# CLIPSep: Learning Text-queried Sound Separation with Noisy Unlabeled Videos

Hao-Wen Dong[1,2] *    Naoya Takahashi[1] †    Yuki Mitsufuji[1]    Julian McAuley[2]    Taylor Berg-Kirkpatrick[2]

[1]Sony Group Corporation    [2]University of California San Diego

* Work done during an internship at Sony    † Corresponding author
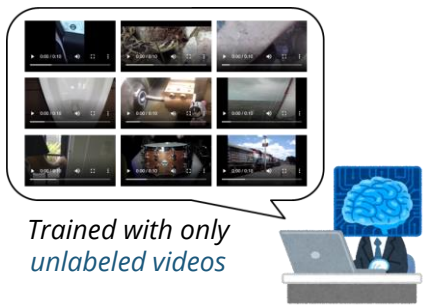
SONY

UC San Diego

ICLR

## Overview

We explore training a sound separation system under a **self-supervised learning** setting. We aim to achieve text-queried universal sound separation by *using only unlabeled data*.
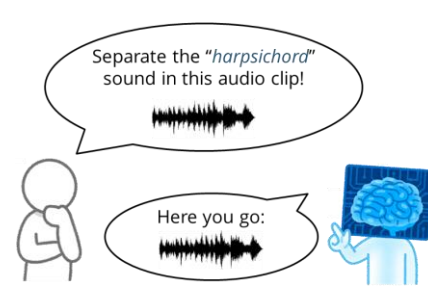
To learn the desired audio-textual correspondence from unlabeled videos, we leverage the visual modality as a bridge using the contrastive image-language pretraining (CLIP) model.

### Training



*Trained with only unlabeled videos*

*Scalable to larger dataset*

### Inference



Separate the "*harpsichord*" sound in this audio clip!

Here you go: 

*Natural text query-based interface*

## Contributions

- We propose the *first text-queried universal sound separation model that can be trained on unlabeled videos*.
- We propose a new approach called *noise invariant training* for *training a query-based sound separation model on noisy data*.

## Data

### MUSIC
(Zhao et al., 2018)



Violin    Acoustic guitar    Accordion

*Music instrument playing videos*

### VGGSound
(Chen et al., 2020)



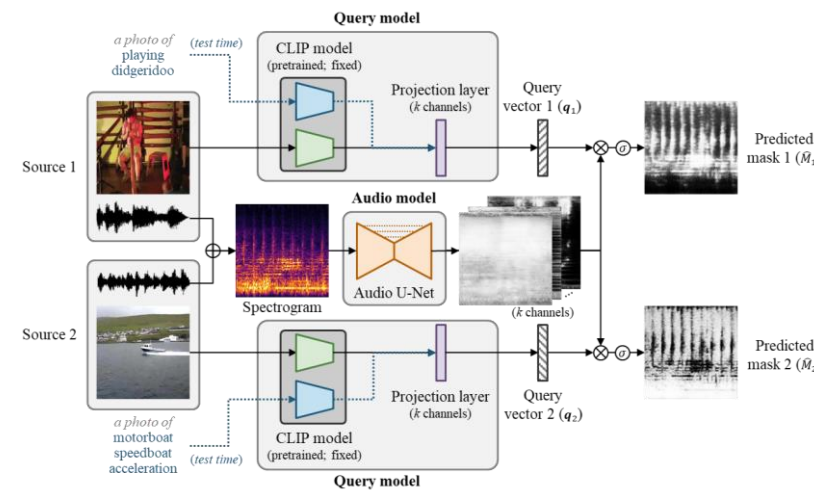Hedge trimmer running    Dog bow-wow    Bird chirping, tweeting

*Noisy videos with diverse sounds*

## CLIPSep

### Training

We mix audio from two videos and train the model to separate each audio source given the corresponding video frame as the query:
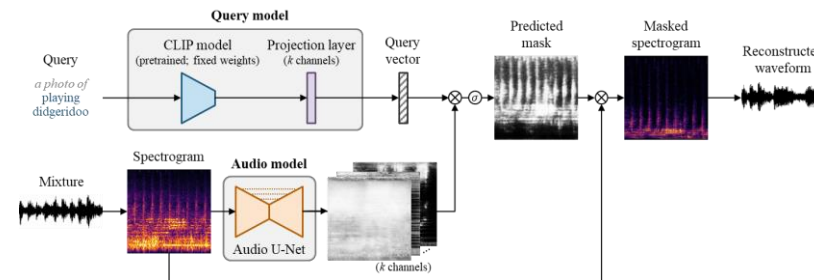
1. An Audio U-Net predicts $k$ intermediate masks $\tilde{M}_1, \ldots, \tilde{M}_k$ from the mixture spectrogram.
2. A pretrained CLIP model encodes the input query into a query vector $q_i$.
3. Construct the predicted masks with $\hat{M}_i = \sum_{j=1}^{k} \sigma(w_{ij}q_{ij}\tilde{M}_j + b_i)$.



### Inference

At test time, we instead use a text query in the form of "*a photo of [user input query]*".
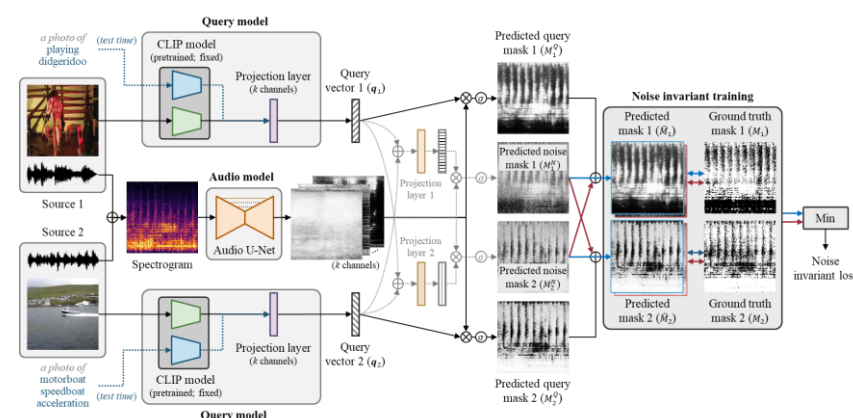
With the pretrained CLIP model, the query vectors we obtain for the image and text queries are expected to be close.



## Noise Invariant Training (NIT)

Videos in the wild may contain off-screen sounds and background noise. We introduce two additional noise masks to capture query-irrelevant sounds.
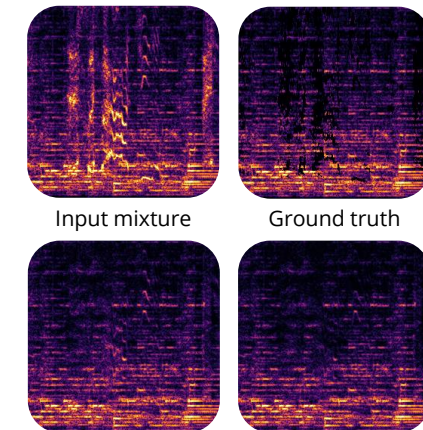
During training, we select the arrangement that has the lowest loss value when combining the noise and query masks. At test time, we discard the noise heads.
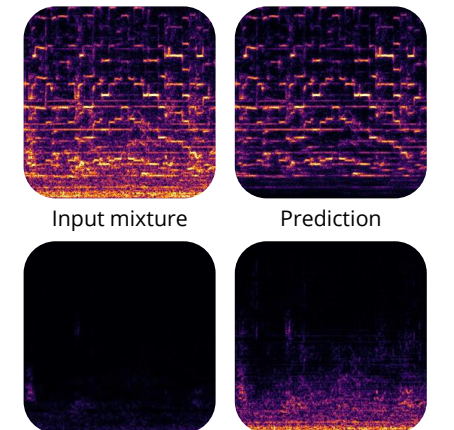


## Results

Audio samples available at sony.github.io/CLIPSep/

### Sound Separation



Input mixture    Ground truth

CLIPSep    CLIPSep-NIT

**Query:** "playing harpsichord"
**Interference:** "people coughing"

### Noise Removal
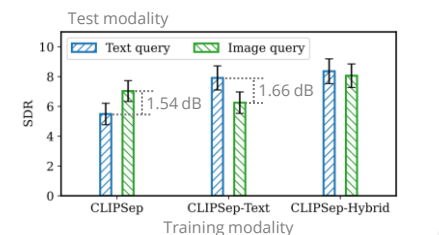


Input mixture    Prediction

Noise head 1    Noise head 2

**Query:** "playing bagpipes"
**Interference:** none

### Results on MUSIC (clean data)

| Model | Unlabeled data | Post-proc. free | Query type Training | Query type Test | SDR [dB] Mean | SDR [dB] Median |
|---|---|---|---|---|---|---|
| Mixture | - | - | | | $0.00 \pm 0.89$ | 0.00 |
| **Text-queried models** | | | | | | |
| CLIPSep | ✓ | ✓ | Image | Text | $5.49 \pm 0.72$ | **4.97** |
| CLIPSep-Text | | ✓ | Text | Text | $7.91 \pm 0.81$ | 7.46 |
| CLIPSep-Hybrid | | ✓ | Text+Image | Text | $8.36 \pm 0.83$ | 8.72 |
| **Image-queried models** | | | | | | |
| SOP (Zhao et al., 2018) | ✓ | | Image | Image | $6.59 \pm 0.85$ | 6.22 |
| CLIPSep | ✓ | ✓ | Image | Image | $7.03 \pm 0.70$ | 5.85 |
| CLIPSep-Text | | ✓ | Text | Image | $6.25 \pm 0.72$ | 6.19 |
| CLIPSep-Hybrid | | ✓ | Text+Image | Image | $8.06 \pm 0.79$ | 8.01 |
| **Nonqueried models** | | | | | | |
| LabelSep | | ✓ | Label | Label | $8.18 \pm 0.80$ | 7.82 |
| PIT (Yu et al., 2017) | ✓ | | × | × | $8.68 \pm 0.76$ | 7.67 |



### Results on VGGSound (noisy data)

| Model | Unlabeled data | Post-proc. free | MUSIC+ Mean SDR | MUSIC+ Median SDR | VGGSound-Clean+ Mean SDR | VGGSound-Clean+ Median SDR |
|---|---|---|---|---|---|---|
| Mixture | - | - | $4.49 \pm 1.41$ | 2.04 | $-0.77 \pm 1.31$ | -0.84 |
| **Text-queried models** | | | | | | |
| CLIPSep | ✓ | ✓ | $9.71 \pm 1.21$ | 8.73 | $2.76 \pm 1.00$ | **3.95** |
| CLIPSep-NIT | ✓ | ✓ | $10.27 \pm 1.04$ | **10.02** | $3.05 \pm 0.73$ | 3.26 |
| BERTSep | | ✓ | $4.67 \pm 0.44$ | 4.41 | $5.09 \pm 0.80$ | 5.49 |
| CLIPSep-Text | | ✓ | $10.73 \pm 0.99$ | 9.93 | $5.49 \pm 0.82$ | 5.06 |
| **Image-queried models** | | | | | | |
| SOP (Zhao et al., 2018) | ✓ | | $11.44 \pm 1.18$ | 11.18 | $2.99 \pm 0.84$ | 3.89 |
| CLIPSep | ✓ | ✓ | $12.20 \pm 1.17$ | 12.42 | $5.46 \pm 0.79$ | **5.35** |
| CLIPSep-NIT | ✓ | ✓ | $11.28 \pm 1.08$ | 10.83 | $4.84 \pm 0.66$ | 3.57 |
| CLIPSep-Text | | ✓ | $9.89 \pm 1.04$ | 8.09 | $2.45 \pm 0.70$ | 1.74 |
| **Nonqueried models** | | | | | | |
| PIT (Yu et al., 2017) | ✓ | | $12.24 \pm 1.20$ | 12.53 | $5.73 \pm 0.79$ | 4.97 |
| LabelSep | | ✓ | - | - | $5.55 \pm 0.81$ | 5.29 |

CLIPSep successfully learned text-queried sound separation on noisy data.

Noise invariant training improves the mean SDRs.

Paper: arxiv.org/abs/2212.07065
Demo: sony.github.io/CLIPSep/
Code: github.com/sony/CLIPSep

Paper    Demo    Code