Mospa: Human motion generation driven by spatial audio. *arXiv preprint arXiv:2507.11949*, 2025.

Zhiwei Zheng, Dongyin Hu, and Mingmin Zhao. Scalable rf simulation in generative 4d worlds. *arXiv preprint arXiv:2508.12176*, 2025.

# A  DATASET CURATION PROCESS

## A.1  DATASET PREPARATION

We construct our training corpus by merging several publicly-available audio dataset, including VGG-Sound, AudioCaps, WavCaps, ESC-50, and FSD50K. Since some of these sets provide audio captions rather than discrete audio label, we first convert every caption to audio labels with GPT-4o-mini API. We only retain single-labeled audio clip and any clip whose caption maps to multiple events is discarded. A CLAP model scores the semantic correspondence between the audio and its new label. Samples with a similarity score below 0.3 are filtered out. The remaining clips from all sources are finally mixed into one large dataset that we use for subsequent data curation.

## A.2  ATOMIC EDIT ACTIONS

We explain the details on creating the single-step atomic edit data pairs. For this single-step audio editing, we have an original audio $a_{i-1}$, a single atomic edit operation $s_i$. Base on atomic edit operation $s_i$, we can generate the edited audio $a_i$.

***Add.*** Assume the original audio is a mix of audio content $A+B+C$. To *add* a new content into this original audio, we sample a new audio content D from the database and mix the D with the original audio $A+B+C+D$. The atomic template is *"Add the sound of {dog barking} at {right} with {3} db"*. The contents inside the {} can be changed to other sound events or sound attributions. We support various sound direction (left, front and right) and dynamic volume adjustment.

***Remove.*** Given an original mix $A+B+C$, let $B$ be the undesired source. We suppress $B$ so the output becomes $A+C$. THe atomic template is Atomic template: *"Remove the sound of {bird chirping} {at right}"*. The directional phrase in braces is optional. If there are similar audio contents in the same clip, the spatial features enables to manipulate it precisely.

***Extract.*** Starting from the same mix $A+B+C$, we isolate one target source $A$ and mute everything else, yielding only $A$. The atomic template is *"Extract the sound of {speaking} {at the right}"*. The direction is optional.

***Turn up/down*** To change loudness of a specific source $B$, we scale it by $\alpha$, where $\alpha > 1$ is for "up" and $0 < \alpha < 1$ is for "down". The resulted audio clip is $A + \alpha B + C$. The atomic template is *"Turn {up / down} the sound of {engine rev} by {2} dB"*. We also support a dynamic range of volume adjustment.

***Change sound direction.*** We alter only the spatial cues of a specific source $C$, producing an edited version $C'$ while leaving the other tracks untouched: $A+B+C'$. The atomic template used is *"Change the sound of {baby crying} {from front} to {right}"*. The "from" clause could also be omitted.

## A.3  COMPLEX AUDIO EDITING DATASET CURATION

In the dataset curation process, we first sample 2-5 audio labels in the database, with LLM randomly assigned volume and sound directions. We then call the GPT-4o batch API with sound sources and attributions. For each API call, we provide 15 sets of sound sources. Through API call, each set of sound sources will return a single data pair containing high-level editing instruction and the corresponding atomic editing steps. We follow these atomic editing steps to manually generate the step-by-step target edited audio with the rules in A.2. We generate 50K data pairs for complex audio editing for training. We generate 1K data pairs from AudioCaps test/validation set for evaluation.

Starting from these 50K complex audio editing pairs, we collect the generated corresponding step-by-step atomic actions and produce about 200K single step audio editing pairs $(s_i, a_{i-1}, a_i)$. We further scale up the this dataset size to 500K to train the LDM audio editor and we also generate another 1K extra audio data pairs to evaluate the performance of single step audio editing. This scaled-up single step audio editing dataset keeps some instructions with audio captions, improving the LDM's robustness to different audio contents in the atomic edit instructions.

We provide the details of our *Base Prompt* for dataset curation as follow:

**Prompt**

You are an expert in spatial audio editing and sound design.

Your task is to generate complex audio editing instructions based on a given list of sound sources (labels). The sound sources will be provided as a list of full sentences (as strings), not character lists. Treat each sentence as a single atomic sound unit. Do not tokenize or split the sound sources into characters.

For example, if you are given: "a baby crying and a man talking; a bird is chirping; dog barking". You should consider "a baby crying and a man talking" as a complete sound source. "a bird is chirping" is another complete sound source and "dog barking" is another sound source. You then generate the step-by-step audio editing instructions based on the given complex instructions.

Task: you need to first brainstorm a complex audio editing instruction

- Imagine a realistic and creative soundscape editing for the given audios.
- You are not limited to the provided audio contents for the editing instruction. And the complex editing instructions should be brief.
- The complex editing instructions could be a soundscape transformation. For example:
  - - Make this sound like it was recorded in a bookstore
  - - Make this sound like a busy coffee shop
  - - Make this sound like a train station
  - - Make this sound like a forest at night
  - - Make this sound like a beach
  - - Make this sound like a sunny day
  - - Craft this sound to feel like a park
  - - Make this audio sound like a quiet farm
  - - Make this audio sound like a firework show
- Your generated complex instructions can be broader than these provided examples. Use your imaginations!
- But remember to keep them brief, the complex editing instructions ideally should not contain the actual sound sources.
- Then, based on the given sound event, give me the detailed editing instructions.

You then generate detailed editing instructions based on the complex audio editing instruction.

- Ensure the editing makes sense (e.g., no waves in a desert, no sheep indoors, rustling leaves in the forest, seas waves in the beach, no raining in the sunny day).
- You are encouraged to remove the original audio contents, but you are required to maintain at least one sound source, not removing all of them.
- Do not split or partially reference a sound source when applying operations.
- Use a combination of simple operations (but NOT necessarily all of them). For example:
  - - Add (e.g., thunderstorm, cat meowing) (for the add operation, you should add up to two additional sources that best align with the complex editing instructions)
  - - Remove (For remove, you should remove at least one sound sources (up to two sound sources), but you must keep at least one sound source!)
  - - Turn up/down (e.g., turn up/turn down the sound of xxx by xxx db (between 0 and 6db))
  - - Change sound direction.
- Each "target" in remove/turn up/turn down/change must exactly match one of the provided "sound sources"
- For the "add" operation:
  - - The "target" must clearly describe the new sound being added (e.g., "crowd chatter", "rain", "footsteps on gravel").
  - - The "effect" should specify the volume and direction (e.g., "at front by 4dB").
  - - Do NOT use "none", "null", or placeholder values as the target. The target must always be a descriptive label of the added sound.
  - - The added sound should not duplicate any original sound source.
- The same operation can be repeated for multiple targets.
- When doing add, you can also have volume and direction attributes
- When editing existing sound sources, you can also have mixed attributes in terms of the volume and sound directions.

- You must ensure each step logically contributes to the final transformation.

Return the output in the following structured JSON format:

```
{
    "sound sources": ["...", "..."], here you should put the sound sources you are given
    "complex editing instruction": "...",
    "atomic editing steps": [
        {"operation": "add", "target": "...", "effect": "at xxx(left, front, right) by xxxdB"},
        {"operation": "remove", "target": "...", "effect": "None"},
        {"operation": "turn up/turn down", "target": "...", "effect": "xxxdB"},
        {"operation": "change", "target": "...", "effect": "to xxx (left, right, or front)"},
    ]
}
```

Do NOT break down sound sources into individual words or characters. Sound sources are complete strings and must remain so in the JSON.

# B IMPLEMENTATION DETAILS

## B.1 BASELINES IMPLEMENTATION

We present the following baseline methods to evaluate the complex audio editing task. One of the baselines is an end-to-end version of Audit without using ALM. All other baseline methods perform multi-step sequential editing with ALM's atomic editing step outputs.

**End-to-End Audit.** We first train an end-to-end version of Audit that directly predicts the final-step edited audio $a_n$ conditioned on the original audio $a_0$ and high-level instruction $\mathcal{P}$ in one step without ALM. We extend the mono-channel Audit to our binaural setting, where we stack the left and right channels of the mel-spectrograms as the model inputs. Follow the original implementation, we convert 10s of audio into mel-spectrograms with a size of $80 \times 624$. using a hop size of 256, a window size of 1024, and mel-bins of size 80. We use the same model configurations in the original paper.

**SDEdit.** SDEdit is a zero-shot method that does not require training a new audio editing model. It uses an off-the-shelf text-to-audio (TTA) generation model, which we use Stable-Audio-Open, as it supports binaural audio generation. We use the default 200 total diffusion steps and start the reverse process from a timestep of 100. We use a classifier-free-guidance (CFG) scale of 7.5. The target caption is composed by concatenating the individual event captions after each editing step.

**DDIM Inversion.** Similar to SDEdit, we also use Stable-Audio-Open as the TTA generation model. We use the default 200 total diffusion steps and start the reverse process from a timestep of 100. We use a CFG scale of 7.5 for both the target and the source. The source text prompt and the target text prompt are composed by concatenating the individual event captions before and after each editing step, respectively.

**ZETA.** Similar to SDEdit and DDIM Inversion, we also use Stable-Audio-Open as the TTA generation model. We use the default 200 total diffusion steps and start the reverse process from a timestep of 100. We use a CFG scale of 7.5 for the target and a CFG scale of 1 for the source. The source text prompt and the target text prompt are composed by concatenating the individual event captions before and after each editing step, respectively.

**AudioEditor.** Specifically in AudioEditor, we replace Affusion with the spatial generator SpatialSonic in BEWO to support binaural editing. We use the default 100 total editing steps with 5 iterations per step to update text embedding for null-text inversion. For the addition task, we use the default punishment ratio (alpha) of -0.001. For the remove or extract task, we use the default punishment ratio (alpha) of 1.

**Audit with ALM.** To evaluate sequential editing with a trainable editor, we use an Audit baseline trained on single-step audio editing data with input audio $a_{i-1}$, atomic instruction $s_i$, and output edited audio $a_i$. The model and data configurations are the same with the end-to-end Audit baseline variant.

## B.2 SMARTDJ IMPLEMENTATION

**ALM.** Our ALM module is initialized from Audio Flamingo 2. This model contains an AF-CLAP audio encoder module to encode the mono-channel audio. The input 10-second audio is first resampled to 16KHz and transformed into a dense audio features $z_a \in \mathbb{R}^{64 \times 2048}$. The mono-channel CLAP encoder understands the audio semantics and sound events, which is enough to reason the atomic steps to make the edited audio semantically align with the high-level instruction. This audio encoder is followed by a representation transformation layers that expand the model capacity. This module has three self-attention layers to the audio feature representation, each with 8 heads and a dimension of 2048. Following this, gated cross-attention layers are used to condition audio representations on the LLM. The LLM uses Qwen2.5-3B, a decoder-only causal LLM with 3B parameters, 36 hidden layers, and 16 attention heads. During training, we keep both the AF-CLAP and LLM frozen during training. The audio representation transformation layers are fully optimized. We apply LoRA only to the gated cross-attention layers with a rank of 16.

**DiT architecture** We conduct experiments at 24KHz audio sample rate. The waveform latent representation operates at 50Hz and consists of 128 channels. Our DiT model has 24 blocks, each with a channel of 1024, 16 attention heads, and an FFN (Feed-Forward Network) dimension of 4096,

resulting in total of 597M parameters. During training, we begin with a 50K linear warm-up phase and the result is training with a cosine learning rate decay. In the inference, we use 100 diffusion steps and a guidance scale of 4, with guidance rescale of 0.8.

## C    MORE RESULTS

In this section, we provide more qualitative results on the ALM-inferenced editing steps and the LDM edited audio results.

### C.1    RESULTS ON ALM-INFERENCED ATOMIC EDITING STEPS

We show multiple examples on ALM-inferenced atomic editing steps. In each example, the "audio labels" are the ground truth audio labels for reference and they are not provided to the ALM module. Only the corresponding audios are provided to the ALM. We then input the high level instruction to the ALM and it output the atomic editing steps.

---

**Example 1: Make this sound like a quiet afternoon in a garden**

Audio labels: clock tick, bird chirp, wind
High-level instructions: Make this sound like a quiet afternoon in a garden
Atomic editing steps:
- Remove the sound of clock tick
- Turn up the the sound of bird tweet by 3dB
- Add the sound of gentle breeze at front by 2dB

---

**Example 2: Make this sound like a busy city street**

Audio labels: engine rev, church bell ring
High-level instructions: Make this sound like a busy city street
Atomic editing steps:
- Remove the sound of bell ring
- Turn down the sound of engine rev by 2dB
- Add the sound of distant siren at left by 2dB
- Add the sound of traffic noise at front by 3dB

---

**Example 3: Make this sound like a bustling park on a sunny day.**

Audio labels: children scream, insect buzz, bird call, chainsaw run
High-level instructions: Make this sound like a bustling park on a sunny day.
Atomic editing steps:
- Remove the sound of chainsaw run
- Add the sound of laughter at left by 3dB
- Turn down children scream by 2dB
- Change the sound of bird call to front

---

**Example 4: Make this sound like a serene beach**

Audio labels: water waves, whistling, music playing
High-level instructions: Make this sound like a serene beach
Atomic editing steps:
- Remove the sound of whistling
- Turn up the sound of wave crash by 4dB
- Add the sound of seagulls calling at front by 3dB

---

> **Example 5: Make this sound like a busy office**
>
> Audio labels: typewriter tap, power tool drill
> High-level instructions: Make this sound like a busy office
> Atomic editing steps:
>  - Remove the sound of drilling
>  - Turn up the sound of typewriter type by 2dB
>  - Add the sound of phone ringing at right by 3dB

> **Example 6: Make this sound like a military base**
>
> Audio labels: fire truck siren, gunshot, vehicle move
> High-level instructions: Make this sound like a military base
> Atomic editing steps:
>  - Remove the sound of fire truck siren
>  - Turn up the sound of gunshot by 3dB
>  - Add the sound of helicopter blades at right by 2dB

> **Example 7: Make this sound like a calm night in the countryside**
>
> Audio labels: snoring, machine run, wind blow
> High-level instructions: Make this sound like a calm night in the countryside
> Atomic editing steps:
>  - Remove the sound of machine run
>  - Turn down the sound of wind blow by 2dB
>  - Add the sound of owl hooting at right by 2dB

> **Example 8: Make this sound like a bustling urban street**
>
> Audio labels: horn honk, cat, siren blare, male laugh
> High-level instructions: Make this sound like a bustling urban street
> Atomic editing steps:
>  - Remove the sound of cat
>  - Turn down the sound of siren blare by 2dB
>  - Add the sound of people chatter at front by 3dB
>  - Add the sound of traffic noise at left by 2dB

## C.2   RESULTS ON ATOMIC EDITING STEPS

We provide more atomic editing results on *add* in Fig 9. As shown by the comparison with the original audio and the edited one, while baseline method tends to replace the original clips and completely generate a new audio clip, SmartDJ can successfully keep the original contents and *add* new sound events into it.

We show more editing examples on *remove* and *extract* in Fig 10 and 11. Compared with baseline method, SmartDJ can effectively either remove unwanted or extract wanted audio contents. The edited results show good alignment with the ground truth edited audios.

Figure 12 presents additional qualitative results for the *change sound direction* task. The y axis in each figure is the sound direction heatmap. SmartDJ consistently relocates the source to the requested spatial direction, and its outputs align well with the ground truth edited audios. By contrast, Audit can not alter spatial effect, showing the limitations of using spectrogram audio encoder for direction-aware editing.
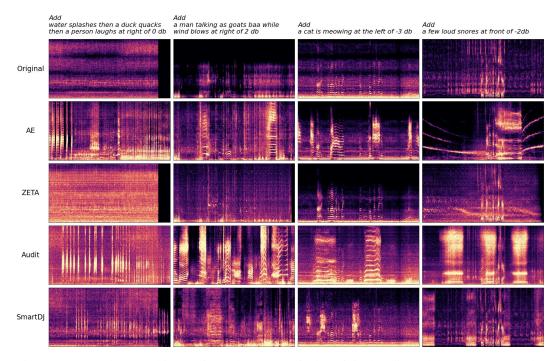
Figure 9: Examples on *add* operation. The top row is the original audio and the rest rows are the edited results. Only SmartDJ can keep the original audio clips while add new audio contents.
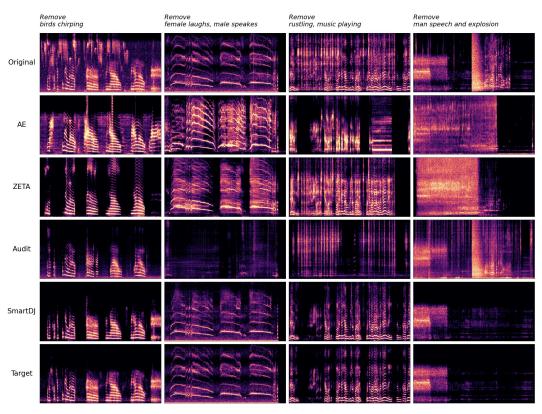


Figure 10: Examples on *remove* operation. The top row is the original audio and the bottom row is the target audio. Only SmartDJ can completely remove unwanted audios parts and keep the remaining part unchanged.
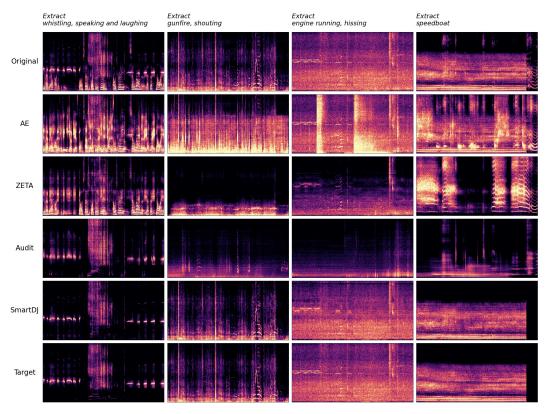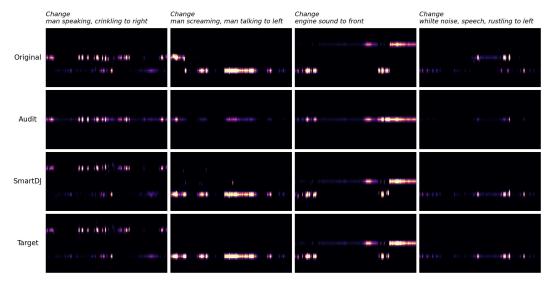
Figure 11: Examples on *extract* operation. The top row is the original audio and the bottom row is the target audio. Only SmartDJ can extract wanted audios that are clean and of high quality.



Figure 12: Examples on *change* operation. The top row is the original audio and the bottom row is the target audio. Only SmartDJ can perfectly edit the sound directions that are matching closely with the target.

## C.3 HUMAN SUBJECTIVE STUDIES DETAILS

We provide more details on the subjective user study. We provide users with data pairs consisting of the original audio, the editing instruction, the SmartDJ edited result, and the edited results from a random competing baseline. In the case of single-step audio editing tasks *remove, extract, turn*

(a) User interface for complex audio editing task     (b) User interface for change direction editing task

Figure 13: The audio pairs, questions, and user interfaces for different audio editing tasks

*up/down, change sound direction* where there exists a ground-truth editing solution, we also provide it as the reference audio as shown in Fig. 13b. We conduct evaluations with 19 participants and 20 (10 complex audio editing, 10 single-step editing) data pairs per participant.

In the 10 complex audio editing question pairs, we ask the user to select between the SmartDJ edited audio and the edited results from a random competing baseline (randomly sampled from AudioEditor, ZETA and Audit) according to the following three questions:

> **Question list for complex audio editing user study**
>
> - Between Audio 1 and Audio 2, which one do you feel aligns with the original audio and the editing instruction better?
> - Between Audio 1 and Audio 2, which one do you feel has the better audio quality?
> - Between Audio 1 and Audio 2, which one do you feel better preserves some of the original audio's contents?

In the single-step editing for *add* operation (3 pairs in total), we compare SmartDJ with a randomly sampled method from three baselines (AudioEditor, ZETA and Audit). We ask the user to answer four questions, with one additional question listed below:

> **Additional question for *add* operation user study**
>
> - Between Audio 1 and Audio 2, which one do you feel the added spatial effect aligns with the text instruction?

For the single-step *remove* and *extract* tasks (4 pairs) we compare SmartDJ with a randomly chosen baseline model. For *turn-up/turn-down* and *change direction* (3 pairs) we benchmark only against AUDIT, since the other baselines cannot perform these operations. Each pair includes a ground-truth reference, and listeners answer the three evaluation questions listed below.

> **Questions list for the *Remove, Extract, Turn up/down, Change sound direction* operation user study**
>
> - Between Audio 1 and Audio 2, which one do you feel aligns with the original audio and the editing instruction better?
> - Between Audio 1 and Audio 2, which one do you feel has the better audio quality?
> - Between Audio 1 and Audio 2, which one do you feel aligns better with the Reference Audio overall?

Across all tasks, SmartDJ is consistently preferred over all baseline methods. For complex audio editing task, SmartDJ receives at least 80% of user votes over the baselines for audio quality, and at least 87% for alignment with the high-level editing instruction and original audio. This shows that

SmartDJ faithfully performs the requested scene transformation while preserving the key elements of the original audio. In the single-step editing task, our method receives more than 77% of user votes for audio quality, and 84% for alignment with the single-step editing text prompt, spatial description, and original or reference audio (when applicable). These results demonstrate that SmartDJ achieves the highest user preference across both quality and alignment metrics, outperforming all competing methods.