

# T2A-FEEDBACK: IMPROVING BASIC CAPABILITIES OF TEXT-TO-AUDIO GENERATION VIA FINE-GRAINED AI FEEDBACK

Anonymous authors

Paper under double-blind review

Suddenly, a man shouted, “Fire!” The man and women joined in. Two children cried together. In no time, thousands of people were shouting, thousands of children were crying, and countless dogs were barking. Amid the chaos, there were sounds of collapsing buildings, explosions, and strong winds, all happening at once. There were also cries for help, the sounds of buildings being dragged, voices of looting, and water splashing everywhere.

(忽一人大呼：“火起”，夫起大呼，妇亦起大呼。两儿齐哭。俄而百千人大呼，百千儿哭，百千犬吠。中间力拉崩倒之声，火爆声，呼呼风声，百千齐作；又夹百千求救声，曳屋许许声，抢夺声，泼水声。)

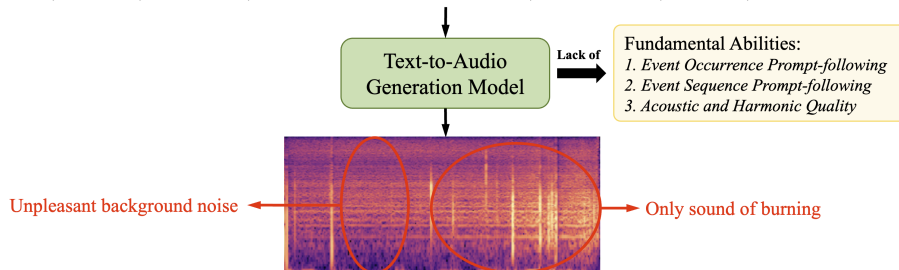


Figure 1: The audio description is from a classic Chinese essay “Kou Ji”, which vividly depicts a performer using only vocal mimicry to recreate an entire dramatic scene. The existing Text-to-Audio generation model struggles to generate such narrative and multi-event audios. The generated audio often fails to contain all events in the described sequence while maintaining acoustic quality and harmony.

## ABSTRACT

Text-to-audio (T2A) generation has achieved remarkable progress in generating a variety of audio outputs from language prompts. However, current state-of-the-art T2A models still struggle to satisfy human preferences for prompt-following and acoustic quality when generating complex multi-event audio. To improve the performance of the model in these high-level applications, we propose to enhance the basic capabilities of the model with AI feedback learning. First, we introduce fine-grained AI audio scoring pipelines to: 1) verify whether each event in the text prompt is present in the audio (**Event Occurrence Score**), 2) detect deviations in event sequences from the language description (**Event Sequence Score**), and 3) assess the overall acoustic and harmonic quality of the generated audio (**Acoustic Harmonic Quality**). We evaluate these three automatic scoring pipelines and find that they correlate significantly better with human preferences than other evaluation metrics. This highlights their value as both feedback signals and evaluation metrics. Utilizing our robust scoring pipelines, we construct a large audio preference dataset, **T2A-FeedBack**, which contains 41k prompts and 249k audios, each accompanied by detailed scores. Moreover, we introduce **T2A-EpicBench**, a benchmark that focuses on long captions, multi-events, and story-telling scenarios, aiming to evaluate the advanced capabilities of T2A models. Finally, we demonstrate how T2A-FeedBack can enhance current state-of-the-art audio model. With simple preference tuning, the audio generation model exhibits significant improvements in both simple (AudioCaps test set) and complex (T2A-EpicBench) scenarios. The project page is available at <https://T2Afeedback.github.io>

## 1 INTRODUCTION

Recent Text-to-Audio (T2A) generation models (Huang et al., 2023b;a; Liu et al., 2023a; 2024; Ghosal et al., 2023; Majumder et al., 2024; Vyas et al., 2023) have made drastic performance improvements. By trained on massive audio-text data (Gemmeke et al., 2017; Fonseca et al., 2021; Chen et al., 2020; Kim et al., 2019), these generative models learn to generate diverse sounds with a given language prompt. For audio generation, generating harmonious multi-event audio or describing a story with audio has important applications in music (Agostinelli et al., 2023), advertising, video-audio generation (Luo et al., 2024; Wang et al., 2024), etc. However, as shown in Figure. 1, existing audio generation models are struggling to generate harmonious and high-quality audio from narrative and complex descriptions, which limits the potential for high-level applications.

The failure of the generated results is often demonstrated in three aspects: 1) cannot fully include all the events described, 2) cannot accurately follow the order of all the events described, and 3) cannot organize all the events harmoniously. Therefore, the model performance in multi-event scenarios is determined by its capabilities in these three fundamental aspects.

To improve the model’s performance across more advanced applications, we focus on strengthening the audio generation model’s fundamental abilities. Inspired by feedback learning in large language models (Ouyang et al., 2022; Bai et al., 2022; Touvron et al., 2023), we propose creating an audio preference dataset centered on three abilities necessary for generating harmonic and complex audio: 1) **Event Occurrence Prompt-Following**, 2) **Event Sequence Prompt-Following**, and 3) **Acoustic&Harmonic Quality**. Based on the preference information, we can refine the model’s core abilities, resulting in better results in both simple and challenging scenarios.

However, due to the scarcity of audio data and the challenges of annotating the scale of user preferences, it is difficult to collect massive audio preference datasets that only rely on human annotators. To fill this void, we explore using AI feedback (Cui et al., 2023; Lee et al., 2023; Yuan et al., 2024; Burns et al., 2023) in text-to-audio generation, utilizing AI models to rank audios instead of relying on human annotators. Compared to manual annotation, automating the data collection and annotation process reduces the cost of obtaining audio preference data and enhances scalability.

Specifically, we develop three AI scoring pipelines to evaluate the generated audio in a fine-grained and holistic manner, corresponding to three core capabilities:

- *Event Occurrence Score*: To specifically check whether each event occurs in, we calculate the audio-text semantic matching score for each described event separately. A lower score suggests that the corresponding event might be absent from the audio.
- *Event Sequence Score*: To verify the correctness of event order, we analyze the start and end times of each event and compare them with the event order outlined in the text prompt. A higher score implies a greater similarity between the event sequences in caption and audio.
- *Acoustics&Harmonic Quality*: Drawing inspiration from aesthetic scoring methods used in image quality scoring, we manually annotate acoustic and harmonic quality for audio samples. These data are then used to train an automatic acoustic&harmonic predictor.

We confirm that our three scoring functions show a stronger correlation with human evaluations compared to existing automatic audio evaluation methods (Wu et al., 2023b; Xie et al., 2024). Consequently, in addition to their application in ranking preference data, these scoring functions can be used as evaluation metrics that more effectively capture human preferences across different aspects.

Leveraging these advanced AI scoring pipelines, we establish a comprehensive data collection and annotation framework. As a result, we construct **T2A-Feedback**, a large audio preference dataset comprising 41,627 captions and 249,762 generated audios, each annotated with detailed scores.

Furthermore, to evaluate the higher-level capabilities of text-to-audio models in multi-event scenarios, we introduce a more challenging benchmark, **T2A-EpicBench**, which features longer, more imaginative, and story-telling captions for audio generation. We enhance the advanced text-to-audio diffusion model, Make-an-Audio 2 (Huang et al., 2023a), with T2A-Feedback. Our results show that using T2A-Feedback not only effectively improves the basic capabilities of the model in simple AudioCaps benchmark, but also emergently improves the performance in complex T2A-EpicBench.

## 2 RELATED WORK

### 2.1 TEXT-TO-AUDIO GENERATION

Text-to-audio generation is an emerging field that aims to convert textual descriptions into corresponding audio outputs. Existing text-to-audio generation methods can be divided into two categories: Diffusion-based and Language model-based. Diffusion-based techniques have gained prominence for generating high-quality, realistic audio by modeling the process of denoising. These methods, like Make-an-Audio (Huang et al., 2023b;a), AudioLDM (Liu et al., 2023a; 2024), Tango (Ghosal et al., 2023; Majumder et al., 2024), start with random noise and iteratively refine it to produce coherent audio over a series of denoising steps. On the other hand, Language model-based methods (Borsos et al., 2023; Agostinelli et al., 2023; Cideron et al., 2024) tokenize audios as acoustic discrete tokens, and predict the tokens within an auto-regressive model conditioned on text inputs.

The above models acquire the ability to generate diverse audio by training on large-scale audio-text datasets. However, current datasets like AudioSet (Gemmeke et al., 2017), AudioCaps (Kim et al., 2019), and FSD50k (Fonseca et al., 2021) only provide tag-level annotations or short captions. As a result, when processing long, detailed language prompts, existing models often produce low-quality, noisy outputs and struggle to accurately follow the text. Due to the difficulty of annotating detailed audio captions, scaling rich and accurate audio descriptions remains a challenge. In this work, we focus on enhancing the model’s basic abilities in event occurrence, sequence, and harmony, thereby improving its performance in both simple scenarios and advanced applications.

### 2.2 PREFERENCE TUNING WITH HUMAN&AI FEEDBACK

Tuning generative models according to human preferences has emerged as a standard practice for improving the quality of outputs. By tuning with feedback information on different aspects, the model can be improved and aligned with human preferences in corresponding aspects. Traditionally, this preference data used for tuning relied heavily on human evaluators who rank multiple generated results, assessing their quality based on various criteria such as relevance, coherence, and aesthetic value (Bai et al., 2022; Touvron et al., 2023; Ouyang et al., 2022; Kirstain et al., 2023; Liang et al., 2024; Wu et al., 2023a; Cideron et al., 2024).

While effective, manual human annotation is costly and time-consuming, which greatly hampers the scalability of preference tuning across more diverse generative tasks. To address the difficulty, more recent developments have focused on leveraging pre-trained AI models to automate the process of scoring generated content (Cui et al., 2023; Lee et al., 2023; Yuan et al., 2024; Burns et al., 2023). Such an AI feedback approach has achieved impressive improvements in large language models.

Recently, some studies have attempted preference fine-tuning in text-to-audio generation models. One recent paper related to our work, Tango2 (Majumder et al., 2024), utilizes contrastive language-audio pre-training (CLAP) (Wu et al., 2023b) to rank audio generated by the Tango model. However, CLAP can only evaluate the global alignment between audio and text but falls short in assessing the fine-grained details, like detailed event occurrence, sequence, and harmony. In this paper, we construct more robust AI audio scoring pipelines with fine-grained recognition ability. Our method shows a much stronger correlation with human preference and the constructed dataset brings significant improvement to the current text-to-audio generation model.

### 2.3 TEXT-TO-AUDIO EVALUATION METRIC

Existing evaluation metrics for audio generation, such as FAD and IS, assess audio distributions but cannot evaluate the quality of individual samples. Additionally, many studies rely on similarity scores from the CLAP model to assess global audio-text semantic alignment. PicoAudio (Xie et al., 2024) uses a text-to-audio grounding model (Xu et al., 2024) to detect audio segments based on language prompts. However, there remains a lack of fine-grained evaluation methods for assessing detailed event occurrence, sequencing, and acoustic quality. Our research fills this gap by creating robust audio AI scoring pipelines, that show a strong correlation with humans, and significantly surpass alternative methods.

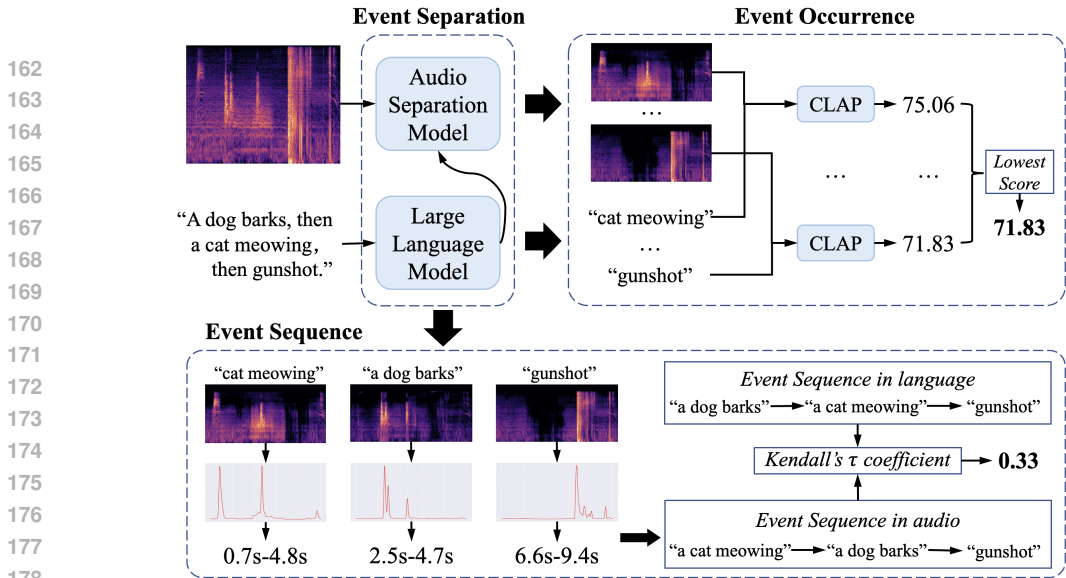


Figure 2: The overview of event occurrence and sequence scoring pipelines.

### 3 T2A-FEEDBACK

In this section, we first dive into the three AI audio scoring pipelines: (i) Event Occurrence Prompt-following, in Section. 3.1; (ii) Event Sequence Prompt-following, in Section. 3.2; (iii) Acoustic Quality, in Section. 3.3. We then describe the specific data generation and sorting method for the T2A-Feedback dataset in Section. 3.4.

#### 3.1 EVENTS OCCURRENCE PROMPT-FOLLOWING

Generating audio that accurately reflects the events described in a given prompt is the fundamental requirement of prompt-following. However, when multiple events are included in the text description, current text-to-audio generation models often struggle to generate each event precisely. To improve the generation model’s event occurrence prompt-following ability, we first build an AI pipeline to determine the occurrence of events in audio.

Previous methods primarily utilize contrastive language-audio pre-training (CLAP) (Wu et al., 2023b) over the audios and language descriptions to assess their semantic relevance. However, in multi-event scenarios, the sentence-level matching score struggles to identify event-level misalignment, and can not pinpoint which specific events are present and which are not, as shown in Figure. 5. To accurately identify misaligned events, we propose to measure the audio-text semantic alignment at the event-level. To this end, we first separate the language description and audio into basic events, as shown in the “Event Separation” part of Figure. 2. Specifically, we utilize a large language model (LLM) (Jiang et al., 2023) to decompose descriptions into event captions according to the described order. Meanwhile, we employ an advanced audio separation model (Liu et al., 2023b) to segment the audio into event-level sub-audios based on these event captions. By calculating the similarity between each event-level description and its corresponding sub-audio in CLAP space, we can gain clearer insights into the specific aligned and misaligned events.

To encourage the models to comprehensively generate all described events, for each audio-text pair, we select the lowest value among all event-level audio-text matching scores as the **Event Occurrence Score**. For audios generated from the same caption, a higher score indicates that the audio is more likely to contain all the described events.

#### 3.2 EVENTS SEQUENCE PROMPT-FOLLOWING

In addition to generating all events, whether these events occur in the temporal order described in the caption is also a crucial aspect of prompt-following. Some recent work attempts to detect the sequence of events in audio. Tango2 (Majumder et al., 2024) computes the CLAP matching score between the temporal description and corresponding audios, but we find the sentence-level



CLAP score is not sensitive to the temporal description in captions, as demonstrated in Figure. 5 and Table. 2. On the other hand, PicoAudio (Xie et al., 2024) employs audio grounding model (Xu et al., 2024) to detect audio segments. However, due to the limitation of the training scale, the generalization performance of the audio grounding model is limited.

To robustly analyze audio event sequences, we propose a new pipeline for event sequence analysis. Similar to event occurrence, we first use the LLM and audio separation model to extract event-level descriptions and their corresponding sub-audios. For each separated audio track, we determine the event’s start and end times based on volume levels. Specifically, we normalize the volume to a range of [0,1], and the period where the normalized volume exceeds a certain threshold is identified as the event’s duration.

In multi-event scenarios, there are multiple complex temporal relationships. To comprehensively assess the temporal alignment between the language prompt and the generated audio, and to specifically identify which temporal relationships are accurate and which are misaligned, we employ Kendall’s  $\tau$  coefficient. This widely used non-parametric statistic measures rank correlation between two variables. Considering  $n$  events and their  $n(n - 1)$  event pairs, we use LLM to analyze the relationships between each event pair in the language description and extract the event sequence in the audio based on the starting time of each event. The **Events Sequence Score** (e.g., Kendall’s  $\tau$  coefficient between event sequences in language and audio) is calculated as:

$$\tau = \frac{C - D}{n(n - 1)} \quad (1)$$

where  $C$  represents the number of concordant event pairs between the description and the audio,  $D$  denotes the number of discordant ones. Higher  $\tau$  indicates a greater alignment of the event sequence in the generated audio with the text description. Specifically,  $\tau = 1$  signifies that the event sequence in the generated audio is identical to the language description, while  $\tau = -1$  indicates that the sequences are completely reversed.

### 3.3 ACOUSTIC&HARMONIC QUALITY

In addition to generating all events accurately following the language prompt, organically integrating different events to create a pleasant-sounding effect is also one of the basic capabilities. However, current audio generation models often produce low-quality and noisy results.

To alleviate this challenge, we first develop an audio acoustic&harmonic quality predictor. Inspired by the image aesthetic predictor in Schuhmann et al. (2022), we first manually score audio samples on a scale from 1 to 4 according to their quality. Two annotators independently score the audio samples according to the same criteria, and samples with consistent scores are accepted as training data. The detailed scoring criteria are as follows:

Annotators need to score the auditory quality of audio from the following four perspectives:

**Acoustic Quality:** Does the generated audio sound realistic and pleasant?

**Harmony:** Do different sound elements integrate well, forming a cohesive auditory scene?

**Background Noise:** Is there noise that disrupts the clarity and naturalness of the audio?

**Dynamic Range:** Are the different audio elements within their reasonable volume range?

The specific standards for each score are as follows:

Score	Standard
1	Poor audio quality; sounds unrealistic with disjointed elements, severe background noise interference, and extremely limited dynamic range.
2	Normal audio quality; some events are natural, but overall harmony is lacking. Background noise affects clarity, and dynamic range is limited.
3	Good audio quality; most events are realistic with harmonious integration. Background noise is minimally disruptive, and dynamic range is reasonable.
4	Excellent audio quality; all events are very realistic with perfect integration, well-managed background noise, and wide dynamic range.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

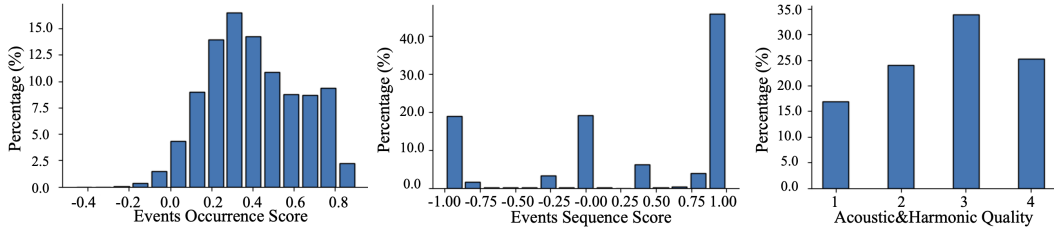


Figure 3: Histograms of three different scores in T2A-Feedback.

Using the human-annotated data, we train a linear predictor on the top of CLAP audio embeddings. With the high-quality pre-trained representation, we find that, akin to aesthetic score predictors for images, a small amount of annotated data can yield a generalized subjective quality predictor. Specifically, we train the acoustic predictor with 1,000 meticulously annotated audio samples using cross-entropy loss. The output of the predictor is termed the **Acoustic&Harmonic Quality**.

### 3.4 PREFERENCE DATA GENERATION

To generate diverse and comprehensive audio, we first augment the text prompts used for audio generation. We begin with the captions from the training set of the large-scale audio-text dataset, AudioCaps. By employing an LLM, we decompose these captions into fundamental event descriptions and calculate their semantic similarity within the CLAP space to filter out non-overlapping, basic event descriptions. Then, we prompt the LLM with randomly selected events to create varied and natural multi-event audio descriptions, with explicit temporal ordering. Finally, we combine the enhanced 3,769 captions with the 37,858 captions from the training set of AudioCaps, serving as the prompt source for audio generation.

As highlighted in Cui et al. (2023), diversity is crucial for preference datasets. To mitigate the potential bias of using a single audio generation model and to enhance the generalization of the generated data, we employ three advanced audio generation models: Make-an-Audio2, AudioLDM2, and Tango2. Each model generates 2 audio per caption, resulting in a total of 6 audio files for each caption. In summary, we produce 249,762 audios from 41,627 descriptions. For audios generated from the same captions, we combine three rankings of each score to derive the overall ranking.

The histogram plots of the scores on all the generated audios are shown in Figure 3. The distribution of Event Occurrence Scores and Acoustic&Harmonic Quality is similar to a Gaussian distribution. Since most descriptions contain one or two sequential events, Event Sequence Scores are concentrated between -1 and 1. As noted in Liang et al. (2024), this discriminative score distribution ensures a balanced ratio of negative to positive samples, enabling effective preference tuning.

## 4 T2A-EPICBENCH

Current text-to-audio generation models are mainly evaluated and compared on the AudioCaps test set. However, the captions in AudioCaps are generally short and simple, averaging 10.3 words per sentence. Specifically, 17% of the captions feature only a single event, and 44% contains two events. This is not enough to assess the model’s capabilities in more advanced applications involving detailed, multi-event, and narrative-style audio generation.

To fill this gap, we propose **T2A-EpicBench**, consisting of 100 detailed, multi-event, and story-telling captions. Each caption averages 54.8 words and 4.2 events, with 86% containing four events and the remainder featuring five or more. Initially, we manually write 10 detailed captions, then used them as in-context examples to prompt LLM for generating the remaining captions. All 100 captions are manually reviewed for accuracy. Several examples from T2A-EpicBench are included in the Appendix.

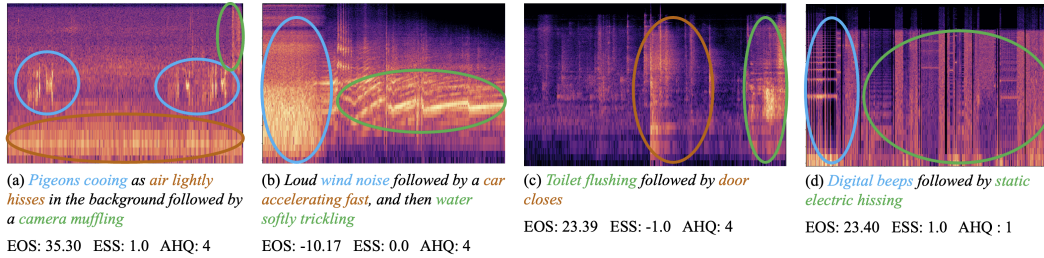


Figure 4: Visualization of the predicted scores from our AI scoring pipeline. We highlight the first, second, and third events described in the captions using blue, brown, and green, respectively.

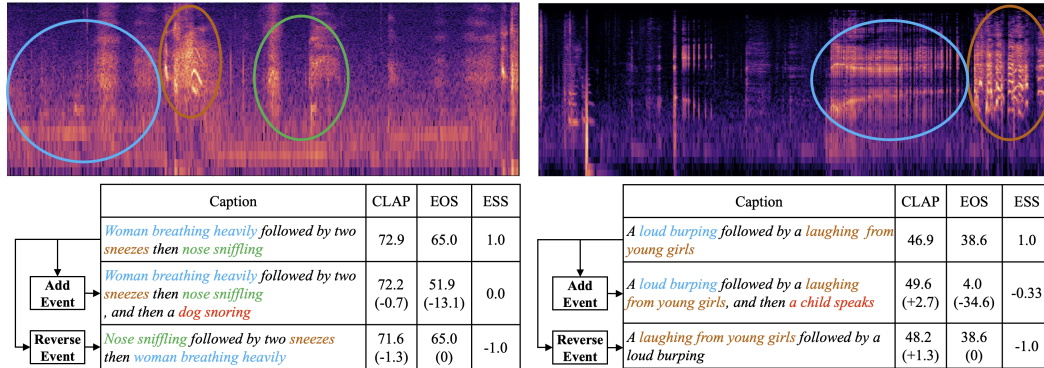


Figure 5: Qualitative comparison between CLAP scores and EOS/ESS scores reveals distinct sensitivities to misalignment. By adding or reversing events in the ground-truth caption, we create captions that are misaligned with the audio in terms of event occurrence and sequence.

## 5 EXPERIMENT

### 5.1 ANALYSIS OF AI SCORING PIPELINES

#### 5.1.1 QUANTITATIVE ANALYSIS

**Evaluation of Event Occurrence Score (EOS)** To evaluate the scoring model’s capability in recognizing whether audios contain all the events described in the text, we propose a missing event recognition task. We construct distracting captions for the AudioCaps test set, by adding random event descriptions to the ground-truth captions. This task challenges models to distinguish the ground-truth caption from the constructed interference captions. There are 3,701 samples in total for this task.

We mainly compare our EOS with sentence-level CLAP score. The caption with the higher matching score to the audio is considered as the prediction. As shown in Table 1, our EOS score showcases a notable advantage over CLAP, demonstrating the superiority of event-level audio-text matching in identifying whether all events are correctly contained in audios.

#### Evaluation of Event Sequence Score (ESS)

To verify the ability to distinguish the alignment of event sequences in text and audio, we collect 450 two-event samples from PicoAudio’s training set, and reverse the events orders in the description as interference caption. Using this dataset, we compare different methods by calculating the accuracy of recognizing the ground-truth description versus the interference description, and by evaluat-

Table 1: Comparison about event occurrence

	Accuracy
Random Guess	50.0%
CLAP	77.5%
EOS(ours)	<b>90.9%</b>

Table 2: Comparison about event sequence.  $ESS_{0,x}$  stands for using  $0.x$  as volume thresholds.

Method	Accuracy	F1 Score	Correlation
CLAP	49.6	-	-
PicoAudio	71.6	0.787	0.30
$ESS_{0.1}$	<b>79.6</b>	0.814	0.43
$ESS_{0.3}$	79.1	<b>0.851</b>	<b>0.52</b>
$ESS_{0.5}$	78.0	0.769	0.52

378 ing the Segment F1 Score (Mesaros et al., 2016) for detecting the start and end times of each audio  
 379 event. Moreover, we manually annotate temporal order alignment for 100 audios generated from our  
 380 temporal-enhanced captions and compute the correlation between different methods and humans.

381  
 382 The results of event sequences are provided in Table. 2. We compare ESS with CLAP score and  
 383 the audio grounding model (Xu et al., 2024) used by PicoAudio (Xie et al., 2024). Compared to  
 384 baselines, our method distinguishes the ground-truth caption from the distracting one more accu-  
 385 rately and achieves higher F1 scores in detecting the start and end times of events in audio. More  
 386 importantly, our method shows a much stronger correlation to human annotations.

387 Additionally, we investigate various volume thresholds used to determine the duration of each event.  
 388 In Table 2, we test thresholds of 0.1, 0.3, and 0.5. ESS consistently performs better than other  
 389 methods across most settings, with 0.3 providing the optimal results and thus chosen as the default  
 390 setting.

391 **Evaluation of Acoustic&Harmonic Quality (AHQ)** To validate our acoustic&harmonic predic-  
 392 tor, we independently annotate 100 additional audios as a test set. The correlation between the model  
 393 predictions and human labels on the test set is 0.786, showing strong generalization ability and high  
 394 consistency with human preferences.  
 395

396 Moreover, we explore building the Acoustic&Harmonic Predic-  
 397 tor on top of various pre-trained audio models and evaluate how  
 398 well each variant correlates with human preferences. The results  
 399 in Table 3 show that the predictor built on CLAP (Wu et al.,  
 400 2023b) outperforms those based on self-supervised models like  
 401 AudioMAE (Huang et al., 2022) and BEAT (Chen et al., 2022).  
 402 Similarly, the image aesthetics predictor (Schuhmann et al., 2022)  
 403 is built on the CLIP model (Iharco et al., 2021). This advan-  
 404 tage may stem from the fact that self-supervised models are task-  
 405 agnostic, whereas CLIP and CLAP align with language, resulting in better semantic discrimination.

Table 3: AHQ Predictor on different base models.

	Correlation
AudioMAE	0.613
BEAT	0.519
CLAP	<b>0.786</b>

406 5.1.2 QUALITATIVE ANALYSIS

407  
 408 We show some example predictions from our scoring  
 409 pipelines in Figure. 4, where our methods can specifi-  
 410 cally identify the misaligned event, the out-of-order event  
 411 order, and the disharmony between events in the audio.  
 412 Moreover, we provide the confusion matrix of acous-  
 413 tic&harmonic predictor on these 100 test samples in Fig-  
 414 ure. 6, which further demonstrates the statistical robust-  
 415 ness of our predictor.

416 Moreover, we provide the qualitative comparison be-  
 417 tween our EOS and ESS with the single CLAP score, in  
 418 Figure. 5. For the ground-truth audio-caption pairs from  
 419 AudioCaps, we perturb the captions by adding an event or  
 420 shuffling the order of events. We find that the CLAP score  
 421 is not sensitive to these perturbations and even yields a  
 422 higher score with the incorrect, perturbed caption. In contrast, our EOS and ESS scores more accu-  
 423 rately reflect the alignment between audio and text regarding event occurrence and event order.

424 5.2 ANALYSIS OF PREFERENCE TUNING

425  
 426 To demonstrate the effect of T2A-Feedback dataset in improving audio generation model, we fine-  
 427 tuning the advanced text-to-audio model, Make-an-Audio 2 (Huang et al., 2023a), with two prefer-  
 428 ence training methods: Direct Preference Optimization (DPO) (Wallace et al., 2024) and Reward  
 429 rAnked FineTuning (RAFT) (Dong et al., 2023). Another audio preference dataset, Audio-Alpaca,  
 430 proposed by Majumder et al. (2024) is the main baseline for comparison. Both the widely-used  
 431 AudioCaps and the new T2A-EpicBench are used as benchmarks, corresponding to applications in  
 simple and complex scenarios respectively.

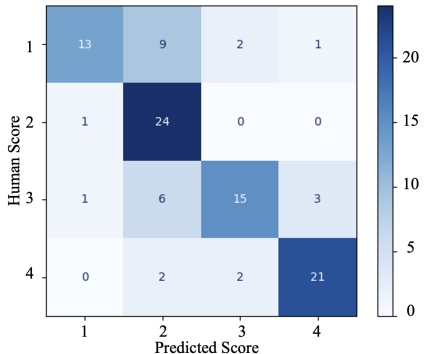


Figure 6: Confusion matrix.

Table 4: Evaluation results on AudioCaps. The EOS, ESS and AHQ represent the Event Occurrence Score, Event Sequence Score and Acoustic&Harmonic Quality, respectively.

		FAD↓	KL↓	IS↑	CLAP↑	EOS.↑	ESS.↑	AQ.↑
Make an Audio 2		<b>1.82</b>	1.44	10.03	69.97	42.05	0.53	2.33
Preference Tuning								
Audio-Alpaca	RAFT	1.93	<b>1.29</b>	10.37	72.23	44.85	0.53	2.45
	DPO	3.20	1.24	<b>12.27</b>	72.36	44.42	0.55	2.14
T2A-Feedback (ours)	RAFT	2.29	1.33	11.66	73.10	45.53	0.51	2.50
	DPO	2.64	1.31	11.35	<b>74.00</b>	<b>49.58</b>	<b>0.57</b>	<b>2.57</b>

Table 5: Evaluation results on T2A-EpicBench. The  $win_{EOS}$ ,  $win_{ESS}$  and  $win_{AHQ}$  represent the win rates of tuned models over the original model in terms of Event Occurrence, Event Sequence and Acoustic&Harmonic Quality, respectively.

		AI Scoring			Human Scoring		
		$win_{EOS}$	$win_{ESS}$	$win_{AHQ}$	$win_{EOS}$	$win_{ESS}$	$win_{AHQ}$
Make an Audio 2		- (14.21)	- (0.03)	- (1.96)	-	-	-
Preference Tuning							
Audio-Alpaca	RAFT	53%(15.73)	51%(0.04)	42%(1.69)	57%	54%	53%
	DPO	55%(16.87)	52%(0.03)	49%(1.96)	65%	<b>64%</b>	59%
T2A-Feedback (ours)	RAFT	52%(15.85)	52%(0.05)	<b>54%(2.14)</b>	61%	57%	61%
	DPO	<b>58%(19.96)</b>	<b>64%(0.13)</b>	52%(2.10)	<b>68%</b>	62%	<b>68%</b>

### 5.2.1 QUANTITATIVE RESULTS ON AUDIOCAPS

The classical automated metrics (FAD, KL, IS, and CLAP), as well as our three new scores (EOS, ESS, and AHQ) are employed to quantitatively evaluate and compare different model variants.

The quantitative results are provided in Table. 4. FAD, KL, and IS assess audio fidelity by evaluating the distribution of the generated audio. For these metrics, both the preference dataset and training methods result in similar overall improvements. CLAP is commonly used to measure the semantic alignment between the input prompt and the generated audio. While both Audio-Alpaca and T2A-Feedback improve the CLAP score, T2A-Feedback yields greater gains.

Moreover, as analyzed in Section. 5.1.1, the proposed EOS and ESS are more accurate than CLAP in judging event occurrence and event sequence, and AHQ shows a strong correlation to human preference in acoustic and harmony. We calculate the three scores for different model variants to evaluate audio generation results more accurately and comprehensively. The significantly better results across these three metrics demonstrate that T2A-Feedback yields far greater improvements compared to Audio-Alpaca, and the DPO method outperforms RAFT in our setting.

### 5.2.2 QUANTITATIVE RESULTS ON T2A-EPICBENCH

Since there are no ground-truth audios for the long and story-telling text prompts in T2A-EpicBench, we primarily measure the win rate of preference-tuned models against the original model outputs across three key areas: event occurrence, event sequence, and acoustic & harmonic quality. In addition to scoring the generated audio with our AI pipeline, we conduct a user study where two human annotators evaluate and select the better output based on each criterion.

The results on T2A-EpicBench, are illustrated in Table. 5, indicate that Audio-Alpaca provides only marginal improvements in handling detailed captions and multi-event scenarios, whereas T2A-Feedback significantly and comprehensively enhances the model’s performance.



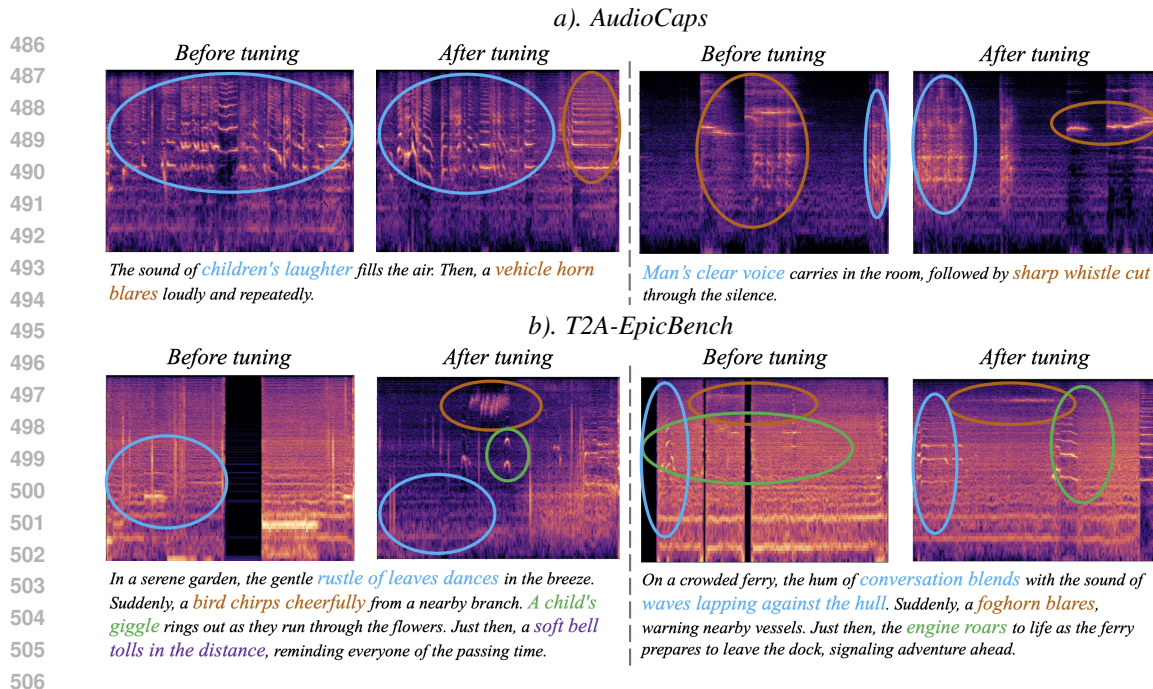


Figure 7: Visualization of the impact of preference tuning with T2A-Feedback.

507  
508  
509  
510  
511  
512  
513  
514  
515  
516

It is worth noting that T2A-Feedback does not include long audio descriptions. The average word count per caption in T2A-Feedback is 9.6, which is considerably shorter than the 54.8 average word number of T2A-EpicBench prompts, and even shorter than Audio-Alpaca’s 10.2 words per caption. T2A-Feedback does not directly provide additional long caption data, and the 65% average win rate in the user study reinforces that by focusing on improving the basic capabilities of short captions, the audio generation model can emergently learn to handle more complex long-text and multi-event scenarios.

### 5.2.3 QUALITATIVE FINDINGS

517  
518  
519  
520  
521  
522  
523  
524  
525  
526

To better demonstrate the effectiveness of preference tuning on T2A-Feedback, we visualize some examples of tuning the original model on our T2A-Feedback with the DPO method in Figure. 7. For the examples of short captions in Figure. 7a, while both models before and after fine-tuning can produce clean audio, the fine-tuned model successfully generates all events in the described order. In the more challenging case from T2A-EpicBench, the original model often generates noisy, low-quality audio, making it difficult to distinguish the events. Preference tuning on T2A-Feedback, as shown in Figure. 7b, significantly reduces background noise and generates audio that more faithfully captures both events and their orders.

## 6 CONCLUSION

527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

In this paper, we build AI scoring pipelines to evaluate three fundamental capabilities of audio generation: Event Occurrence Prompt-following, Event Sequence Prompt-following, and Acoustics&Harmonic Quality. Using these automatic evaluation metrics, which are highly correlated with human preferences, we build a large-scale audio preference dataset, **T2A-Feedback**. Experimentally, we extensively demonstrate the accuracy and robustness of our AI scoring pipelines. The three scores demonstrate a strong correlation to human preferences, which highlights its potential to better evaluate text-to-audio generation models. To assess the model’s ability in complex multi-event scenarios, we propose a new challenging benchmark, **T2A-EpicBench**, which requires models to generate detailed and narrative audios. Using our T2A-Feedback to tune the audio generation model effectively improves its capabilities in the three core aspects and achieves better performance in both simple (AudioCaps) and complex (T2A-EpicBench) scenarios.

## REPRODUCIBILITY STATEMENT

The newly proposed audio AI scoring pipeline, preference dataset (T2A-Feedback) and benchmark (T2A-EpicBench) will be open-sourced. In addition, in Section 3, 5 and A, we describe our pipelines, evaluation tasks and data, and other implementation details in detail to ensure the reproducibility of our method.

## REFERENCES

- Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325*, 2023.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, et al. Audioldm: a language modeling approach to audio generation. *IEEE/ACM transactions on audio, speech, and language processing*, 31:2523–2533, 2023.
- Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, et al. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. *arXiv preprint arXiv:2312.09390*, 2023.
- 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.
- Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, and Furu Wei. Beats: Audio pre-training with acoustic tokenizers. *arXiv preprint arXiv:2212.09058*, 2022.
- Geoffrey Cideron, Sertan Girgin, Mauro Verzetti, Damien Vincent, Matej Kastelic, Zalán Borsos, Brian McWilliams, Victor Ungureanu, Olivier Bachem, Olivier Pietquin, et al. Musicrl: Aligning music generation to human preferences. *arXiv preprint arXiv:2402.04229*, 2024.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*, 2023.
- Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023.
- 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.
- Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 776–780. IEEE, 2017.
- Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction-tuned llm and latent diffusion model. *arXiv preprint arXiv:2304.13731*, 2023.
- 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.



- 594 Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian  
595 Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. *Advances in Neural Infor-*  
596 *mation Processing Systems*, 35:28708–28720, 2022.
- 597 Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin  
598 Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced  
599 diffusion models. In *International Conference on Machine Learning*, pp. 13916–13932. PMLR,  
600 2023b.
- 601 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori,  
602 Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali  
603 Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL [https://doi.org/10.5281/](https://doi.org/10.5281/zenodo.5143773)  
604 [zenodo.5143773](https://doi.org/10.5281/zenodo.5143773). If you use this software, please cite it as below.
- 605 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
606 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
607 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 608 Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating  
609 captions for audios in the wild. In *Proceedings of the 2019 Conference of the North American*  
610 *Chapter of the Association for Computational Linguistics: Human Language Technologies, Vol-*  
611 *ume 1 (Long and Short Papers)*, pp. 119–132, 2019.
- 612 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-  
613 a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural*  
614 *Information Processing Systems*, 36:36652–36663, 2023.
- 615 Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton  
616 Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. Rlaif: Scaling reinforcement learning  
617 from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*, 2023.
- 618 Youwei Liang, Junfeng He, Gang Li, Peizhao Li, Arseniy Klimovskiy, Nicholas Carolan, Jiao Sun,  
619 Jordi Pont-Tuset, Sarah Young, Feng Yang, et al. Rich human feedback for text-to-image genera-  
620 tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
621 pp. 19401–19411, 2024.
- 622 Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and  
623 Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv*  
624 *preprint arXiv:2301.12503*, 2023a.
- 625 Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu  
626 Wang, Yuxuan Wang, and Mark D Plumbley. Audioldm 2: Learning holistic audio generation  
627 with self-supervised pretraining. *IEEE/ACM Transactions on Audio, Speech, and Language Pro-*  
628 *cessing*, 2024.
- 629 Xubo Liu, Qiuqiang Kong, Yan Zhao, Haohe Liu, Yi Yuan, Yuzhuo Liu, Rui Xia, Yuxuan  
630 Wang, Mark D Plumbley, and Wenwu Wang. Separate anything you describe. *arXiv preprint*  
631 *arXiv:2308.05037*, 2023b.
- 632 Simian Luo, Chuanhao Yan, Chenxu Hu, and Hang Zhao. Diff-foley: Synchronized video-to-audio  
633 synthesis with latent diffusion models. *Advances in Neural Information Processing Systems*, 36,  
634 2024.
- 635 Navonil Majumder, Chia-Yu Hung, Deepanway Ghosal, Wei-Ning Hsu, Rada Mihalcea, and Sou-  
636 janya Poria. Tango 2: Aligning diffusion-based text-to-audio generations through direct prefer-  
637 ence optimization. *arXiv preprint arXiv:2404.09956*, 2024.
- 638 Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen. Metrics for polyphonic sound event  
639 detection. *Applied Sciences*, 6(6):162, 2016.
- 640 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
641 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-  
642 low instructions with human feedback. *Advances in neural information processing systems*, 35:  
643 27730–27744, 2022.

- 648 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi  
649 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An  
650 open large-scale dataset for training next generation image-text models. *Advances in Neural  
651 Information Processing Systems*, 35:25278–25294, 2022.
- 652 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
653 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
654 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 655 Apoorv Vyas, Bowen Shi, Matthew Le, Andros Tjandra, Yi-Chiao Wu, Baishan Guo, Jiemin Zhang,  
656 Xinyue Zhang, Robert Adkins, William Ngan, et al. Audiobox: Unified audio generation with  
657 natural language prompts. *arXiv preprint arXiv:2312.15821*, 2023.
- 658 Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,  
659 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using  
660 direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
661 and Pattern Recognition*, pp. 8228–8238, 2024.
- 662 Yongqi Wang, Wenxiang Guo, Rongjie Huang, Jiawei Huang, Zehan Wang, Fuming You, Ruiqi Li,  
663 and Zhou Zhao. Frieren: Efficient video-to-audio generation with rectified flow matching. *arXiv  
664 preprint arXiv:2406.00320*, 2024.
- 665 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.  
666 Human preference score v2: A solid benchmark for evaluating human preferences of text-to-  
667 image synthesis. *arXiv preprint arXiv:2306.09341*, 2023a.
- 668 Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov.  
669 Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption  
670 augmentation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and  
671 Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023b.
- 672 Zeyu Xie, Xuenan Xu, Zhizheng Wu, and Mengyue Wu. Picoaudio: Enabling precise times-  
673 tamp and frequency controllability of audio events in text-to-audio generation. *arXiv preprint  
674 arXiv:2407.02869*, 2024.
- 675 Xuenan Xu, Ziyang Ma, Mengyue Wu, and Kai Yu. Towards weakly supervised text-to-audio  
676 grounding. *arXiv preprint arXiv:2401.02584*, 2024.
- 677 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason  
678 Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024.
- 683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A IMPLEMENTATION DETAILS

**Audio Generation** During the audio generation process in T2A-Feedback, all models are set to 100 denoising steps with the DDIM scheduler, and classifier-free guidance is configured at 4.0.

**Training Details** For Acoustic&Harmonic Predictor, we train an extra two-layer MLP projector on the top of CLAP audio representations using Cross Entropy(CE) loss. The predictor is trained using the Adam optimizer with a learning rate of  $1e-2.5$  for 6 epochs on 1,000 manually annotated data. For preference tuning, we employ the AdamW optimizer with a learning rate of  $1e-5$  for both DPO and RAFT strategy, and train one epoch for both Audio-Alpaca and T2A-Feedback.

## B EXAMPLES FROM T2A-EPICBENCH

1. At a lively beach, the waves crash rhythmically against the shore, providing a soothing melody. Suddenly, a seagull caws overhead, drawing attention from sunbathers. Children’s giggles fill the air as they splash in the water. Just then, a distant drumbeat starts, adding a festive atmosphere to the scene.

2. In a vibrant classroom, the teacher’s voice resonates as she explains a lesson. Suddenly, a pencil rolls off a desk and clatters to the floor, causing a brief distraction. A student whispers a joke, provoking a wave of giggles. Just then, the school bell rings, signaling the end of the period and the excitement of freedom.

3. In a busy city street, the honking of cars creates a chaotic symphony. Suddenly, a bicycle bell rings sharply as a cyclist weaves through traffic. The murmur of pedestrians chatting fills the air, blending with the distant sound of street performers playing music. Just then, the sound of footsteps approaches, adding to the urban rhythm.

4. At a busy construction site, the sound of drills and saws fills the air, creating a symphony of labor. Suddenly, a heavy beam falls with a loud thud, causing workers to pause. A whistle blows, signaling a break, and conversations buzz among the crew. Just then, a truck backs up, beeping insistently as it arrives.

5. In a vibrant downtown area, the honking of cars creates a chaotic symphony. Suddenly, a street vendor shouts out their specials, trying to attract customers. The laughter of people enjoying a nearby café adds warmth to the urban sounds. Just then, a bus rumbles past, its engine growling as it continues on its route.

6. In a vibrant market, the chatter of vendors fills the air as they hawk their goods. Suddenly, a loud crash echoes as a stack of crates falls over, causing startled gasps. A nearby musician strums a guitar, trying to restore the upbeat mood. Just then, a child squeals with delight, tugging at their parent’s hand to explore further.

7. In a sunlit meadow, the buzzing of bees fills the air as they flit from flower to flower. Suddenly, a cow moos softly from a nearby barn, adding a pastoral charm. A couple of children giggle as they run through the grass, their joyful sounds mingling with nature. Just then, a breeze stirs, causing the wind chimes to tinkle gently.

8. In a tranquil village square, the chirping of crickets fills the evening air. Suddenly, a family gathers around a fire pit, laughter and chatter rising in the dusk. The crackle of flames adds warmth to the scene. Just then, the distant call of an owl echoes, signaling the approach of night.

9. In a dense forest, the soft rustle of leaves whispers through the trees as a gentle breeze blows. Suddenly, a twig snaps underfoot, startling a nearby deer, which bounds away with a soft thud. A bird sings a cheerful melody, filling the air with life. Just then, the distant sound of a waterfall cascades, creating a peaceful backdrop to the vibrant sounds of nature.

## C LLM PROMPTS

The LLM used in Section 1 to separate basic events in the audio description, and in Section 2 for caption augmentation, is Mistral-7B-instruct-v0.2. The respective prompts are provided below:

```
messages = [ "role": "user", "content": "'I will provide a description of an audio, you need
to break down the sentence and figure out the single sound elements. Use the words that
appears in the sentence and appropriately replace the demonstrative pronouns if possible,
such as it, she, him. The output should only include the decomposed sub-events. Here are
some examples:
Description: A man speaks while ducks honk then birds vocalize. '",
"role": "assistant", "content": "Answers: (a man speaks while duck honk)@(birds vocal-
ize)",
"role": "user", "content": "Description: Rain falls on a hard surface.",
"role": "assistant", "content": "Answers: (rain falls on a hard surface)",
"role": "user", "content": "Description: A female makes a speech into a microphone and it
is very loud.",
"role": "assistant", "content": "Answers: (a female makes a speech into a microphone)",
"role": "user", "content": "Description: {new caption}" ]
```

```
messages = [ "role": "user", "content": "'Generate a sentence that contains several different
sounds to make a relative whole story, organize them by words indicating time order, don't
describe the things irrelevant to sound. First print out the generated sentence and later list
the single sound events in the time order they occur. Here are some examples:
Events: honking of a toy trumpet; dog's howl; child's laughter'",
"role": "assistant", "content": "Answer: Child's laughter rings out, followed by the obnox-
ious honking of a toy trumpet. After that, a dog's howl reverberates.
(child's laughter)@(honking of a toy trumpet)@(dog's howl)",
"role": "user", "content": "Events: {new caption}" ]
```

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809