

T2A-Feedback: Improving Basic Capabilities of Text-to-Audio Generation via Fine-grained AI Feedback

Anonymous ACL submission

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。)

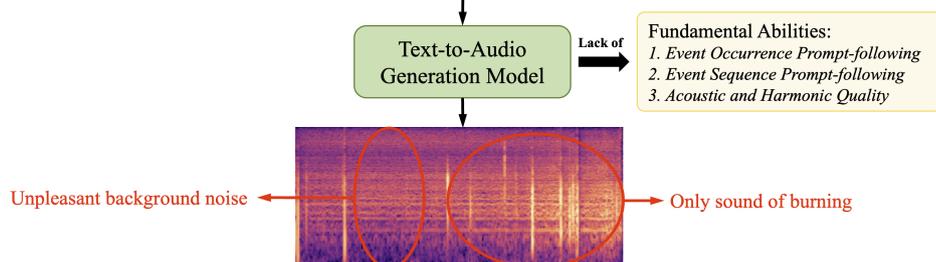


Figure 1: The audio description is from a classic Chinese essay “Kou Ji”, which vividly depicts a performer using only vocal mimicry to recreate an entire dramatic scene. The existing Text-to-Audio generation model struggles to generate such narrative and multi-event audios. The generated audio often 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

057	narrative and complex descriptions, which limits	prompt. A higher score implies a greater simi-	107
058	the potential for high-level applications.	larity between the event sequences in caption	108
059	The failure of the generated results is often	and audio.	109
060	demonstrated in three aspects: 1) cannot fully in-		
061	clude all the events described, 2) cannot accurately	• <i>Acoustics&Harmonic Quality</i> : Drawing inspi-	110
062	follow the order of all the events described, and 3)	ration from aesthetic scoring methods used	111
063	cannot organize all the events harmoniously. There-	in image quality scoring, we manually anno-	112
064	fore, the model performance in multi-event scenar-	tate acoustic and harmonic quality for audio	113
065	ios is determined by its capabilities in these three	samples. These data are then used to train an	114
066	fundamental aspects.	automatic acoustic&harmonic predictor.	115
067	To improve the model’s performance across		116
068	more advanced applications, we focus on strength-	We confirm that our three scoring functions show	117
069	ening the audio generation model’s fundamental	a stronger correlation with human evaluations com-	118
070	abilities. Inspired by feedback learning in large	pared to existing automatic audio evaluation meth-	119
071	language models (Ouyang et al., 2022; Bai et al.,	ods (Wu et al., 2023b; Xie et al., 2024). Conse-	120
072	2022; Touvron et al., 2023), we propose creating an	quently, in addition to their application in ranking	121
073	audio preference dataset centered on three abilities	preference data, these scoring functions can be used	122
074	necessary for generating harmonic and complex	as evaluation metrics that more effectively capture	123
075	audio: 1) Event Occurrence Prompt-Following ,	human preferences across different aspects.	124
076	2) Event Sequence Prompt-Following , and 3)	Leveraging these advanced AI scoring pipelines,	125
077	Acoustic&Harmonic Quality . Based on the pref-	we establish a comprehensive data collection and	126
078	erence information, we can refine the model’s core	annotation framework. As a result, we construct	127
079	abilities, resulting in better results in both simple	T2A-Feedback , a large audio preference dataset	128
080	and challenging scenarios.	comprising 41,627 captions and 249,762 generated	129
081	However, due to the scarcity of audio data and	audios, each annotated with detailed scores.	130
082	the challenges of annotating the scale of user pref-	Furthermore, to evaluate the higher-level capa-	131
083	erences, it is difficult to collect massive audio pref-	abilities of text-to-audio models in multi-event sce-	132
084	erence datasets that only rely on human annotators.	narios, we introduce a more challenging bench-	133
085	To fill this void, we explore using AI feedback (Cui	mark, T2A-EpicBench , which features longer,	134
086	et al., 2023; Lee et al., 2023; Yuan et al., 2024;	more imaginative, and story-telling captions for	135
087	Burns et al., 2023) in text-to-audio generation, uti-	audio generation. We enhance the advanced text-to-	136
088	lizing AI models to rank audios instead of relying	audio diffusion model, Make-an-Audio 2 (Huang	137
089	on human annotators. Compared to manual annota-	et al., 2023a), with T2A-Feedback. Our results	138
090	tion, automating the data collection and annotation	show that using T2A-Feedback not only effecti-	139
091	process reduces the cost of obtaining audio prefer-	vely improves the basic capabilities of the model	140
092	ence data and enhances scalability.	in simple AudioCaps benchmark, but also emer-	141
093	Specifically, we develop three AI scoring	gently improves the performance in complex T2A-	142
094	pipelines to evaluate the generated audio in a fine-	EpicBench.	143
095	grained and holistic manner, corresponding to three		
096	core capabilities:	2 Related Work	144
097	• <i>Event Occurrence Score</i> : To specifically	2.1 Text-to-Audio Generation	144
098	check whether each event occurs in, we calcu-	Text-to-audio generation is an emerging field that	145
099	late the audio-text semantic matching score	aims to convert textual descriptions into corre-	146
100	for each described event separately. A lower	sponding audio outputs. Existing text-to-audio	147
101	score suggests that the corresponding event	generation methods can be divided into two cat-	148
102	might be absent from the audio.	egories: Diffusion-based and Language model-	149
103	• <i>Event Sequence Score</i> : To verify the cor-	-based. Diffusion-based techniques have gained	150
104	rectness of event order, we analyze the start	prominence for generating high-quality, realis-	151
105	and end times of each event and compare	tic audio by modeling the process of denoising.	152
106	them with the event order outlined in the text	These methods, like Make-an-Audio (Huang et al.,	153
		2023b,a), AudioLDM (Liu et al., 2023a, 2024),	154
		Tango (Ghosal et al., 2023; Majumder et al., 2024),	155

156	start with random noise and iteratively refine it to	One recent paper related to our work, Tango2 (Ma-	207
157	produce coherent audio over a series of denois-	jumder et al., 2024), utilizes contrastive language-	208
158	ing steps. On the other hand, Language model-	audio pre-training (CLAP) (Wu et al., 2023b) to	209
159	based methods (Borsos et al., 2023; Agostinelli	rank audio generated by the Tango model. How-	210
160	et al., 2023; Cideron et al., 2024) tokenize audios	ever, CLAP can only evaluate the global alignment	211
161	as acoustic discrete tokens, and predict the tokens	between audio and text but falls short in assessing	212
162	within an auto-regressive model conditioned on	the fine-grained details, like detailed event occur-	213
163	text inputs.	rence, sequence, and harmony. In this paper, we	214
164	The above models acquire the ability to gener-	construct more robust AI audio scoring pipelines	215
165	ate diverse audio by training on large-scale audio-	with fine-grained recognition ability. Our method	216
166	text datasets. However, current datasets like Au-	shows a much stronger correlation with human	217
167	dioSet (Gemmeke et al., 2017), AudioCaps (Kim	preference and the constructed dataset brings sig-	218
168	et al., 2019), and FSD50k (Fonseca et al., 2021)	nificant improvement to the current text-to-audio	219
169	only provide tag-level annotations or short captions.	generation model.	220
170	As a result, when processing long, detailed lan-		
171	guage prompts, existing models often produce low-	2.3 Text-to-Audio Evaluation Metric	221
172	quality, noisy outputs and struggle to accurately	Existing evaluation metrics for audio generation,	222
173	follow the text. Due to the difficulty of annotating	such as FAD and IS, assess audio distributions but	223
174	detailed audio captions, scaling rich and accurate	cannot evaluate the quality of individual samples.	224
175	audio descriptions remains a challenge. In this	Additionally, many studies rely on similarity scores	225
176	work, we focus on enhancing the model’s basic	from the CLAP model to assess global audio-text	226
177	abilities in event occurrence, sequence, and har-	semantic alignment. PicoAudio (Xie et al., 2024)	227
178	mony, thereby improving its performance in both	uses a text-to-audio grounding model (Xu et al.,	228
179	simple scenarios and advanced applications.	2024) to detect audio segments based on language	229
180	2.2 Preference Tuning with Human&AI	prompts. However, there remains a lack of fine-	230
181	Feedback	grained evaluation methods for assessing detailed	231
182	Tuning generative models according to human pref-	event occurrence, sequencing, and acoustic quality.	232
183	erences has emerged as a standard practice for im-	Our research fills this gap by creating robust audio	233
184	proving the quality of outputs. By tuning with feed-	AI scoring pipelines, that show a strong correlation	234
185	back information on different aspects, the model	with humans, and significantly surpass alternative	235
186	can be improved and aligned with human pref-	methods.	236
187	erences in corresponding aspects. Traditionally,	3 T2A-Feedback	237
188	this preference data used for tuning relied heav-	In this section, we first dive into the three AI audio	238
189	ily on human evaluators who rank multiple gener-	scoring pipelines: (i) Event Occurrence Prompt-	239
190	ated results, assessing their quality based on vari-	following, in Section. 3.1; (ii) Event Sequence	240
191	ous criteria such as relevance, coherence, and aes-	Prompt-following, in Section. 3.2; (iii) Acoustic	241
192	thetic value (Bai et al., 2022; Touvron et al., 2023;	Quality, in Section. 3.3. We then describe the spe-	242
193	Ouyang et al., 2022; Kirstain et al., 2023; Liang	cific data generation and sorting method for the	243
194	et al., 2024; Wu et al., 2023a; Cideron et al., 2024).	T2A-Feedback dataset in Section. 3.4.	244
195	While effective, manual human annotation is	3.1 Events Occurrence Prompt-following	245
196	costly and time-consuming, which greatly hampers	Generating audio that accurately reflects the events	246
197	the scalability of preference tuning across more	described in a given prompt is the fundamental re-	247
198	diverse generative tasks. To address the difficulty,	quirement of prompt-following. However, when	248
199	more recent developments have focused on leverag-	multiple events are included in the text description,	249
200	ing pre-trained AI models to automate the process	current text-to-audio generation models often strug-	250
201	of scoring generated content (Cui et al., 2023; Lee	gle to generate each event precisely. To improve	251
202	et al., 2023; Yuan et al., 2024; Burns et al., 2023).	the generation model’s event occurrence prompt-	252
203	Such an AI feedback approach has achieved im-	following ability, we first build an AI pipeline to	253
204	pressive improvements in large language models.	determine the occurrence of events in audio.	254
205	Recently, some studies have attempted prefer-		
206	ence fine-tuning in text-to-audio generation models.		

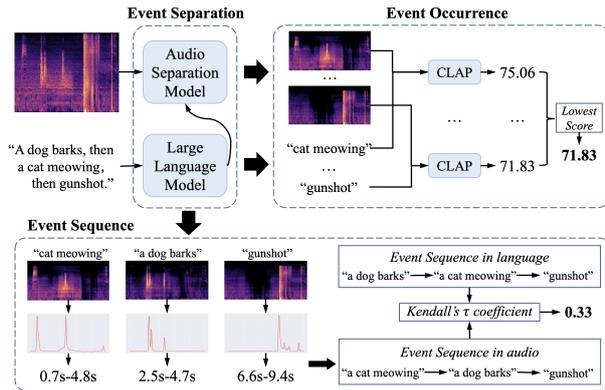


Figure 2: The overview of event occurrence and sequence scoring pipelines.

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. 4. 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. 4 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 τ 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 τ coefficient between event sequences in language and audio) is calculated as:

$$\tau = C - Dn(n-1) \quad (1)$$

where C represents the number of concordant event pairs between the description and the audio, D denotes the number of discordant ones. Higher τ 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. Three annotators independently score the audio samples according to the same criteria, and samples with consistent scores are accepted as training data. Detailed scoring criteria is provided in Appendix.

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 2,000 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 cap-

tion. 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 Appendix. 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 storytelling 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.

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 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. The test sets of AudioCaps (3,701 samples), Clotho (5,225 samples),

	AudioCaps	Clotho	MusicCaps
Random Guess	50.0%	50.0%	50.0%
CLAP	77.5%	86.4%	69.4%
PANNs	82.0%	79.9%	56.1%
EOS(ours)	90.9%	90.4%	99.8%

Table 1: Results of event occurrence recognition

	Two Events		Three Events		Correlation
	Acc	F1	Acc	F1	
CLAP	49.6%	-	53.7%	-	-
PicoAudio	71.6%	0.787	51.3	0.574	0.30
ESS _{0.1}	79.6%	0.814	54.2	0.606	0.43
ESS _{0.3}	79.1%	0.851	57.6	0.587	0.52
ESS _{0.5}	78.0%	0.769	56.7	0.535	0.52

Table 2: Results of event sequence recognition. ESS_{0.x} stands for using 0.x as volume thresholds.

and MusicCaps (4,434 samples) are employed for evaluation.

We mainly compare our EOS with CLAP and PANNs. The caption with the higher matching score to the audio is considered as the prediction. For the audio tagging model, PANNs, we match the top 5 recognized audio categories with open-domain descriptions. As shown in Table 1, our EOS score showcases a notable advantage over baselines on all the benchmarks, 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 and 200 three-event samples from PicoAudio’s data, 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 evaluating the Segment F1 Score (Mesaros et al., 2016) for detecting the start and end times of each audio event. Moreover, we manually annotate temporal order alignment for 100 audios generated from our temporal-enhanced captions and compute the correlation between different methods and humans.

The results of event sequences are provided in Table. 2. We compare ESS with CLAP score and the audio grounding model (Xu et al., 2024) used by PicoAudio (Xie et al., 2024). Compared to base-

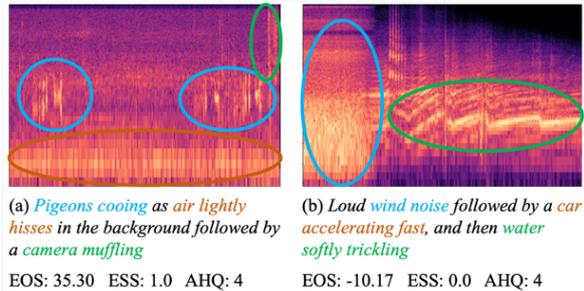


Figure 3: Visualization of the predicted scores from our AI scoring pipeline. We highlight the first, second, and third events described in the captions using blue, brown, and green, respectively.

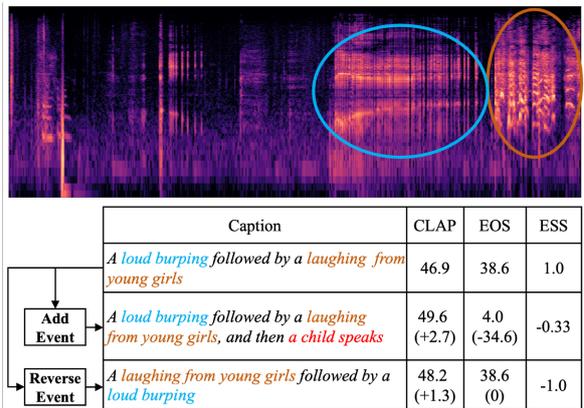


Figure 4: Qualitative comparison between CLAP scores and EOS/ESS scores reveals distinct sensitivities to misalignment. By adding or reversing events in the ground-truth caption, the captions become misaligned with the audio in terms of event occurrence and sequence.

lines, our method distinguishes the ground-truth caption from the distracting one more accurately and achieves higher F1 scores in detecting the start and end times of events in audio. More importantly, our method shows a much stronger correlation to human annotations.

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

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

Moreover, we explore building the Acoustic&Harmonic Predictor on top of various pre-

		FAD↓	KL↓	IS↑	CLAP↑	EOS↑	ESS↑	AQ↑
Make an Audio 2		1.82	1.44	10.03	69.97	42.05	0.53	2.33
Preference Tuning								
Audio-Alpaca	RAFT	1.93	1.29	10.37	72.23	44.85	0.53	2.45
	DPO	3.20	1.24	12.27	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	74.00	49.58	0.57	2.57

Table 3: Evaluation results on AudioCaps. The EOS, ESS and AHQ represent the Event Occurrence Score, Event Sequence Score 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%	64%	59%
T2A-Feedback (ours)	RAFT	52%(15.85)	52%(0.05)	54%(2.14)	61%	57%	61%
	DPO	58%(19.96)	64%(0.13)	52%(2.10)	68%	62%	68%

Table 4: 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.

trained audio models and evaluate how well each variant correlates with human preferences. The correlation of the predictor built on CLAP (Wu et al., 2023b) (0.79) outperforms those based on self-supervised models like AudioMAE (Huang et al., 2022) (0.61) and BEAT (Chen et al., 2022) (0.52). Similarly, the image aesthetics predictor (Schuhmann et al., 2022) is built on the CLIP model (Ilharco et al., 2021). This advantage may stem from the alignment with language, resulting in better semantic discrimination.

5.1.2 Qualitative Analysis

We show some example predictions from our scoring pipelines in Figure. 3, where our methods can specifically identify the misaligned event, the out-of-order event order, and the disharmony between events in the audio. Moreover, we provide the qualitative comparison between our EOS and ESS with the single CLAP score, in Figure. 4. For the ground-truth audio-caption pairs from AudioCaps, we perturb the captions by adding an event or shuffling the order of events. We find that the CLAP score is not sensitive to these perturbations and even yields a higher score with the incorrect, perturbed caption. In contrast, our EOS and ESS scores more accurately reflect the alignment between audio and text regarding event occurrence and event order.

5.2 Analysis of Preference Tuning

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

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. 3. 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

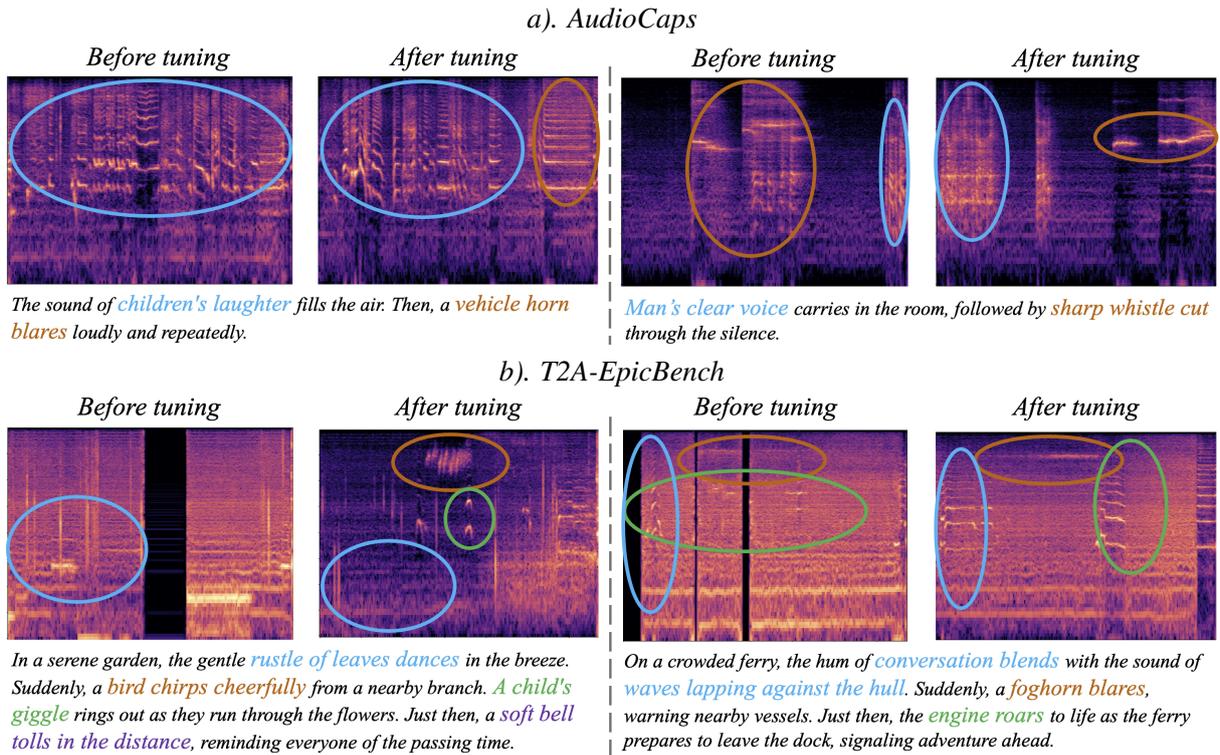


Figure 5: Visualization of the impact of preference tuning with T2A-Feedback.

human annotators evaluate and select the better output based on each criterion.

The results on T2A-EpicBench, are illustrated in Table. 4, 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.

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, the audio generation model can emergently learn to handle more complex long-text and multi-event scenarios.

5.2.3 Qualitative Findings

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. 5. For the examples of short captions in Figure. 5a, 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. 5b, significantly reduces background noise and generates audio that more faithfully captures both events and their orders.

6 Conclusion

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, 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 basic capabilities and achieves better performance in both simple (AudioCaps) and complex (T2A-EpicBench) scenarios.

626 Limitation

627 Automatically generating high-quality and harmo-
628 nious audio from detailed, narrative, and multi-
629 event scenarios remains a long-term goal. The per-
630 formance of the audio generation model depends
631 on both the pre-training phase and the post-training
632 phase (fine-tuning and feedback learning). To fully
633 address the challenge of generating coherent au-
634 dio for long narrative prompts, improvements are
635 needed across the entire process.

636 References

637 Andrea Agostinelli, Timo I Denk, Zalán Borsos,
638 Jesse Engel, Mauro Verzetti, Antoine Caillon,
639 Qingqing Huang, Aren Jansen, Adam Roberts, Marco
640 Tagliasacchi, et al. 2023. Musiclm: Generating mu-
641 sic from text. *arXiv preprint arXiv:2301.11325*.

642 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda
643 Askeell, Anna Chen, Nova DasSarma, Dawn Drain,
644 Stanislav Fort, Deep Ganguli, Tom Henighan, et al.
645 2022. Training a helpful and harmless assistant with
646 reinforcement learning from human feedback. *arXiv*
647 *preprint arXiv:2204.05862*.

648 Zalán Borsos, Raphaël Marinier, Damien Vincent,
649 Eugene Kharitonov, Olivier Pietquin, Matt Shar-
650 ifi, Dominik Roblek, Olivier Teboul, David Grang-
651 ier, Marco Tagliasacchi, et al. 2023. Audioldm: a
652 language modeling approach to audio generation.
653 *IEEE/ACM transactions on audio, speech, and lan-*
654 *guage processing*, 31:2523–2533.

655 Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner,
656 Bowen Baker, Leo Gao, Leopold Aschenbrenner,
657 Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan
658 Leike, et al. 2023. Weak-to-strong generalization:
659 Eliciting strong capabilities with weak supervision.
660 *arXiv preprint arXiv:2312.09390*.

661 Honglie Chen, Weidi Xie, Andrea Vedaldi, and An-
662 drew Zisserman. 2020. Vggsound: A large-scale
663 audio-visual dataset. In *ICASSP 2020-2020 IEEE*
664 *International Conference on Acoustics, Speech and*
665 *Signal Processing (ICASSP)*, pages 721–725. IEEE.

666 Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu,
667 Daniel Tompkins, Zhuo Chen, and Furu Wei. 2022.
668 Beats: Audio pre-training with acoustic tokenizers.
669 *arXiv preprint arXiv:2212.09058*.

670 Geoffrey Cideron, Sertan Girgin, Mauro Verzetti,
671 Damien Vincent, Matej Kastelic, Zalán Borsos, Brian
672 McWilliams, Victor Ungureanu, Olivier Bachem,
673 Olivier Pietquin, et al. 2024. Musicrl: Aligning mu-
674 sic generation to human preferences. *arXiv preprint*
675 *arXiv:2402.04229*.

676 Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao,
677 Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and

Maosong Sun. 2023. Ultrafeedback: Boosting lan-
guage models with high-quality feedback. *arXiv*
preprint arXiv:2310.01377. 678
679
680

Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan
Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng
Zhang, Kashun Shum, and Tong Zhang. 2023. Raft:
Reward ranked finetuning for generative foundation
model alignment. *arXiv preprint arXiv:2304.06767*. 681
682
683
684
685

Eduardo Fonseca, Xavier Favory, Jordi Pons, Frederic
Font, and Xavier Serra. 2021. Fsd50k: an open
dataset of human-labeled sound events. *IEEE/ACM*
Transactions on Audio, Speech, and Language Pro-
cessing, 30:829–852. 686
687
688
689
690

Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman,
Aren Jansen, Wade Lawrence, R Channing Moore,
Manoj Plakal, and Marvin Ritter. 2017. Audio set:
An ontology and human-labeled dataset for audio
events. In *2017 IEEE international conference on*
acoustics, speech and signal processing (ICASSP),
pages 776–780. IEEE. 691
692
693
694
695
696
697

Deepanway Ghosal, Navonil Majumder, Ambuj
Mehrish, and Soujanya Poria. 2023. Text-to-audio
generation using instruction-tuned llm and latent dif-
fusion model. *arXiv preprint arXiv:2304.13731*. 698
699
700
701

Jiawei Huang, Yi Ren, Rongjie Huang, Dongchao Yang,
Zhenhui Ye, Chen Zhang, Jinglin Liu, Xiang Yin,
Zejun Ma, and Zhou Zhao. 2023a. Make-an-audio 2:
Temporal-enhanced text-to-audio generation. *arXiv*
preprint arXiv:2305.18474. 702
703
704
705
706

Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski,
Michael Auli, Wojciech Galuba, Florian Metze, and
Christoph Feichtenhofer. 2022. Masked autoen-
coders that listen. *Advances in Neural Information*
Processing Systems, 35:28708–28720. 707
708
709
710
711

Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren,
Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xi-
ang Yin, and Zhou Zhao. 2023b. Make-an-audio:
Text-to-audio generation with prompt-enhanced dif-
fusion models. In *International Conference on Ma-*
chine Learning, pages 13916–13932. PMLR. 712
713
714
715
716
717

Gabriel Ilharco, Mitchell Wortsman, Ross Wightman,
Cade Gordon, Nicholas Carlini, Rohan Taori, Achal
Dave, Vaishal Shankar, Hongseok Namkoong, John
Miller, Hannaneh Hajishirzi, Ali Farhadi, and Lud-
wig Schmidt. 2021. [Openclip](#). If you use this soft-
ware, please cite it as below. 718
719
720
721
722
723

Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-
sch, Chris Bamford, Devendra Singh Chaplot, Diego
de las Casas, Florian Bressand, Gianna Lengyel, Guil-
laume Lample, Lucile Saulnier, et al. 2023. Mistral
7b. *arXiv preprint arXiv:2310.06825*. 724
725
726
727
728

Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee,
and Gunhee Kim. 2019. Audiocaps: Generating cap-
tions for audios in the wild. In *Proceedings of the*
2019 Conference of the North American Chapter of
the Association for Computational Linguistics. 729
730
731
732

733		Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. <i>Advances in Neural Information Processing Systems</i> , 35:25278–25294.	789
734			790
735			791
736	Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. 2023. Pick-a-pic: An open dataset of user preferences for text-to-image generation. <i>Advances in Neural Information Processing Systems</i> , 36:36652–36663.		792
737			793
738			794
739			795
740			
741	Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. <i>arXiv preprint arXiv:2309.00267</i> .	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	796
742			797
743			798
744			799
745			800
746			801
747	Youwei Liang, Junfeng He, Gang Li, Peizhao Li, Arseniy Klimovskiy, Nicholas Carolan, Jiao Sun, Jordi Pont-Tuset, Sarah Young, Feng Yang, et al. 2024. Rich human feedback for text-to-image generation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 19401–19411.	Apoorv Vyas, Bowen Shi, Matthew Le, Andros Tjandra, Yi-Chiao Wu, Baishan Guo, Jiemin Zhang, Xinyue Zhang, Robert Adkins, William Ngan, et al. 2023. Audiobox: Unified audio generation with natural language prompts. <i>arXiv preprint arXiv:2312.15821</i> .	802
748			803
749			804
750			805
751			806
752			
753			
754	Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. 2023a. Audioldm: Text-to-audio generation with latent diffusion models. <i>arXiv preprint arXiv:2301.12503</i> .	Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. 2024. Diffusion model alignment using direct preference optimization. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 8228–8238.	807
755			808
756			809
757			810
758			811
759	Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley. 2024. Audioldm 2: Learning holistic audio generation with self-supervised pretraining. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> .	Yongqi Wang, Wenxiang Guo, Rongjie Huang, Jiawei Huang, Zehan Wang, Fuming You, Ruiqi Li, and Zhou Zhao. 2024. Frieren: Efficient video-to-audio generation with rectified flow matching. <i>arXiv preprint arXiv:2406.00320</i> .	812
760			813
761			
762			
763			
764			
765	Xubo Liu, Qiuqiang Kong, Yan Zhao, Haohe Liu, Yi Yuan, Yuzhuo Liu, Rui Xia, Yuxuan Wang, Mark D Plumbley, and Wenwu Wang. 2023b. Separate anything you describe. <i>arXiv preprint arXiv:2308.05037</i> .	Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. 2023a. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. <i>arXiv preprint arXiv:2306.09341</i> .	814
766			815
767			816
768			817
769			818
770	Simian Luo, Chuanhao Yan, Chenxu Hu, and Hang Zhao. 2024. Diff-foley: Synchronized video-to-audio synthesis with latent diffusion models. <i>Advances in Neural Information Processing Systems</i> , 36.	Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. 2023b. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In <i>ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 1–5. IEEE.	819
771			820
772			821
773			822
774			823
775	Navonil Majumder, Chia-Yu Hung, Deepanway Ghosal, Wei-Ning Hsu, Rada Mihalcea, and Soujanya Poria. 2024. Tango 2: Aligning diffusion-based text-to-audio generations through direct preference optimization. <i>arXiv preprint arXiv:2404.09956</i> .	Zeyu Xie, Xuenan Xu, Zhizheng Wu, and Mengyue Wu. 2024. Picoaudio: Enabling precise timestamp and frequency controllability of audio events in text-to-audio generation. <i>arXiv preprint arXiv:2407.02869</i> .	824
776			825
777			826
778			827
779			828
780	Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen. 2016. Metrics for polyphonic sound event detection. <i>Applied Sciences</i> , 6(6):162.	Xuenan Xu, Ziyang Ma, Mengyue Wu, and Kai Yu. 2024. Towards weakly supervised text-to-audio grounding. <i>arXiv preprint arXiv:2401.02584</i> .	829
781			830
782			
783	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-rewarding language models. <i>arXiv preprint arXiv:2401.10020</i> .	831
784			832
785			833
786			834
787			
788			
		Wendi Zheng, Jiayan Teng, Zhuoyi Yang, Weihang Wang, Jidong Chen, Xiaotao Gu, Yuxiao Dong, Ming Ding, and Jie Tang. 2024. Cogview3: Finer and faster	835
			836
			837
			838
			839
			840
			841
			842
			843
			844

A 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.

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 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.

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.

B Scoring Criteria for Acoustic&Harmonic Quality

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 harmony is lacking; Background noise affects clarity; and dynamic range is limited.
3	Good audio quality; most events are realistic with good integration; Background noise is minimally disruptive; and dynamic range is reasonable.
4	Excellent audio quality; all events are realistic with perfect integration; well-managed background noise; and wide dynamic range.

C 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.

D Other models on T2A-EpicBench

The performance of AudioLDM 2 and Tango 2 on T2A-EpicBench is as follows:

	EOS	ESS	AHQ
Make an Audio 2	14.21	0.03	1.96
AudioLDM2	16.35	0.04	1.98
Tango2	19.42	0.07	2.11
Make an Audio 2 + T2A-Feedback (DPO)	19.96	0.13	2.10

Table 5: Results of AudioLDM2 and Tango2 on our T2A-EpicBench.

	AudioCaps	Clotho	MusicCaps
Average	89.3	88.8	99.8
Lowest	90.9	90.4	99.8

Table 6: Comparison between selecting lowest or average score for event occurrence score

As shown in Table. 5, the improvements observed across Make-an-audio 2, AudioLDM2, and Tango2 on EpicBench align with their inherent capabilities, with newer and more advanced models achieving better results. This indirectly validates the robustness and effectiveness of our benchmark and AI-based scoring pipeline.

Moreover, we observed that although the Make-an-audio 2 model does not perform well on EpicBench initially, it achieves the best performance after feedback alignment with T2A-Feedback. This highlights the practicality and significance of our dataset.

E Study about Lowest Score for EOS

We tested the effect of selecting the average score and the lowest score among all matching scores for event occurrence judgment in Table. 6. We find that using the lowest score can better distinguish the caption with extra events for current audio-text datasets. According to the statistical results, we empirically select the lowest score for event occurrence.

F Negative Effect to FAD Score

FAD and FID estimate the mean and covariance of two sample groups in a high-dimensional feature space and calculate their similarity. A negative correlation between FAD (FID) and subjective metrics is widely observed in the text-to-image and text-to-audio generations. The study Pick-a-Pic (Kirstain et al., 2023) for text-to-image feedback learning has discussed this phenomenon, suggesting that it

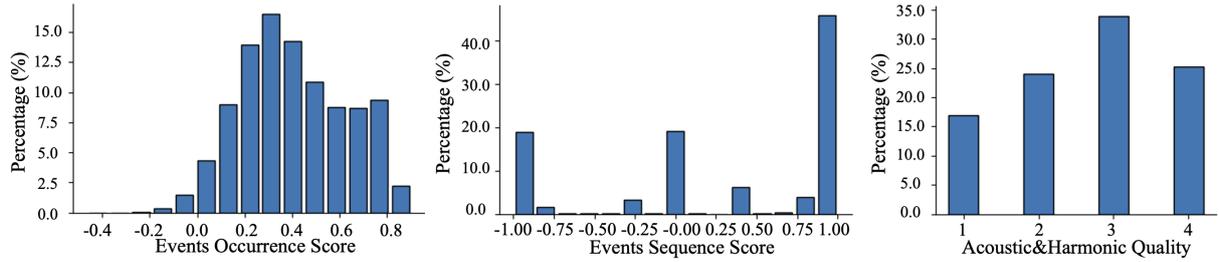


Figure 6: Histograms of three different scores in T2A-Feedback.

	A-1	A-2	A-3	Majority
Predictor	70.31%	68.75%	62.50%	73.44%
A-1	-	64.06%	68.75%	74.33%
A-2	64.06%	-	65.63%	71.88%
A-3	68.75%	65.63%	-	70.31%

Table 7: Agreement between AHQ annotations and predictions on 100 testing samples. A-1, A-2 and A-3 are three human annotators. “Majority” stands for the agreement between each judge and the other three judges’s majority votes.

I Potential Risks

Since our audio data is generated by the model based on the text, its content is mainly determined by the provided text. Therefore, if using generation models without safety checkers, offensive and unsafe content may be generated. In our work, we checked the content of the text prompt to ensure that the generated data does not contain offensive content.

may be correlated to the classifier-free guidance scale mechanism. Larger classifier-free guidance scales tend to produce more vivid samples, which humans generally prefer, but deviate from the distribution of ground truth samples in the test set, resulting in worse (higher) FID (FAD) scores.

More specifically, this phenomenon is witnessed in Tables 1 and 2 of CogView3 (Zheng et al., 2024) (text-to-image method) and Table 3 of Tango2 (Majumder et al., 2024) (text-to-audio method), where models achieve higher human preference scores but worse FID (FAD) scores. The negative correlation between FID (FAD) and subjective scores, as consistently shown by previous methods, appears to be an expected outcome when aligning generative models with human preferences.

G Statistic of Each Score

We provide the histogram maps of three different scores in Figure. 6.

H Agreement between AHQ Annotations

We provide the agreement between Acoustic&Harmonic Quality (AHQ) annotations and predictions in Table. 7. All the annotators exhibit an agreement rate of over 70% with the majority vote, which demonstrates the reliability of our annotation process.