

# NATURAL LANGUAGE SUPERVISION FOR GENERAL-PURPOSE AUDIO REPRESENTATIONS



Benjamin Elizalde\*, Soham Deshmukh\*, Huaming Wang  
Emails: {benjaminm, sdeshmukh, huawang}@microsoft.com  
Microsoft

## POSTER IN 3 SENTENCES

- 1 Contrastive Language-Audio Pretraining (CLAP) connects audio and language with two encoders and a contrastive learning loss.
- 2 Pretrained with 4.6M audio-text pairs; the audio encoder (HTSAT) is trained on 22 audio tasks; the text encoder is a modified autoregressive decoder-only model (GPT2).
- 3 Our model achieves SoTA in about 26 tasks outperforming 4 different models.

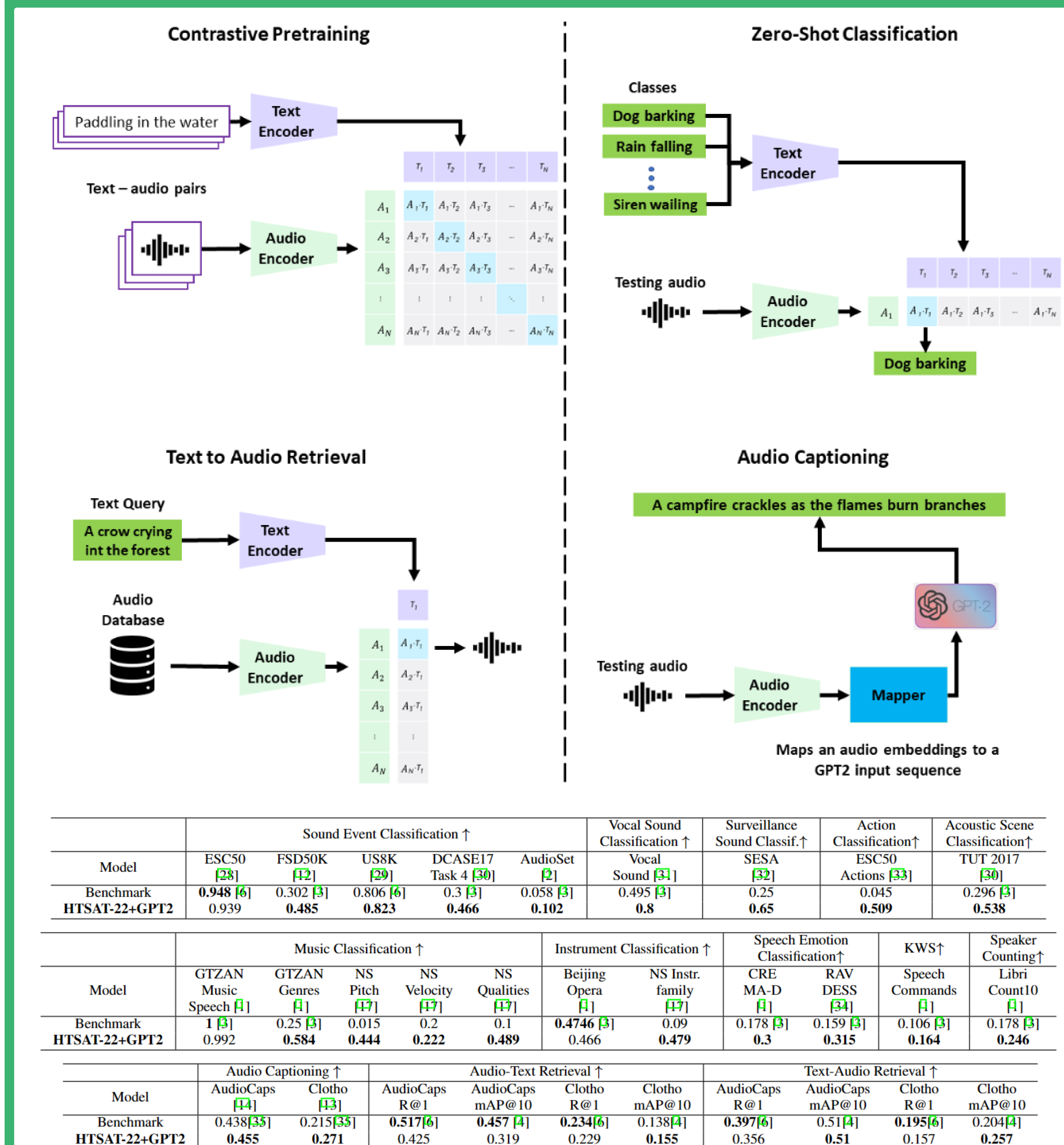
## CODE AVAILABLE ON:

[github.com/microsoft/CLAP](https://github.com/microsoft/CLAP)

## KEYWORDS

- **Contrastive Pretraining Data:** 4.6M pairs from WavCaps, AudioSet, FSD50K, Clotho, AudioCaps, MACS, WavText5k, SoundDesc, NSynth, FMA, Mosi, Meld, Iemocap, Mosei, MSP-Podcast, CochScene, LJspeech, EpicKitchen, Kinetics700, findsounds.com.
- **Text Encoder (GPT2)** is adapted to embed sentences by appending the special token  $\langle \text{endof text} \rangle$  at the end of each input text.
- **Audio Encoder (HTSAT-22)** is independently trained on 22 audio tasks.
- **Contrastive Loss:** We used cosine similarity to compare text and audio; multimodal embeddings are 1024 dims: batch size 1536.
- **Audio Captioning:** At training time, only the weights of the mapper network are learned with a captioning loss.
- **Benchmark** - SoTA in the literature.

## CLAP 🍌 PRETRAINING AND RESULTS ON 26 DOWNSTREAM TASKS



**Table 2:** Performance on 26 downstream tasks using our proposed encoders and 4.6M training pairs. As the benchmark, we used the best numbers in the literature, when no number was available we used random performance. Higher is better for all tasks. The evaluation metrics are mAP for FSD50k, ESC50-Actions, AudioSet, and NS Qualities; F1-score for DCASE17; and SPIDER for Captioning; all others use Accuracy.

## INSIGHTS

- 1 Our proposed encoders HTSAT-22+GPT2 outperformed the best combination of encoders in the literature (Figure 1).
- 2 Our HTSAT-22 audio encoder is the major contributor to performance improvement.
- 3 Adding diversity and scaling the audio-text pairs in training presents a trade-off that increases performance in some tasks but decreases it in others.
- 4 Zero-Shot models should be evaluated across different domains and tasks with focus on generalization rather than on over fitting to specific tasks.
- 5 We only fine-tuned the audio encoder for ESC50 and achieved SoTA with 98.25% acc.

	Zero-Shot Score ↑
Model	Average
CNN14+BERT	0.428
HTSAT+CLIP	0.430
HTSAT+RoBERTa	0.431
HTSAT+GPT2	0.435
HTSAT-22+RoBERTa	0.454
HTSAT-22+CLIP	0.469
HTSAT-22+GPT2	0.480

Figure 1: Comparison of different audio and text encoders in the literature. Zero-Shot score is the average of the metrics on 16 downstream tasks (higher is better). All models are trained on 119k training pairs. Our proposed encoders (HTSAT-22+GPT2) outperformed them all.