
SemanticAudio: Audio Generation and Editing in Semantic Space

Zheqi Dai¹, Guangyan Zhang², Haolin He¹, Xiquan Li³, Jingyu Li², Chunyat Wu¹,
Yiwen Guo^{4,*}, Qiuqiang Kong^{1,*}

¹The Chinese University of Hong Kong

²LIGHTSPEED, ³Shanghai Jiao Tong University

⁴Independent Researcher

Abstract

In recent years, Text-to-Audio Generation has achieved remarkable progress, offering sound creators powerful tools to transform textual inspirations into vivid audio. However, existing models predominantly operate directly in the acoustic latent space of a Variational Autoencoder (VAE), often leading to suboptimal alignment between generated audio and textual descriptions. In this paper, we introduce **SemanticAudio**, a novel framework that conducts both audio generation and editing directly in a high-level semantic space. We define this semantic space as a compact representation capturing the global identity and temporal sequence of sound events, distinct from fine-grained acoustic details. SemanticAudio employs a two-stage Flow Matching architecture: the **Semantic Planner** first generates these compact semantic features to sketch the global semantic layout, and the **Acoustic Synthesizer** subsequently produces high-fidelity acoustic latents conditioned on this semantic plan. Leveraging this decoupled design, we further introduce a training-free text-guided editing mechanism that enables precise attribute-level modifications on general audio without retraining. Specifically, this is achieved by steering the semantic generation trajectory via the difference of velocity fields derived from source and target text prompts. Extensive experiments demonstrate that SemanticAudio surpasses existing mainstream approaches in semantic alignment. Demo available at: <https://semanticaudio1.github.io/>

1 Introduction

Text-to-Audio (TTA) Generation [5, 20, 21] aims to synthesize diverse and high-fidelity auditory content directly from natural language textual prompts. This technology serves as a pivotal creative tool for applications including virtual reality, gaming, film post-production, and human-computer interaction. Recent years have witnessed a paradigm shift in this field, fueled by the scaling of data and model parameters alongside architectural innovations. In particular, the adoption of continuous generative objectives, exemplified by Diffusion Models and Flow Matching, has elevated the fidelity and controllability of synthesized audio.

Most mainstream TTA models perform modeling directly in the acoustic latent space, typically utilizing compressed representations from a Variational Autoencoder (VAE) [20, 21]. While this design excels at preserving low-level acoustic fidelity, it often falls short in high-level semantic understanding. These models frequently struggle to precisely capture the intent in textual prompts, resulting in insufficient *alignment*—defined here as the accurate correspondence between the presence and sequence of auditory events and the text description.

Addressing this limitation requires a clear distinction between the *semantic* and *acoustic* levels of audio. In this work, we define semantics as the high-level conceptual content—specifically the

identity, occurrence, and temporal sequence of sound events—as distinct from fine-grained acoustic details. Audio signals exhibit significant semantic redundancy: high-level semantics are relatively compact and abstract compared to dense acoustic details. Drawing inspiration from two-stage semantic planning approaches in video generation, we hypothesize that directly modeling dense low-level representations in a high-dimensional acoustic latent space is suboptimal for achieving semantic alignment. Instead, the generation process should be decomposed: first accomplishing global content planning in a compact high-level semantic space, followed by the progressive synthesis of acoustic details.

Motivated by this insight, we propose **SemanticAudio**, a novel two-stage Flow Matching-based framework. The core innovation lies in performing the audio generation process via a high-level semantic space. First, a **Semantic Planner** generates compact semantic features from text to sketch the global event layout. Second, conditioned on these features, an **Acoustic Synthesizer** produces high-fidelity VAE latent representations. This design effectively addresses the limitations in high-level semantic modeling inherent in conventional acoustic-space approaches.

Beyond generation, we demonstrate that this decoupled architecture naturally extends to audio editing tasks. While attempting training-free text-guided editing [14, 27] with standard acoustic-space models, we observed unsatisfactory results due to the substantial semantic gap between text and acoustic latents. Leveraging SemanticAudio, we introduce a training-free editing mechanism that operates directly in the semantic space. By steering the generation trajectory via the difference of velocity fields derived from source and target prompts, we achieve precise attribute-level modifications. This stands in contrast to traditional audio editing methods [18, 25], which are typically limited to predefined operations such as addition or deletion. Our mechanism, by fully capitalizing on the advantages of semantic space, enables flexible, text-driven manipulation of high-level semantics on general audio without additional training.

The main contributions of this work are summarized as follows:

SemanticAudio Framework: We propose a two-stage framework comprising a **Semantic Planner** and an **Acoustic Synthesizer**. This architecture performs audio generation directly in a high-level semantic space, effectively decoupling content planning from acoustic synthesis.

Superior Semantic Consistency: By first **sketching the global event layout** in the semantic space, our method achieves substantial outperformance over existing mainstream methods in high-level semantic alignment between generated audio and textual prompts.

Training-free Audio Editing: We introduce a training-free mechanism that enables **flexible, text-driven manipulation of high-level semantics on general audio**. By directly steering the semantic ODE trajectory, this approach achieves versatile attribute-level modifications without requiring additional training or complex inversion steps.

2 Related Work

Text-to-Audio Generation Recent advances in TTA generation have been driven by the scaling of latent diffusion models and Flow Matching frameworks. The prevailing paradigm involves compressing audio into an acoustic latent space via a Variational Autoencoder (VAE) trained on mel-spectrograms, followed by modeling the noise-to-data distribution within this space. Prominent approaches include AudioLDM [20], Make-An-Audio [9], AudioGen [13], and Tango [22]. More recently, Flow Matching-based models such as MeanAudio [17] and LAFMA [7] have demonstrated improved training stability and sampling efficiency. Despite achieving high acoustic fidelity, these models predominantly operate directly in the high-dimensional acoustic latent space. This design conflates fine-grained acoustic details with high-level event logic, often leading to suboptimal semantic alignment, particularly regarding the temporal sequence and structure of sound events described in complex textual prompts.

Semantic Representations in Audio To bridge the semantic gap, prior works have explored various high-level audio representations. Early efforts utilized discrete semantic tokens, as seen in AudioLM [2], or continuous embeddings from contrastive models like CLAP [26] and AudioMAE [8]. However, these representations have largely served as auxiliary conditioning signals rather than the primary generation target. Furthermore, global descriptors like CLAP aggregate information into a single vector, losing the temporal granularity required for detailed event planning. In contrast, the

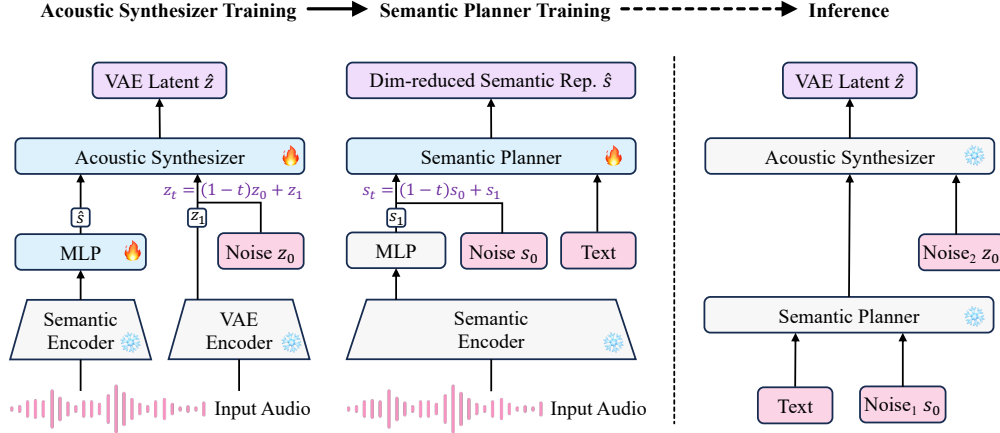


Figure 1: Overview of the SemanticAudio framework. The model employs a two-stage Flow Matching architecture: the **Semantic Planner** first generates low-dimensional semantic latents conditioned on text, followed by the **Acoustic Synthesizer** which produces high-fidelity acoustic latents for VAE decoding.

recent Perception Encoder series, specifically PE-A-Frame [24], provides frame-level semantic embeddings trained with fine-grained audiovisual objectives. By capturing precise temporal alignment between audio frames and textual descriptions, PE-A-Frame offers a temporally rich semantic space suitable for the decoupled planning strategy we propose in this work.

Audio Editing Audio editing approaches typically fall into training-based or training-free categories. Training-based methods, such as Audit [25] and RFM-Editing [6], rely on supervised learning with paired data (e.g., original/edited pairs) to learn specific instruction-following capabilities. While precise, they suffer from high data annotation costs and limited generalization to unseen instructions. Conversely, training-free methods leverage the inherent priors of pre-trained generative models. These often follow an inversion-based paradigm—exemplified by AudioMorphix [18]—where the input audio is inverted to a noise latent and resampled with modified text conditions. However, these approaches are susceptible to inversion errors and struggle to disentangle semantic content from acoustic texture. While inversion-free editing via vector field composition (e.g., FlowEdit [14]) has proven effective in the image domain, its application to audio, particularly within a high-level semantic space, remains underexplored.

Inspirations from Video and Image Generation The concept of decoupling semantic planning from low-level synthesis has gained traction in visual generation. In video generation, SemanticGen [1] demonstrated that generating global layouts in a compact semantic space prior to pixel-level refinement significantly improves coherence in long sequences. Similar "coarse-to-fine" paradigms have been applied to image generation (e.g., RCG [16] and TokensGen [23]). SemanticAudio adapts this insight to the auditory domain, being the first framework to perform audio generation and editing directly within a continuous, high-level semantic space, effectively decoupling content planning from acoustic rendering.

3 SemanticAudio Framework

In this section, we present the detailed architecture of the SemanticAudio framework. As illustrated in Figure 1, our framework effectively decouples text-to-audio generation into two distinct stages: (1) a **Semantic Planner** that sketches the global event layout in a compact semantic space, and (2) an **Acoustic Synthesizer** that produces high-fidelity acoustic details conditioned on the semantic plan. We first detail the representation spaces, followed by the design of the two generative stages.

3.1 Pre-trained VAE and Semantic Representation

SemanticAudio builds upon a pre-trained variational autoencoder (VAE) and a semantic encoder to bridge raw audio waveforms and high-level semantics.

Acoustic Representation. SemanticAudio leverages a variational autoencoder (VAE) to compress a raw audio waveform into a compact acoustic latent space $z \in \mathbb{R}^{T \times C}$. Formally, the encoder E_{VAE} maps the input waveform a to a latent representation $z = E_{\text{VAE}}(a)$, where T denotes the number of acoustic time steps and C represents the channel dimension. The decoder D_{VAE} reconstructs the audio from this latent, $\hat{a} = D_{\text{VAE}}(z)$, ensuring high perceptual fidelity. In this work, we adopt the pre-trained Descript Audio Codec (DAC) [15] as our acoustic VAE.

Semantic Representation. To enable precise control over the temporal layout and content of sound events, we require a semantic encoder E_{sem} capable of extracting continuous, frame-level embeddings $s \in \mathbb{R}^{N \times D}$. Here, N corresponds to the number of semantic frames (determined by the frame rate of the encoder) and D is the embedding dimension. Unlike global descriptors that aggregate information into a single vector (e.g., CLAP [4]), frame-level representations are essential for preserving the fine-grained temporal structure required for event planning.

In this work, we adopt the pre-trained **Perception Encoder** [24]. This model is trained via **fine-grained supervised contrastive learning** on large-scale audio-text datasets. By explicitly aligning audio frames with their corresponding textual descriptions, it excels at **capturing precise semantic-temporal correspondences**. This makes it uniquely capable of tasks requiring detailed event sequencing and distinguishing overlapping sound concepts, providing a robust foundation for our Semantic Planner. To enable tractable modeling in the generative process, we introduce a lightweight MLP projection head P_θ that reduces these high-dimensional embeddings ($D = 1024$) to a compact low-dimensional space:

$$\hat{s} = P_\theta(s) \in \mathbb{R}^{N \times d}, \quad d \ll D. \quad (1)$$

The projection head P_θ , which is randomly initialized, is trained jointly with the Acoustic Synthesizer and remains fixed during the subsequent training of the Semantic Planner. This design ensures that the reduced representations \hat{s} preserve essential semantic content necessary for accurate acoustic synthesis.

3.2 Semantic Planner: Text-to-Semantic Generation

The **Semantic Planner** is responsible for high-level content planning. It learns to generate low-dimensional semantic representations directly conditioned on text prompts, effectively sketching the global event layout.

Given a text prompt y , we extract complementary semantic conditions using two distinct **pre-trained** encoders. We employ the text encoder from CLAP [4] to extract a global sentence embedding c_g , capturing the high-level atmosphere. Simultaneously, we use the Flan-T5 [3] encoder to extract a sequence of token-level embeddings c_d , preserving fine-grained syntactic structures and dynamic instructions. To simplify notation, we denote the full textual conditioning set as $C = \{c_g, c_d\}$. These representations serve as dual inputs to ensure both global coherence and local precision.

The Semantic Planner is a Flow Matching model $\mathcal{F}_{\text{plan}}$ that learns a velocity field $v_\theta^{\text{plan}}(t, \hat{s}_t, C)$ to transport noise $\hat{s}_0 \sim \mathcal{N}(0, I)$ to the target semantic latent \hat{s}_1 . The training objective follows the Flow Matching [19] loss:

$$\mathcal{L}_{\text{FM}}^{\text{plan}} = \mathbb{E}_{t, \hat{s}_t, C} \left\| v_\theta^{\text{plan}}(t, \hat{s}_t, C) - (\hat{s}_1 - \hat{s}_0) \right\|^2, \quad (2)$$

where $\hat{s}_t = (1 - t)\hat{s}_0 + t\hat{s}_1$, and the target $\hat{s}_1 = P_\theta(E_{\text{sem}}(a_{\text{gt}}))$ is obtained using the fixed projection head derived from the Acoustic Synthesizer training.

During inference, we sample $\hat{s}_0 \sim \mathcal{N}(0, I)$ and integrate the ODE $d\hat{s}^t = v_\theta^{\text{plan}}(t, \hat{s}^t, C) dt$ to obtain the planned semantic features \hat{s}_1 .

3.3 Acoustic Synthesizer: Semantic-to-Acoustic Synthesis

The **Acoustic Synthesizer** bridges abstract semantic plans and concrete auditory signals. Conditioned on the semantic features \hat{s}_1 , it learns to synthesize high-fidelity acoustic latents $z_1 \in \mathbb{R}^{T \times C}$.

Training Strategy. A critical aspect of our framework is that **the Acoustic Synthesizer is trained prior to the Semantic Planner**. We jointly optimize the synthesizer and the projection head P_θ . This ensures that the projected semantic features $\hat{s} = P_\theta(E_{\text{sem}}(a_{\text{gt}}))$ retain sufficient information for

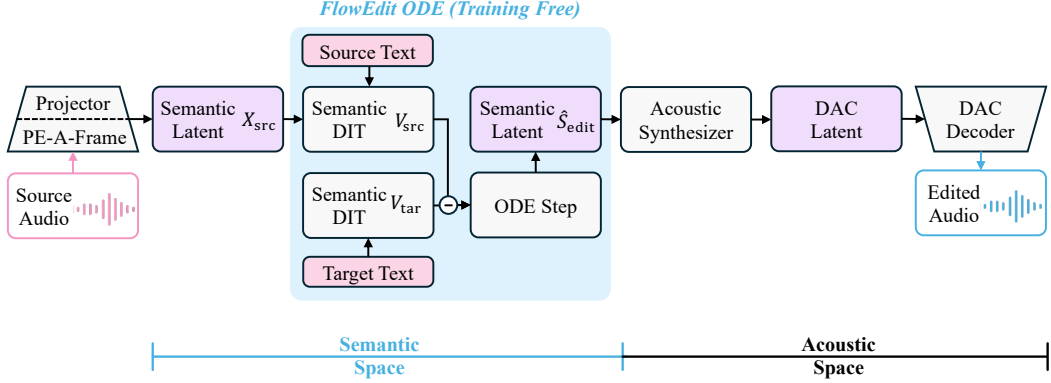


Figure 2: Overview of our training-free text-guided audio editing method. The process leverages the pre-trained velocity fields of the **Semantic Planner** to perform semantic-level editing in the low-dimensional latent space via difference velocity integration, followed by high-fidelity reconstruction using the **Acoustic Synthesizer**. The method requires no additional training, inversion, or optimization.

reconstruction while discarding redundant noise. Once trained, P_θ is frozen to provide target labels for the Semantic Planner.

Modeling. The synthesizer adopts the same Flow Matching formulation as the Semantic Planner (Equation 2). It learns a velocity field $v_\theta^{\text{syn}}(t, z_t, \hat{s}_1)$ to map noise z_0 to the ground-truth acoustic latents $z_1 = E_{\text{VAE}}(a_{\text{gt}})$, conditioned on \hat{s}_1 .

Inference. The full generation pipeline is executed sequentially: we first generate the semantic plan \hat{s}_1 using the Semantic Planner, which then serves as the condition for the Acoustic Synthesizer to generate z_1 . Finally, the waveform is reconstructed via the VAE decoder $\hat{a} = D_{\text{VAE}}(z_1)$.

3.4 Training-Free Text-Guided Audio Editing

A key advantage of our decoupled SemanticAudio framework is its inherent capability for **training-free audio editing**. Unlike pixel- or acoustic-space editing methods that often struggle to disentangle semantic content from background noise, our approach operates directly on the high-level semantic layout. This allows users to modify specific auditory events while preserving the underlying temporal structure, all without requiring model fine-tuning.

Building upon this insight, we introduce a mechanism inspired by FlowEdit [14], as shown in Figure 2. It directly leverages the velocity fields learned by the **Semantic Planner** to perform precise semantic-level modifications, while the **Acoustic Synthesizer** ensures high-fidelity acoustic reconstruction.

Given a source audio a_{src} and its semantic latent \hat{s}_{src} , we define the editing trajectory using a **Delta Velocity Field** v_Δ^t . This field represents the directional difference between the Semantic Planner’s velocity fields conditioned on the target (C_{tgt}) and source (C_{src}) prompts:

$$v_\Delta^t(\hat{s}^t, t) = v_\theta^{\text{plan}}(\hat{s}^t, t, C_{\text{tgt}}) - v_\theta^{\text{plan}}(\hat{s}^t, t, C_{\text{src}}). \quad (3)$$

where C_{src} can be the source text embedding or **null conditioning** if the source text is unavailable.

In practice, to ensure stability against stochastic variations, we approximate v_Δ^t by averaging over N_{avg} noisy realizations at each timestep:

$$v_\Delta^t \approx \frac{1}{N_{\text{avg}}} \sum_{i=1}^{N_{\text{avg}}} \left[v_\theta^{\text{plan}}(\hat{s}_{\text{tgt},i}^t, t, C_{\text{tgt}}) - v_\theta^{\text{plan}}(\hat{s}_{\text{src},i}^t, t, C_{\text{src}}) \right]. \quad (4)$$

Starting from the source semantic latent $\hat{s}^1 = \hat{s}_{\text{src}}$, we integrate this delta field backward to $t = 0$ using standard discrete steps (e.g., Euler method) to obtain the edited semantic latent \hat{s}_{edit} . Finally, \hat{s}_{edit} is decoded by the **Acoustic Synthesizer** into the final audio.

4 Experiments

In this section, we empirically evaluate SemanticAudio on two primary tasks: text-to-audio generation and training-free semantic editing. We aim to verify our core hypothesis: decoupling global semantic planning from acoustic synthesis leads to superior semantic alignment without compromising audio fidelity.

4.1 Datasets and Evaluation Protocols

Training Data. We train both SemanticAudio and all baseline models exclusively on AudioCaps [11], a benchmark dataset containing approximately 46k audio-text pairs (128 hours). Unlike larger weakly-supervised datasets, AudioCaps offers high-quality, human-annotated captions rich in semantic detail. All audio clips are standardized to a 10-second duration via silence padding or truncation.

Test Set for Generation. For standard text-to-audio generation, we utilize the official AudioCaps [11] test split (957 clips). Following standard protocols [17], we randomly select one caption per clip as the generation prompt.

Protocol for Training-Free Editing. Since no standard benchmark exists for open-domain semantic audio editing, we construct a rigorous evaluation set derived from the AudioCaps test split:

Source Selection: We select 50 representative clips as source audio.

Instruction Generation: Using GPT-4, we generate 10 diverse editing instructions per clip (e.g., timbre modification, event replacement) based on the original caption.

Semantic Filtering: To ensure the editing task is non-trivial, we compute CLAP similarity between the original audio and the new instructions. We filter out the top 400 pairs with high similarity (which imply trivial changes) and retain the 100 "hard" prompts that require substantial semantic alteration.

4.2 Implementation Details

Architecture and Conditioning. We implement SemanticAudio using PyTorch, comprising two decoupled DiT-based Flow Matching modules: the **Semantic Planner** and the **Acoustic Synthesizer**. To ensure a rigorous comparison, the control baseline (**Base Model**) shares the exact same backbone configuration: 28 transformer layers, 16 attention heads, and a hidden dimension of 1152 (~610M parameters). Conditioning signals are processed by a dual-encoder setup: **FLAN-T5-small**¹ for text prompts and **PE-A-Frame-small**² for frame-level audio-text alignment. For the acoustic target, we utilize the **DAC-VAE**³ continuous latent space. We set the semantic latent dimension to $d = 64$ based on ablation results. The critical distinction is that the Base Model operates directly in the high-dimensional acoustic space, whereas our method decouples semantic planning from synthesis.

Training Protocol. All models are trained on AudioCaps for 200k iterations with a batch size of 32. We use the AdamW optimizer with a learning of 10^{-4} and a linear warm-up (1k steps) and step decay schedule. Time steps are sampled from a logit-normal distribution ($\mu = 0.4, \sigma = 1.0$).

Inference and Editing. We adopt a differentiated sampling strategy to balance alignment and fidelity: the Semantic Planner utilizes CFG (scale 3.0, 50 steps) to ensure semantic adherence, while the Acoustic Synthesizer uses unguided sampling (scale 1.0, 25 steps). Editing is performed via the training-free ODE path construction using $N_{\text{avg}} = 8$ noisy realizations for robust velocity approximation.

4.3 Evaluation Metrics

We adopt a multi-faceted evaluation protocol to assess the model across three distinct dimensions: acoustic reconstruction, semantic generation, and instruction-guided editing.

Reconstruction Quality. To verify the **Acoustic Synthesizer**'s ability to decode semantic plans into high-fidelity waveforms, we rely on signal-level distance metrics. Following standard DAC protocols, we report the **Mel-spectrogram loss** and **Multi-Scale STFT loss** to measure the precision of spectral reconstruction.

¹FLAN-T5-small: <https://huggingface.co/google/flan-t5-small>

²PE-A-Frame-small: <https://huggingface.co/facebook/pe-a-frame-small>

³DACVAE: <https://huggingface.co/facebook/dacvae-watermarked>

Table 1: Quantitative comparison on AudioCaps test set for text-to-audio generation, including baselines and ablation on semantic latent dimension (last checkpoint only). Lower is better (\downarrow) for FD and KL; higher is better (\uparrow) for IS and LAION-CLAP. **Best** values are bolded.

Model / Config	Dim	FD \downarrow	KL \downarrow	IS \uparrow	CLAP \uparrow
AudioLDM-L-Full	-	29.50	1.68	8.17	0.208
Tango-Full-FT	-	15.64	1.24	8.78	0.291
TangoFlux	-	20.65	1.27	12.81	0.318
Base (Ours)	-	20.64	1.58	9.02	0.338
SemanticAudio (Ours)	128	22.26	1.61	7.38	0.348
SemanticAudio (Ours)	64	21.43	1.54	7.61	0.354
SemanticAudio (Ours)	32	22.90	1.57	7.37	0.340

Table 2: Reconstruction metrics for the **Acoustic Synthesizer** on AudioCaps. We utilize the DAC decoder for waveform reconstruction. Lower is better (\downarrow) for both Mel Loss and STFT Loss. **Best** values are bolded.

Config	Mel Loss \downarrow	STFT Loss \downarrow
Ours ($d = 128$)	0.813	0.794
Ours ($d = 64$)	0.850	0.815
Ours ($d = 32$)	0.926	0.884

Table 3: LAION-CLAP scores (\uparrow) on the "Hard" editing set (50 source clips, 100 filtered prompts). We compare Conditional vs. Unconditional settings. The **Original Source Audio score is 0.2635**. **Best** values are bolded.

Method	Editing CLAP Score (\uparrow)	
	Conditional	Unconditional
Base Model	0.2956	0.2936
Ours ($d = 32$)	0.3191	0.3246
Ours ($d = 128$)	0.3495	0.3493
Ours ($d = 64$)	0.3539	0.3557

Text-to-Audio Generation. We employ standard objective metrics on the full AudioCaps test set to evaluate generation performance. To assess *Semantic Alignment*, we compute **CLAP scores** (using both LAION and MS variants), which quantify the semantic similarity between the generated audio and the input text. For *Fidelity and Diversity*, we measure the **Fréchet Distance (FD)** across feature space of PANNs [12] to evaluate the distributional distance to real audio, alongside the **Inception Score (IS)** to quantify sample quality and diversity. Additionally, **Kullback–Leibler (KL) Divergence** is computed on classifier outputs (PANNs) to ensure the generated event distribution matches the ground truth.

Editing Protocol and Metrics. Since open-domain editing lacks paired ground-truth references, we establish a rigorous *Zero-Shot Editing Benchmark*. We construct a dataset of 100 "hard" editing instances by selecting 50 source clips and generating diverse modification instructions via GPT. Crucially, we filter these instructions based on low CLAP similarity to the source audio, ensuring that the task requires substantial semantic alteration rather than trivial changes. For evaluation, we prioritize **CLAP** (consistency with the edit instruction) and **IS** (overall quality). We also report **FD** to measure distributional adherence to the real audio manifold, using a balanced reference set of source and disjoint real clips. Note that KL Divergence is omitted for editing due to the lack of a defined reference class distribution for open-ended instructions. *Remark:* As the sample size for editing ($N=100$) differs from generation ($N=957$), metric scales (especially FD) are not directly comparable across tasks.

4.4 Results and Analysis

We evaluate SemanticAudio on text-to-audio generation and training-free editing using the AudioCaps test set. We strictly follow standard evaluation protocols, comparing against state-of-the-art baselines and our controlled **Base Model**.

4.4.1 Text-to-Audio Generation

Superior Semantic Alignment. As detailed in Table 1, **SemanticAudio** (with the optimal configuration of $d = 64$) establishes a new state-of-the-art in semantic alignment, achieving a **CLAP**

score of 0.354. This performance significantly surpasses strong baselines such as TangoFlux [10] (0.318) and Tango-Full-FT (0.291). Crucially, our model outperforms the **Base Model** (0.338) by a clear margin. Since the Base Model shares the exact same backbone and parameter count but operates in the acoustic latent space, this result directly validates our core hypothesis: *decoupling semantic planning from acoustic synthesis enables the model to capture high-level textual intent more effectively than modeling directly in a noisy, high-dimensional acoustic space.*

Fidelity vs. Alignment Trade-off. While our Fréchet Distance (FD) is slightly higher than models optimized purely for texture reconstruction (like Tango), it remains highly competitive. We argue that this is a worthwhile trade-off, prioritizing the structural and semantic correctness of the audio (reflected in high CLAP scores) over pixel-perfect acoustic texture matching.

Impact of Semantic Dimension. Our ablation study reveals that a semantic dimension of $d = 64$ strikes the optimal balance. Lower dimensions ($d = 32$) lead to information bottlenecks that degrade alignment, while higher dimensions ($d = 128$) introduce redundancy without proportional performance gains.

4.4.2 Acoustic Reconstruction Quality

Table 2 isolates the performance of the **Acoustic Synthesizer**. The model exhibits excellent reconstruction fidelity (Mel Loss 0.813, STFT Loss 0.794 at $d = 128$). The performance degrades gracefully as the dimension decreases, confirming that our lightweight MLP projection successfully compresses semantic information while retaining sufficient cues for the synthesizer to reconstruct high-fidelity waveforms.

4.4.3 Training-Free Semantic Editing

We evaluate editing performance on the "hard" subset of 100 attribute-level prompts (Table 3).

Zero-Shot Editing Capability. **SemanticAudio** demonstrates remarkable editing control, achieving a Conditional CLAP score of **0.3539**, a substantial leap from the original source audio (0.2635). This confirms the model’s ability to precisely modify audio attributes (e.g., timbre, atmosphere) to match new text instructions without any task-specific tuning.

Robustness to Missing Source Text. A key finding is that the **Unconditional** setting (where source text is null) performs on par with the Conditional setting (0.3557 vs. 0.3539). This suggests our flow-matching-based editing mechanism is highly robust: the model can infer the transformation trajectory purely from the semantic difference between the source audio and target text, even without explicit knowledge of the original caption.

Comparison with Baselines. **SemanticAudio** consistently outperforms the Base Model (approx. 0.29) in editing tasks. The Base Model’s entangled latent space struggles to isolate specific attributes for modification, whereas our decoupled semantic space allows for targeted, composition-aware editing.

5 Conclusion

In this work, we presented **SemanticAudio**, a novel two-stage Flow Matching framework that fundamentally rethinks text-to-audio generation by prioritizing semantic planning over direct acoustic synthesis. By explicitly modeling a compact, high-level semantic space, our **Semantic Planner** captures global event structures and textual intent with superior precision, as evidenced by state-of-the-art CLAP scores on AudioCaps. We further leveraged this decoupled design to introduce a training-free editing mechanism. Inspired by differential flow compositions, this method enables intuitive, inversion-free attribute modification, demonstrating exceptional robustness even in the absence of source text. Our results quantitatively confirm that separating semantic reasoning from acoustic realization not only enhances generation alignment but also provides a flexible, unified foundation for controllable audio editing. Future work will explore scaling this paradigm to longer-form audio and integrating multi-modal controls.

Limitations

Data Scale and Temporal Constraints. Our current implementation prioritizes high-quality semantic alignment by training exclusively on AudioCaps. While this rigorous supervision ensures precise text-audio correspondence, the dataset’s limited scale (~128 hours) and standardized 10-second duration impose constraints on generalization. Consequently, **SemanticAudio** may exhibit reduced robustness when handling long-form audio generation or highly complex, overlapping acoustic scenes compared to models trained on massive, weakly-supervised datasets (e.g., WavCaps or AudioSet). Future work will focus on scaling the semantic-space framework to larger, diverse corpora to capture long-tail acoustic distributions and extend temporal consistency beyond short clips.

Evaluation Challenges in Generative Editing. Standardizing the evaluation of open-domain audio editing remains an industry-wide challenge due to the absence of paired ground-truth references. While our constructed "Hard-Negative" benchmark allows for quantitative measurement via proxy metrics (CLAP, FD, IS), these automated scores may not fully capture human perceptual nuances in attribute modification. Our current study relies on objective metrics to validate the methodology; however, rigorous subjective evaluations (e.g., large-scale Mean Opinion Scores or AB preference tests) are necessary to further validate the practical utility of the editing features. We aim to contribute to the establishment of more comprehensive, paired source-target editing benchmarks in future iterations.

References

- [1] Jianhong Bai, Xiaoshi Wu, Xintao Wang, Xiao Fu, Yuanxing Zhang, Qinghe Wang, Xiaoyu Shi, Menghan Xia, Zuozhu Liu, Haoji Hu, Pengfei Wan, and Kun Gai. Semanticgen: Video generation in semantic space, 2025. URL <https://arxiv.org/abs/2512.20619>.
- [2] Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, and Neil Zeghidour. Audioldm: a language modeling approach to audio generation, 2023. URL <https://arxiv.org/abs/2209.03143>.
- [3] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models, 2022. URL <https://arxiv.org/abs/2210.11416>.
- [4] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning audio concepts from natural language supervision. In *Proc. ICASSP*, pages 1–5, 2023.
- [5] Zach Evans, Julian D. Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Stable audio open, 2024. URL <https://arxiv.org/abs/2407.14358>.
- [6] Liting Gao, Yi Yuan, Yaru Chen, Yuelan Cheng, Zhenbo Li, Juan Wen, Shubin Zhang, and Wenwu Wang. Rfm-editing: Rectified flow matching for text-guided audio editing, 2025. URL <https://arxiv.org/abs/2509.14003>.
- [7] Wenhao Guan, Kaidi Wang, Wangjin Zhou, Yang Wang, Feng Deng, Hui Wang, Lin Li, Qingyang Hong, and Yong Qin. Lafma: A latent flow matching model for text-to-audio generation. In *Proc. Interspeech*, pages 4813–4817, 2024.
- [8] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen, 2023. URL <https://arxiv.org/abs/2207.06405>.
- [9] Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models, 2023. URL <https://arxiv.org/abs/2301.12661>.

- [10] Chia-Yu Hung, Navonil Majumder, Zhifeng Kong, Ambuj Mehrish, Amir Ali Bagherzadeh, Chuan Li, Rafael Valle, Bryan Catanzaro, and Soujanya Poria. Tangoflux: Super fast and faithful text to audio generation with flow matching and clap-ranked preference optimization, 2025. URL <https://arxiv.org/abs/2412.21037>.
- [11] 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)*, pages 119–132, 2019.
- [12] Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D. Plumbley. Panns: Large-scale pretrained audio neural networks for audio pattern recognition, 2020. URL <https://arxiv.org/abs/1912.10211>.
- [13] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation, 2023. URL <https://arxiv.org/abs/2209.15352>.
- [14] Vladimir Kulikov, Matan Kleiner, Inbar Huberman-Spiegelglas, and Tomer Michaeli. Flowedit: Inversion-free text-based editing using pre-trained flow models, 2025. URL <https://arxiv.org/abs/2412.08629>.
- [15] 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.
- [16] Tianhong Li, Dina Katabi, and Kaiming He. Return of unconditional generation: A self-supervised representation generation method, 2024. URL <https://arxiv.org/abs/2312.03701>.
- [17] Xiquan Li, Junxi Liu, Yuzhe Liang, Zhikang Niu, Wenxi Chen, and Xie Chen. Meanaudio: Fast and faithful text-to-audio generation with mean flows, 2025. URL <https://arxiv.org/abs/2508.06098>.
- [18] Jinhua Liang, Yuanzhe Chen, Yi Yuan, Dongya Jia, Xiaobin Zhuang, Zhuo Chen, Yuping Wang, and Yuxuan Wang. Audiomorphix: Training-free audio editing with diffusion probabilistic models, 2025. URL <https://arxiv.org/abs/2505.16076>.
- [19] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. In *11th International Conference on Learning Representations, ICLR 2023*, 2023.
- [20] 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. *Proceedings of the International Conference on Machine Learning*, pages 21450–21474, 2023.
- [21] 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*, 32:2871–2883, 2024. doi: 10.1109/TASLP.2024.3399607.
- [22] Navonil Majumder, Chia-Yu Hung, Deepanway Ghosal, Wei-Ning Hsu, Rada Mihalcea, and Soujanya Poria. Tango 2: Aligning diffusion-based text-to-audio generations through direct preference optimization, 2024. URL <https://arxiv.org/abs/2404.09956>.
- [23] Wenqi Ouyang, Zeqi Xiao, Danni Yang, Yifan Zhou, Shuai Yang, Lei Yang, Jianlou Si, and Xingang Pan. Tokensgen: Harnessing condensed tokens for long video generation, 2025. URL <https://arxiv.org/abs/2507.15728>.
- [24] Apoorv Vyas, Heng-Jui Chang, Cheng-Fu Yang, Po-Yao Huang, Luya Gao, Julius Richter, Sanyuan Chen, Matt Le, Piotr Dollár, Christoph Feichtenhofer, Ann Lee, and Wei-Ning Hsu. Pushing the frontier of audiovisual perception with large-scale multimodal correspondence learning, 2025. URL <https://arxiv.org/abs/2512.19687>.

- [25] Yuancheng Wang, Zeqian Ju, Xu Tan, Lei He, Zhizheng Wu, Jiang Bian, and Sheng Zhao. Audit: Audio editing by following instructions with latent diffusion models. *NeurIPS 2023*, 2023.
- [26] Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Marianna Nezhurina, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation, 2024. URL <https://arxiv.org/abs/2211.06687>.
- [27] Sihan Xu, Yidong Huang, Jiayi Pan, Ziqiao Ma, and Joyce Chai. Inversion-free image editing with natural language, 2023. URL <https://arxiv.org/abs/2312.04965>.