

# CLAP With Me

Step by Step Semantic Search on Audio Sources

# Content

## Introduction

Who/What/Why/How

## Contrastive What Now?

Brief Overview of CLAP

## Building The App

Recording/Searching/Classifying/Similar Sounds

## Closing

Summary/Next Steps



# Introduction

# About Me

[ajwallace.com](http://ajwallace.com)

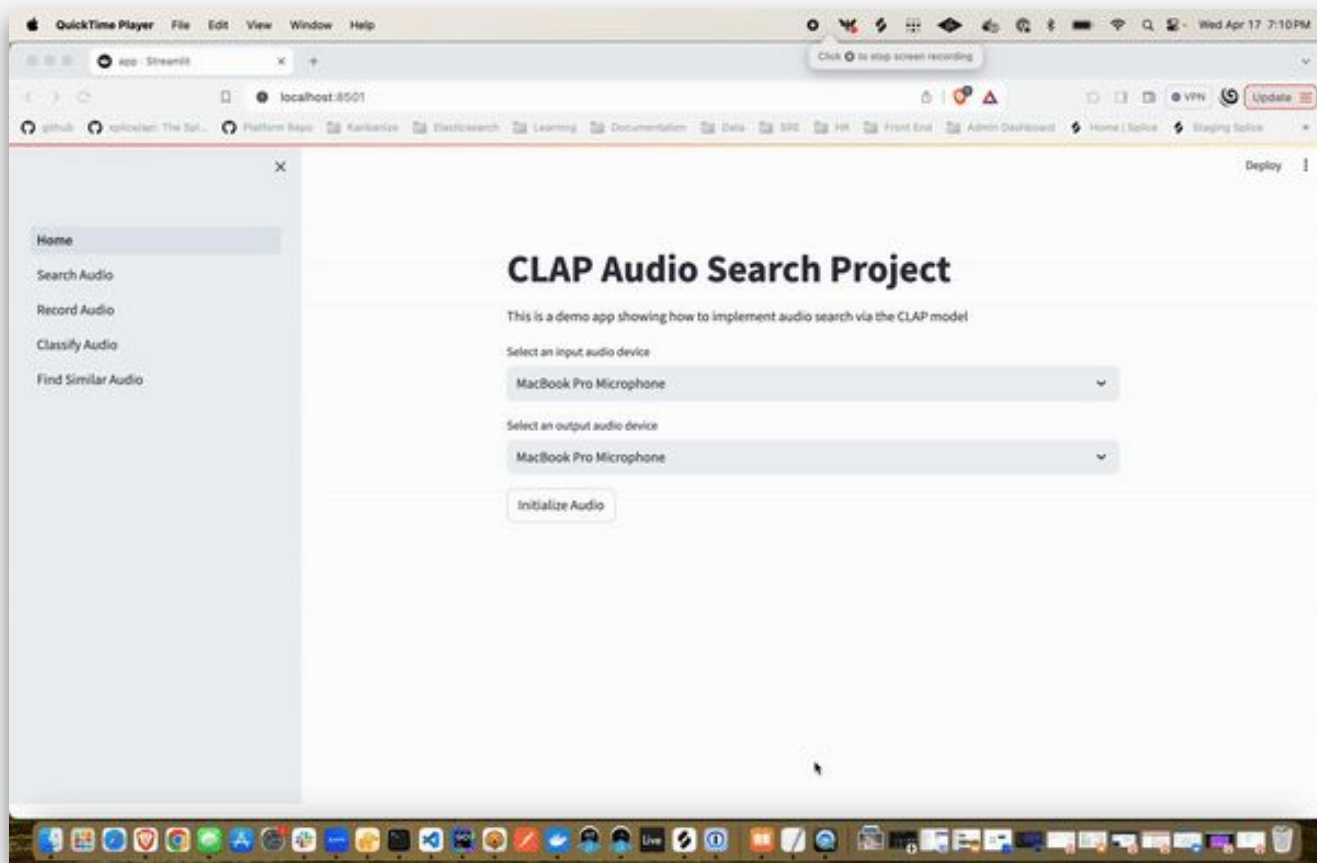
<https://github.com/ajwallacemusic>

Sr. Software Engineer, [splice.com](http://splice.com)

Search Consultant



# What We're Gonna Build



# Why?

Music and Audio are Awesome.  
So is Search.

- $\frac{4}{5}$  People listen to music using a streaming service.
- Nearly 60% of musicians use AI in production.
- 50% of US uses voice search.
- Audiobooks revenue has grown double digits each year this decade.
- Multimodal Search is the future.

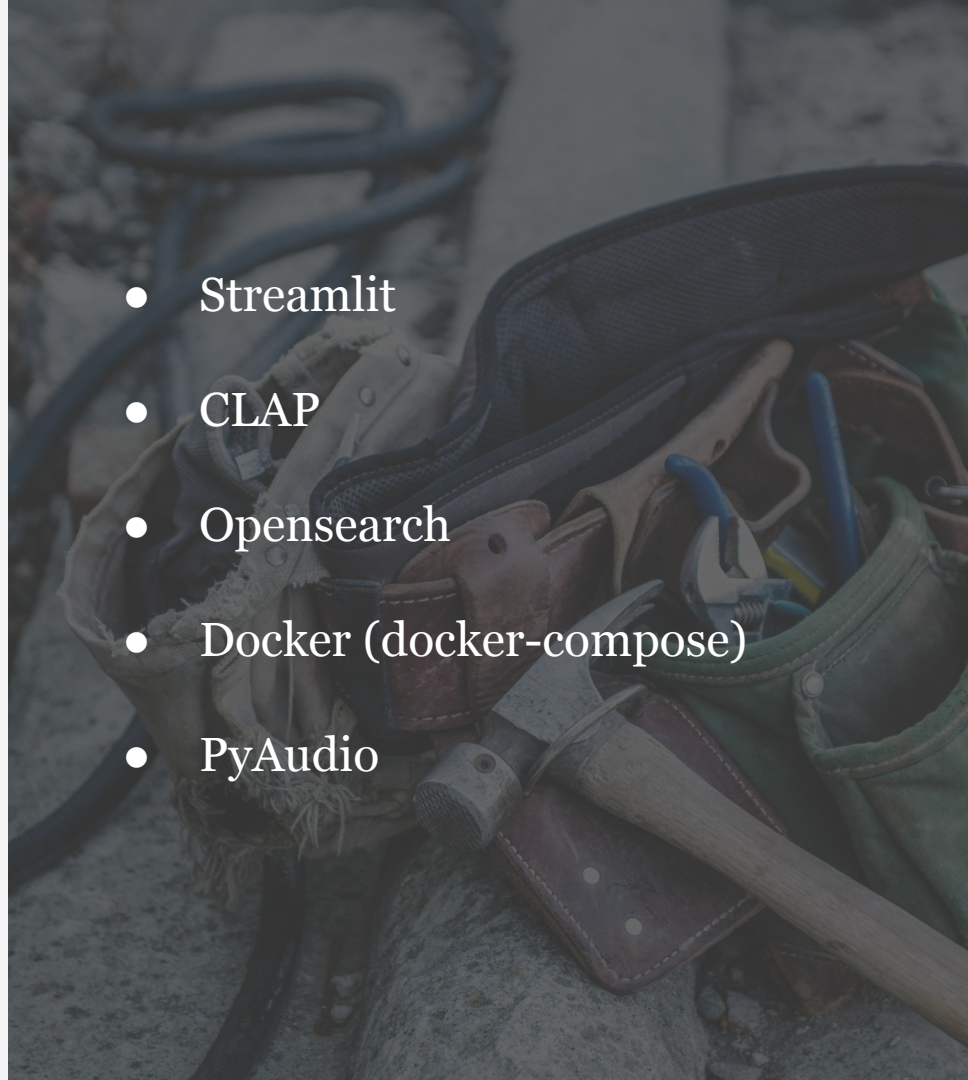


If I can do it...

# How Are We Gonna Build It?

With code, duh.

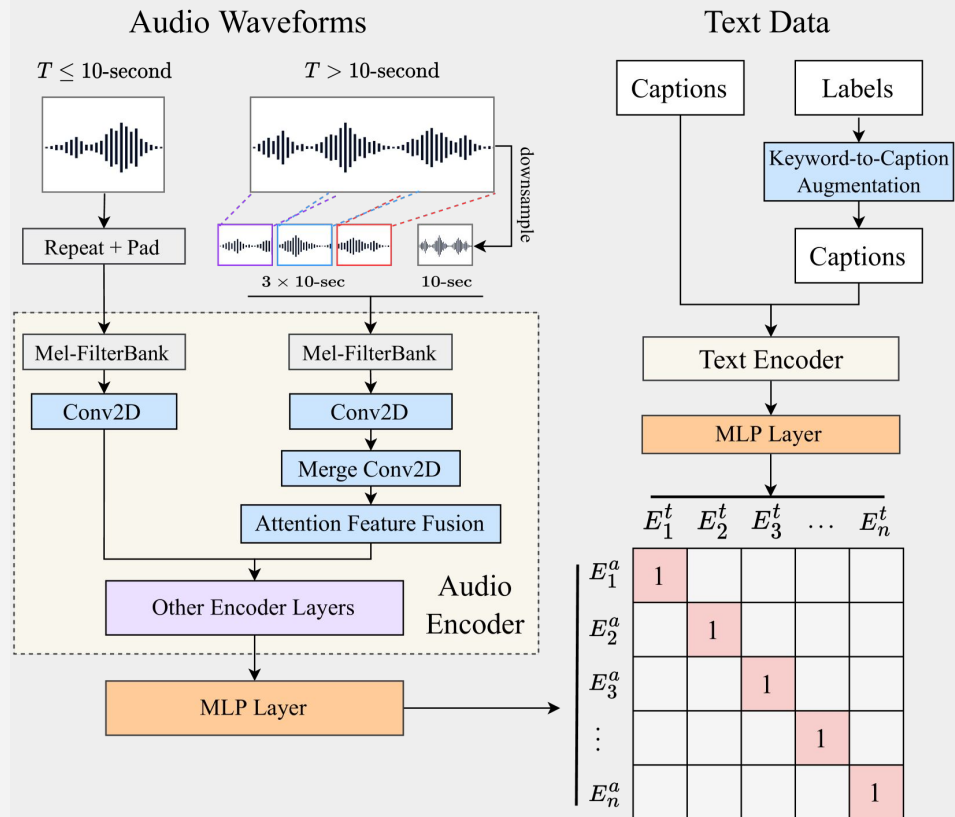
- Streamlit
- CLAP
- Opensearch
- Docker (docker-compose)
- PyAudio

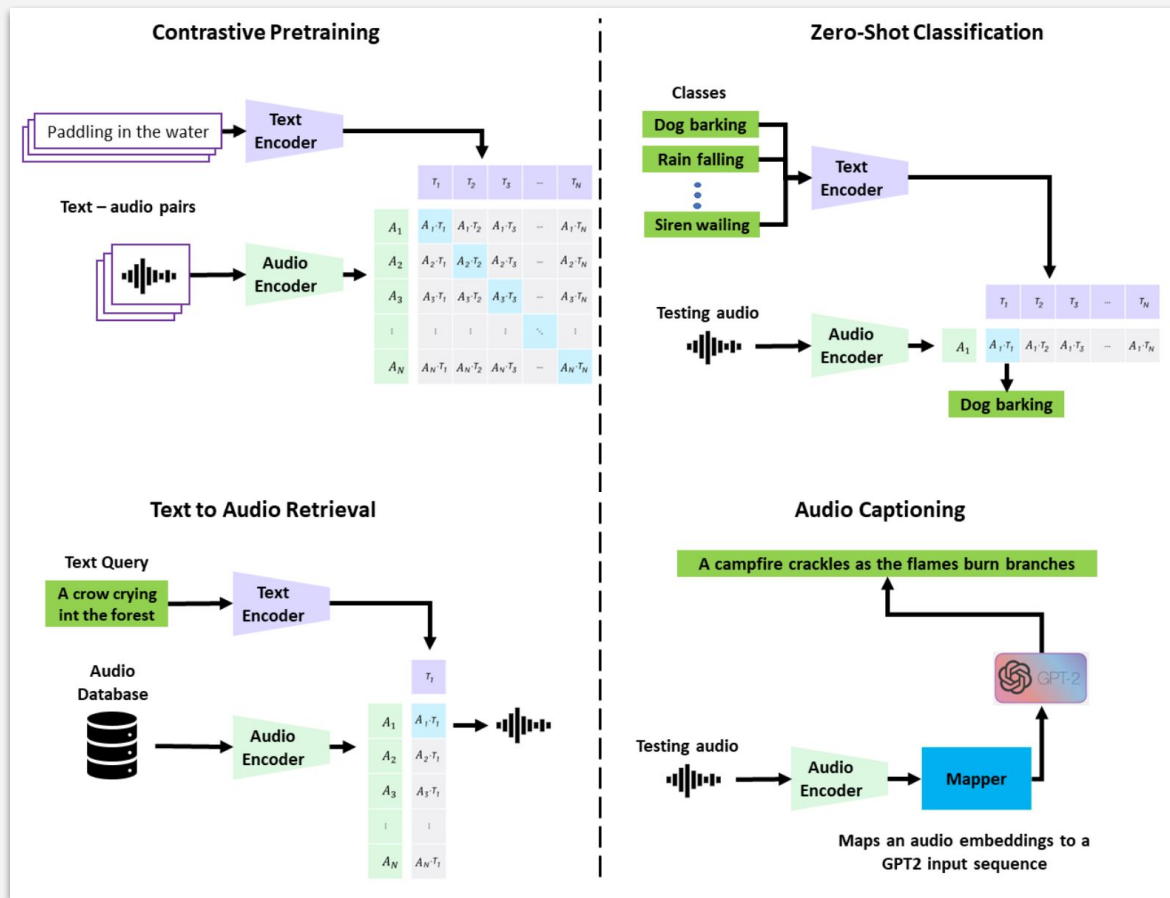


Contrastive What Now?

# CLAP

## Contrastive Language-Audio Pretraining





Audio and text pairs are used to create a joint multimodal space.

# Research Papers

[CLAP: Learning Audio Concepts From Natural Language Supervision](#)

[Large-scale Contrastive Language-Audio Pretraining with Feature Fusion and Keyword-to-Caption Augmentation](#)

<https://github.com/LAION-AI/CLAP>

[https://huggingface.co/docs/transformers/v4.27.2/en/model\\_doc/clap](https://huggingface.co/docs/transformers/v4.27.2/en/model_doc/clap)

## LARGE-SCALE CONTRASTIVE LANGUAGE-AUDIO PRETRAINING FEATURE FUSION AND KEYWORD-TO-CAPTION AUGMENTATION

Yusong Wu<sup>1\*</sup>, Ke Chen<sup>2\*</sup>, Tianyu Zhang<sup>1\*</sup>, Yuchen Hui<sup>1,3\*</sup>, Marianna Nezhurina<sup>1</sup>,  
Taylor Berg-Kirkpatrick<sup>2</sup>, Shlomo Dubnov<sup>2</sup>

<sup>1</sup>Mila, Quebec Artificial Intelligence Institute, Université de Montréal

<sup>2</sup>University of California San Diego <sup>3</sup>LAION

### ABSTRACT

Contrastive learning has shown remarkable success in the field of multimodal representation learning. In this paper, we propose a pipeline of contrastive language-audio pretraining to develop an audio representation by combining audio data with natural language descriptions. To accomplish this target, we first release LAION-Audio-630K, a large collection of 633,526 audio-text pairs from different data sources. Second, we construct a contrastive language-audio pretraining model by considering different audio encoders and text encoders. We incorporate the feature fusion mechanism and keyword-to-caption augmentation into the model design to further enable the model to process audio inputs of variable lengths and enhance the performance. Third, we perform comprehensive experiments to evaluate our model across three tasks: text-to-audio retrieval, zero-shot audio classification, and supervised audio classification. The results demonstrate that our model achieves superior performance in text-to-audio retrieval task. In audio classification tasks, the model achieves state-of-the-art performance in the zero-shot setting and is able to obtain performance comparable to models' results in the non-zero-shot setting. LAION-Audio-630K<sup>1</sup> and the proposed model<sup>2</sup> are both available to the public.

**Index Terms**— Contrastive Learning, Representation Learning, Text-to-Audio Retrieval, Audio Classification, Audio Dataset

### 1. INTRODUCTION

Audio is one of the most common information types in the world alongside text and image data. However, different audio tasks typically require finely-annotated data, which limits the amount of available audio data due to the labor-intensive collection procedure. Consequently, designing an effective audio representation for many audio tasks without requiring a lot of supervision remains a challenge.

The contrastive learning paradigm is a successful solution for training a model on large-scale noisy data collected from internet. The recently proposed Contrastive Language-Image Pretraining (CLIP) [1] learns the cross-modal relationship between images by

variations of ImageNet dataset [2]. Additionally, CLIP shows success in downstream tasks such as text-to-image retrieval, image-guided captioning. Similar to vision, audio data also contains overlapping information. In audio tasks, for instance, some text descriptions are related to the corresponding audio. These text descriptions contain meaning that could be learned together with audio to form an audio representation of cross-modal information. Additionally, training such a model requires simply audio data, which is easy to collect.

Several recent studies [3–9] have presented different contrastive language-audio pretraining models for audio retrieval task. [6] utilizes Pretrained Audio Neural Network [10] as the audio encoder, BERT [11] as the text encoder, and loss functions to evaluate the text-to-audio retrieval task. [5] further ensemble HTSAT [12] and RoBERTa [13] encoder list to further enhance performance. The effectiveness of the learned representation is evaluated in audio classification. Some other studies, such as CLAP [7] and WaveCLIP [9], focus more on the contrastive language-audio (language) pretraining model. The results show a great potential for contrastive learning in the audio domain.

Nonetheless, current studies have not fully explored the power of the language-audio contrastive learning. The models mentioned above are trained on relatively small-scale audio data. Large-scale data collection and augmentation are essential. Second, prior work lacks a full investigation of the model parameter settings of audio/text encoders, which is essential for determining the basic contrastive language-audio pretraining model. The model struggles to accommodate various audio inputs. Third, for the transformer-based audio encoder, the model's resolution to handle audio inputs of variable lengths is limited. Fourth, language-audio model studies focus on supervised learning without assessing their audio representation quality on audio tasks. As a representation model, we expect the model to have a generalization ability to more downstream tasks.

In this paper, we make contributions

arXiv:2211.06687v4 [cs.LG] 21 Mar 2024

# Using It

```
# init ML processors/models/tokenizers
audio_classifier = pipeline(task="zero-shot-audio-classification",
model="laion/larger_clap_music_and_speech")
model = ClapModel.from_pretrained("laion/larger_clap_music_and_speech")
processor = ClapProcessor.from_pretrained("laion/larger_clap_music_and_speech")
tokenizer = AutoTokenizer.from_pretrained("laion/larger_clap_music_and_speech")

# function to create audio embedding
def embed_audio(filepath):
    y, sr = librosa.load(filepath)
    inputs = processor(audios=y, return_tensors="pt", sampling_rate=48000)
    audio_embed = model.get_audio_features(**inputs)
    arr = audio_embed.detach().numpy()
    return arr
```

# Building The App

# Opensearch

```
"mappings": {  
  "properties": {  
    "audio_embedding": {  
      "type": "knn_vector",  
      "dimension": 512,  
      "method": {  
        "name": "hnsw",  
        "space_type": "l2",  
        "engine": "lucene",  
        ...  
      }  
    }  
  }  
}
```

# Streamlit Homepage

```
st.title("CLAP Audio Search Project")
desc = "This is a demo app showing how to implement audio search via the CLAP model"
st.write(desc)

# initialize the things
init_session_state()
menu()
set_up_audio()
init_clap()
create_search_client()
```

# Recording & Indexing

```
if start:  
    st.session_state.audio_data = []  
    record_audio()  
  
if stop:  
    save_audio(filename)
```

# Searching

```
if st.button("Search"):  
    if st.session_state.search_type == "Vector":  
        response = vector_search(query)  
        display_results(response)  
    elif st.session_state.search_type == "BM25":  
        response = bm25_search(query)  
        display_results(response)  
    elif st.session_state.search_type ==  
"Hybrid":  
        response = hybrid_search(query)  
        display_results(response)
```

# Classifying

```
def classify_audio(text1, text2, audiofile):  
    y, _ = librosa.load(audiofile)  
    audio_classifier = st.session_state.clap_audio_classifier  
    output = audio_classifier(y, candidate_labels=[text1, text2])  
  
    return get_label_with_highest_score(output)
```

# Audio to Audio Search

```
if uploaded_file is not None:
    source = uploaded_file
    st.audio(source)
else:
    source = audiofile

if st.button("Search"):
    response = audio_to_audio_search(source)
    display_results(response)
```

Closing

# Summary

Pretty good out of the gate

Hybrid helps

# Next Steps

Opensearch hosted ML model

<https://opensearch.org/docs/latest/search-plugins/multimodal-search/>

Measure Relevance

<https://quepid.com/>

Fine Tuning

Sujit Pal: <https://haystackconf.com/us2022/talk-1>

Thank You

# Sources

<https://github.com/ajwallacemusic/clap-demo>

[Large-scale Contrastive Language-Audio Pretraining with Feature Fusion and Keyword-to-Caption Augmentation](#)

[CLAP: Learning Audio Concepts From Natural Language Supervision](#)

Images: <https://unsplash.com/>

Statistics: <https://www.ifpi.org/>, <https://www.musicradar.com/>, <https://upcity.com/>,  
<https://www.statista.com/>