

KAUST Rising Stars in AI Symposium 2024

Learning Text-to-Audio Synthesis from Videos

Hao-Wen (Herman) Dong

UC San Diego

February 20, 2024

AI for Music & Audio

New technology creates new art form



AI

**Empowering music and audio creation
with machine learning**



Music & Audio

Music & Audio for AI

New art form inspires new technology



Generative AI for Music & Audio

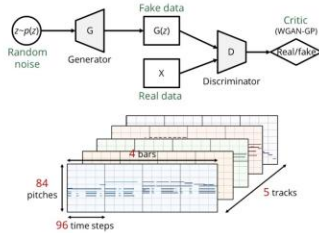
Empowering music and audio creation with machine learning

Multitrack Music Generation

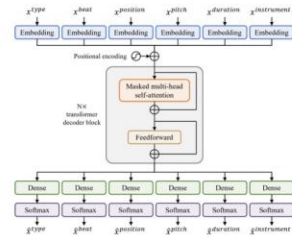
Advancing deep generative models for multitrack music



MuseGAN (AAAI 2018)



MMT (ICASSP 2023)

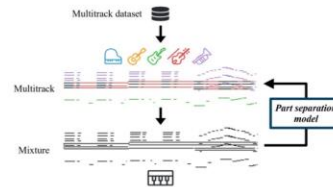


Assistive Music Creation Tools

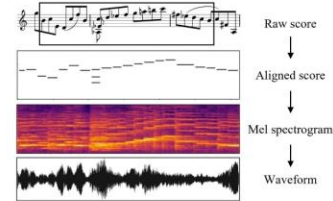
Developing AI-augmented assistive music creation tools



Arranger (ISMIR 2021)



Deep Performer (ICASSP 2022)

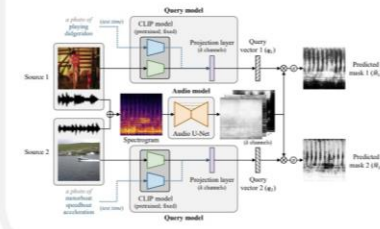


Multimodal Learning for Audio & Music

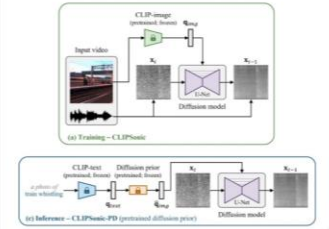
Learning sound separation and synthesis from videos



CLIPSep (ICLR 2023)



CLIPsonic (WASPAA 2023)



🧠 Generative AI for Music & Audio 🎵

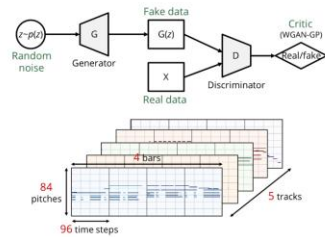
Empowering music and audio creation with machine learning

Multitrack Music Generation

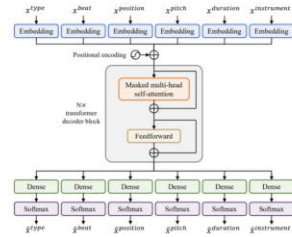
Advancing deep generative models for multitrack music



MuseGAN (AAAI 2018)



MMT (ICASSP 2023)



Assistive Music Creation Tools

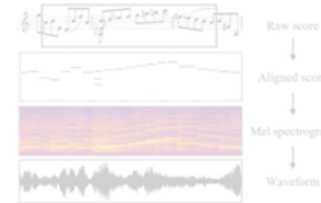
Developing AI-augmented assistive music creation tools



Arranger (ISMIR 2021)



Deep Performer (ICASSP 2022)

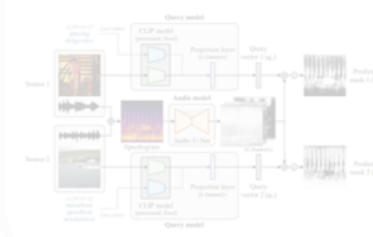


Multimodal Learning for Audio & Music

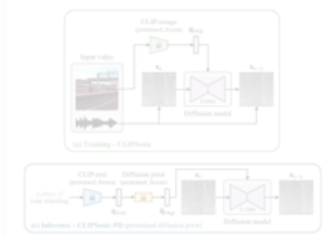
Learning sound separation and synthesis from videos



CLIPSep (ICLR 2023)



CLIPsonic (WASPAA 2023)



🧠 Generative AI for Music & Audio 🎵

Empo

Multitrack Music Generation

arning

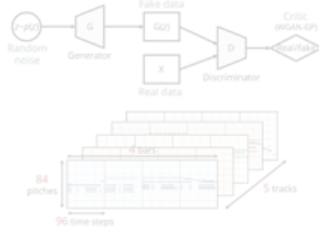
Advancing deep generative models for multitrack music



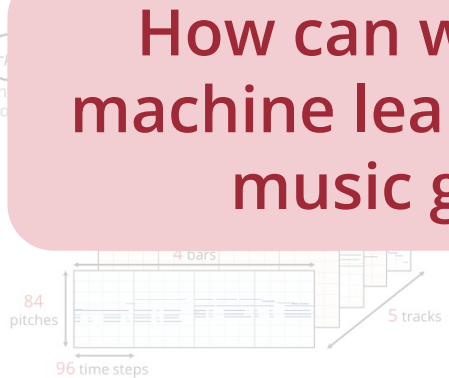
Multitrack Music Genera

Advancing deep generative models for multitrack music

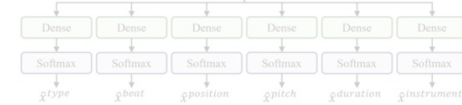
MuseGAN (AAAI 2018)



MuseGAN (AAAI 2018)



MMT (ICASSP 2023)



How can we build better machine learning models for music generation?

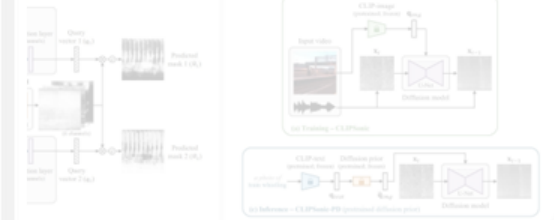
Supervised Learning for Audio & Music

Sound separation analysis from videos



Sep 2023)

CLIPsonic (WASPAA 2023)



🧠 Generative AI for Music & Audio 🎵

Empo

Multitrack Music Generation

arning

Advancing deep generative models for multitrack music



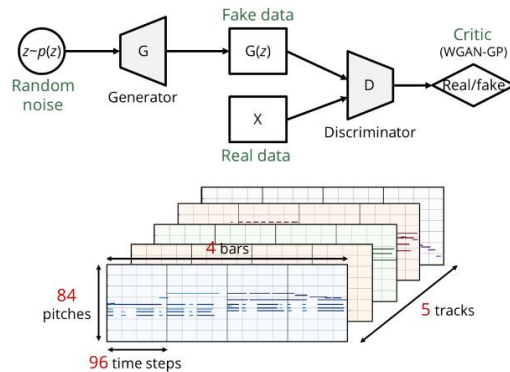
Multitrack Music Genera

dal Learning for Audio & Music

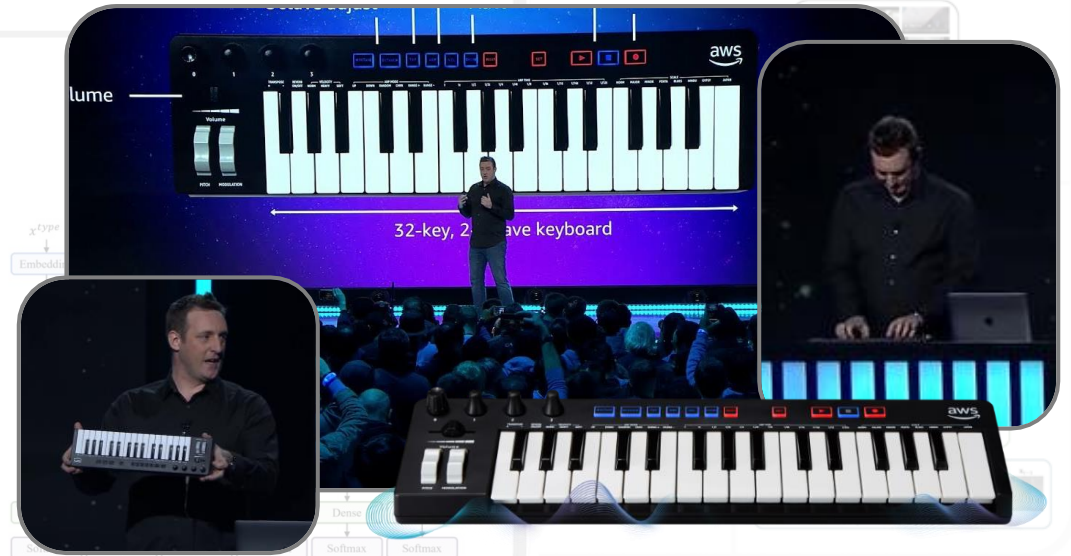
Advancing deep generative models for multitrack music



MuseGAN (AAAI 2018)



Pop Music Generation



Featured in Amazon AWS DeepComposer

First neural net that can generate multi-instrument music from scratch

🧠 Generative AI for Music & Audio 🎵

Empo

arning

Multitrack Music Generation

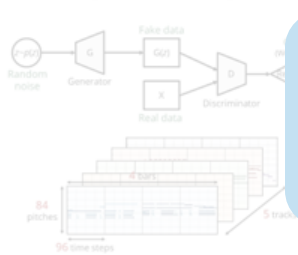
Advancing deep generative models for multitrack music



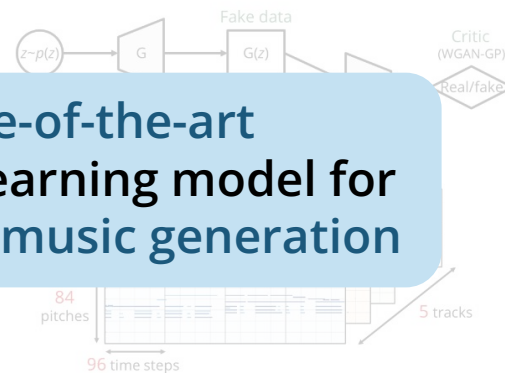
Multitrack Music Genera

Advancing deep generative models for multitrack music

MuseGAN (AAAI 2018)

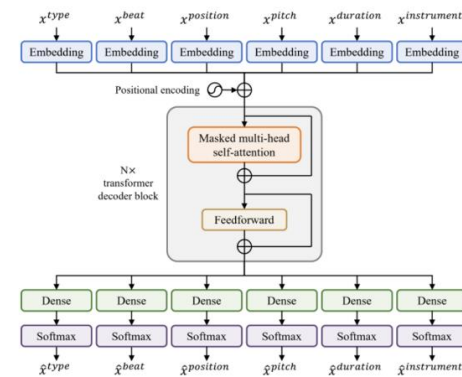


MuseGAN (AAAI 2018)



State-of-the-art machine learning model for orchestral music generation

MMT (ICASSP 2023)

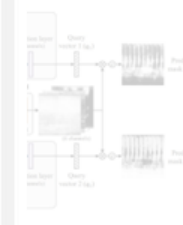


Self-supervised Learning for Audio & Music

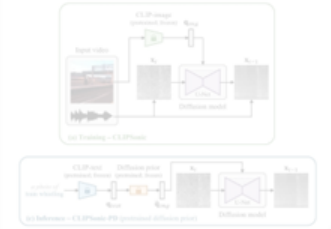
Sound separation analysis from videos



Sep (2023)



CLIPsonic (WASPAA 2023)



Orchestral Music Generation

🧠 Generative AI for Music & Audio 🎵

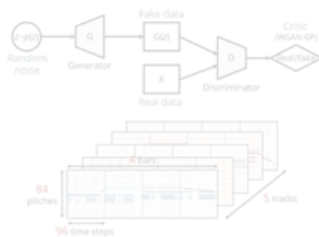
Empowering music and audio creation with machine learning

Multitrack Music Generation

Advancing deep generative models for multitrack music



MuseGAN (AAAI 2018)



MMT (ICASSP 2023)

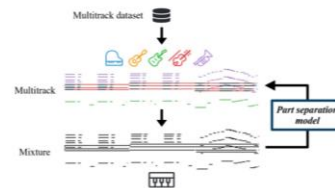


Assistive Music Creation Tools

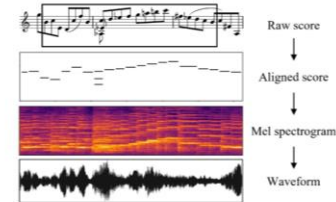
Developing AI-augmented assistive music creation tools



Arranger (ISMIR 2021)



Deep Performer (ICASSP 2022)

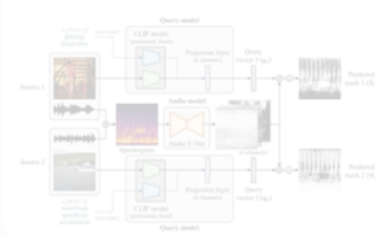


Multimodal Learning for Audio & Music

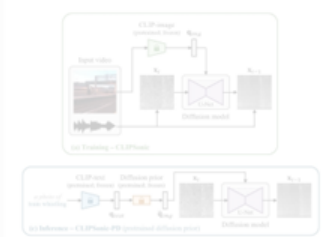
Learning sound separation and synthesis from videos



CLIPSep (ICLR 2023)



CLIPsonic (WASPAA 2023)



🧠 Generative AI for Music & Audio 🎵

Empowering

Assistive Music Creation Tools

Learning

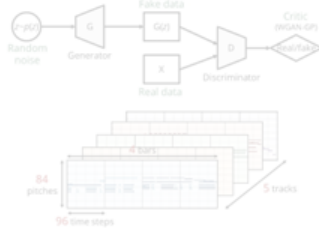
Developing AI-augmented assistive music creation tools



Multitrack Music Generation

Advancing deep generative models for multitrack music

MuseGAN (AAAI 2018)



(ICASSP 2021)



Arranger (ISMIR 2021)

Deep Performer (ICASSP 2022)

How can AI help professionals and amateurs create music?

Multitrack

Mixture



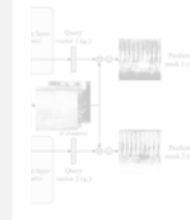
Mel spectrogram
↓
Waveform

Deep Learning for Audio & Music

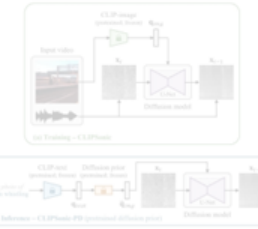
Sound separation and source separation from videos



Deep Learning (2023)



CLIPsonic (WASPAA 2023)



🧠 Generative AI for Music & Audio 🎵

Empirical

Learning

Assistive Music Creation Tools

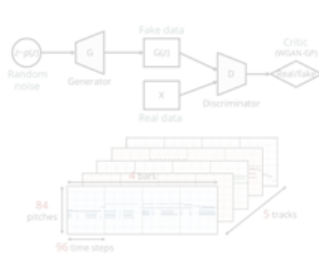
Developing AI-augmented assistive music creation tools



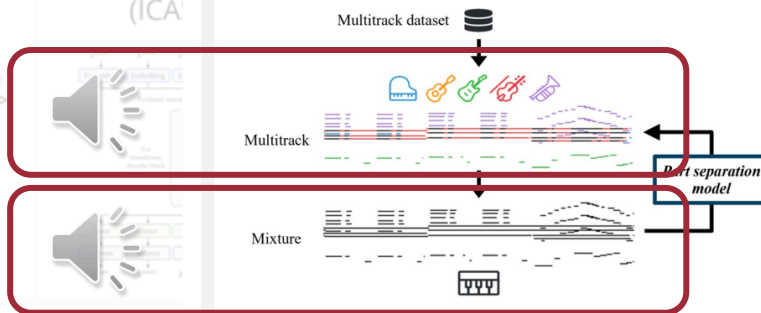
Multitrack Music Generation

Advancing deep generative models for multitrack music

MuseGAN (AAAI 2018)



Arranger (ISMIR 2021)



Deep Performer (ICASSP 2022)



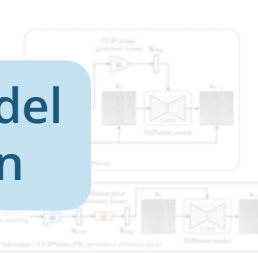
Deep Learning for Audio & Music

Audio and separation analysis from videos



Deep Learning for Audio & Music

CLIPsonic (WASPAA 2023)



First ever machine learning model for automatic instrumentation

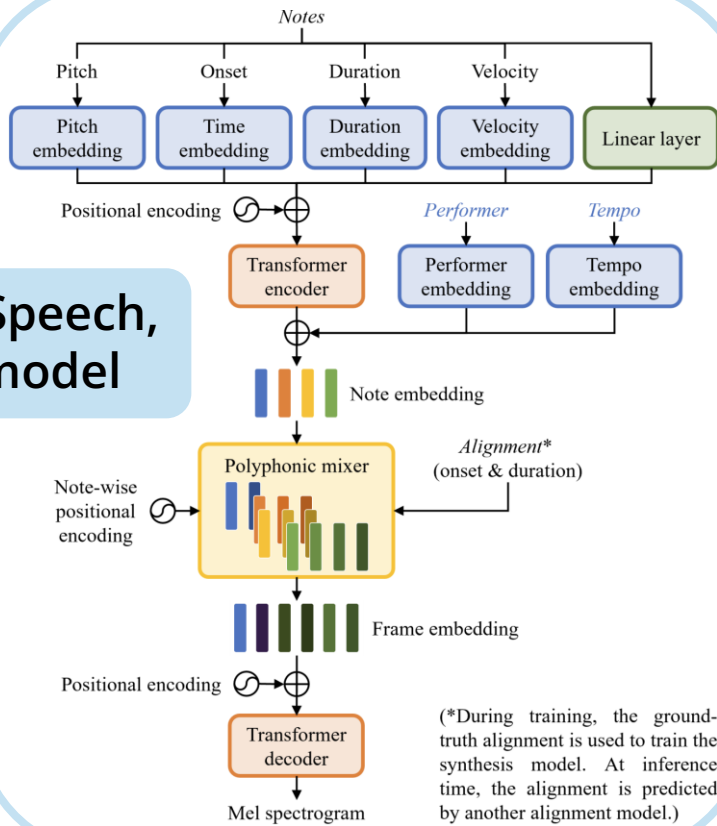
Automatic Instrumentation

Generative AI for Music & Audio 🎵

Empirical

Assistive Music Creation Tools

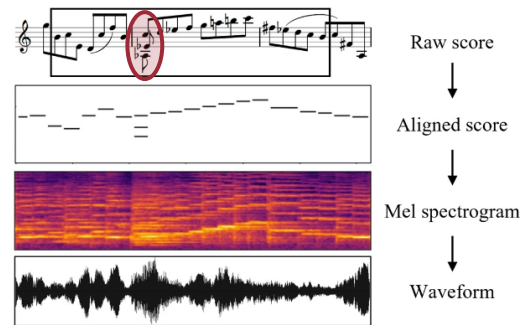
Learning



Adapted from FastSpeech, a text-to-speech model



Deep Performer (ICASSP 2022)



Score-to-audio synthesis

Multitrack Music

ed tools

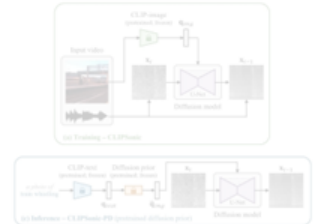
Deep Learning for Audio & Music

Sound separation and source separation from videos



Step 1 (2023)

CLIPsonic (WASPAA 2023)



🧠 Generative AI for Music & Audio 🎵

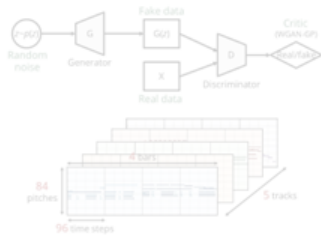
Empowering music and audio creation with machine learning

Multitrack Music Generation

Advancing deep generative models for multitrack music



MuseGAN (AAAI 2018)



MMT (ICASSP 2023)



Assistive Music Creation Tools

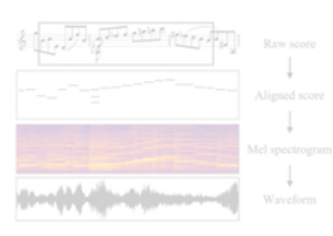
Developing AI-augmented assistive music creation tools



Arranger (ISMIR 2021)



Deep Performer (ICASSP 2022)

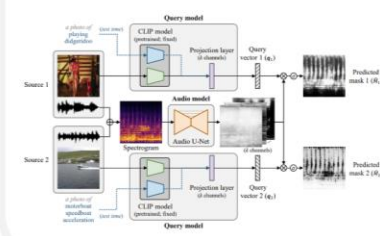


Multimodal Learning for Audio & Music

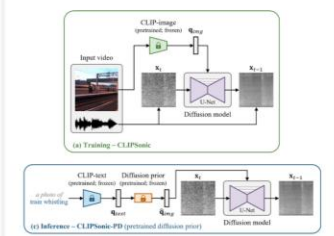
Learning sound separation and synthesis from videos



CLIPSep (ICLR 2023)



CLIPsonic (WASPAA 2023)



🧠 Generative AI for Music & Audio 🎵

Empc

arning

Multimodal Learning for Audio & Music

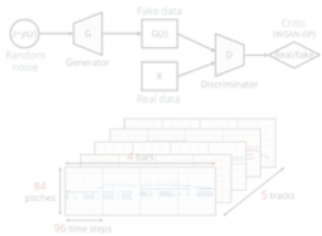
Learning sound separation and synthesis from videos



Multitrack Music Genera

Advancing deep generative models for multitrack music

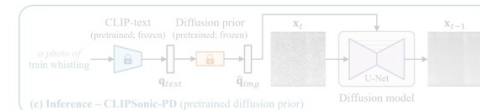
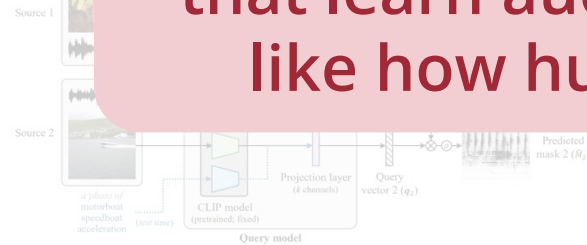
MuseGAN (AAAI 2018)



CLIPSep (ICLR 2023)

CLIPsonic (WASPAA 2023)

How can we build AI systems that learn audio concepts like how humans do?



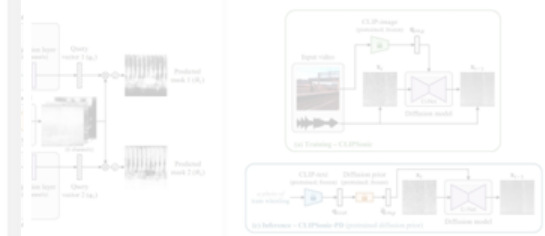
Multimodal Learning for Audio & Music

Sound separation and synthesis from videos

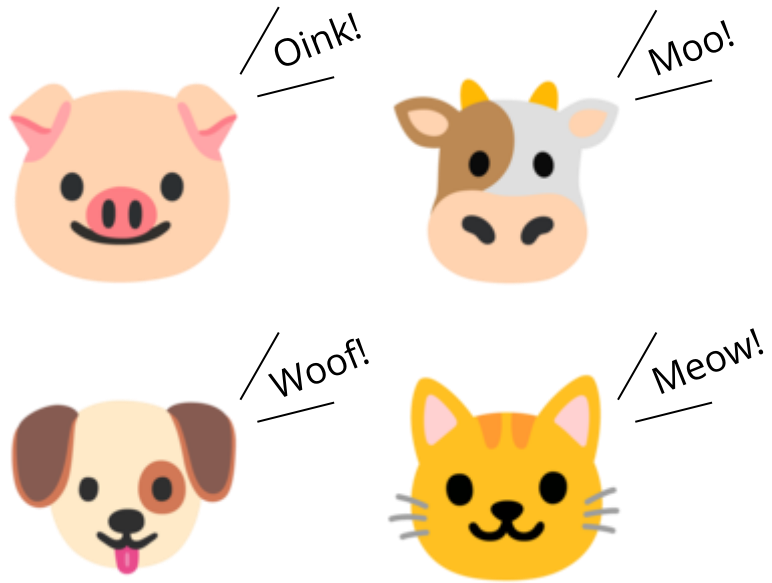


CLIPSep (ICLR 2023)

CLIPsonic (WASPAA 2023)

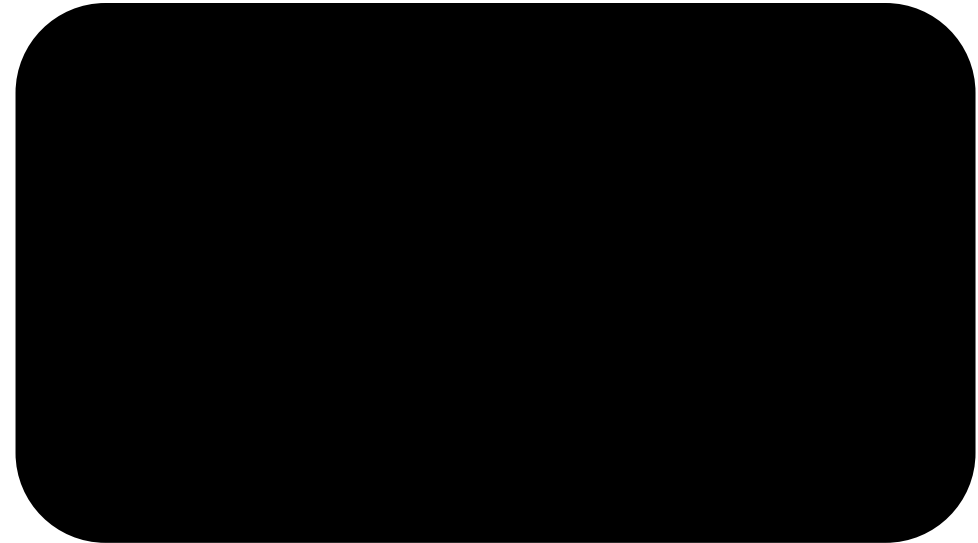
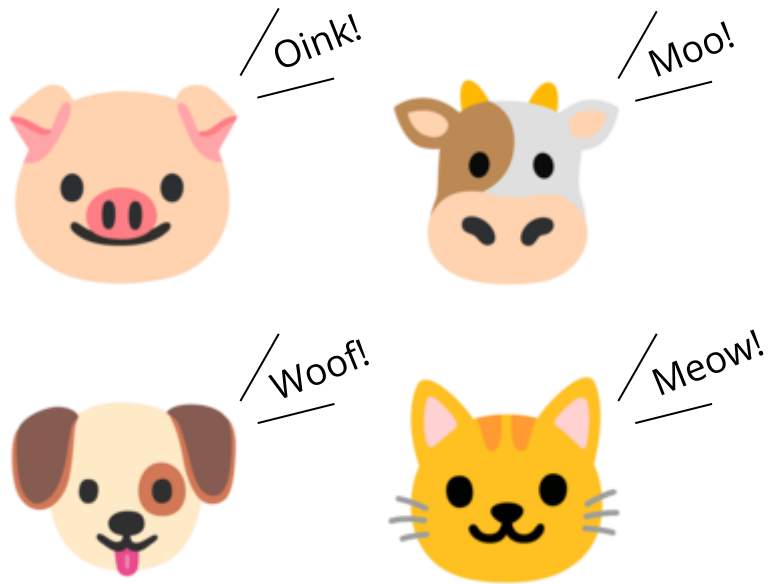


Learning Sounds from Observations



What does the fox say?

Learning Sounds from Observations



Can machines learn to synthesize sounds from watching *noisy* videos?

🧠 Generative AI for Music & Audio 🎵

Empc

arning

Multimodal Learning for Audio & Music

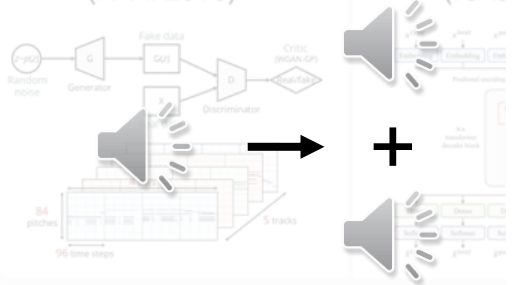
Learning sound separation and synthesis from videos



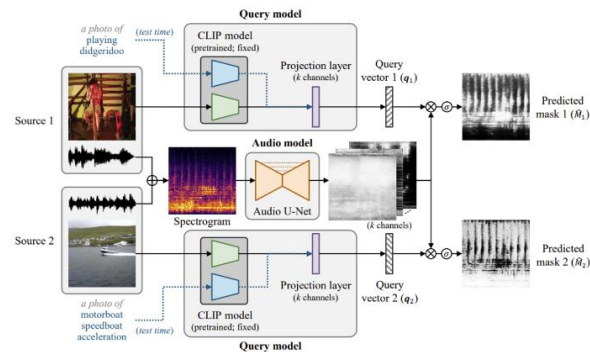
Multitrack Music Genera

Advancing deep generative models for multitrack music

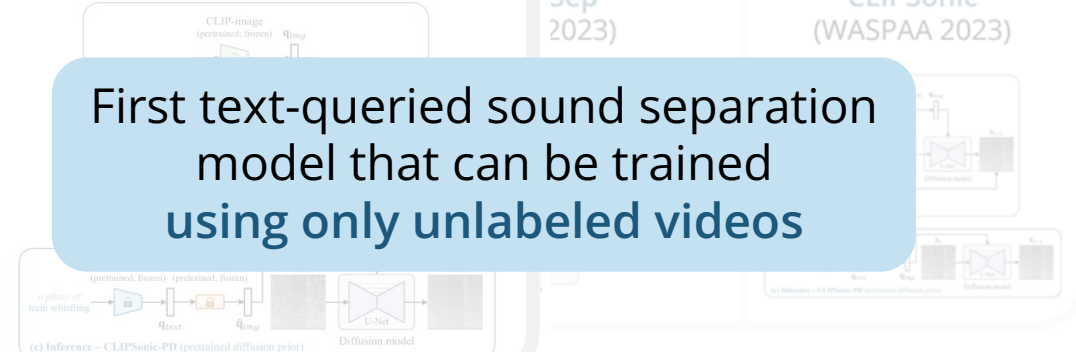
Query: "playing harpsichord" (AAAI 2018)



CLIPSep (ICLR 2023)



CLIP Sonic (WASPAA 2023)



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

Text-queried sound separation

🧠 Generative AI for Music & Audio 🎵

Empc

arning

Multimodal Learning for Audio & Music

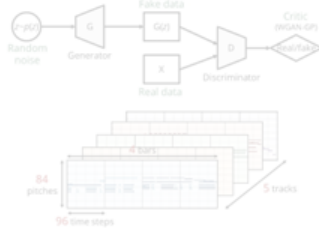
Learning sound separation and synthesis from videos



Multitrack Music Genera

Advancing deep generative models for multitrack music

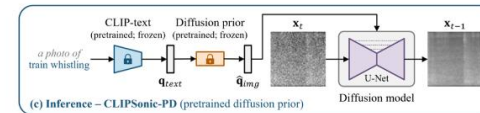
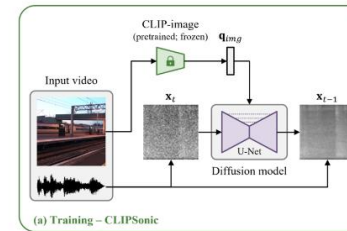
MuseGAN (AAAI 2018)



CLIPSep (ICLR 2023)

First text-to-sound synthesis model that can be trained using only unlabeled videos

CLIPsonic (WASPAA 2023)



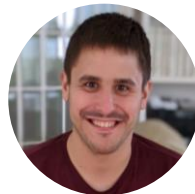
Text-to-audio synthesis

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

Hao-Wen Dong^{1,2*} Xiaoyu Liu¹ Jordi Pons¹ Gautam Bhattacharya¹
Santiago Pascual¹ Joan Serrà¹ Taylor Berg-Kirkpatrick² Julian McAuley²

¹ Dolby Laboratories ² University of California San Diego

* Work done during an internship at Dolby

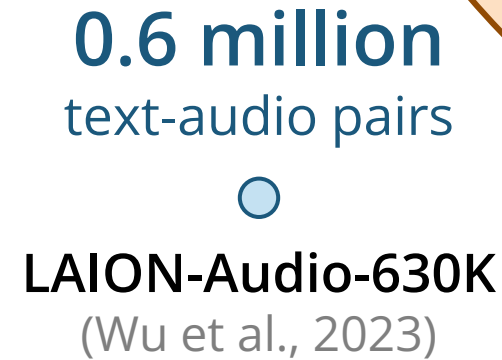
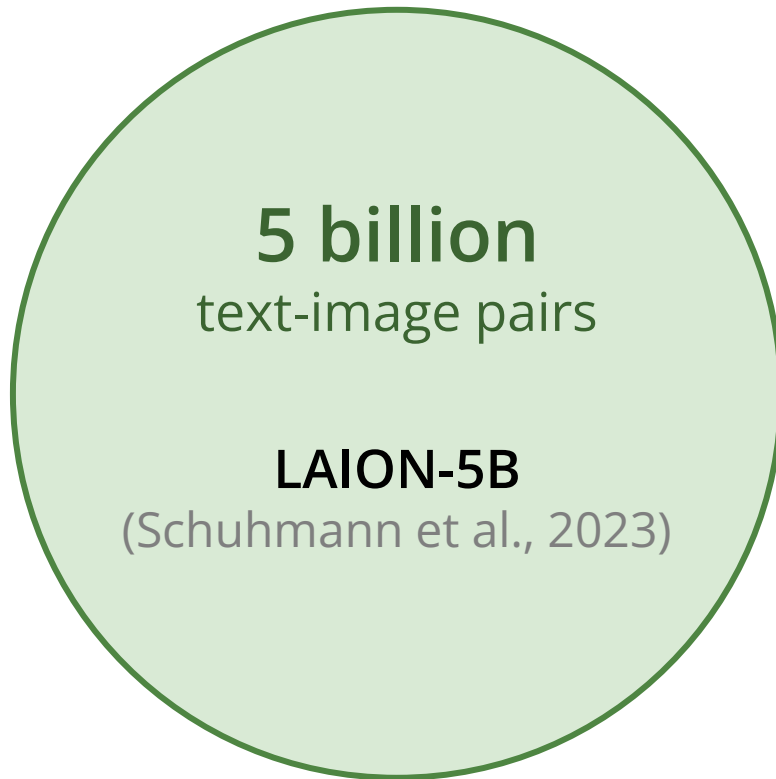


What is Text-to-Audio Synthesis?

- Goal: Given a text query, generate the corresponding sounds

(These samples are generated by our proposed model.)

Why NOT Text-audio Pairs?

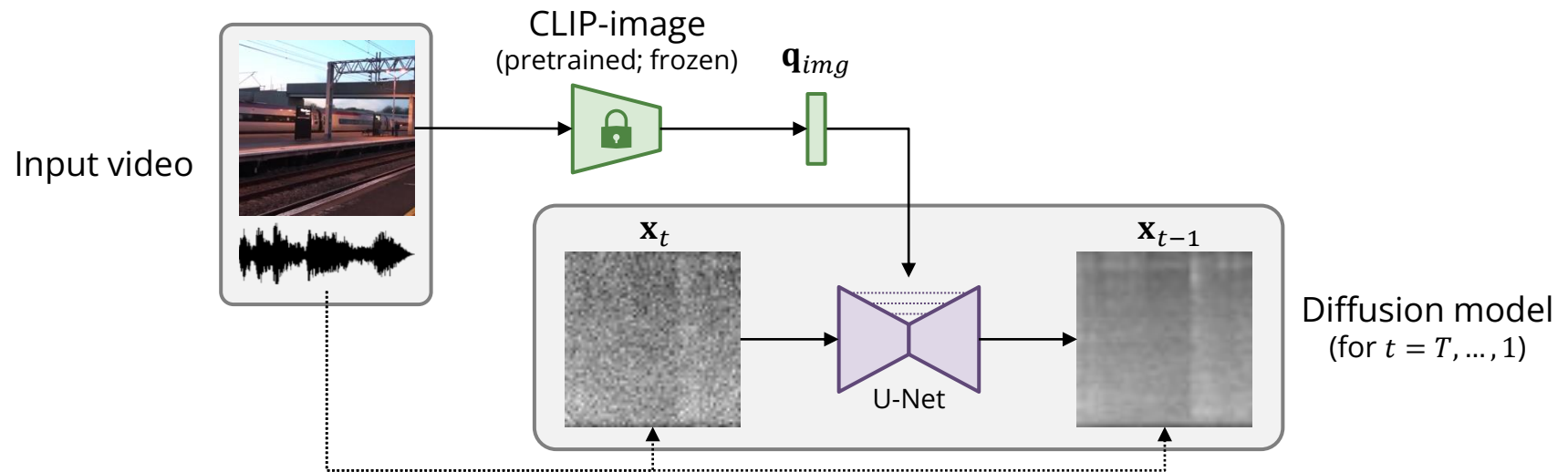


YouTube videos!

500 hours of videos
uploaded per minute

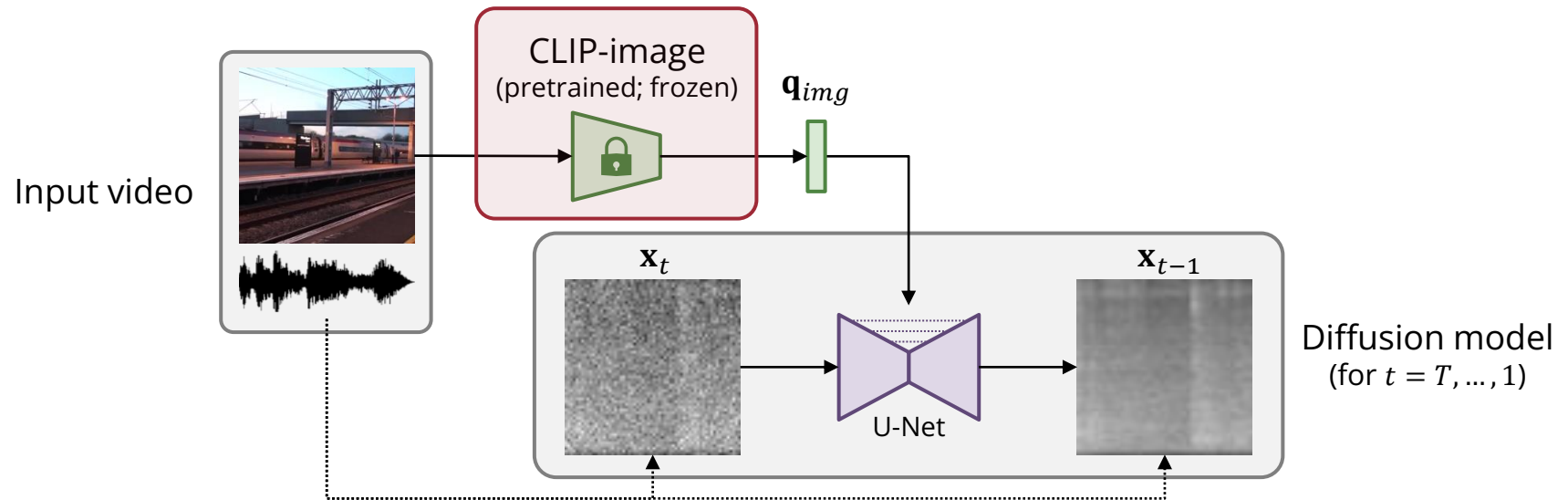
Training an Image-to-Audio Synthesis Model

- We start by training an image-to-audio synthesis model



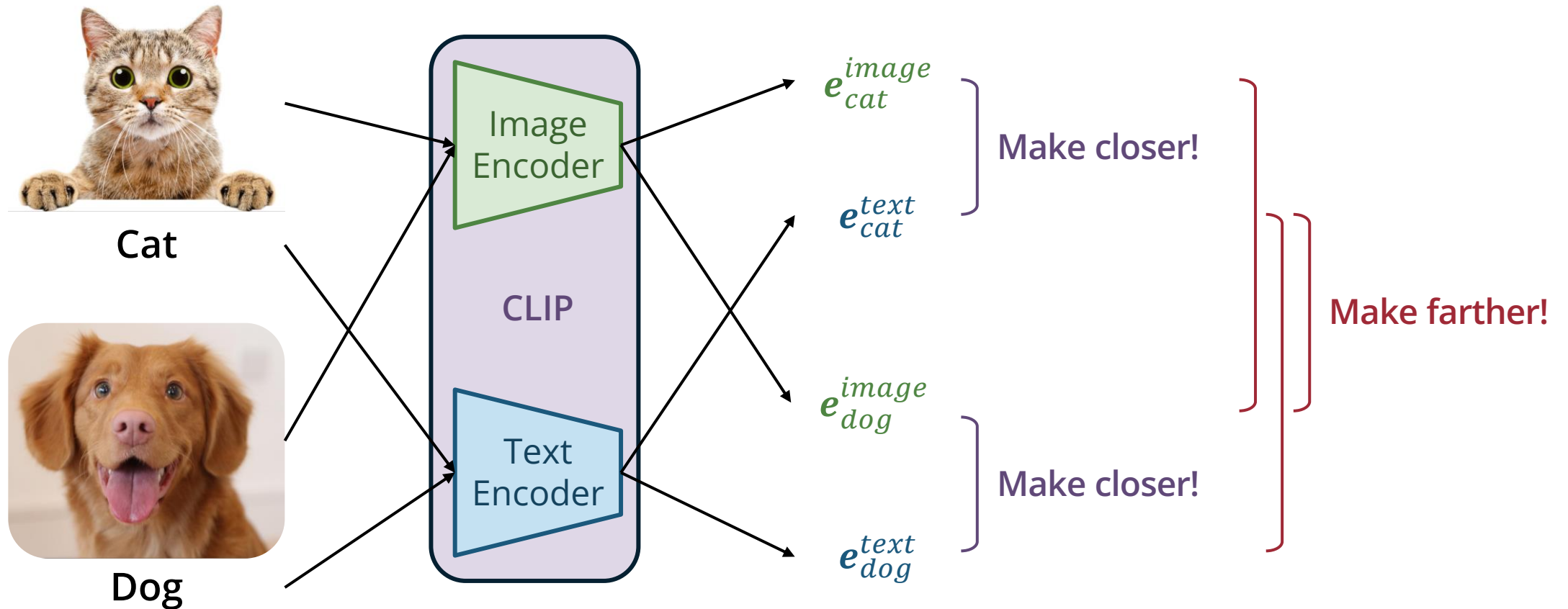
Training an Image-to-Audio Synthesis Model

- We start by training an image-to-audio synthesis model



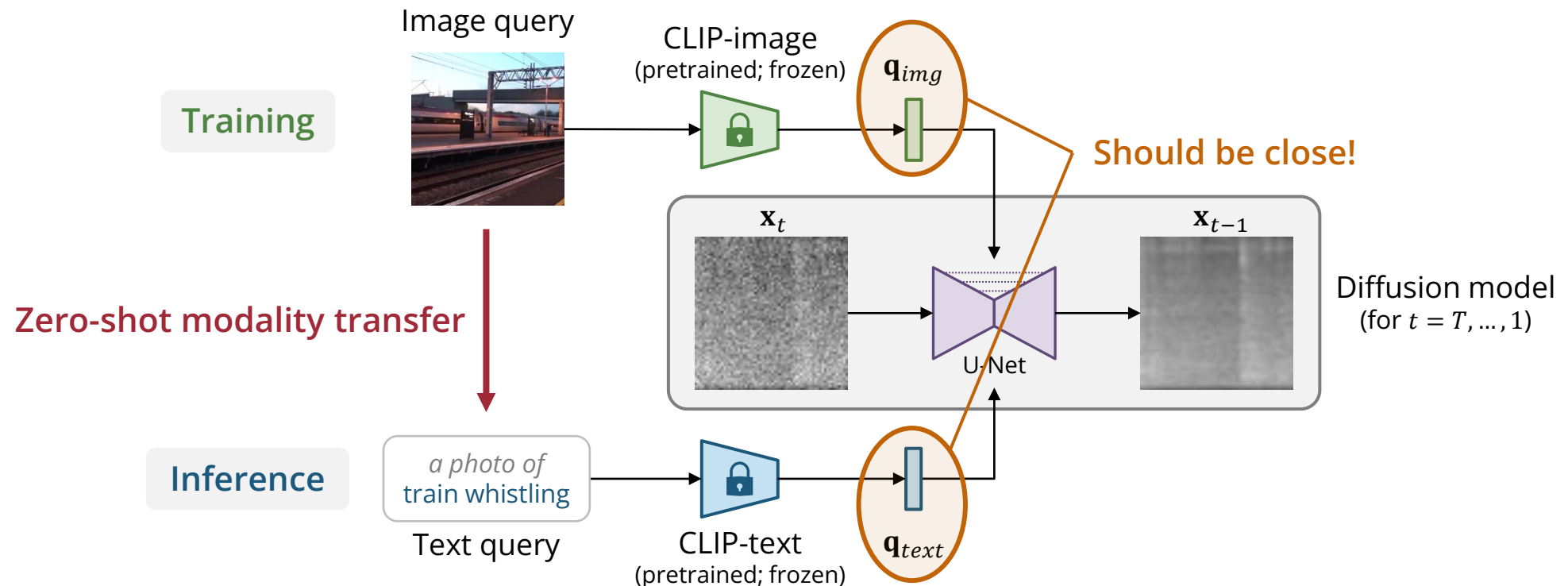
CLIP (Contrastive Language-Image Pretraining)

- Learn a **shared embedding space** for images and texts via *contrastive learning*



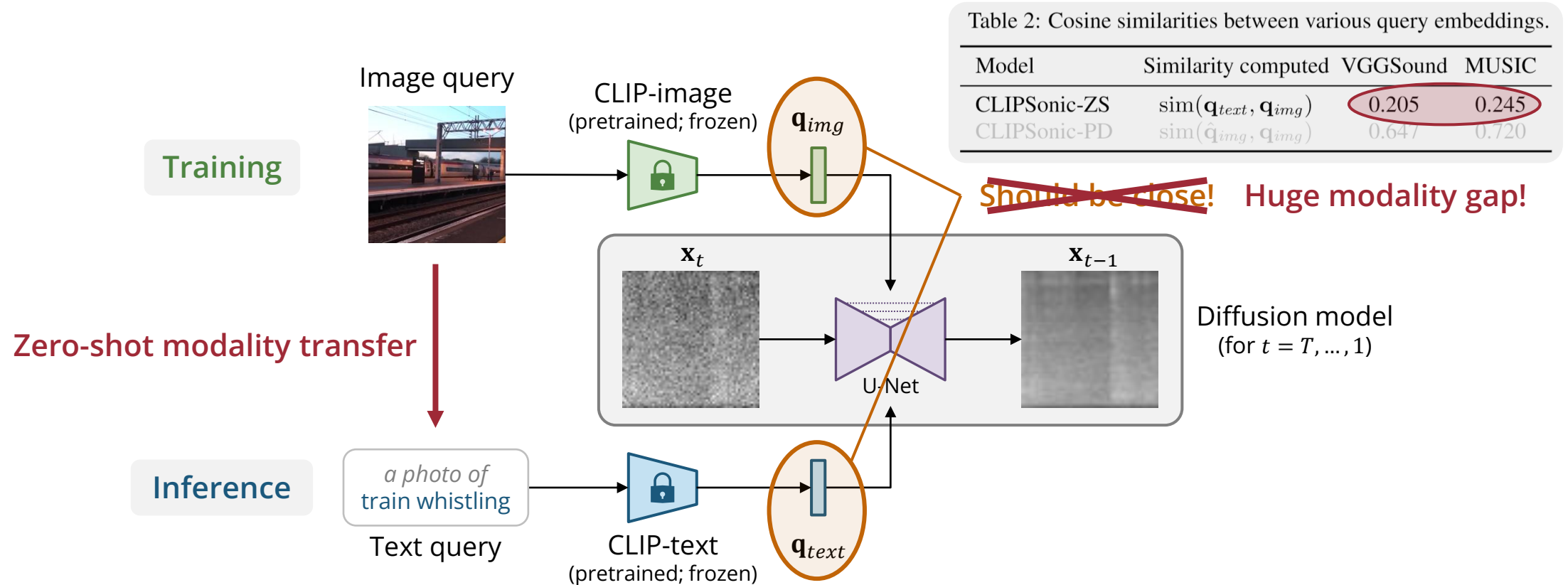
Inference – Zero-shot Modality Transfer

- We switch to a pretrained CLIP-text encoder for text-to-sound synthesis



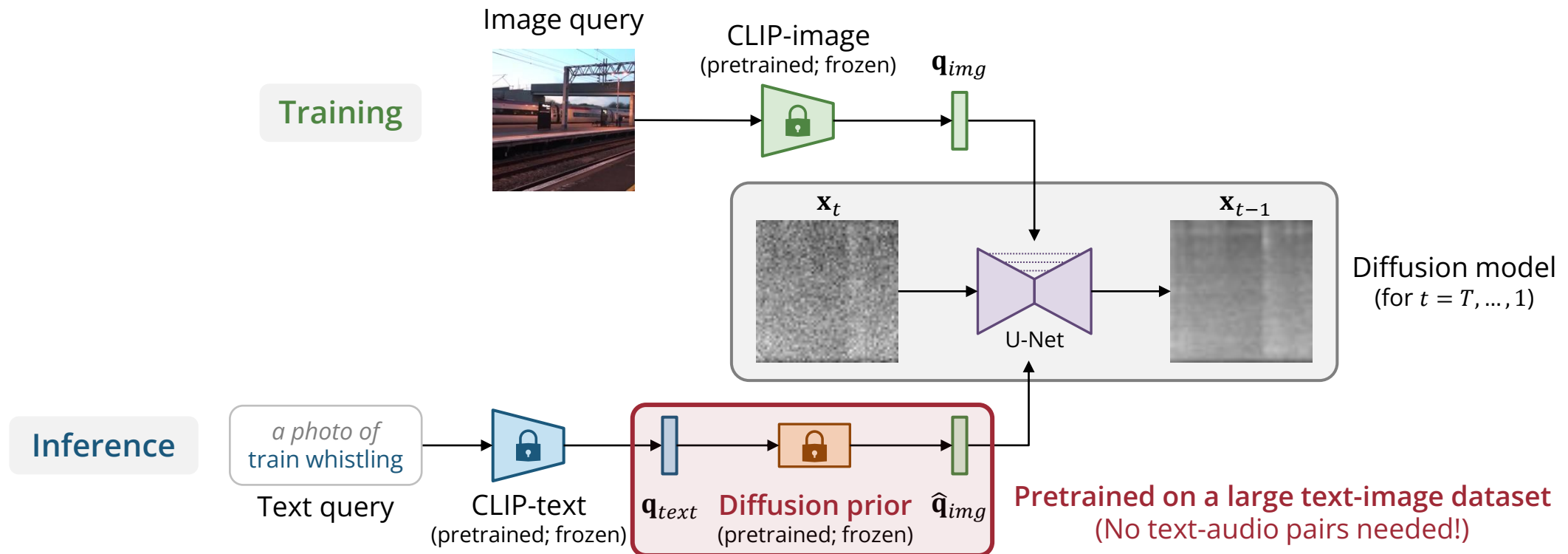
Inference – Zero-shot Modality Transfer

- We switch to a pretrained CLIP-text encoder for text-to-sound synthesis



Leveraging Diffusion Prior to Close the Modality Gap

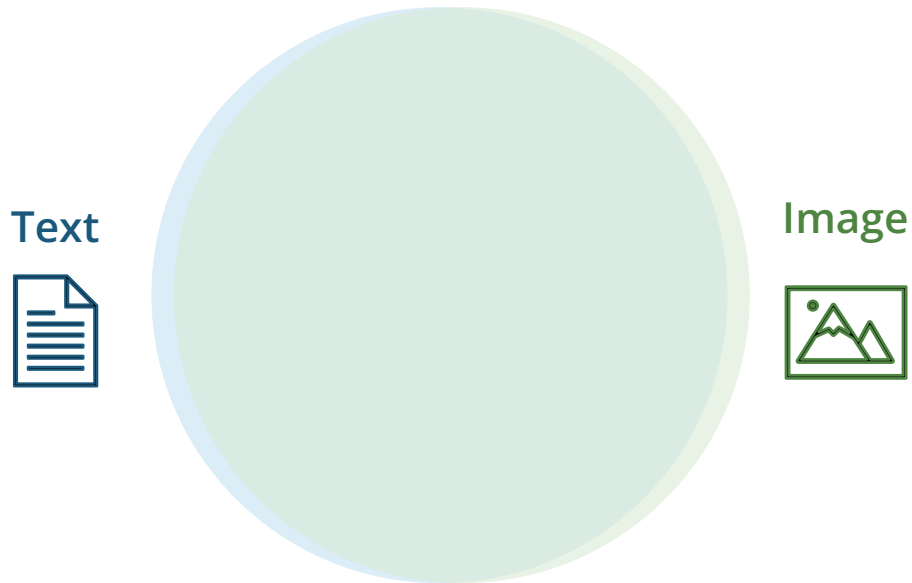
- We adopt a pretrained diffusion prior model to reduce the modality gap



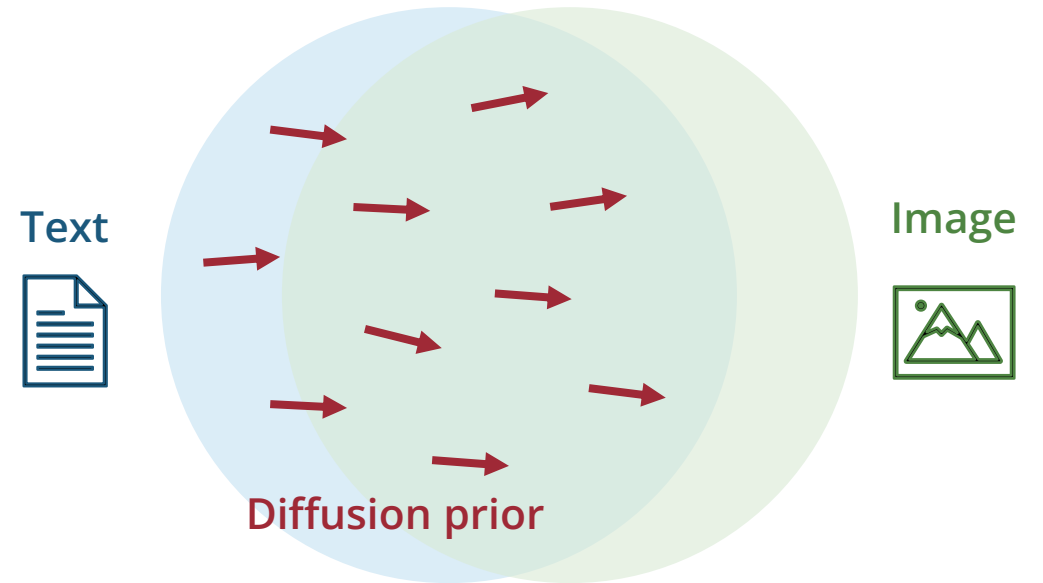
Diffusion Prior (Ramesh et al., 2022)

CLIP embedding spaces

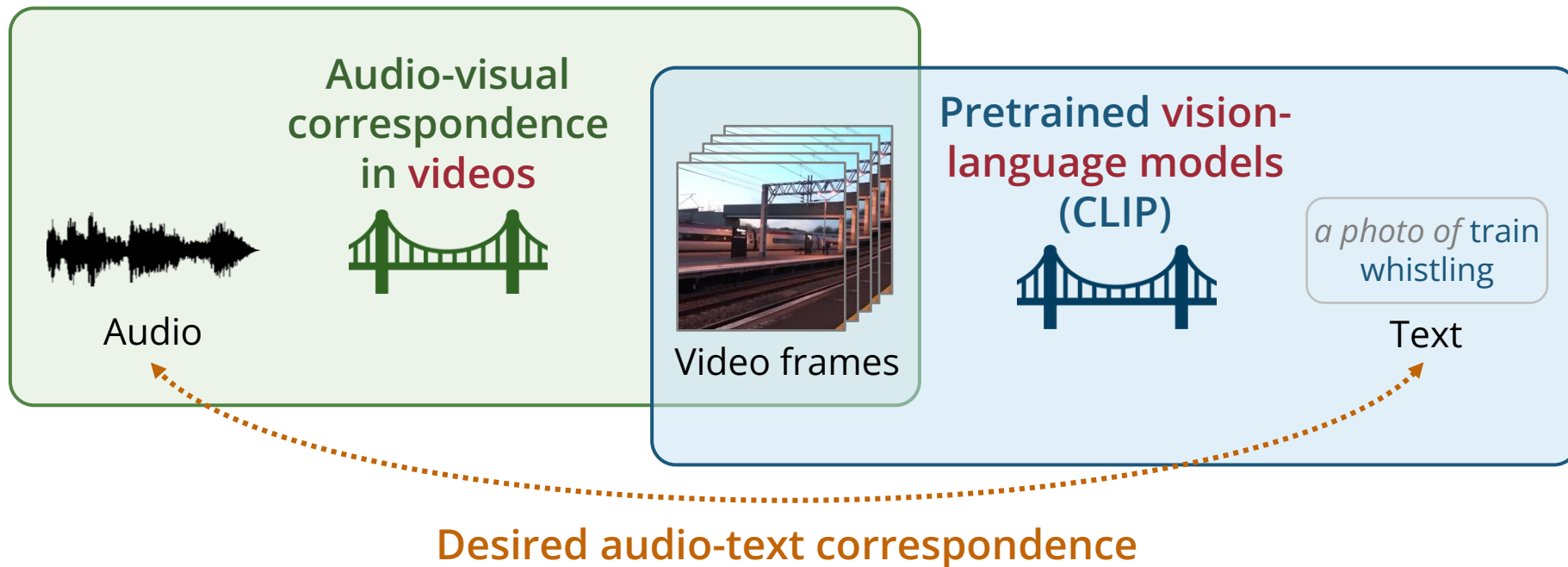
Ideal case



In practice



Leveraging the Visual Domain as a Bridge



No text-audio pairs required!

Scalable to large video datasets!

Data

MUSIC

(Zhao et al., 2018)



Violin



Acoustic guitar



Accordion

Music instrument playing videos

(1,055 videos, 21 instruments)

VGGSound

(Chen et al., 2020)



Hedge trimmer
running



Dog bow-wow



Bird chirping,
tweeting

Noisy videos with diverse sounds

(172K videos, 310 classes)

Example Text-to-Audio Synthesis Results

Rapping



Sea waves



Thunder



Smoke detector beeping



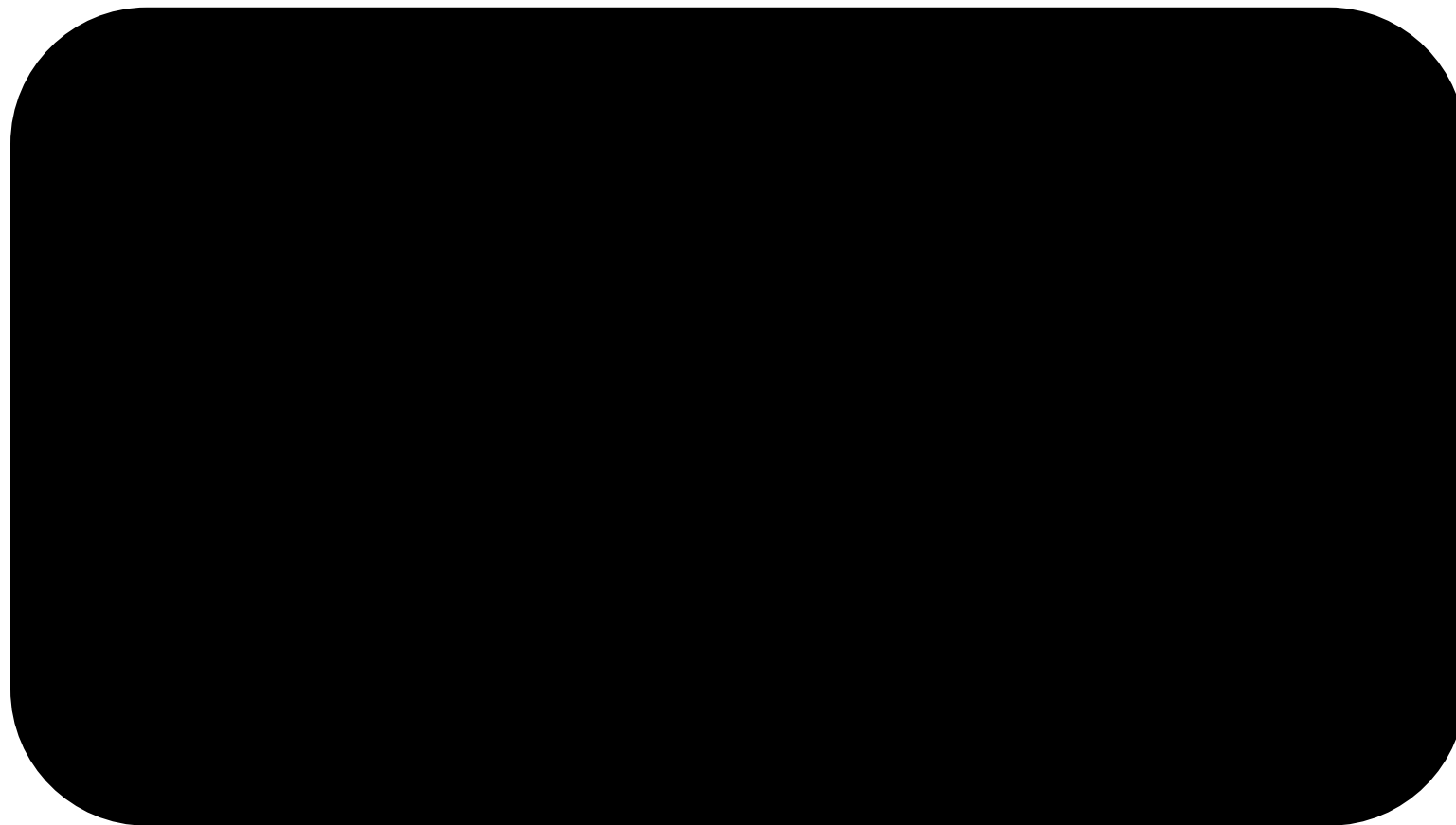
Playing table tennis



Playing violin fiddle



Example Image-to-Audio Synthesis Results (Out-of-distribution)



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

Subjective & Objective Evaluation Results

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

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

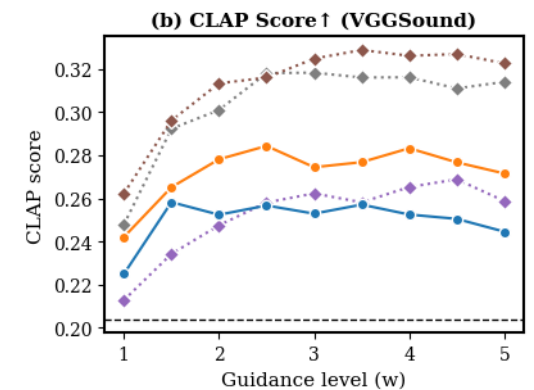
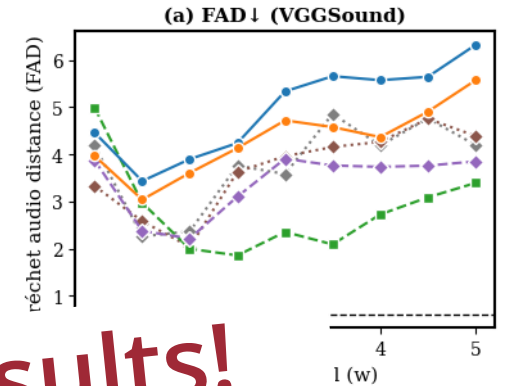
Table 4: Listening test results for image-to-audio synthesis (MOS).

Model	Fidelity	Relevance
CLIPSonic-IQ (image-queried)	3.29 ± 0.16	3.80 ± 0.19
SpecVQGAN [20]	2.15 ± 0.17	2.54 ± 0.23
im2wav [21]	2.19 ± 0.15	3.90 ± 0.22

Table 1: Evaluation results on VGGSound and MUSIC.

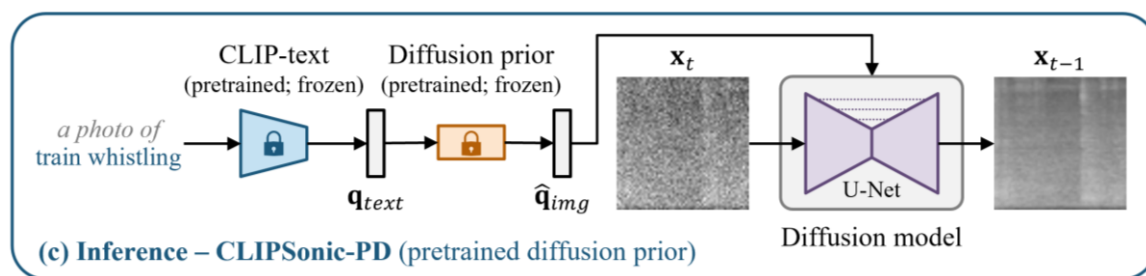
Model				VGGSound		MUSIC	
		Image	Text	FAD ↓	CLAP score ↑	FAD ↓	CLAP score ↑
CLIPSonic-IQ	-	Image	Image	2.97	-	4.71	-
CLIPSonic-ZS (zero-shot modality transfer)	✓	Image	Text	3.43	0.258	19.30	0.284
CLIPSonic-PD (pretrained diffusion prior)	✓	Image	Text	3.04	0.265	13.51	0.254
CLIPSonic-SD (supervised diffusion prior)	✗	Image	Text	2.37	0.234	12.13	0.299
CLIP-TTA	✗	Text	Text	2.26	0.292	9.39	0.298
CLAP-TTA	✗	Text	Text	2.58	0.296	10.92	0.303
BigVGAN mel spectrogram reconstruction	-	-	-	0.60	0.204	6.21	0.272

Check out our paper for more results!

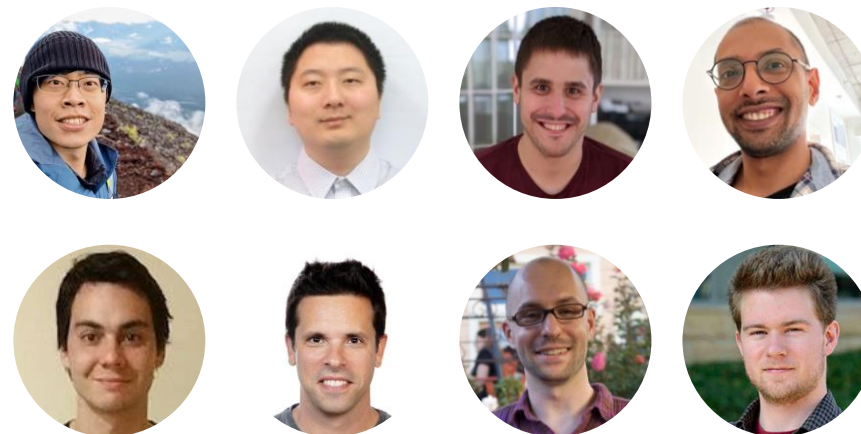


Summary

- First text-to-audio synthesis model that **requires *no* text-audio pairs**
- **Strong text-to-audio** synthesis performance without text-audio data
- **State-of-the-art image-to-audio** synthesis performance



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



What's Next?




Video → Music & sound effects
Text → Video with music & sound effects

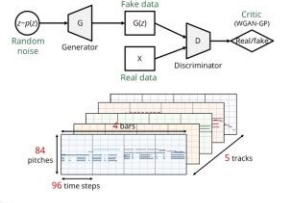
Generative AI for Music & Audio

Empowering music and audio creation with machine learning

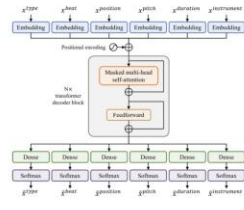
Multitrack Music Generation

Advancing deep generative models for multitrack music 


MuseGAN (AAAI 2018)



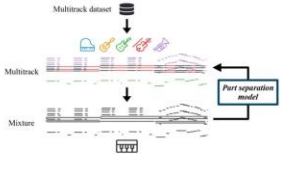
MMT (ICASSP 2023)



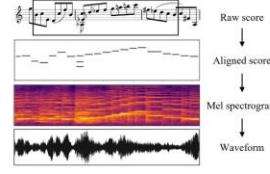
Assistive Music Creation Tools

Developing AI-augmented assistive music creation tools 


Arranger (ISMIR 2021)



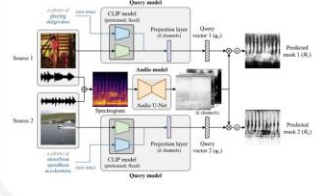
Deep Performer (ICASSP 2022)



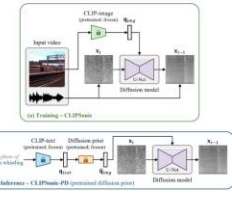
Multimodal Learning for Audio & Music

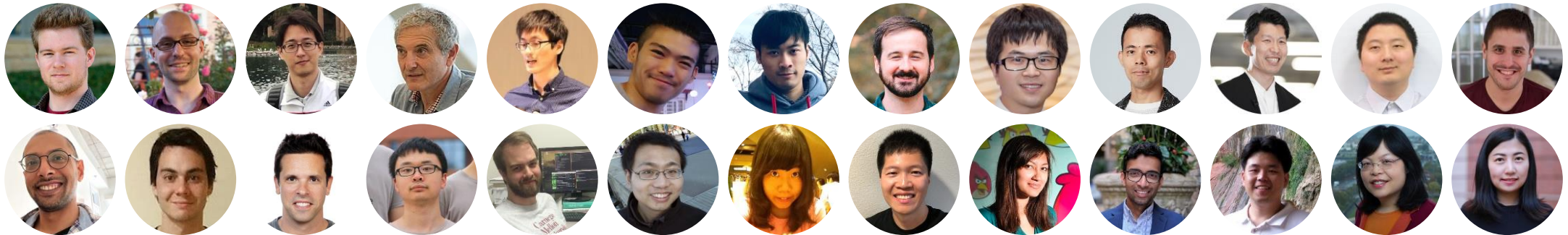
Learning sound separation and synthesis from videos 

CLIPSep (ICLR 2023)



CLIPsonic (WASPAA 2023)





UC San Diego

中央研究院
ACADEMIA SINICA

Dolby

SONY

amazon

