

Synthio: Augmenting Small-Scale Audio Classification Datasets with Synthetic Data

Sreyan Ghosh^{♦♦*}, Sonal Kumar^{♦*}, Zhifeng Kong[♦], Rafael Valle[♦], Bryan Catanzaro[♦], Dinesh Manocha[♦]
 ♦NVIDIA, CA, USA, ♪University of Maryland, College Park, USA

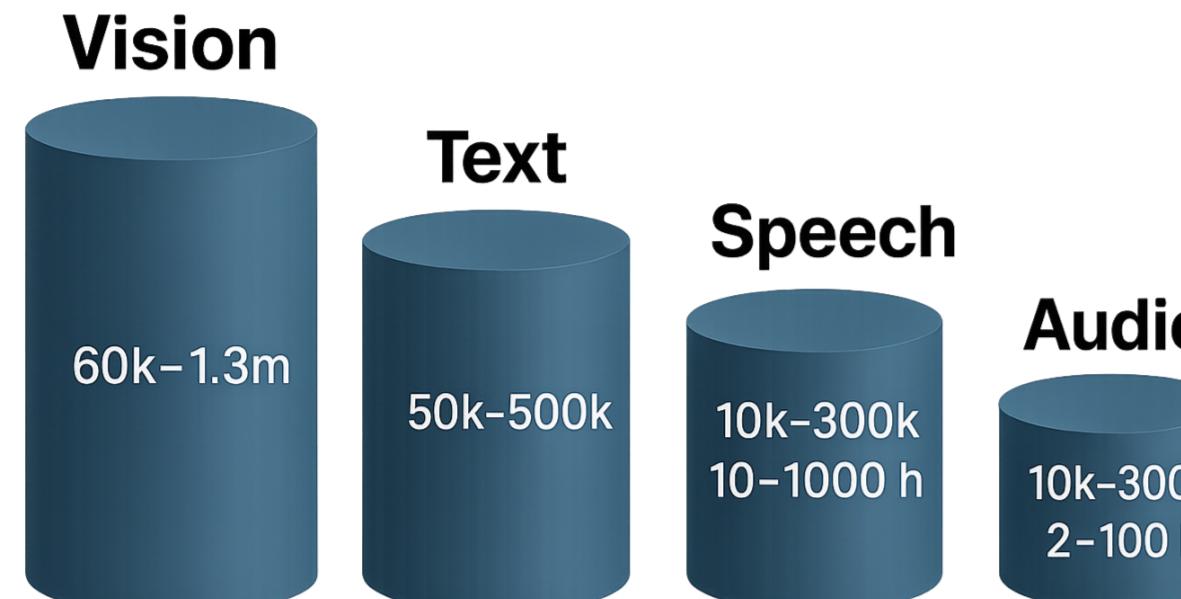
Modular framework for
Synthetic Data Generation

T2A model trained on ~8M
data points



Audio Understanding Suffers from the Lack of Enough Data

- (1) Audio (sounds & music) is low-resource compared to text, vision, and speech.
- (2) Clean, diverse audio is hard to scrape or collect.
- (3) Models trained on benchmarks generalize poorly to audio beyond the specific settings.
- (4) Traditional augmentations alter data statistically—not semantically—limiting diversity and realism. They focus on surface-level patterns in the data rather than capturing the fundamental mechanisms that drive real-world data generation.



Problems with Generating Synthetic Data for Audio

Synthetic data can help—but comes with hurdles:

- Lack of strong open-source text-to-audio (T2A) models to generate **diverse** data at scale.
- Shortage of large-scale training data for training T2A models.
- T2A models trained on internet scale data struggle to generate audios **acoustically consistent** with the downstream dataset.
- T2A models struggle with prompt-following and fine-grained control, making tailored generation hard.



Figure 1. (Left:) Acoustic characteristics for the label **bus** vary by region or dataset (e.g., loud engines vs. quiet ones). This hurts consistency. (Right:) T2A models spuriously associate the label **airport** with announcements, which vary across regions. This hurts diversity.

Synthio: A Novel Synthetic Data Generation Pipeline for Audio

We propose Synthio, a new pipeline for generating high-quality synthetic audio. The primary aim is to generate synthetic data to expand small-scale audio classification datasets. It has 3 key parts:

- (1) **A T2A model trained on large-scale data:** We train a DiT-based T2A model ~8M cleaned audio-caption pairs (using Stable Audio Open), ensuring no leakage with benchmark datasets.
- (2) **Aligning the Text-to-Audio Models with Preference Optimization:** We align T2A outputs with dataset-specific acoustic traits using **preference optimization**—so the synthetic audio “sounds like” examples in the downstream dataset.
- (3) **Generating Diverse Synthetic Augmentations:** We propose **MixCap**, where we diverse synthetic augmentations. Figure 2 shows Synthio outperforms other methods on ESC-50 (100 samples) by generating consistent and prompt the LLM to remix acoustic elements and imagine new soundscapes. Additionally, a self-reflection module filters and improves captions.

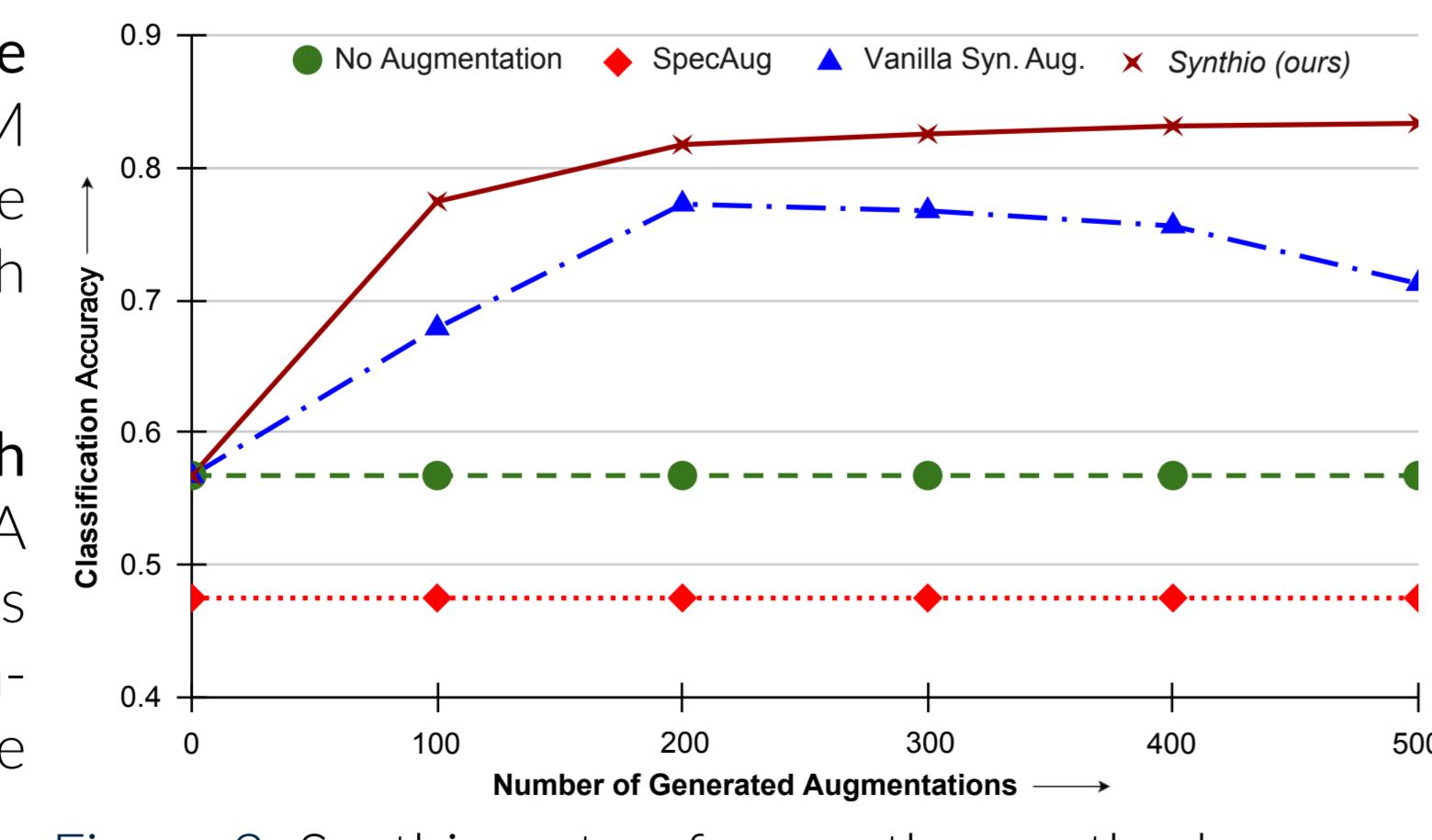
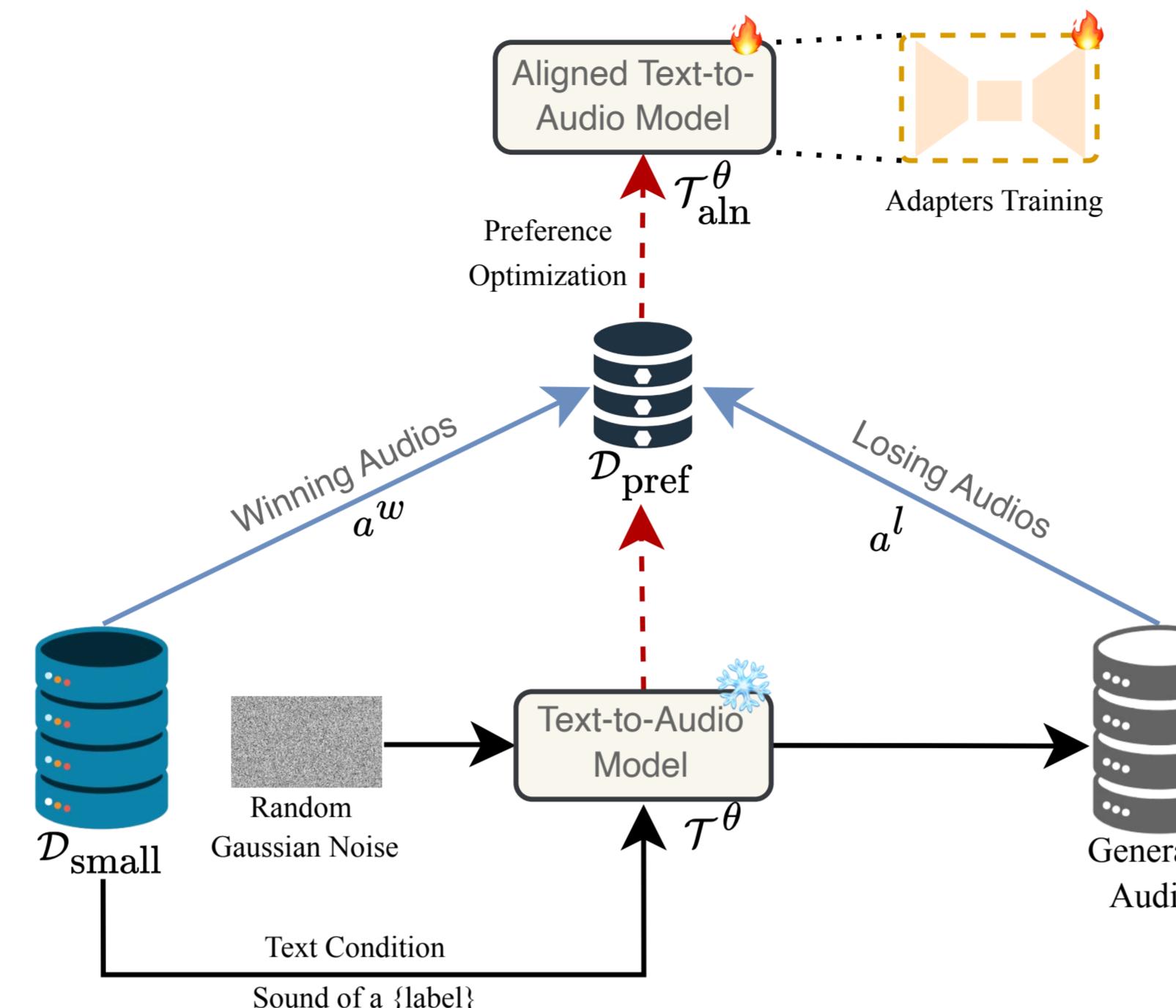


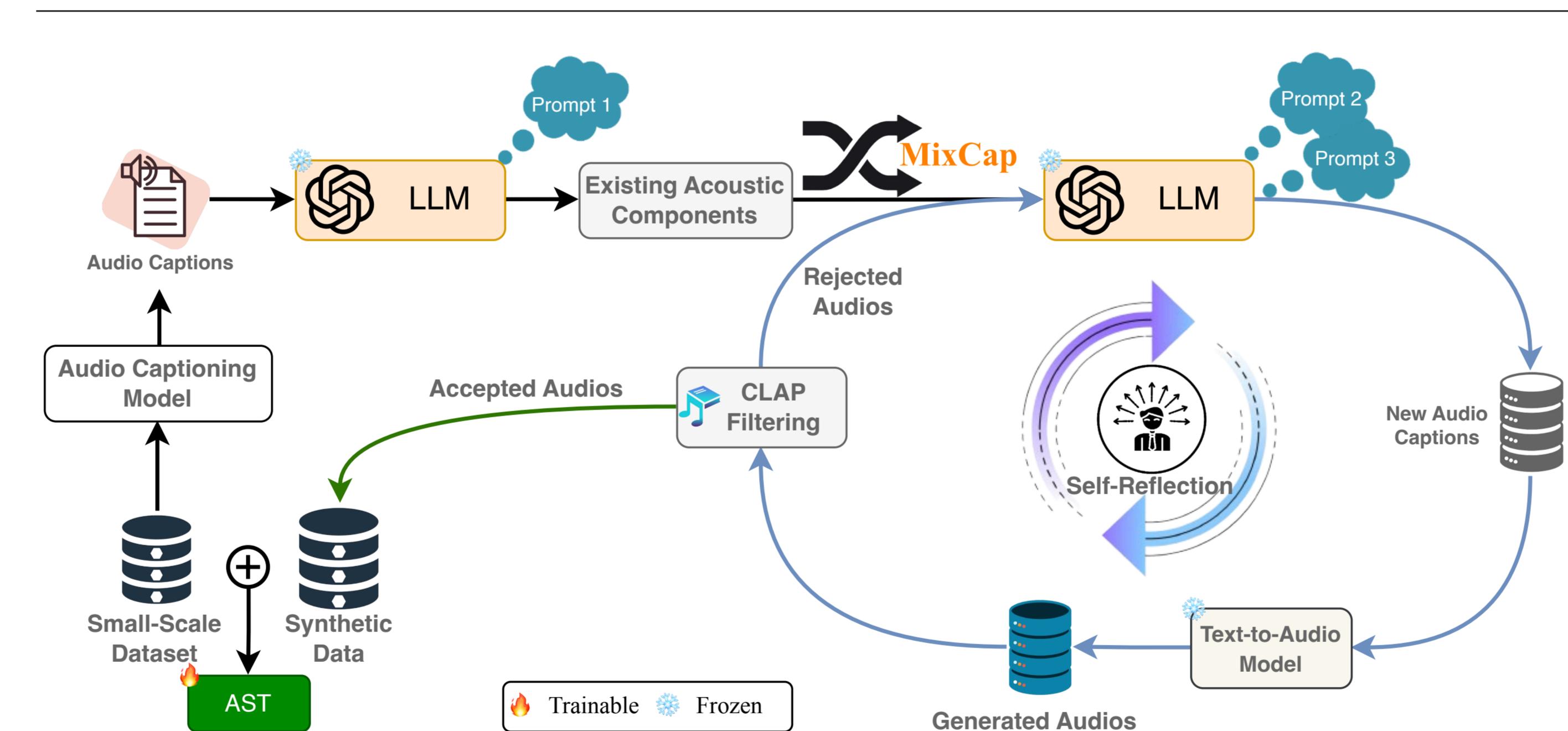
Figure 2. Synthio outperforms other methods on ESC-50 (100 samples) by generating consistent and prompt the LLM to remix acoustic elements and imagine new soundscapes. Additionally, a self-reflection module filters and improves captions.

1. Aligning the T2A Model with Preference Optimization



- We propose aligning the student (T2A model) with the teacher (human-labeled dataset) using preference optimization (DPO) to ensure generated audio sounds similar to desired dataset characteristics
- In **Step 1**, we create a preference dataset by generating multiple synthetic audios per caption and pairing them with gold-standard audio, labeling them as “winner” or “loser”.
- In **Step 2**, we train the T2A model using DPO so it prefers generating audio that matches the labeled “winner” characteristics.
- Standard fine-tuning (ERM) encourages “similar” audios and overfits on small-scale datasets. Our approach encourages diverse, dataset-aligned augmentations that “sound similar” and outperform ERM.

2. Generating Diverse Synthetic Augmentations



- We propose **language-guided audio imagination**, where we use LLMs to generate meaningful and controllable captions for synthetic audio generation.
- Our method, **MixCap**, extracts acoustic components (e.g., background sounds, events) using a captioning model and prompts the LLM to remix these with ground-truth labels for richer, more diverse captions.
- We filter generated audio-caption pairs using **CLAP-based similarity**: only audios closely matching the target label are accepted.
- For rejected samples, we use **LLM-based self-reflection** to revise captions and re-generate audio, iterating until quality thresholds are met or no rejected samples remain.
- The final accepted set of audios and captions is used to **expand the small-scale dataset and fine-tune an audio classification model** for downstream tasks, helping it learn both global structure and compositional acoustic detail.

Main Results

#	Method	ESC-50	USD8K	GTZAN	Medley	TUT	NSynth	VS	MSDB	DCASE	FSD50K
50	Gold-only (No Aug.)	22.25	55.09	47.05	47.23	37.60	33.32	77.49	56.85	12.09	7.16
	Random Noise	18.50	57.42	45.20	46.55	35.86	32.42	76.41	52.55	13.21	8.04
	Pitch Shifting	20.55	59.32	46.80	48.17	37.22	34.34	78.17	54.50	12.93	10.04
	SpecAugment	19.50	58.36	46.00	47.18	36.73	27.32	77.27	53.25	12.81	7.93
	Audiomentations	20.35	60.13	47.25	48.30	38.24	28.15	79.12	54.51	13.28	10.17
	Retrieval	19.20	37.14	42.55	43.65	35.80	31.27	71.42	51.35	10.53	7.28
	Vanilla Syn. Aug. + LLM Caps.	40.75	63.54	55.35	47.23	41.50	33.17	78.37	54.10	15.89	10.63
	Synthio (ours)	46.80	65.84	63.74	55.36	40.90	38.17	78.77	57.05	13.07	10.70
		49.50_{+12%}	76.12_{+38%}	68.20_{+44%}	60.58_{+28%}	43.84_{+17%}	40.83_{+22%}	80.67_{+4%}	60.15_{+5%}	17.23_{+42%}	13.91_{+94%}
100	Gold-only (No Aug.)	56.75	72.89	64.15	57.81	47.14	39.11	84.32	65.60	12.50	10.53
	Random Noise	58.50	71.54	65.50	56.98	46.21	38.20	83.33	66.15	13.35	13.71
	Pitch Shifting	59.55	73.52	66.75	58.46	47.50	39.53	85.07	68.25	12.19	13.11
	SpecAugment	47.50	72.43	69.75	58.06	50.07	41.96	85.14	66.40	14.17	14.80
	Audiomentations	48.50	73.82	71.05	59.32	51.14	42.15	85.24	68.40	16.93	13.55
	Retrieval	52.45	68.24	61.55	54.83	45.39	37.84	83.27	58.55	10.93	10.05
	Vanilla Syn. Aug. + LLM Caps.	77.25	77.31	68.25	49.96	42.31	48.78	63.55	15.73	12.63	12.63
	Synthio (ours)	86.05 _{+47%}	85.00 _{+17%}	71.20 _{+11%}	71.23 _{+23%}	52.42 _{+11%}	44.92 _{+15%}	86.70 _{+3%}	68.80 _{+5%}	19.38 _{+55%}	16.35 _{+53%}
200	Gold-only (No Aug.)	84.75	74.80	77.00	67.41	55.32	48.77	87.38	68.80	23.15	13.59
	Random Noise	83.55	75.15	75.50	66.71	54.42	47.83	86.45	65.45	24.82	15.32
	Pitch Shifting	84.90	74.48	78.55	67.74	55.44	48.12	87.47	69.80	23.11	17.51
	SpecAugment	85.10	76.46	76.25	65.70	55.72	54.80	87.42	69.25	27.36	17.93
	Audiomentations	85.25	75.80	77.30	67.00	55.21	53.15	86.08	70.50	26.29	18.36
	Retrieval	82.55	71.20	73.65	65.80	53.25	47.63	86.28	63.55	19.51	15.36
	Vanilla Syn. Aug. + LLM Caps.	85.40	77.96	77.10	78.97	55.51	55.20	86.49	72.95	28.55	19.04
	Synthio (ours)	86.10 _{+2%}	82.81 _{+11%}	82.05 _{+7%}	79.40 _{+18%}	56.83 _{+3%}	57.10 _{+17%}	87.52 _{+0.2%}	80.40 _{+17%}	32.81 _{+42%}	20.85 _{+53%}

Result Analysis

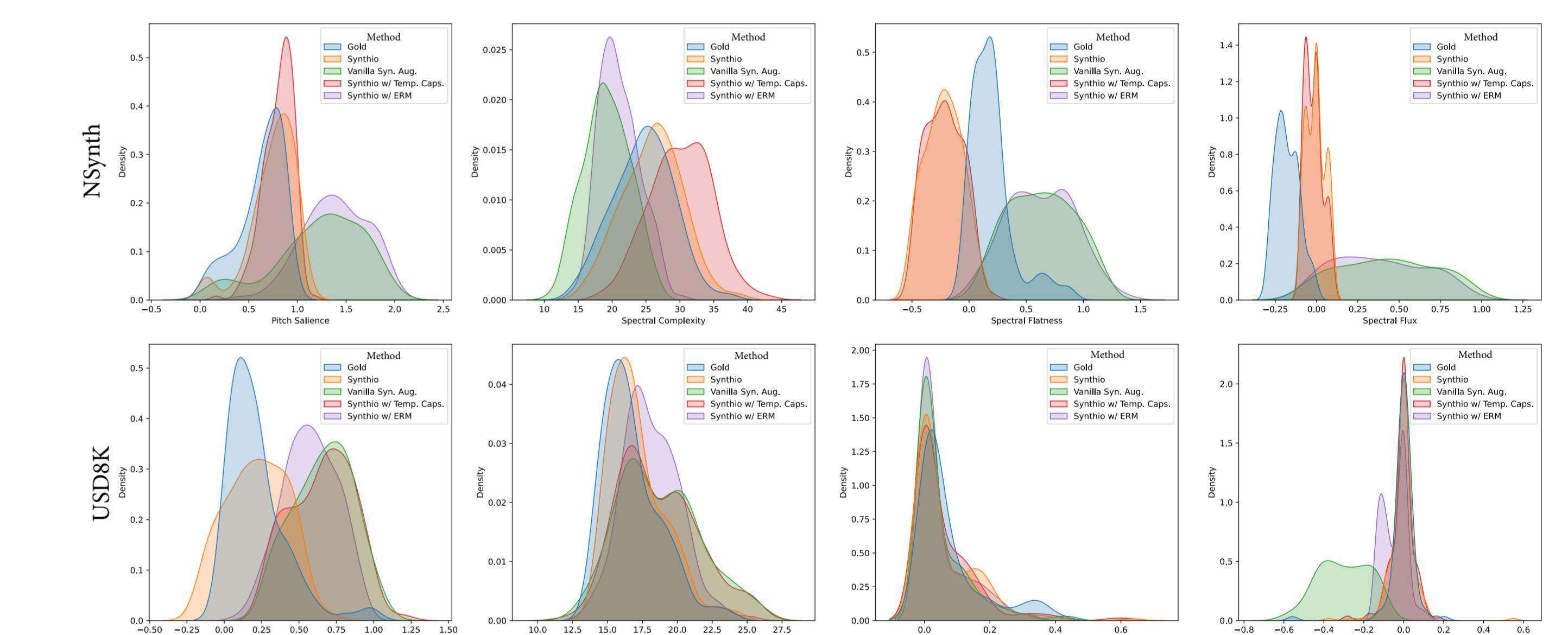


Figure 4. Comparison of spectral and pitch features between generated audios in D_{syn} and real audios in D_{small} ($n = 100$). Synthio-generated audios closely replicate the features of real data.

#	Method	USD8K(↓)	NSynth(↓)
100	Gold-only (No Aug.)	64.15	84.32
	Vanilla Syn. Aug.	29.05	34.13
	Synthio (ours)	35.10	29.20
	w/ Template Captions	24	