# Generative AI for Music and Audio

**Hao-Wen (Herman) Dong**

董皓文

UC San Diego

# About me



Hi, I'm Herman.
I do AI x Music research.
I love music and movies!

國立臺灣大學 National Taiwan University — 2013 – 2017
*B.S. in Electrical Engineering*

中央研究院 ACADEMIA SINICA — 2017 – 2019
*Research Assistant*

UC San Diego — 2019 – 2021
*M.S. in Computer Science*

UC San Diego — 2019 – present
*Ph.D. in Computer Science (expected)*

Summer 2019 — YAMAHA
*Research Intern*

Summer 2021 — Dolby
*Deep Learning Audio Intern*

Summer 2022 — SONY
*Student Intern*

Fall 2022 — amazon
*Applied Scientist Intern*

Winter 2023 — Dolby
*Speech/Audio Deep Learning Intern*

Summer 2023 — Adobe
*Research Scientist/Engineer Intern*

Fall 2023 — NVIDIA
*Research Intern*

# About me

### EE



a female cat engineer making
an electric chip in a classroom

### Music



a cat playing heavy metal

### CS



a cat engineer debugging on laptop

# Introduction

Mumbai, the city of dreams.

# Multimodal generative AI for Films



Now it's the cybernetic heart of India,

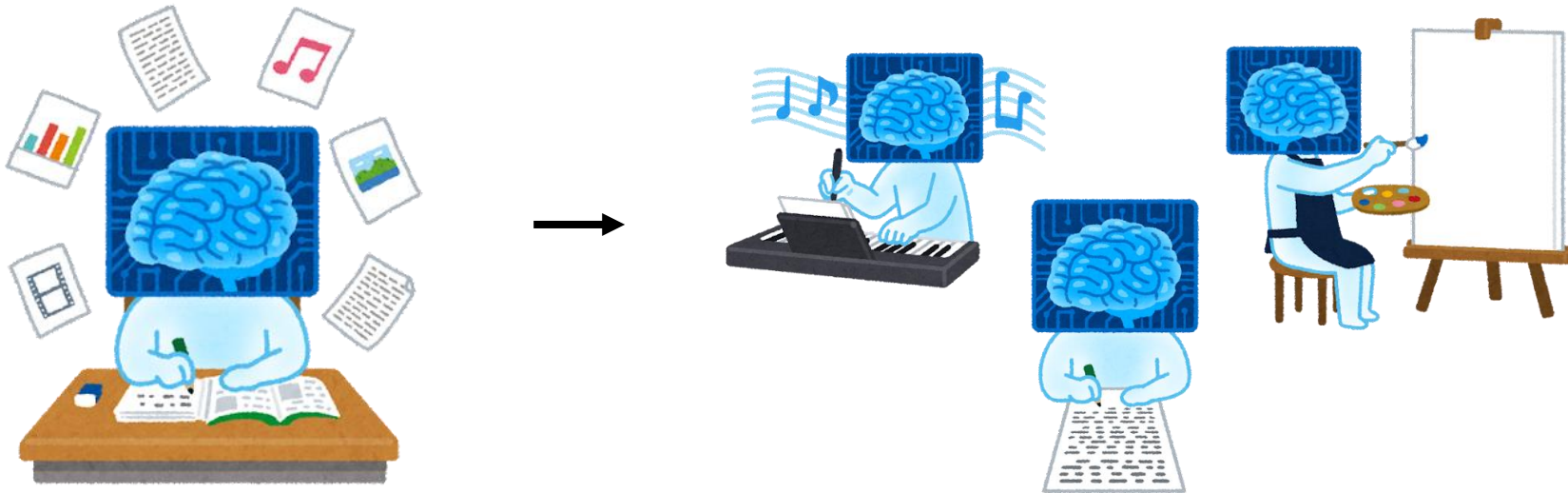| | |
|---:|:---|
| Visuals | **Midjourney** |
| Video | **Runway** |
| Narration (script) | **ChatGPT** |
| Narration (voice) | **ElevenLabs** |
| Sound effects | **Audiocraft** |

# What is Generative AI?

- Generative AI is **AI capable of generating text, images, or other media.**

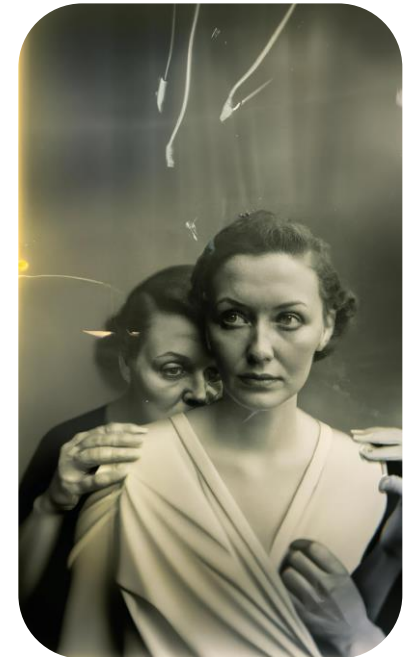# Generative AI for Visual Arts

**AI made a magazine cover**



(Source: Cosmopolitan)

**AI won an art contest**



(Source: CNN Business)

**AI won a photography contest**



(Source: CNN)

Gloria Liu, "The World's Smartest Artificial Intelligence Just Made Its First Magazine Cover," *Cosmopolitan*, June 21, 2022.
Rachel Metz, "AI won an art contest, and artists are furious," *CNN Business*, September 3, 2022.
Lianne Kolirin, "Artist rejects photo prize after AI-generated image wins award," *CNN*, April 18, 2023.
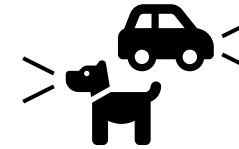
8

# Types of Audio

Speech

Music

Sound effects

(Source: Wikimedia Commons)

(Source: Wikimedia Commons)

(Source: Wikimedia Commons)

# Generative AI for Music

**Prompt**: relaxing and smooth jazz played in a stylish cafe

**Prompt**: delightful country music with acoustic guitars

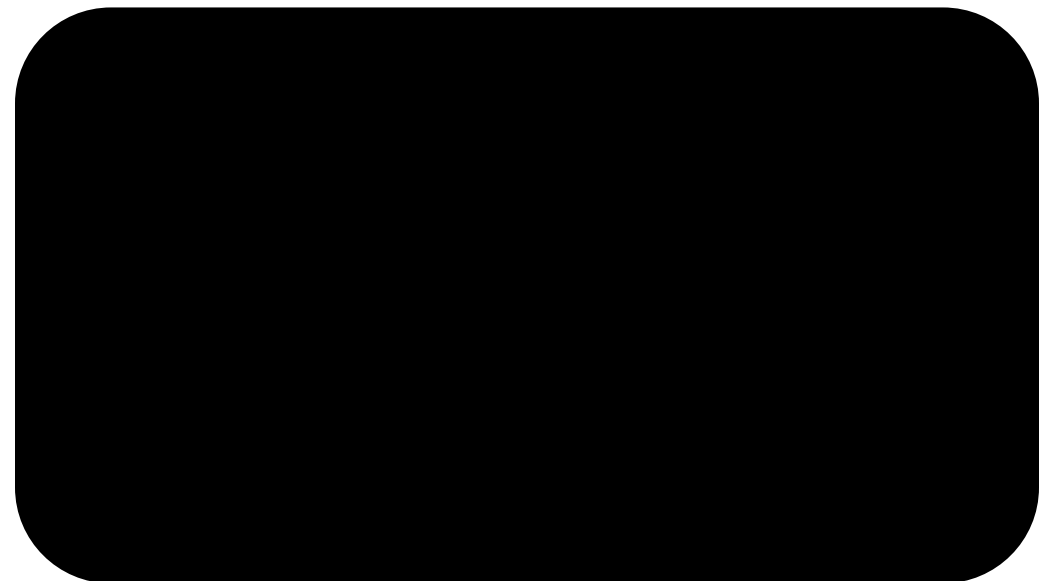**Prompt**: cinematic and suspenseful orchestral music

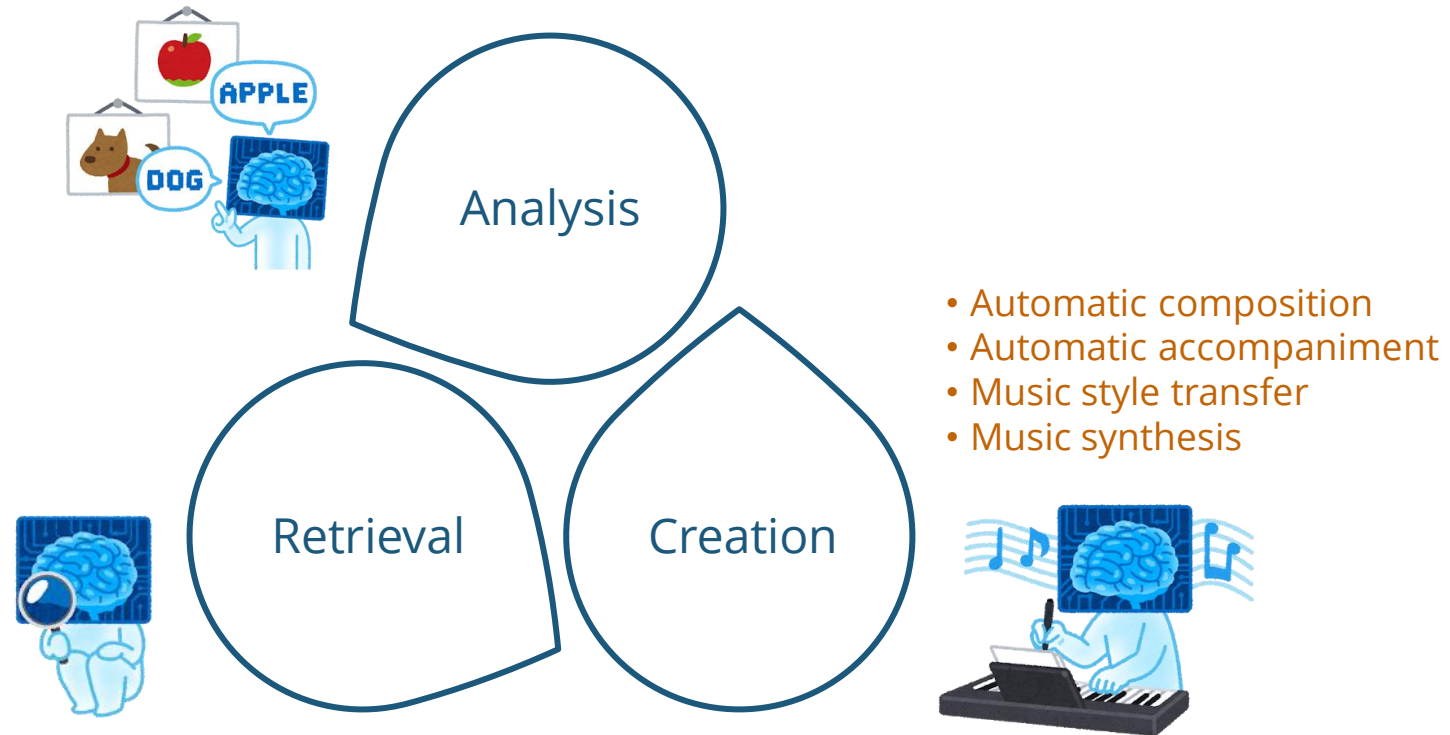# Generative AI for Sound Effects

**Text-to-audio Synthesis**



whistling with wind blowing

**Image-to-audio Synthesis**

Liu et al., "AudioLDM: Text-to-Audio Generation with Latent Diffusion Models," *ICML*, 2023.
Dong et al., "CLIPSonic: Text-to-Audio Synthesis with Unlabeled Videos and Pretrained Language-Vision Models," *WASPAA*, 2023.

# Music Information Research (MIR)

- *"Intelligent ways to analyze, retrieve and create music"* (Yang 2018)



- Automatic composition
- Automatic accompaniment
- Music style transfer
- Music synthesis

Analysis

Retrieval

Creation

Yang, "Music Information Research," *SNHCC*, TIGP, lecture notes, April 2018.

12

# My Research



AI × Music

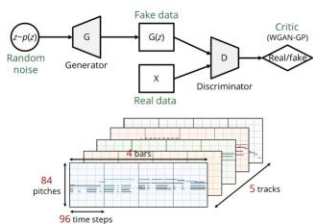**Multitrack Music Generation**

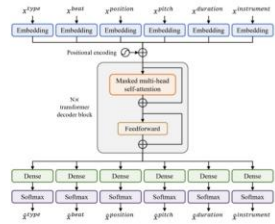Generating new music contents automatically

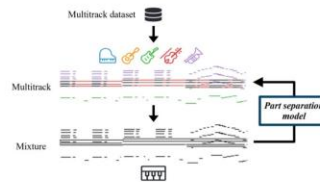MuseGAN (AAAI 2018)

Multitrack Music Transformer (ICASSP 2023)

**Assistive Music Creation Tools**

Assisting humans to create and perform music

Arranger (ISMIR 2021)

Deep Performer (ICASSP 2022)

**Multimodal Learning for Audio & Music**

Learning sound separation and synthesis from videos

CLIPSep (ICLR 2023)

CLIPSonic (WASPAA 2023)

# My Research

AI × Music 🎵

**Multitrack Music Generation**

Generating new music contents automatically

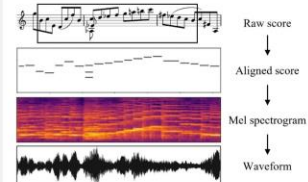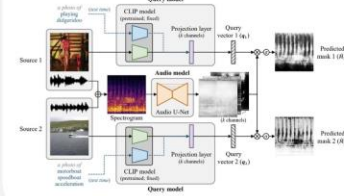| MuseGAN (AAAI 2018) | Multitrack Music Transformer (ICASSP 2023) |

**Assistive Music Creation Tools**

Assisting humans to create and perform music

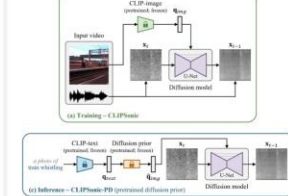| Arranger (ISMIR 2021) | Deep Performer (ICASSP 2022) |

**Multimodal Learning for Audio & Music**

Learning sound separation and synthesis from videos

| CLIPSep (ICLR 2023) | CLIPSonic (WASPAA 2023) |

**Featured in
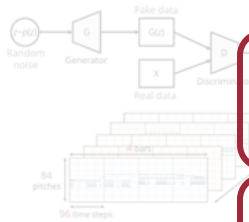Amazon AWS DeepComposer**

# My Research



AI × Music

**Multitrack Music Generation**
Generating new music contents automatically

MuseGAN (AAAI 2018)

Multitrack Music Transformer (ICASSP 2023)

**Assistive Music Creation Tools**
Assisting humans to create and perform music

Arranger (ISMIR 2021)

Deep Performer (ICASSP 2022)

Raw score
Aligned score
Mel spectrogram
Waveform

**Multimodal Learning for Audio & Music**
Learning sound separation and synthesis from videos

CLIPSep (ICLR 2023)

CLIPSonic (WASPAA 2023)

# My Research

AI × Music ♫

**Multitrack Music Generation**

Generating new music contents automatically

MuseGAN
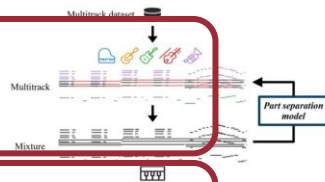(AAAI 2018)

Multitrack Music Transformer
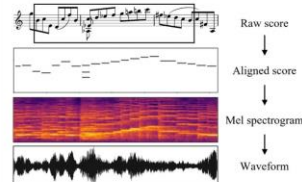(ICASSP 2023)

**Assistive Music Creation Tools**

Assisting humans to create and perform music

Arranger
(ISMIR 2021)

Deep Performer
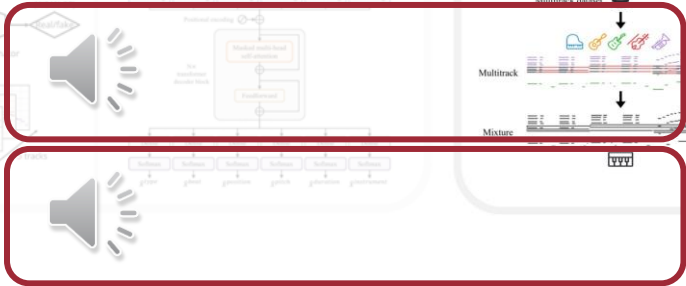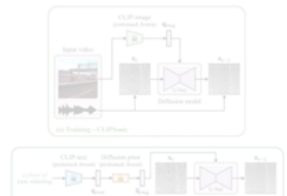(ICASSP 2022)

**Multimodal Learning for Audio & Music**

Learning sound separation and synthesis from videos
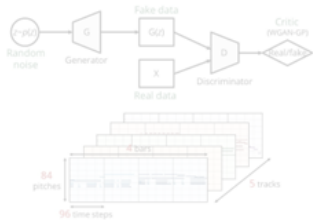
CLIPSep
(ICLR 2023)

CLIPSonic
(WASPAA 2023)

# 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 Corporation    [2]University of California San Diego

* Work done during an internship at Sony    † Corresponding author

# CLIP (Contrastive Language-Image Pretraining)
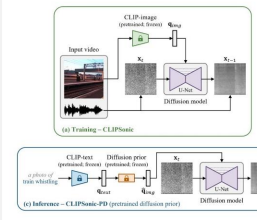
- Learn a shared embedding space for images and texts via *contrastive learning*



Radford et al., "Learning Transferable Visual Models From Natural Language Supervision," *Proc. ICML*, 2021.

# Leveraging the Visual Domain as a Bridge

Pretrained **vision-language models**

Audio-visual correspondence in **videos**

**Text**

**Image**

**Audio**

Desired text-audio correspondence

No text-audio pairs required!

Scalable to large video datasets!

# Why NOT Text-audio Pairs?

**YouTube videos!**

500 hours of videos
uploaded per minute

**5 billion**
text-image pairs

**LAION-5B**
(Schuhmann et al., 2023)

**0.6 million**
text-audio pairs

**LAION-Audio-630K**
(Wu et al., 2023)

Schuhmann et al., "LAION-5B: An open large-scale dataset for training next generation image-text models," *NeurIPS, Datasets and Benchmarks Track*, 2023.
Wu et al., "Large-scale Contrastive Language-Audio Pretraining with Feature Fusion and Keyword-to-Caption Augmentation," *ICASSP*, 2023.

# CLIPSep

# 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**

Zhao et al., "The Sound of Pixels," *ECCV*, 2018. (dataset)
Chen et al., "VGGSound: A Large-Scale Audio-Visual Dataset," *ICASSP*, 2020. (dataset)

# Demo – CLIPSep

Query: "*playing harpsichord*"

# Noise Invariant Training (NIT)

# Demo – CLIPSep-NIT

Query: "*playing harpsichord*"



| Mixture | CLIPSep | CLIPSep-NIT | Ground truth |

# Quantitative Results

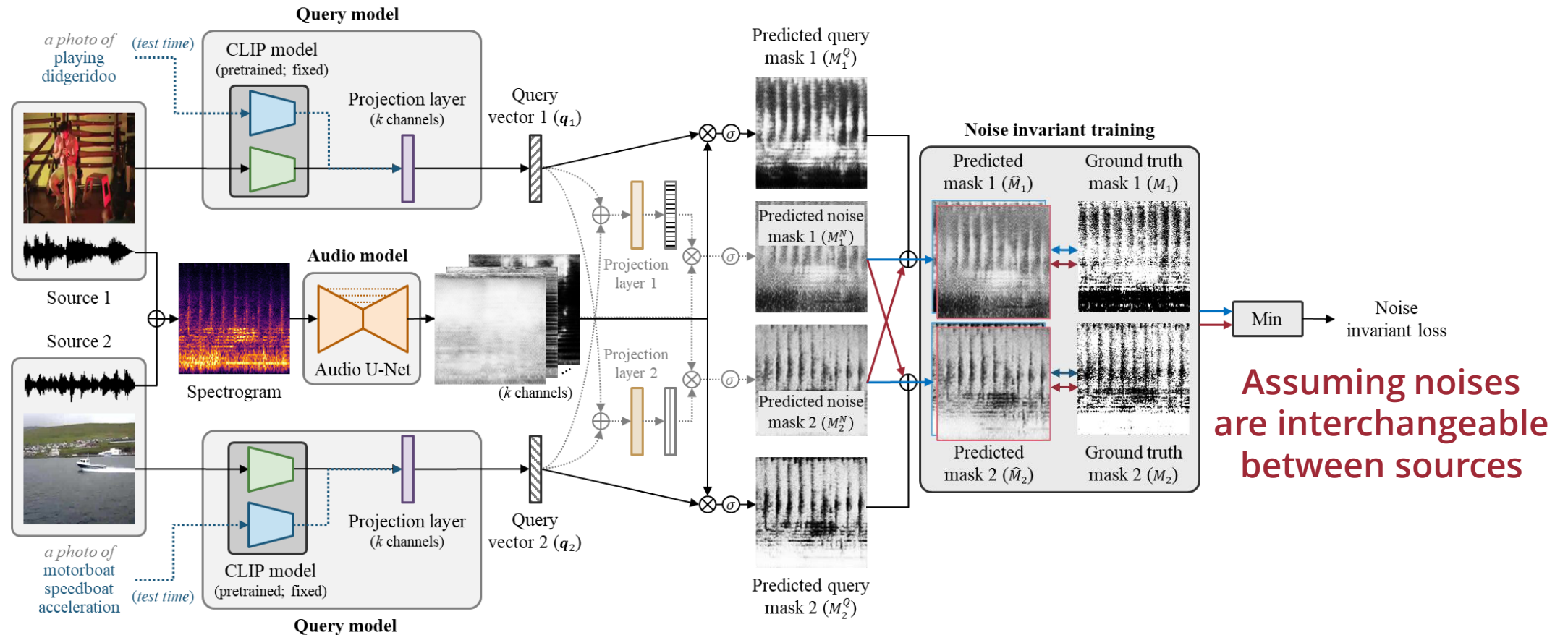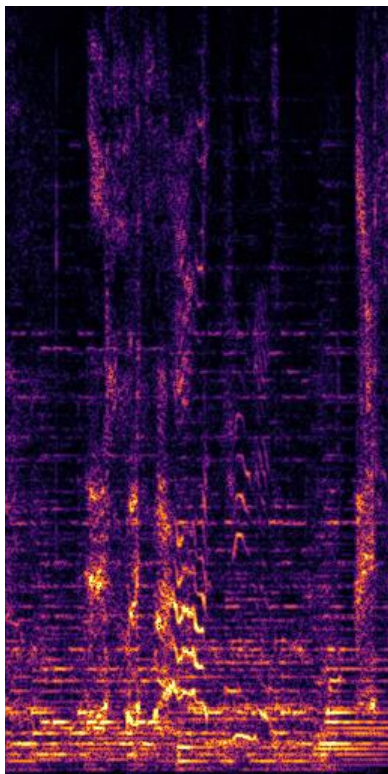| Model | Unlabeled data | Post-proc. free | MUSIC$^+$ | | VGGSound-Clean$^+$ | |
|---|---|---|---|---|---|---|
| | | | Mean SDR | Median SDR | Mean SDR | Median SDR |
| Mixture | - | - | $4.49 \pm 1.41$ | 2.04 | $-0.77 \pm 1.31$ | -0.84 |
| **Text-queried models** | | | | | | |
| CLIPSep | ✓ | ✓ | $9.71 \pm 1.21$ | 8.73 | $2.76 \pm 1.00$ | **3.95** |
| CLIPSep-NIT | ✓ | ✓ | $\mathbf{10.27 \pm 1.04}$ | **10.02** | $\mathbf{3.05 \pm 0.73}$ | 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 |

**Significant performance improvement** against the baseline!

# Demo – Noise Removal

Query: "*playing bagpipe*"



**Mixture**  **Prediction**  **Noise head 1**  **Noise head 2**

# Summary

## CLIPSep

First text-queried universal sound separation model that can be trained **using only unlabeled videos**

## Noise Invariant Training

A new approach for training a query-based sound separation model with **noisy data in the wild**



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

# My Research



**Multimodal Learning for Audio & Music**

Learning sound separation and synthesis from videos

Multitrack Music Ge...

Generating new music contents automatically

MuseGAN
(AAAI 2018)

Multit

CLIPSep
(ICLR 2023)

CLIPSonic
(WASPAA 2023)

...earning for Audio & Music

...separation
...rom videos

CLIPSonic
(WASPAA 2023)

# CLIPSonic: Text-to-Audio Synthesis with Unlabeled Videos and Pretrained Language-Vision Models

**Hao-Wen Dong**[1,2]*    Xiaoyu Liu[1]    Jordi Pons[1]    Gautam Bhattacharya[1]

Santiago Pascual[1]    Joan Serrà[1]    Taylor Berg-Kirkpatrick[2]    Julian McAuley[2]

[1] Dolby Laboratories    [2] University of California San Diego

* Work done during an internship at Dolby

**Dolby**          **UC San Diego**

# Diffusion model

**Add noise** gradually
(Forward diffusion process)



$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

$\mathbf{x}_T \longrightarrow \cdots \longrightarrow \mathbf{x}_t \longrightarrow \mathbf{x}_{t-1} \longrightarrow \cdots \longrightarrow \mathbf{x}_0$

**Remove noise** gradually
(Backward diffusion process)

Input



Ho et al., "Denoising Diffusion Probabilistic Models," *Proc. NeurIPS*, 2020.

# CLIPSonic – Training (Image-queried)

- We train the model to perform image-to-audio synthesis
  - Encode a video frame using a pretrained CLIP-image encoder (Radford et al., 2021)

Radford et al., "Learning Transferable Visual Models From Natural Language Supervision," *Proc. ICML*, 2021.

# CLIPSonic – Inference (Text-queried)

- We use a pretrained diffusion prior model (Ramesh et al., 2022)
  - To generate a CLIP-image embedding given a CLIP-text embedding

Ramesh et al., "Hierarchical Text-Conditional Image Generation with CLIP Latents," *arXiv preprint arXiv:2204.06125*, 2022.

# CLIPSonic – Inference Examples

# 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**

Zhao et al., "The Sound of Pixels," *ECCV*, 2018. (dataset)
Chen et al., "VGGSound: A Large-Scale Audio-Visual Dataset," *ICASSP*, 2020. (dataset)

# Text-to-Audio Synthesis – Demo

Rapping

Sea waves

Thunder

Smoke detector beeping

Playing table tennis

Playing violin fiddle

# Text-to-Audio Synthesis – Listening Test

Table 3: Listening test results for text-to-audio synthesis (MOS).

| Model | VGGSound | | MUSIC | |
|---|---|---|---|---|
| | Fidelity | Relevance | Fidelity | Relevance |
| CLIPSonic-ZS | $2.55 \pm 0.22$ | $2.01 \pm 0.27$ | $2.98 \pm 0.23$ | $3.87 \pm 0.24$ |
| CLIPSonic-PD | $\mathbf{3.04 \pm 0.20}$ | $2.86 \pm 0.25$ | $\mathbf{3.67 \pm 0.18}$ | $3.91 \pm 0.24$ |
| Ground truth | $3.78 \pm 0.19$ | $3.54 \pm 0.29$ | $3.90 \pm 0.17$ | $4.34 \pm 0.18$ |

**Significant performance improvement** against the baseline!

# Image-to-Audio Synthesis – Demo (Out-of-distribution)

# Image-to-Audio Synthesis – Listening Test

Table 4: Listening test results for image-to-audio synthesis (MOS).

| Model | Fidelity | Relevance |
|---|---|---|
| CLIPSonic-IQ (image-queried) | $\mathbf{3.29 \pm 0.16}$ | $3.80 \pm 0.19$ |
| SpecVQGAN [20] | $2.15 \pm 0.17$ | $2.54 \pm 0.23$ |
| im2wav [21] | $2.19 \pm 0.15$ | $\mathbf{3.90 \pm 0.22}$ |

**State-of-the-art** image-to-audio performance!

Sheffer and Adi, "I Hear Your True Colors: Image Guided Audio Generation," *ICASSP*, 2023.
Iashin and Rahtu, "Taming Visually Guided Sound Generation," *BMVC*, 2021.

# Summary

- Proposed a text-to-audio synthesis model that requires *no* text-audio pairs

- Achieves strong performance in objective and subjective evaluations

- Achieves state-of-the-art performance in image-to-audio synthesis



(c) Inference – CLIPSonic-PD (pretrained diffusion prior)

Paper: arxiv.org/abs/2306.09635
Demo: salu133445.github.io/clipsonic

# Conclusion

# Leveraging the Visual Domain as a Bridge



Pretrained **vision-language models**

Audio-visual correspondence in **videos**

Text

Image

Audio

Desired text-audio correspondence

No text-audio pairs required!

Scalable to large video datasets!

# A Lot More to Learn from Videos

- Free audio-visual correspondence

- Rich context information

- Rich temporal dynamics



THAT'S TOO MUCH INFORMATION!

# Future Directions

# Challenges

**Multimodality**

Text

Image

Video

Emotion

**Usability**

**Licensing**

COPYRIGHT

WEBSITE

# Multimodal Generative AI



Text

Text-to-image generation
Text-guided image editing

Text-to-audio generation
Text-guided audio editing

?

Image

Audio

Image-to-audio generation
Audio-to-image generation

# Multimodal generative AI for Ads



Video **Runway Gen-2**

Music **MusicGen**

進捗共有チャンネル, "Runway Gen-2 Sample / Tiger Whiskey," *YouTube*, July 9, 2023.
Kaoru Naito, *Twitter*, https://twitter.com/ka0ru_1620/status/1678313226453520385, July 10, 2023.

# Generative AI for News



*Generate an audio in Science Fiction theme: Mars News reporting that Humans send light-speed probe to Alpha Centauri.*
*Start with news anchor, followed by a reporter interviewing a chief engineer from an organization that built this probe, founded by United Earth and Mars Government, and end with the news anchor again.*

| | |
|---:|:---|
| Script | **GPT-4** |
| Music | **MusicGen** |
| Narration | **Bark** |
| Sound effects | **AudioLDM** |

Liu et al., "WavJourney: Compositional Audio Creation with Large Language Models," *arXiv preprint arXiv:2307.14335*, 2023.

# Controllable Generative AI

**Large language models**
(GPT-4)

**Pretrained generative audio models**
(MusicGen, AudioLDM, Bark)

Instructions    ➝    Audio Script    ➝    Audio

| Audio Type | Layout | ID | Character | Volume | Action | Content Description | Duration |
|---|---|---|---|---|---|---|---|
| Music | Background | 1 | N/A | -30 | Begin | Dramatic orchestral news theme. | Auto |
| Speech | Foreground | N/A | Host | -15 | N/A | Welcome to Mars News ... | Auto |
| Music | Background | 1 | N/A | N/A | End | N/A | Auto |
| Speech | Foreground | N/A | Host | -15 | N/A | Now let's connect with our on-site reporter ... | Auto |
| Sound effect | Foreground | N/A | N/A | -35 | N/A | Transition swoosh. | 1 |
| Sound effect | Background | 2 | N/A | -30 | Begin | Background noise of busy engineering office. | Auto |
| Speech | Foreground | N/A | Reporter | -15 | N/A | We're here at the headquarters of ... | Auto |
| Speech | Foreground | N/A | Director | -15 | N/A | Thank you, so it's a fantastic ... | Auto |
| Speech | Foreground | N/A | Reporter | -15 | N/A | This is truly an impressive feat ... | Auto |

Liu et al., "WavJourney: Compositional Audio Creation with Large Language Models," *arXiv preprint arXiv:2307.14335*, 2023.

# Controllable Generative AI

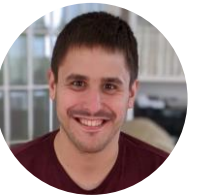| Audio Type | Layout | ID | Character | Volume | Action | Content Description | Duration |
|---|---|---|---|---|---|---|---|
| Music | Background | 1 | N/A | -30 | Begin | Dramatic orchestral news theme. | Auto |
| Speech | Foreground | N/A | Host | -15 | N/A | Welcome to Mars News … | Auto |
| Music | Background | 1 | N/A | N/A | End | N/A | |
| Speech | Foreground | N/A | Host | -15 | N/A | Now let's connect with our on-site reporter … | |
| Sound effect | Foreground | N/A | N/A | -35 | N/A | Transition swoosh. | |
| Sound effect | Background | 2 | N/A | -30 | Begin | Background noise of busy engineering office. | |
| Speech | Foreground | N/A | Reporter | -15 | N/A | We're here at the headquarters of … | |
| Speech | Foreground | N/A | Director | -15 | N/A | Thank you, so it's a fantastic … | |
| Speech | Foreground | N/A | Reporter | -15 | N/A | This is truly an impressive feat … | |

**Integration into professional creative workflow**

# Licensing Example – Adobe Firefly



**Trained with royalty-free Adobe Stock images**

# Acknowledgements

Thank you!