

Human-Centered Generative AI for Content Creation

Generating Music and Audio with Machine Learning

Hao-Wen (Herman) Dong



January 14, 2025

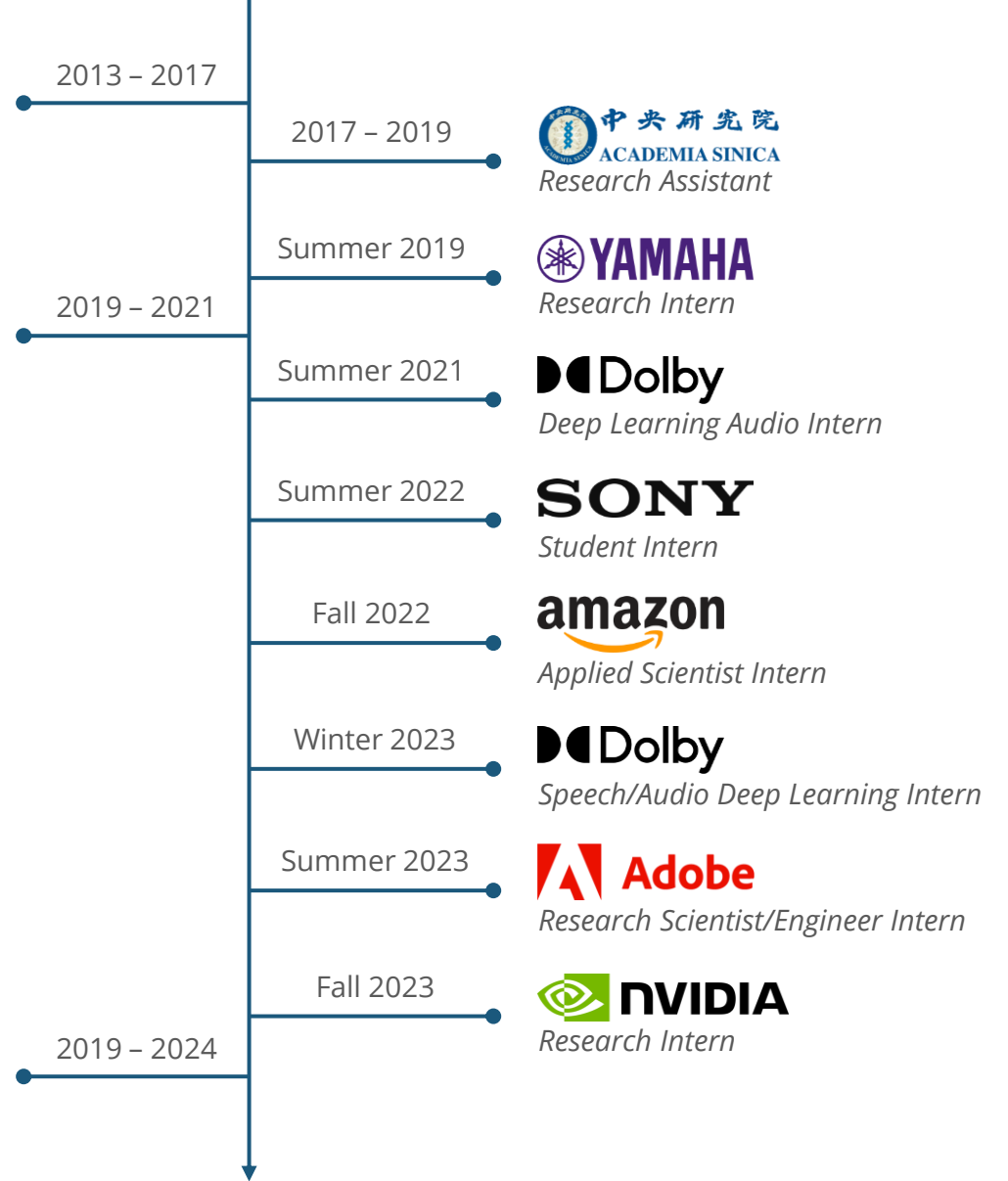
About Me



 國立臺灣大學
National Taiwan University
B.S. in Electrical Engineering

UC San Diego
M.S. in Computer Science

UC San Diego
Ph.D. in Computer Science



SCHOOL OF MUSIC, THEATRE & DANCE
PERFORMING ARTS TECHNOLOGY
UNIVERSITY OF MICHIGAN

My Background

Electrical Engineering



a female cat engineer making an electric chip in a classroom

Music



a cat playing heavy metal

Computer Science



a cat engineer debugging on laptop

Made in September 2023!

My Background

Electrical Engineering



a cat engineer making an electric chip in a classroom

Music



a cat playing heavy metal

January 2025
Made in ~~September 2023!~~

Computer Science



a cat engineer debugging Python programs on a laptop

Analytic AI vs Generative AI



Generative AI for Visual Arts



(Source: Cosmopolitan)

First Prize in Digital arts
at Colorado State Fair
Fine Arts Competition

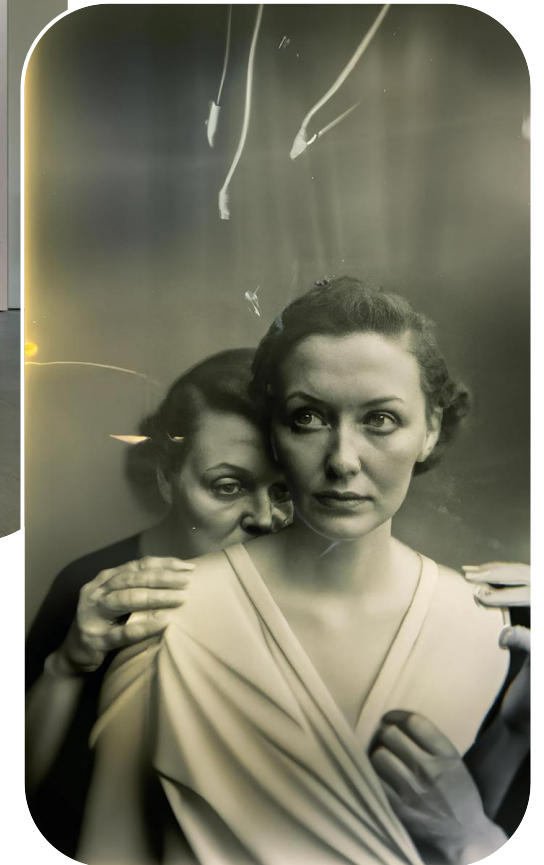


(Source: CNN Business)

(Source: MoMA Magazine)



Sony World
Photography Award in
Creative Open Category



(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.
Refik Anadol, "[Refik Anadol on AI, Algorithms, and the Machine as Witness](#)," *MoMA Magazine*, December 20, 2022.
Lianne Kolirin, "[Artist rejects photo prize after AI-generated image wins award](#)," *CNN*, April 18, 2023.

Generative AI for Video Generation

Sora released in February 2024!



Generative AI for Video Generation

Sora released in February 2024!



Challenges in Generative AI for Content Generation

- **New application domains may require new generative models**
- **Professionals need assistive tools that augment their creativity and productivity** in addition to fully automated tools
- **Certain media require handling multimodal data streams** at the same time

Human-Centered Generative AI for Content Creation

Augmenting human creativity with machine learning

- **Novel Generative Models for New Domains**
 - **Multitrack music generation** (AAAI 2018, ISMIR 2018, ISMIR 2020, ICASSP 2023, ISMIR 2024), **controllable music generation** (AIMG 2024, arXiv 2024), **documentary teaser generation** (arXiv 2024)
- **AI-Assisted Tools for Content Creation**
 - **Violin performance synthesis** (ICASSP 2022, arXiv 2024), **music instrumentation** (ISMIR 2021), **music arrangement** (AAAI 2018), **music harmonization** (JNMR 2020)
- **Multimodal Generative Models for Content Creation**
 - **Queried sound separation** (ICLR 2023), **text-to-audio synthesis** (WSS 2023, WASPAA 2023), **text-to-music generation** (ISMIR LBD 2024, arXiv 2024), **documentary teaser generation** (arXiv 2024)

Generating Music & Audio with Machine Learning



Coevolution of Music & Technology



Hildegard Dodel, Public domain, via Wikimedia Commons.
Taken at Hamamatsu Museum of Musical Instruments, August 2019.
yan, [CC BY-SA 4.0](https://creativecommons.org/licenses/by-sa/4.0/), via Wikimedia Commons.

Music & AI

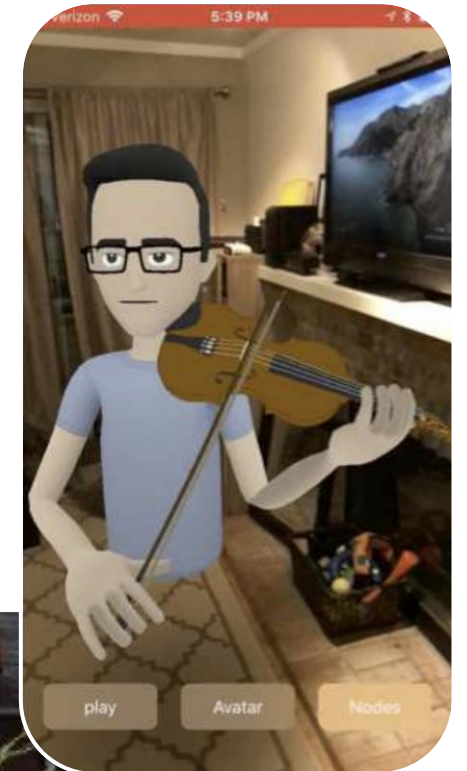
(Source: Yamaha)



(Source: Sankei Shimbun)



(Shlizerman et al., 2019)



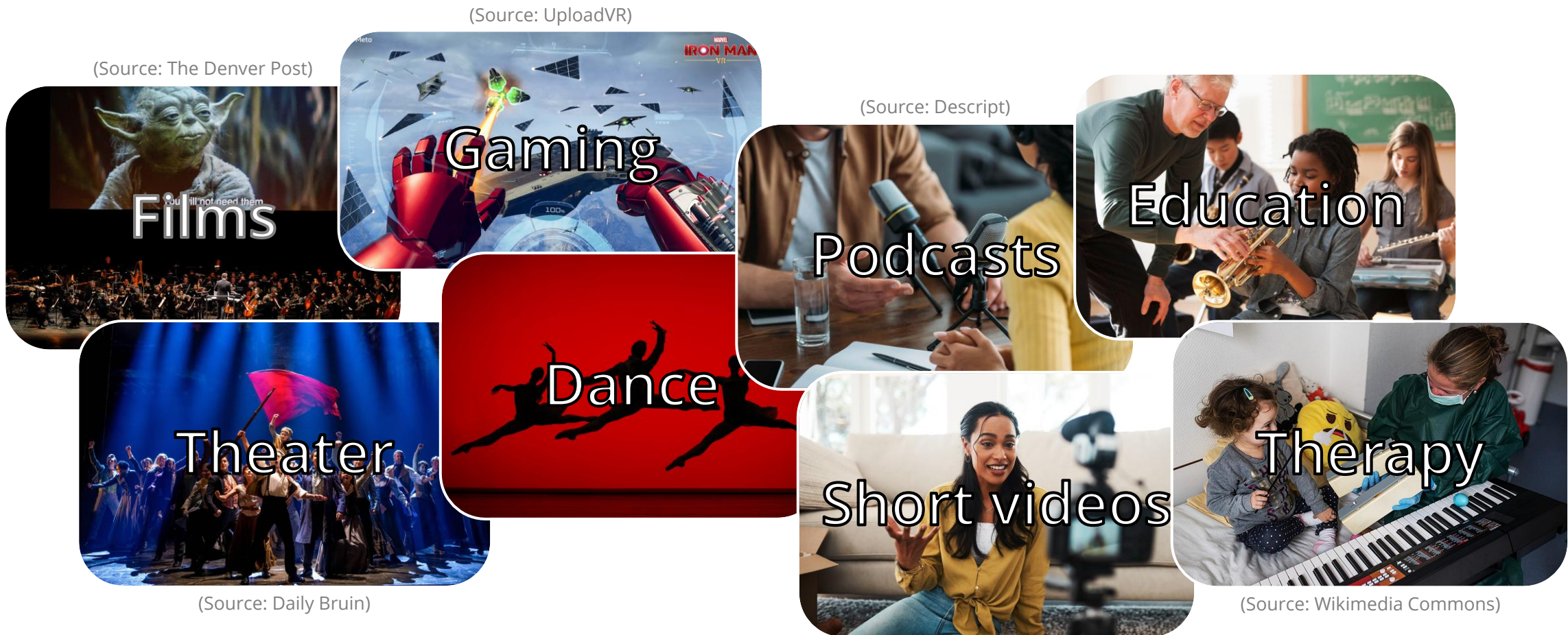
(Source: Robot Gizmos)



(Source: NBC DFW)

yamaha.com/en/news_release/2018/18013101/
sankei.com/article/20240113-CQCOSQHJWFIYPJJKZDCITRTRVI/
roboticgizmos.com/shimon-musical-robot-deep-learning/
nbcdfw.com/entertainment/the-scene/how-verdigris-ensemble-is-using-ai-to-create-a-new-concert-experience/3366031/
Shlizerman et al., "Audio to Body Dynamics," *Proc. CVPR*, 2018.

Generative AI for Content Creation



Universitaetsmedizin, [CC BY-SA 4.0](https://commons.wikimedia.org/wiki/File:Universitaetsmedizin), via Wikimedia Commons
uploadvr.com/iron-man-vr-breaks-free-from-cords-load-screens-on-quest-2/
descript.com/blog/article/what-is-the-best-audio-interface-for-recording-a-podcast
denverpost.com/2019/08/02/colorado-symphony-movie-scores-harry-potter-star-wars/
dailybruin.com/2023/08/04/theater-review-the-musical-les-misrables-offers-stellar-displays-and-impassioned-vocals

Part 1: Generating Music

Music Generation – Four Paradigms



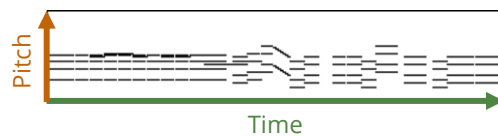
Symbolic music generation

Text-based

```
Program_change_0,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_76, Time_shift_2, Note_off_67,  
Note_on_67, Time_shift_2, Note_off_67,  
...
```

MIDI

Image-based



Piano roll



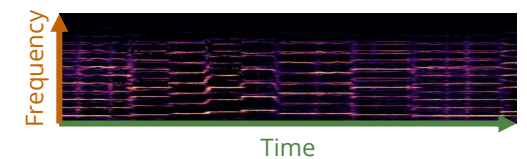
Audio-domain music generation

Time series-based



Waveform

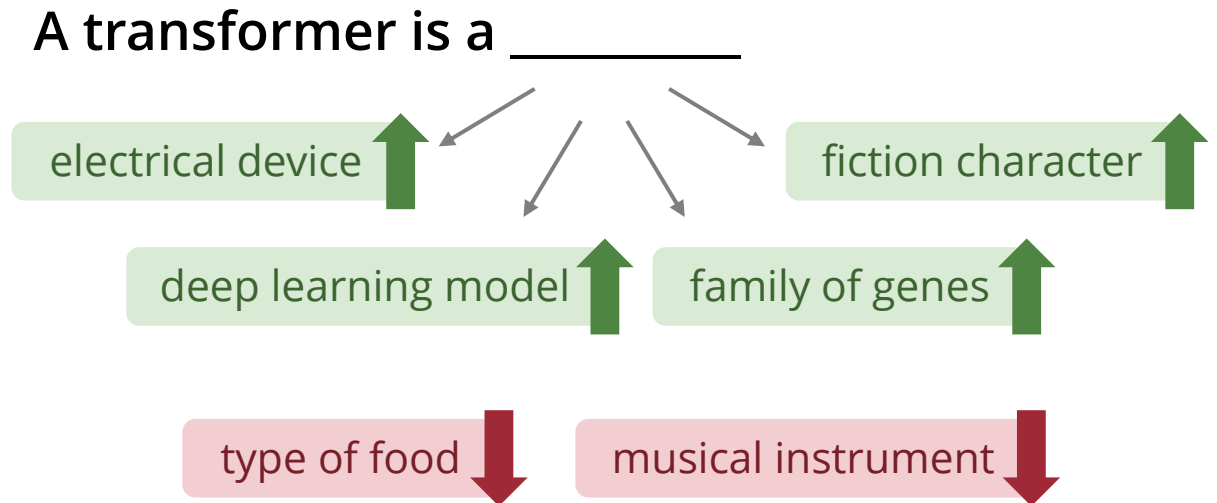
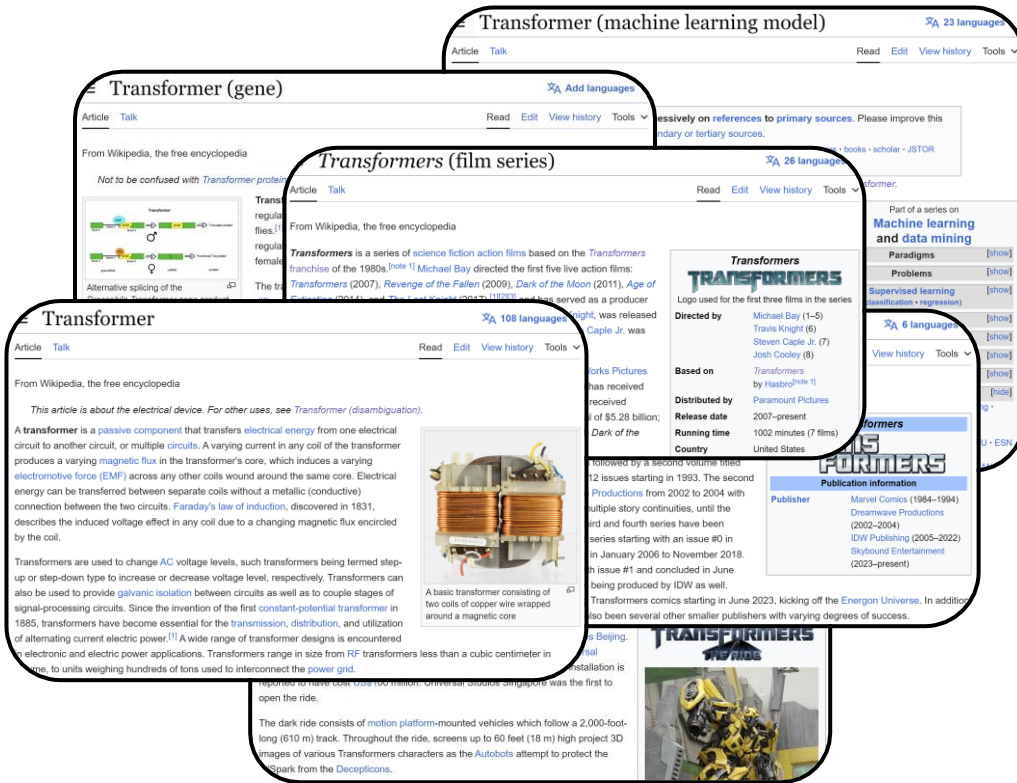
Image-based



Spectrogram

Generating Text using Language Models

- Predicting the next word given the past sequence of words



Generating Text using Language Models

- How do we generate a new sentence with a language model?

A transformer is a

→ Model → deep

A transformer is a deep

→ Model → learning

A transformer is a deep learning

→ Model → model

A transformer is a deep learning model

→ Model → introduced

A transformer is a deep learning model introduced

→ Model → in

A transformer is a deep learning model introduced in

→ Model → 2017



Multitrack Music Transformer

Hao-Wen Dong Ke Chen Shlomo Dubnov Julian McAuley Taylor Berg-Kirkpatrick

University of California San Diego



UC San Diego

Designing a Machine-readable Music Language

- We represent a music piece as a sequence of “**super words**”

$$\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$$

- Each super word \mathbf{x}_i encodes:

$$\mathbf{x}_i = (x_i^{\text{type}}, x_i^{\text{beat}}, x_i^{\text{position}}, x_i^{\text{pitch}}, x_i^{\text{duration}}, x_i^{\text{instrument}})$$

Specify note & instrument information

Structural

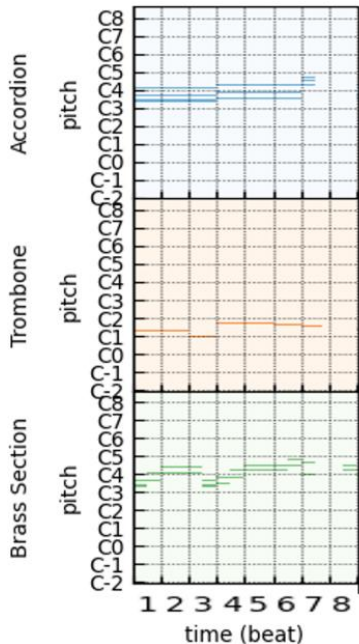
Start of song
Start of notes
End of song

Data

Instrument
Note



An Example of the Proposed Representation



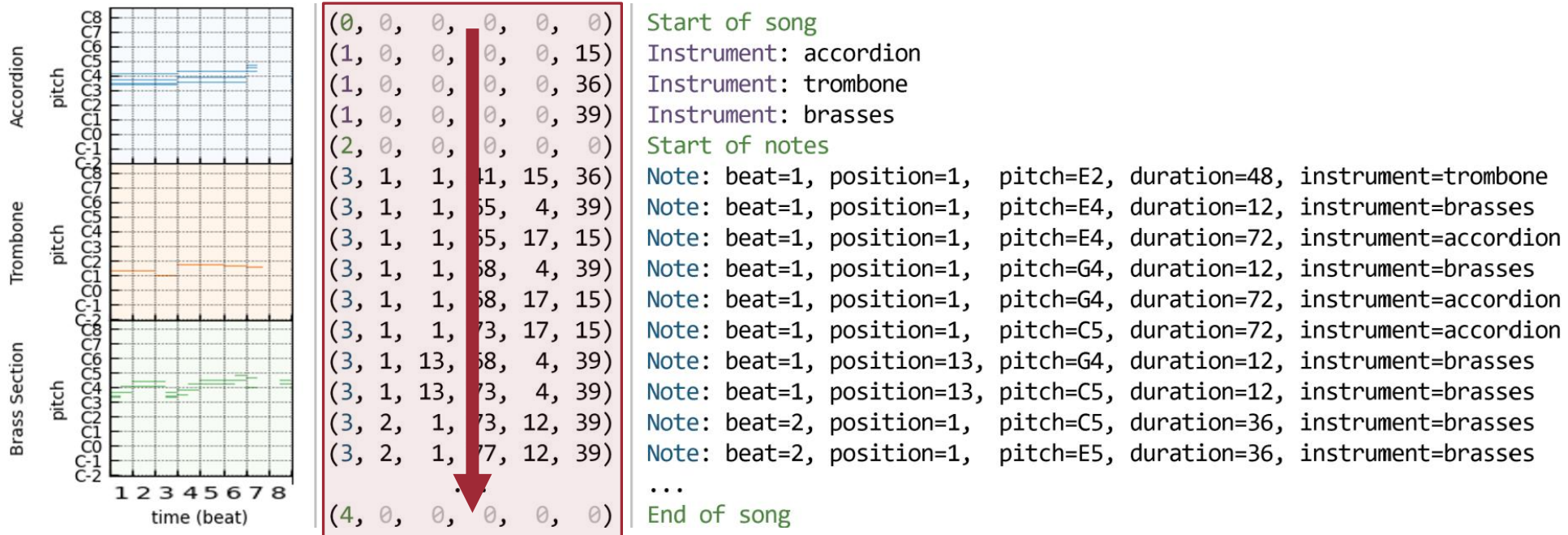
Structural events

(0, 0, 0, 0, 0, 0)	Start of song
(1, 0, 0, 0, 0, 15)	Instrument: accordion
(1, 0, 0, 0, 0, 36)	Instrument: trombone
(1, 0, 0, 0, 0, 39)	Instrument: brasses
(2, 0, 0, 0, 0, 0)	Start of notes
(3, 1, 1, 41, 15, 36)	Note: beat=1, position=1, pitch=E2, duration=48, instrument=trombone
(3, 1, 1, 65, 4, 39)	Note: beat=1, position=1, pitch=E4, duration=12, instrument=brasses
(3, 1, 1, 65, 17, 15)	Note: beat=1, position=1, pitch=E4, duration=72, instrument=accordion
(3, 1, 1, 68, 4, 39)	Note: beat=1, position=1, pitch=G4, duration=12, instrument=brasses
(3, 1, 1, 68, 17, 15)	Note: beat=1, position=1, pitch=G4, duration=72, instrument=accordion
(3, 1, 1, 73, 17, 15)	Note: beat=1, position=1, pitch=C5, duration=72, instrument=accordion
(3, 1, 13, 68, 4, 39)	Note: beat=1, position=13, pitch=G4, duration=12, instrument=brasses
(3, 1, 13, 73, 4, 39)	Note: beat=1, position=13, pitch=C5, duration=12, instrument=brasses
(3, 2, 1, 73, 12, 39)	Note: beat=2, position=1, pitch=C5, duration=36, instrument=brasses
(3, 2, 1, 77, 12, 39)	Note: beat=2, position=1, pitch=E5, duration=36, instrument=brasses
...	...
(4, 0, 0, 0, 0, 0)	End of song

Instrument events

Note events

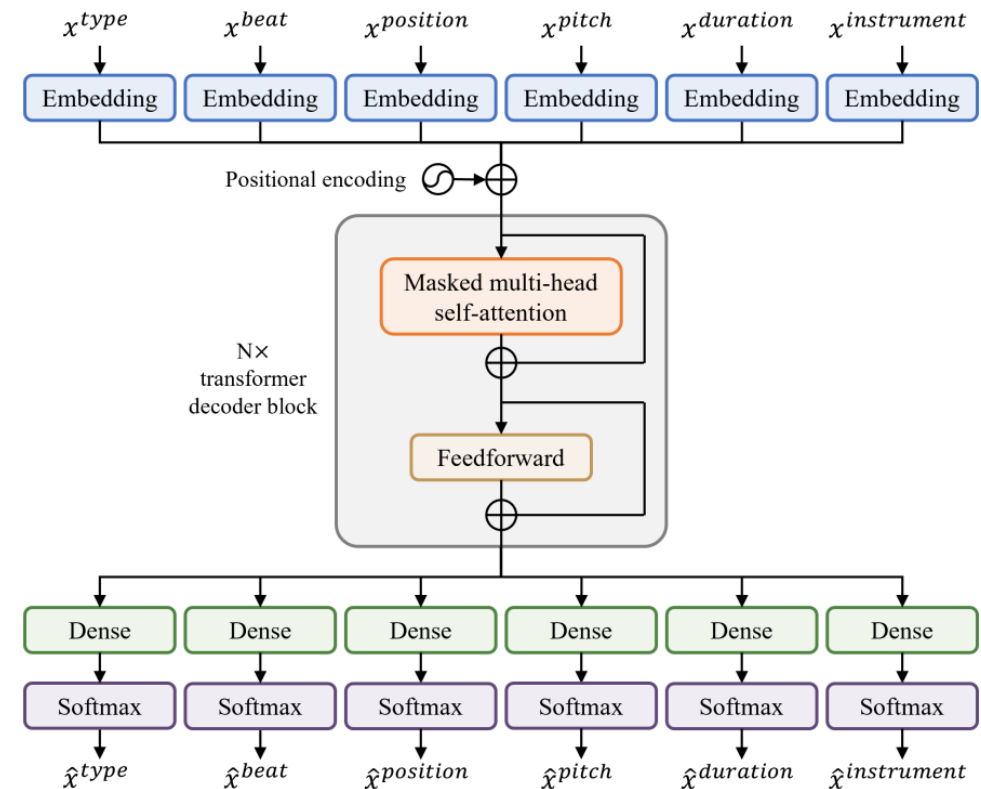
An Example of the Proposed Representation



Multitrack Music Transformer

- A decoder-only transformer model
- Predicts six fields at the same time
- Trained autoregressively

Word-by-word



Symbolic Orchestral Database (SOD)

- 5,743 orchestral pieces (**357 hours** in total)
- Contains various ensembles: choir, string quartet, symphony, etc.



Example Results

Unconditional
generation



Three Sampling Modes

Unconditional generation

Input

(0, 0, 0, 0, 0, 0)	Start of song
(1, 0, 0, 0, 0, 15)	Instrument: accordion
(1, 0, 0, 0, 0, 36)	Instrument: trombone
(1, 0, 0, 0, 0, 39)	Instrument: brasses
(2, 0, 0, 0, 0, 0)	Start of notes
(3, 1, 1, 41, 15, 36)	Note: beat=1, position=1, pitch=E2, duration=48, instrument=trombone
(3, 1, 1, 65, 4, 39)	Note: beat=1, position=1, pitch=E4, duration=12, instrument=brasses
(3, 1, 1, 65, 17, 15)	Note: beat=1, position=1, pitch=E4, duration=72, instrument=accordion
(3, 1, 1, 68, 4, 39)	Note: beat=1, position=1, pitch=G4, duration=12, instrument=brasses
(3, 1, 1, 68, 17, 15)	Note: beat=1, position=1, pitch=G4, duration=72, instrument=accordion
(3, 1, 1, 73, 17, 15)	Note: beat=1, position=1, pitch=C5, duration=72, instrument=accordion
(3, 1, 13, 68, 4, 39)	Note: beat=1, position=13, pitch=G4, duration=12, instrument=brasses
(3, 1, 13, 73, 4, 39)	Note: beat=1, position=13, pitch=C5, duration=12, instrument=brasses
(3, 2, 1, 73, 12, 39)	Note: beat=2, position=1, pitch=C5, duration=36, instrument=brasses
(3, 2, 1, 77, 12, 39)	Note: beat=2, position=1, pitch=E5, duration=36, instrument=brasses
...	...
(4, 0, 0, 0, 0, 0)	End of song

Instrument-informed generation

Input

(0, 0, 0, 0, 0, 0)	Start of song
(1, 0, 0, 0, 0, 15)	Instrument: accordion
(1, 0, 0, 0, 0, 36)	Instrument: trombone
(1, 0, 0, 0, 0, 39)	Instrument: brasses
(2, 0, 0, 0, 0, 0)	Start of notes
(3, 1, 1, 41, 15, 36)	Note: beat=1, position=1, pitch=E2, duration=48, instrument=trombone
(3, 1, 1, 65, 4, 39)	Note: beat=1, position=1, pitch=E4, duration=12, instrument=brasses
(3, 1, 1, 65, 17, 15)	Note: beat=1, position=1, pitch=E4, duration=72, instrument=accordion
(3, 1, 1, 68, 4, 39)	Note: beat=1, position=1, pitch=G4, duration=12, instrument=brasses
(3, 1, 1, 68, 17, 15)	Note: beat=1, position=1, pitch=G4, duration=72, instrument=accordion
(3, 1, 1, 73, 17, 15)	Note: beat=1, position=1, pitch=C5, duration=72, instrument=accordion
(3, 1, 13, 68, 4, 39)	Note: beat=1, position=13, pitch=G4, duration=12, instrument=brasses
(3, 1, 13, 73, 4, 39)	Note: beat=1, position=13, pitch=C5, duration=12, instrument=brasses
(3, 2, 1, 73, 12, 39)	Note: beat=2, position=1, pitch=C5, duration=36, instrument=brasses
(3, 2, 1, 77, 12, 39)	Note: beat=2, position=1, pitch=E5, duration=36, instrument=brasses
...	...
(4, 0, 0, 0, 0, 0)	End of song

N-beat continuation

Input

(0, 0, 0, 0, 0, 0)	Start of song
(1, 0, 0, 0, 0, 15)	Instrument: accordion
(1, 0, 0, 0, 0, 36)	Instrument: trombone
(1, 0, 0, 0, 0, 39)	Instrument: brasses
(2, 0, 0, 0, 0, 0)	Start of notes
(3, 1, 1, 41, 15, 36)	Note: beat=1, position=1, pitch=E2, duration=48, instrument=trombone
(3, 1, 1, 65, 4, 39)	Note: beat=1, position=1, pitch=E4, duration=12, instrument=brasses
(3, 1, 1, 65, 17, 15)	Note: beat=1, position=1, pitch=E4, duration=72, instrument=accordion
(3, 1, 1, 68, 4, 39)	Note: beat=1, position=1, pitch=G4, duration=12, instrument=brasses
(3, 1, 1, 68, 17, 15)	Note: beat=1, position=1, pitch=G4, duration=72, instrument=accordion
(3, 1, 1, 73, 17, 15)	Note: beat=1, position=1, pitch=C5, duration=72, instrument=accordion
(3, 1, 13, 68, 4, 39)	Note: beat=1, position=13, pitch=G4, duration=12, instrument=brasses
(3, 1, 13, 73, 4, 39)	Note: beat=1, position=13, pitch=C5, duration=12, instrument=brasses
(3, 2, 1, 73, 12, 39)	Note: beat=2, position=1, pitch=C5, duration=36, instrument=brasses
(3, 2, 1, 77, 12, 39)	Note: beat=2, position=1, pitch=E5, duration=36, instrument=brasses
...	...
(4, 0, 0, 0, 0, 0)	End of song

Only needs to train ONE model!

Example Results

Unconditional generation



Instrument-informed generation



church-organ, viola,
contrabass, strings,
voices, horn, oboe

4-beat continuation

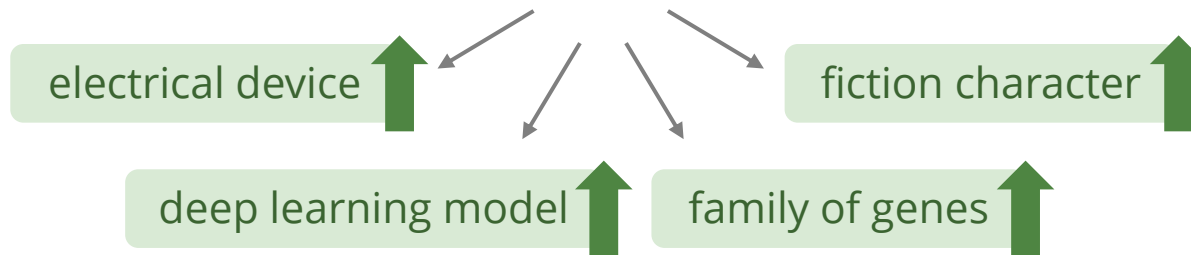


Mozart's
Eine kleine Nachtmusik

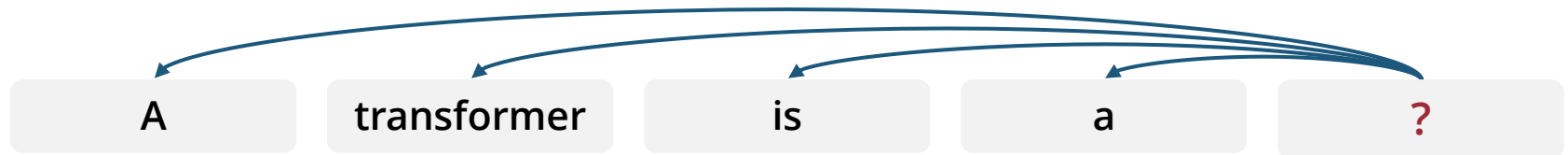


The Magic of Transformers – Self-attention Mechanism

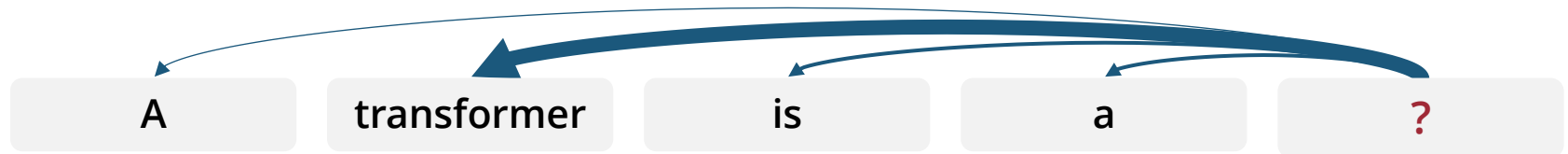
A transformer is a _____



Uniform attention



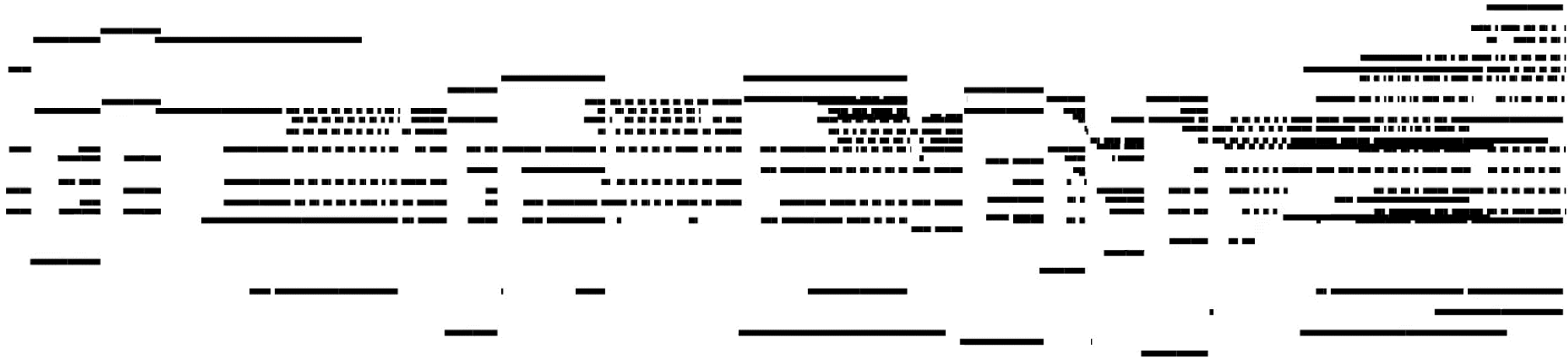
Variable attention



Transformers learn what to attend to from big data!

Visualizing Musical Self-attention (Huang et al., 2018)

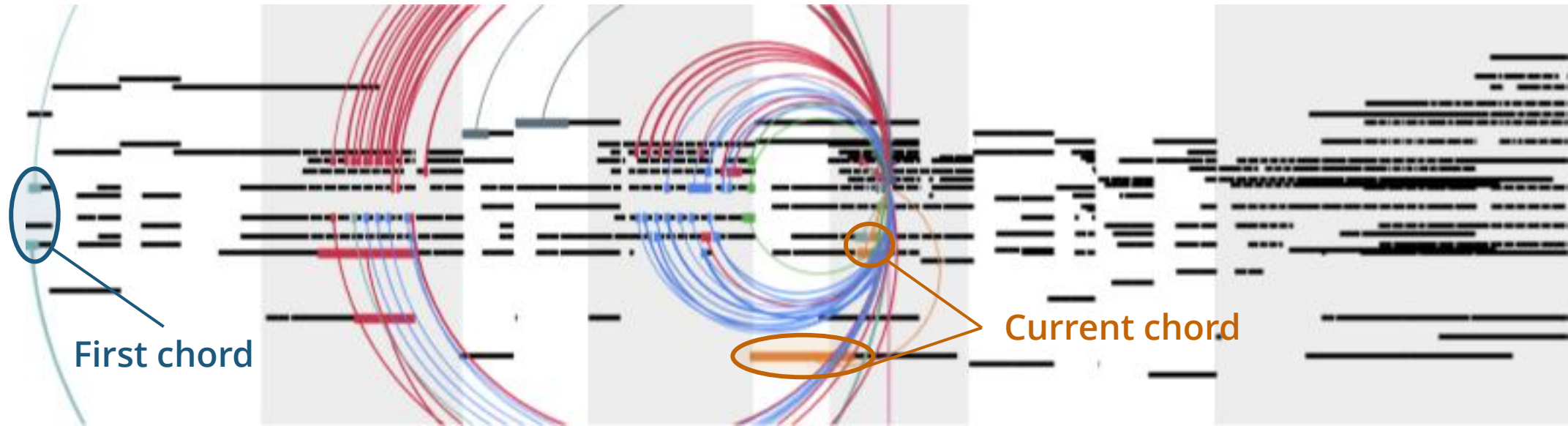
(Each color represents an attention head)



(Source: Huang et al., 2018)

Visualizing Musical Self-attention (Huang et al., 2018)

(Each color represents an attention head)



(Source: Huang et al., 2018)

Can we go beyond case studies?

Systematically Analyzing Musical Self-attention

- We proposed two new quantities for measuring **mean relative attention**

$$\gamma_k^{(d)} = \frac{\sum_{\mathbf{x} \in \mathcal{D}} \sum_{s > t} a_{s,t}(\mathbf{x}) \mathbb{1}_{x_t^{(d)} - x_s^{(d)} = k}}{\sum_{\mathbf{x} \in \mathcal{D}} \sum_{s > t} a_{s,t}(\mathbf{x})}$$

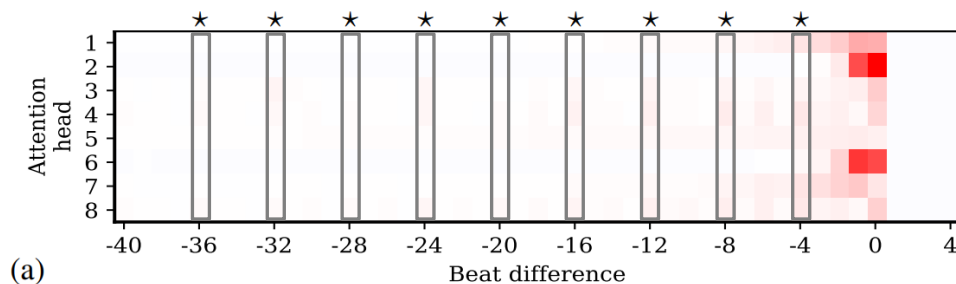
$$\tilde{\gamma}_k^{(d)} = \gamma_k^{(d)} - \frac{\sum_{\mathbf{x} \in \mathcal{D}} \sum_{s > t} \mathbb{1}_{x_t^{(d)} - x_s^{(d)} = k}}{\sum_{\mathbf{x} \in \mathcal{D}} \sum_{s > t} 1}$$

- The MMT model attends more to notes

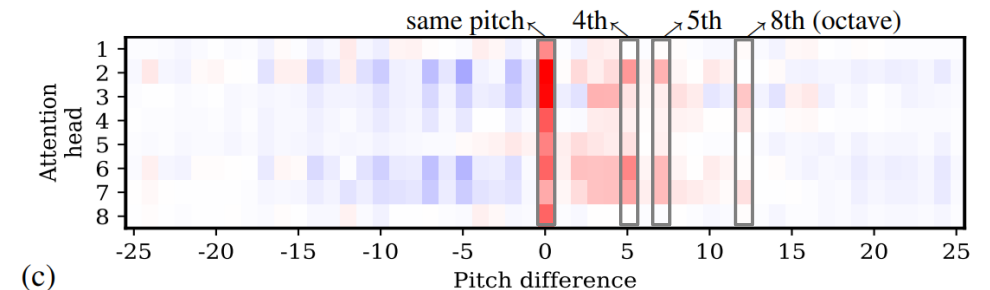
that are $4N$ beats away in the past

that has a pitch in an octave above which forms a consonant interval

Positive and negative mean relative attention gain



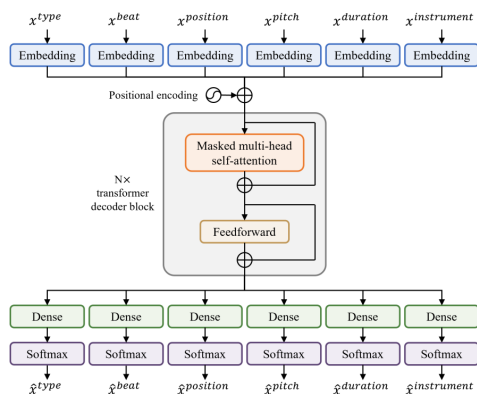
Positive and negative mean relative attention gain



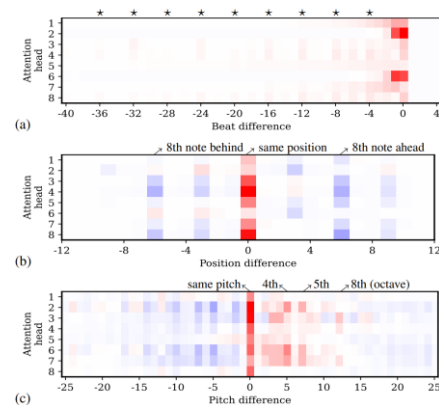
Contributions

- **State-of-the-art orchestral music generation model**
- Presented the **first systematic analysis of musical self-attention**
- Showed that MMT learns a **relative self-attention for beat and pitch**

Multitrack Music Transformer



Musical Self-attention



Paper: arxiv.org/abs/2207.06983
Demo: salu133445.github.io/mmt/
Code: github.com/salu133445/mmt



UC San Diego

Part 2: Generating Audio

Types of Audio

Speech



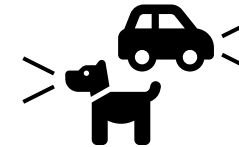
(Source: Wikimedia Commons)

Music



(Source: Wikimedia Commons)

Sound effects



(Source: Wikimedia Commons)

BPJ Media Inc, [CC BY-SA 3.0](#), via Wikimedia Commons.
Vancouver Film School Retouched version by User:Quenhitrn., [CC BY 2.0](#), via Wikimedia Commons.
The Blackbird Academy, [CC BY-SA 2.0](#), via Wikimedia Commons.
One Man Films, ["One Shot - WAR ACTION SHORT FILM," YouTube](#), September 11, 2022.

Music Generation – Four Paradigms



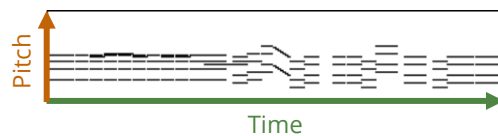
Symbolic music generation

Text-based

```
Program_change_0,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_76, Time_shift_2, Note_off_67,  
Note_on_67, Time_shift_2, Note_off_67,  
...
```

MIDI

Image-based



Piano roll



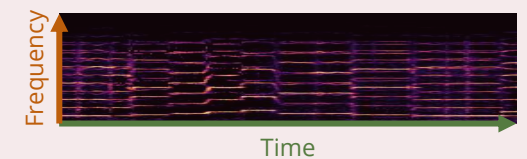
Audio-domain music generation

Time series-based



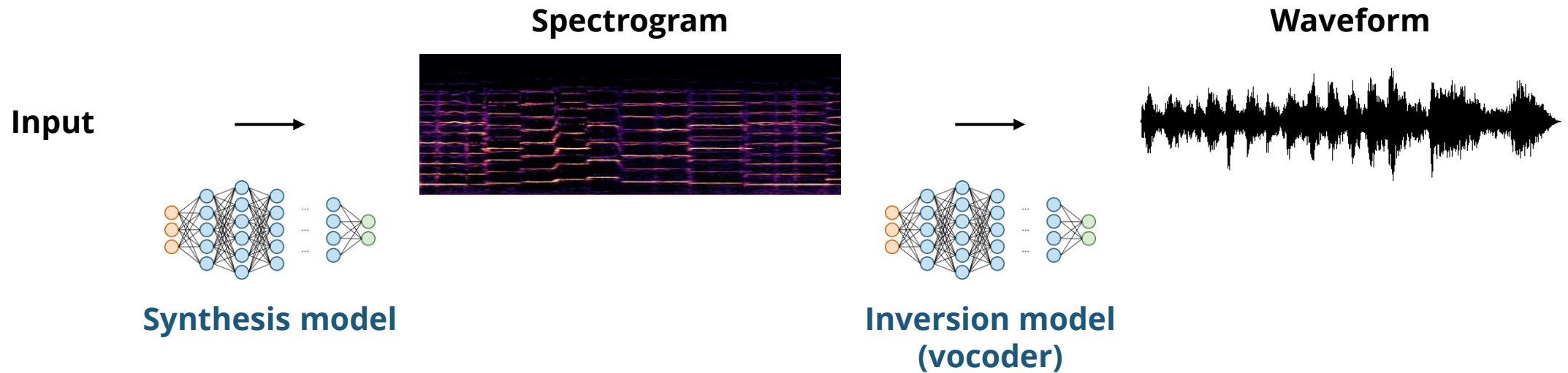
Waveform

Image-based



Spectrogram

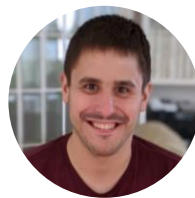
Frequency-domain Audio Synthesis



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

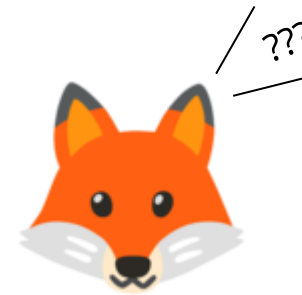
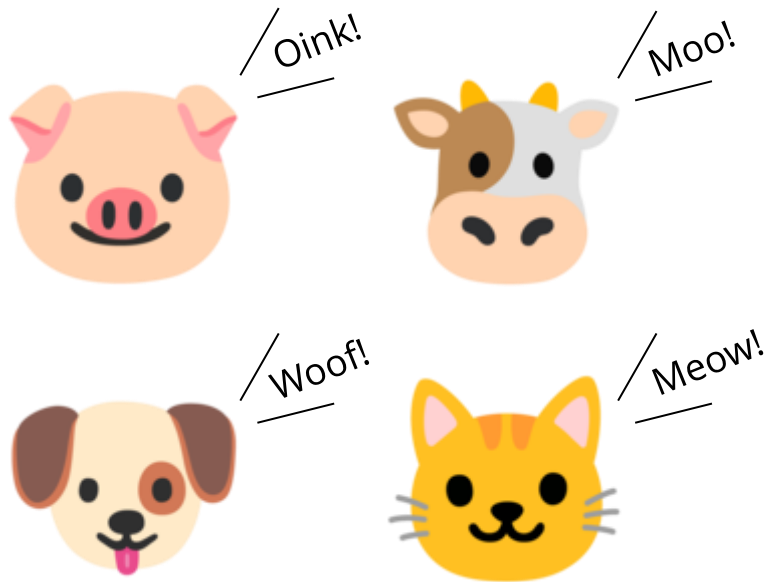
Hao-Wen Dong^{1,2*} Xiaoyu Liu¹ Jordi Pons¹ Gautam Bhattacharya¹
Santiago Pascual¹ Joan Serrà¹ Taylor Berg-Kirkpatrick² Julian McAuley²

¹ Dolby Laboratories ² University of California San Diego



Learning Sounds from Observations

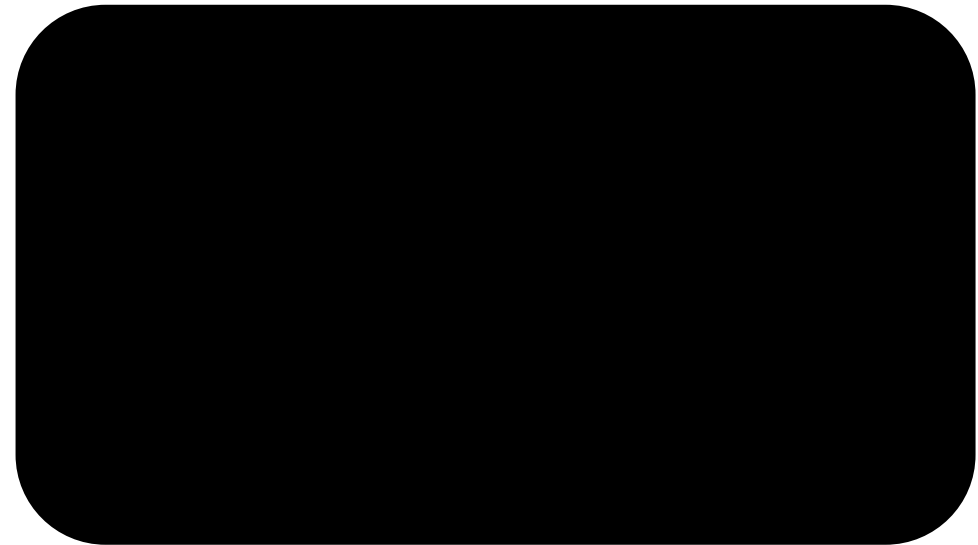
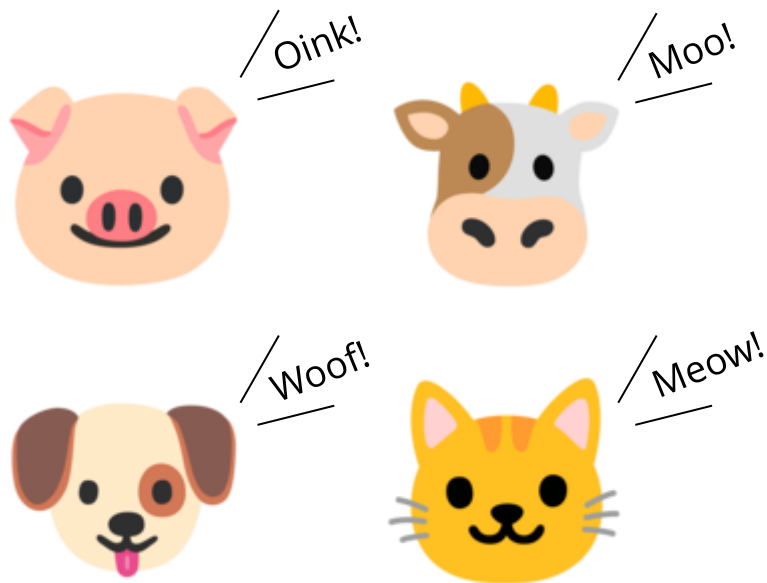
- Watching a dog barking, humans can *associate the barking sound to the dog*



What does the fox say?

Learning Sounds from **Noisy Videos**

- Watching a dog barking, humans can *associate the barking sound to the dog*



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

Data

VGGSound

(Chen et al., 2020)



Hedge trimmer
running



Dog bow-wow



Bird chirping,
tweeting

Noisy videos with diverse sounds

(172K videos, 310 classes)

MUSIC

(Zhao et al., 2018)



Violin



Acoustic guitar



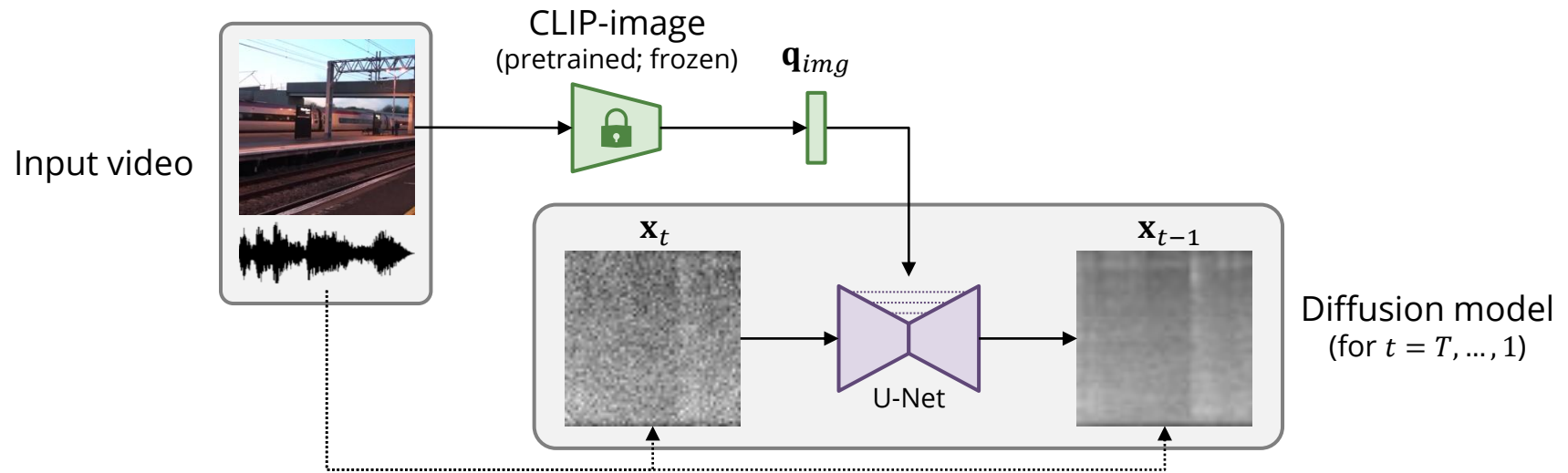
Accordion

Music instrument playing videos

(1,055 videos, 21 instruments)

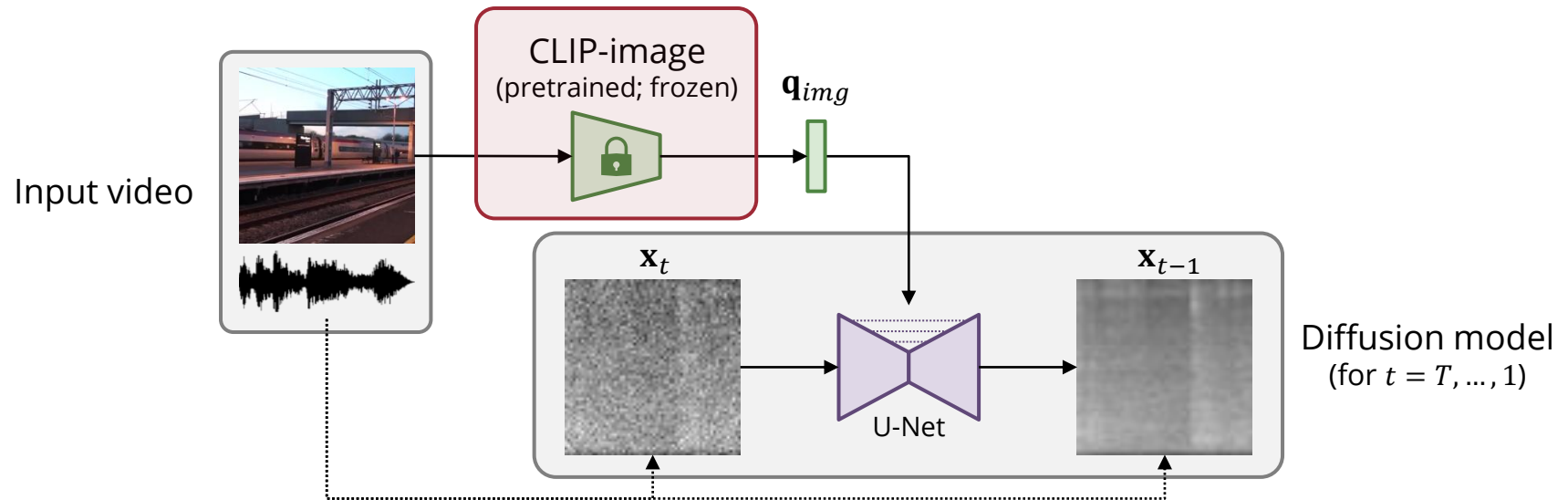
Training an Image-to-Audio Synthesis Model

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



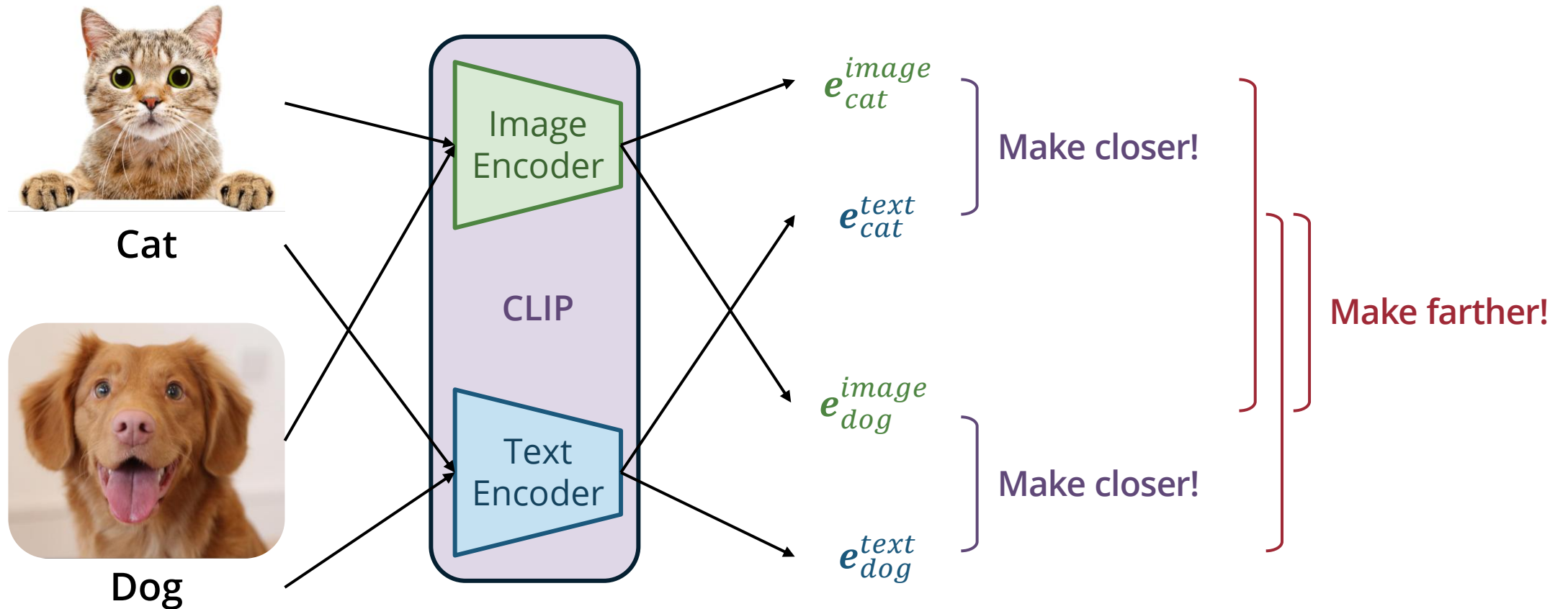
Training an Image-to-Audio Synthesis Model

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



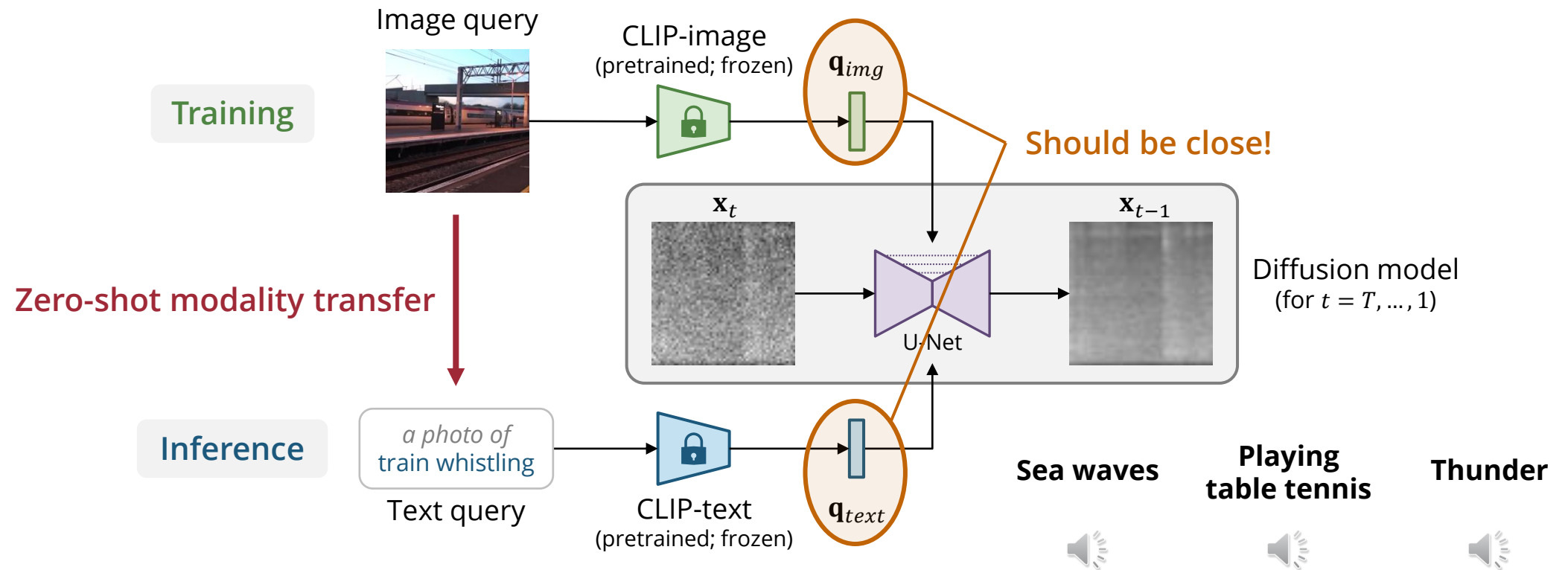
CLIP (Contrastive Language-Image Pretraining)

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



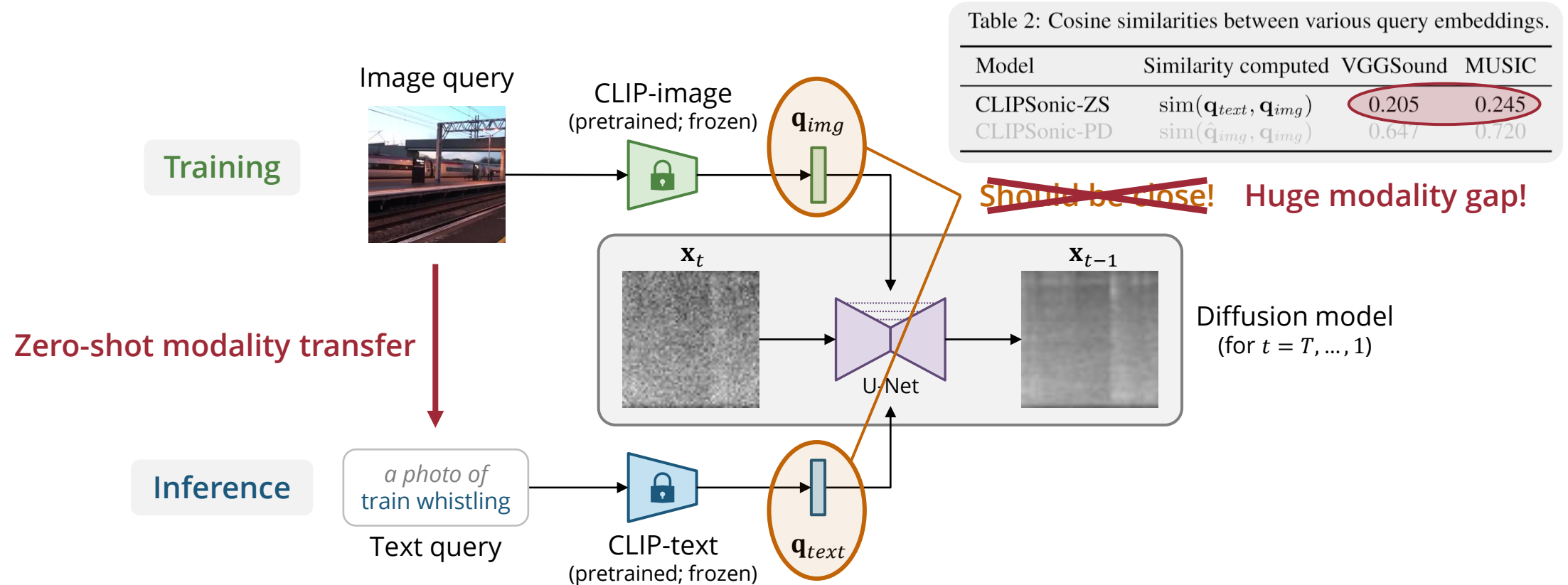
Inference – Zero-shot Modality Transfer

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



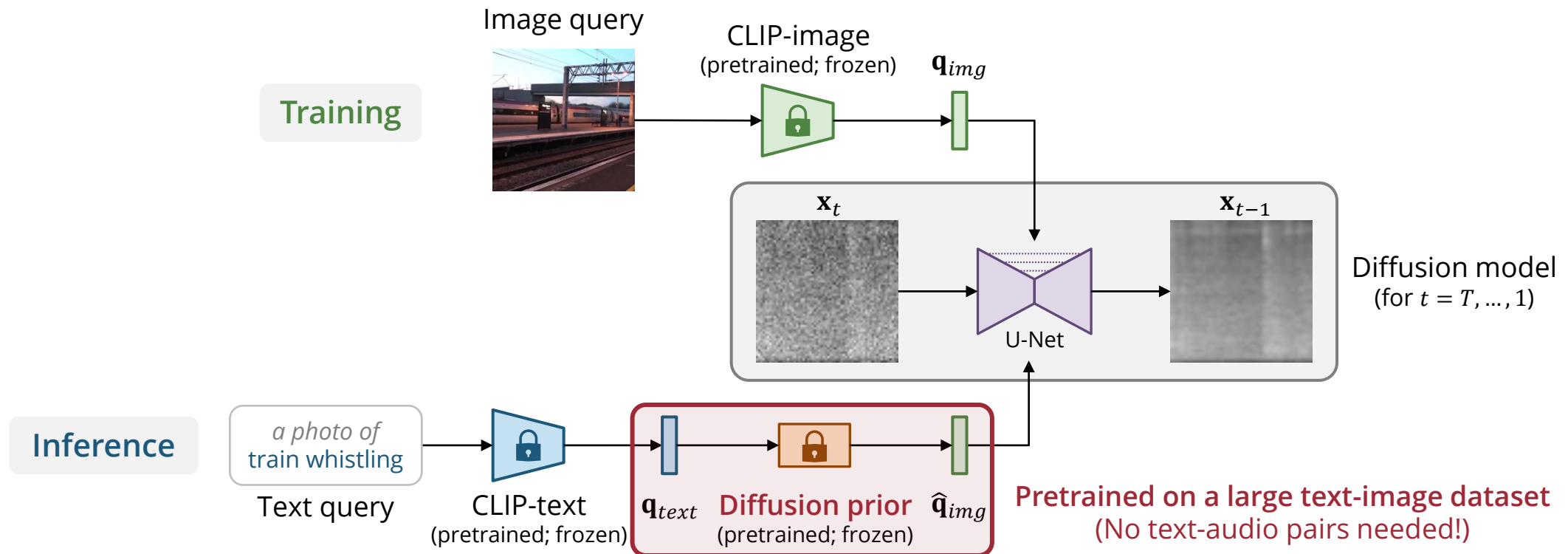
Inference – Zero-shot Modality Transfer

- 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



Example Text-to-Audio Synthesis Results

Rapping



Sea waves



Thunder



Smoke detector beeping



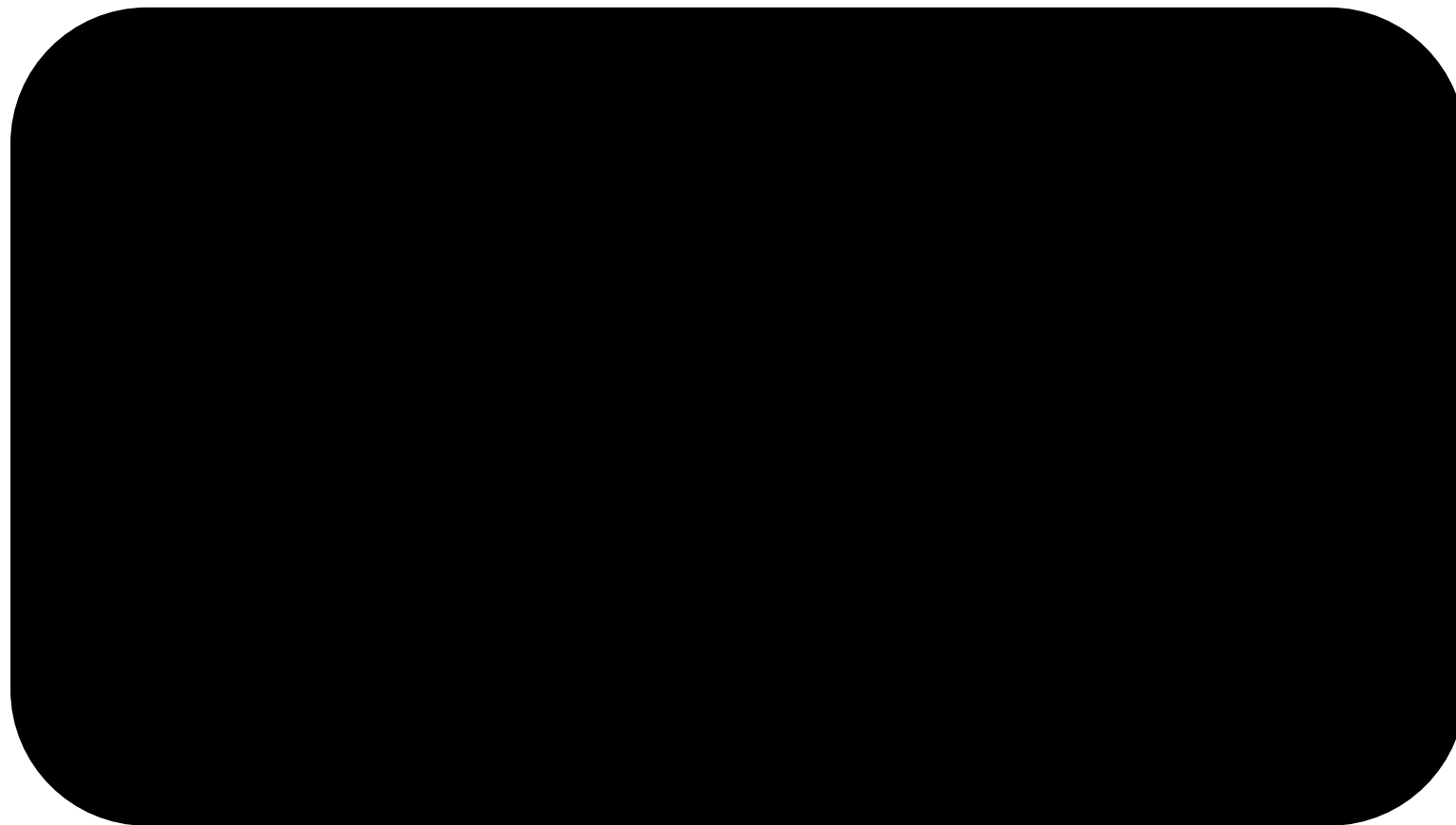
Playing table tennis



Playing violin fiddle



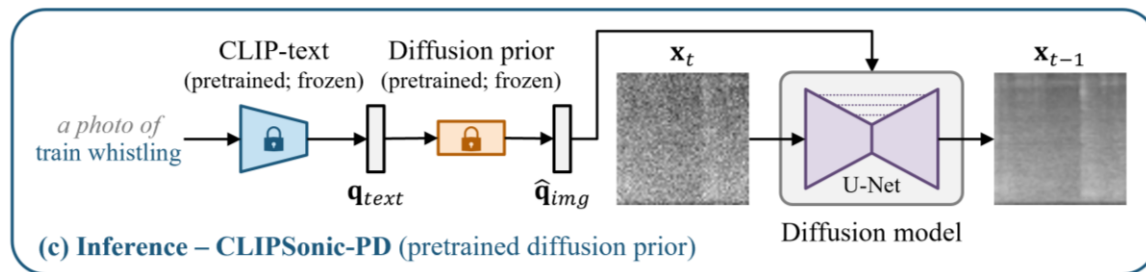
Example Image-to-Audio Synthesis Results (Out-of-distribution)



(Then!) State-of-the-art image-to-audio synthesis performance!

Contributions

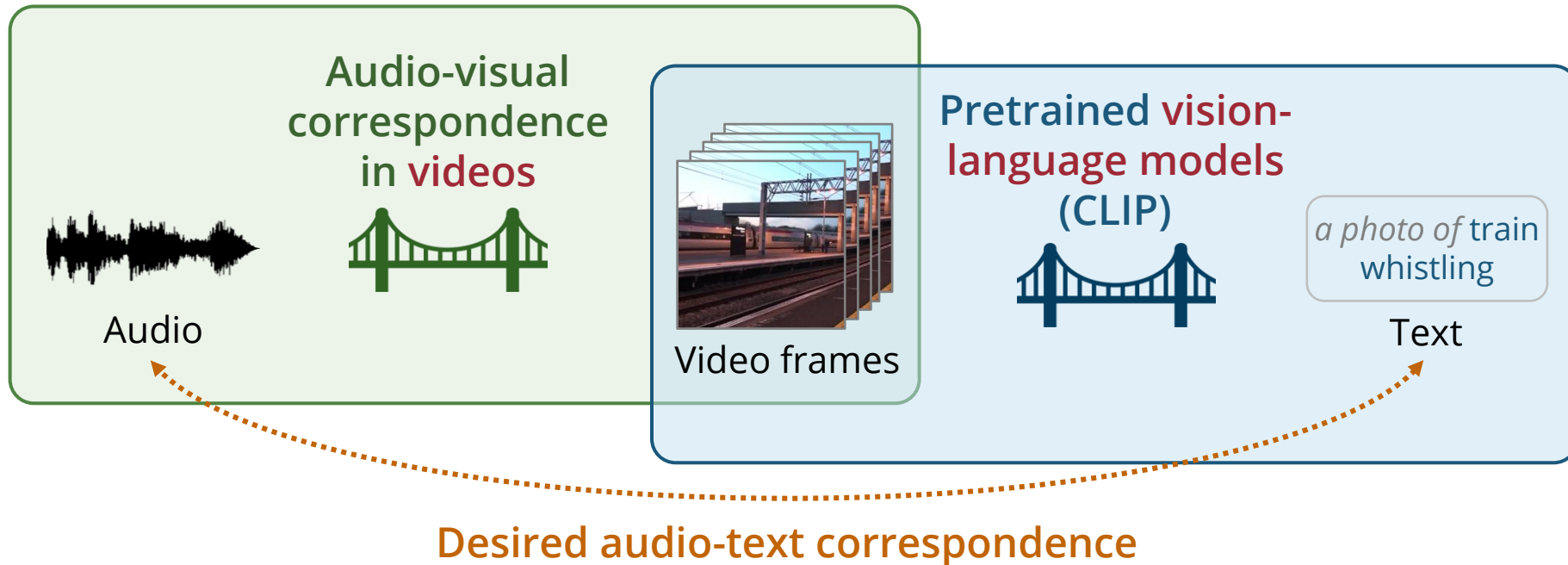
- First text-to-audio synthesis model that **requires *no* text-audio pairs**
- **Strong text-to-audio synthesis performance** without text-audio data
- (Then!) **State-of-the-art image-to-audio synthesis performance**



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



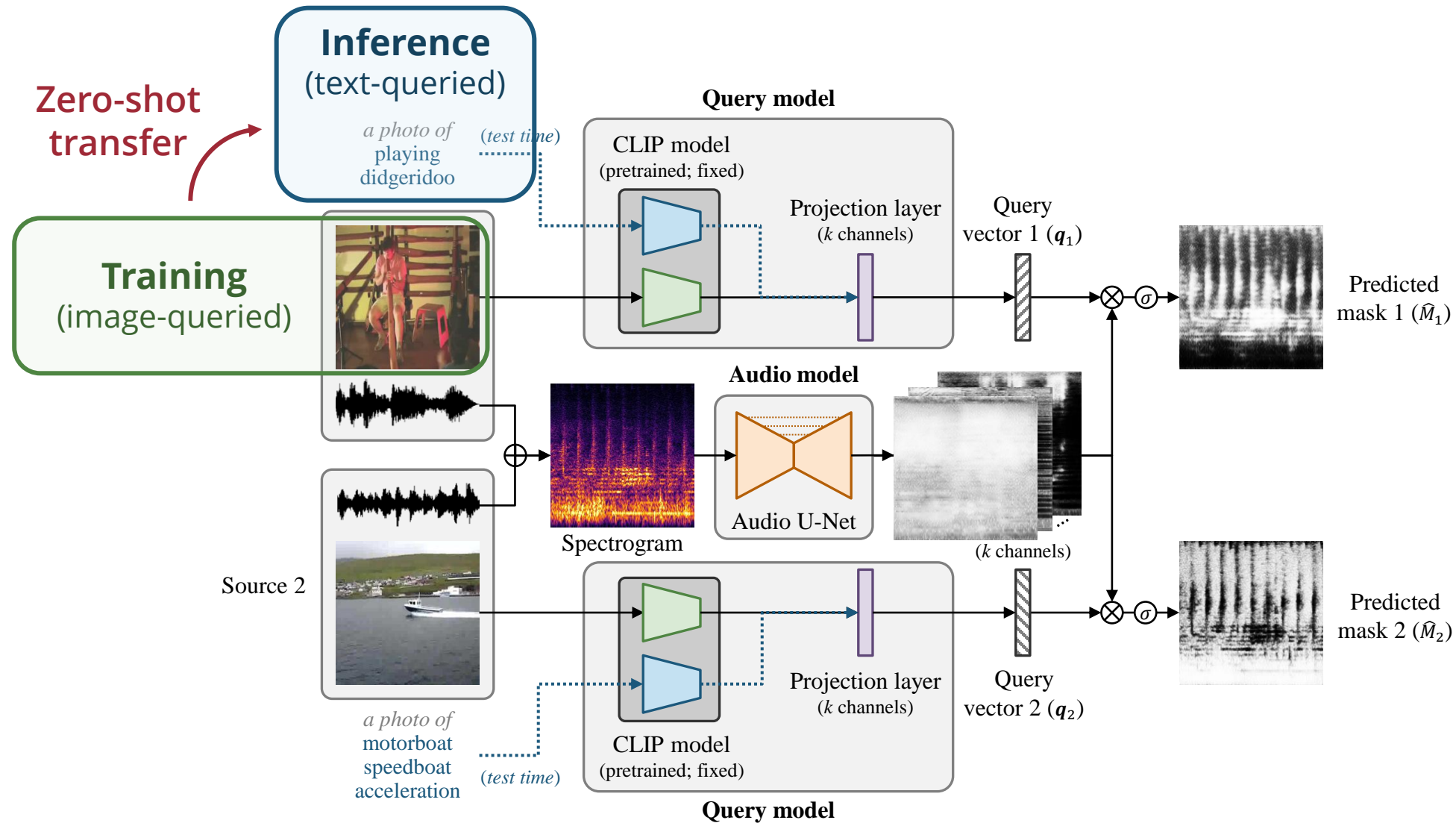
Leveraging the Visual Domain as a Bridge



No text-audio pairs required!

Scalable to large video datasets!

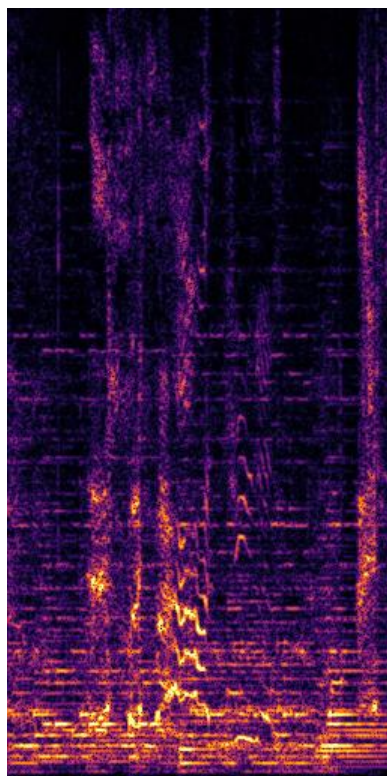
CLIPSep: Text-queried Sound Separation



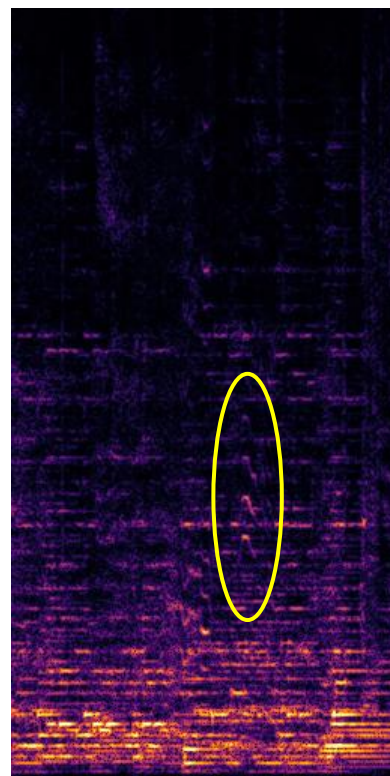
CLIPSep: Text-queried Sound Separation

Query: *"playing harpsichord"*

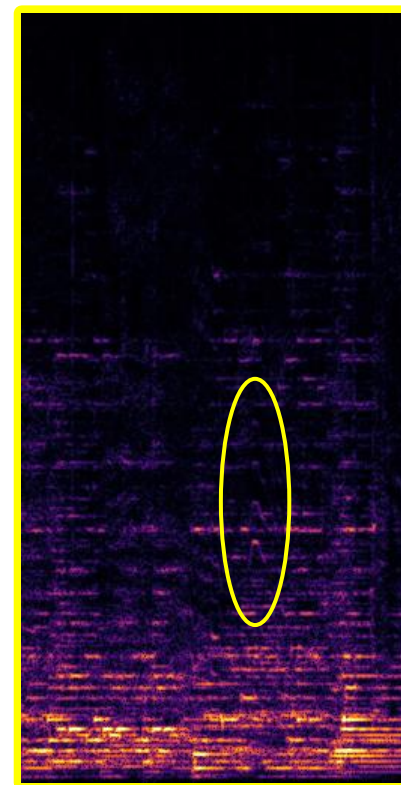
Mixture



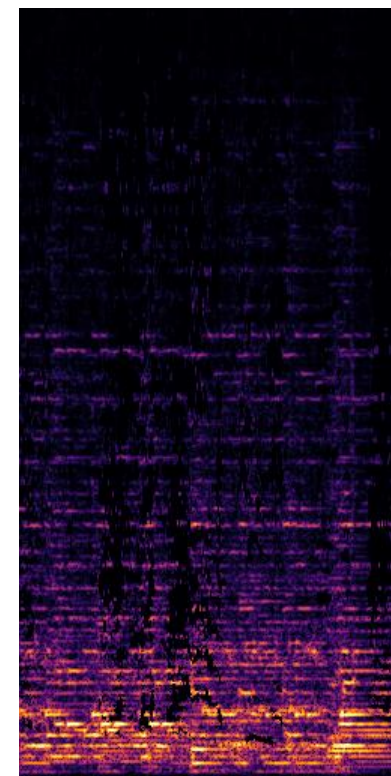
CLIPSep



CLIPSep-NIT



Ground truth



CLIPSep: Noise Removal

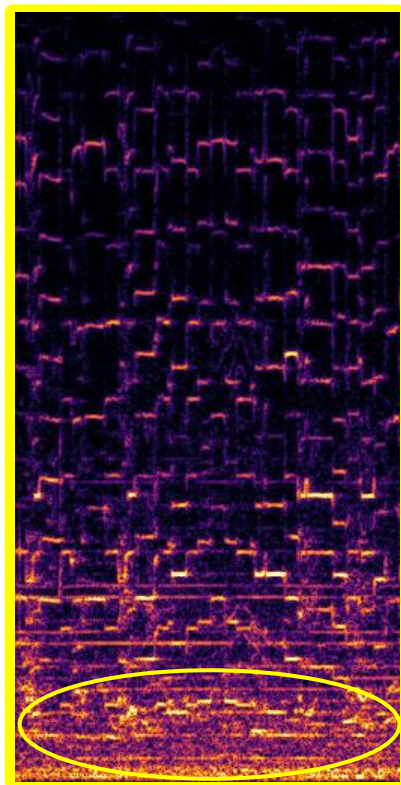


SONY

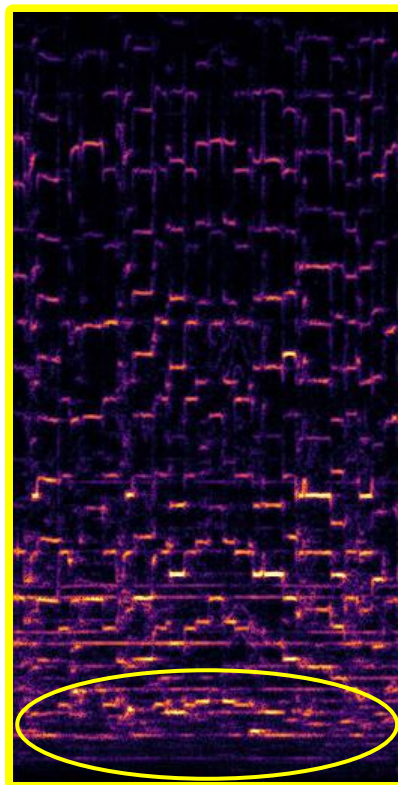
UC San Diego

Query: *"playing bagpipe"*

Mixture



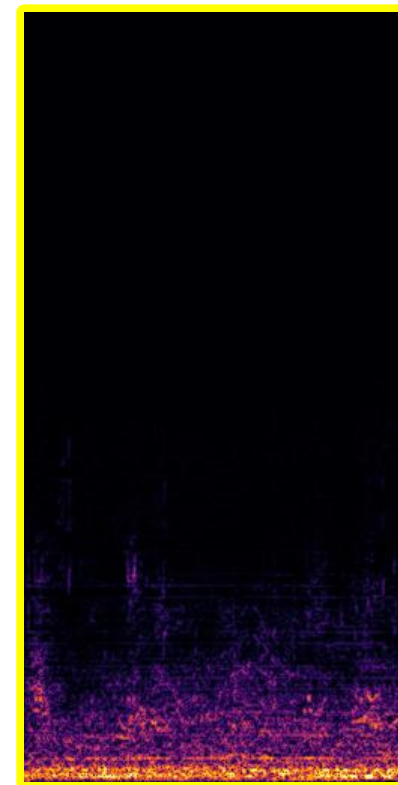
Prediction



Noise head 1



Noise head 2



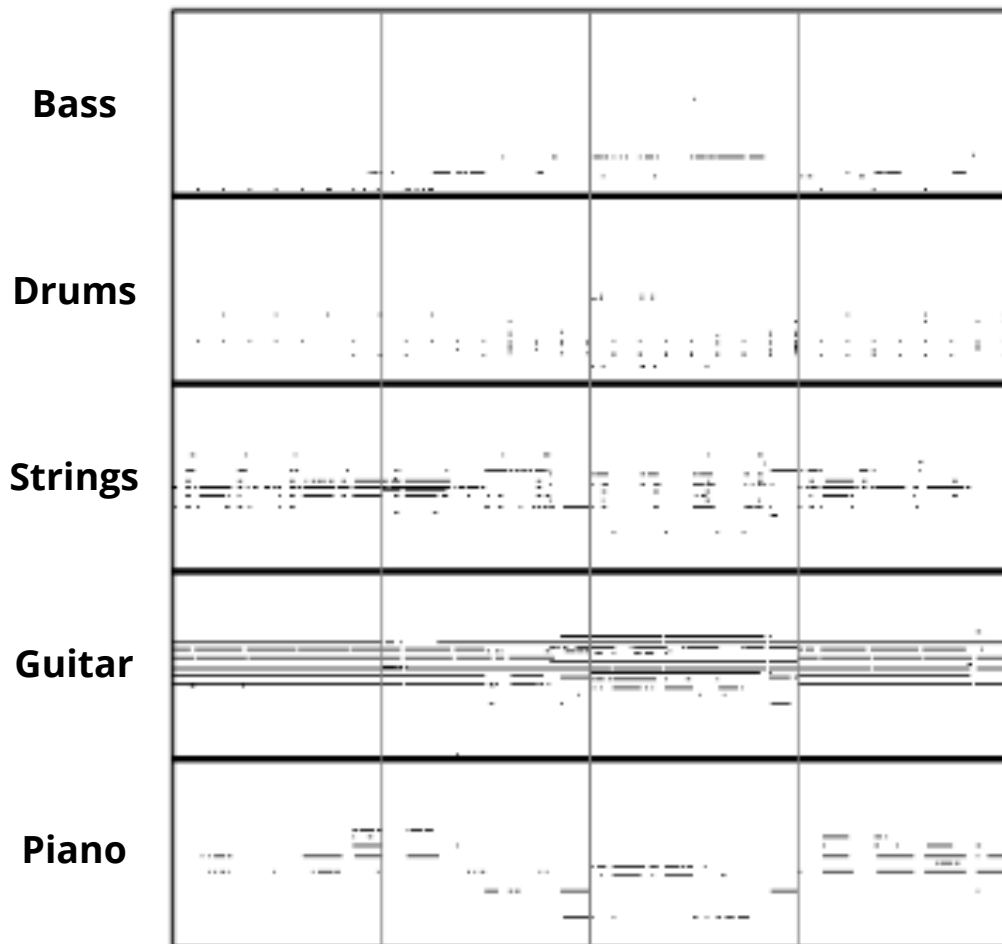
Human-Centered Generative AI for Content Creation

Augmenting human creativity with machine learning

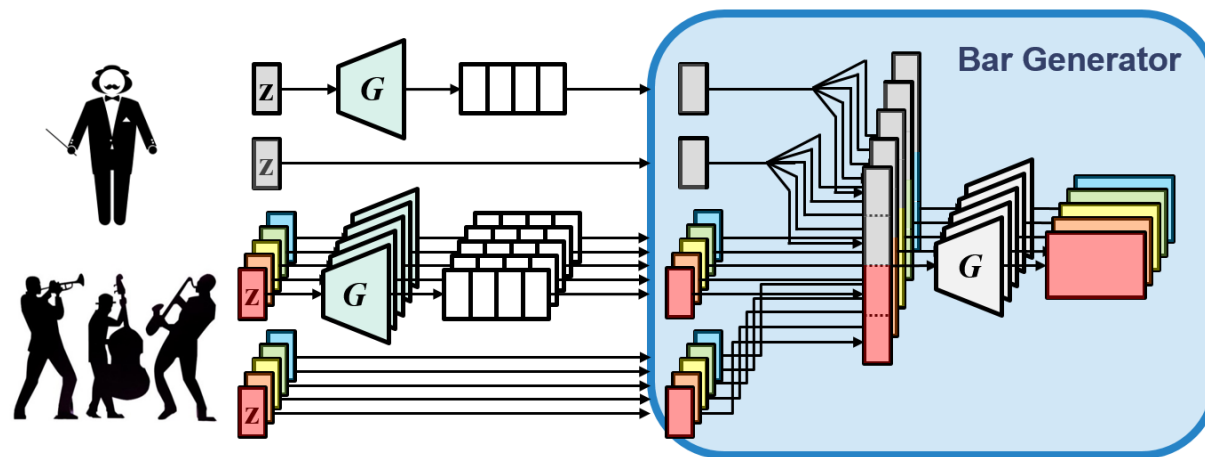
- **Novel Generative Models for New Domains**
 - **Multitrack music generation** (AAAI 2018, ISMIR 2018, ISMIR 2020, ICASSP 2023, ISMIR 2024), **controllable music generation** (AIMG 2024, arXiv 2024), **documentary teaser generation** (arXiv 2024)
- **AI-Assisted Tools for Content Creation**
 - **Violin performance synthesis** (ICASSP 2022, arXiv 2024), **music instrumentation** (ISMIR 2021), **music arrangement** (AAAI 2018), **music harmonization** (JNMR 2020)
- **Multimodal Generative Models for Content Creation**
 - **Queried sound separation** (ICLR 2023), **text-to-audio synthesis** (WSS 2023, WASPAA 2023), **text-to-music generation** (ISMIR LBD 2024, arXiv 2024), **documentary teaser generation** (arXiv 2024)

Generating Multi-instrument Music using GANs (AAAI 2018)

Multitrack Piano Roll



MuseGAN Generator



MuseGAN Features in AWS DeepComposer (2020)

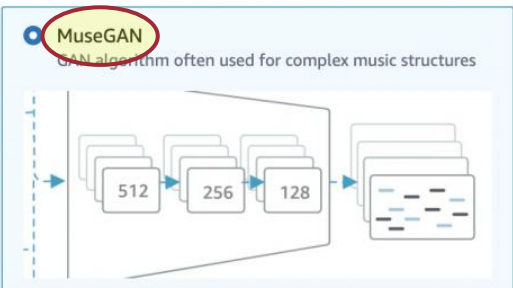
AWS DeepComposer > Models > Train a model

Train a model

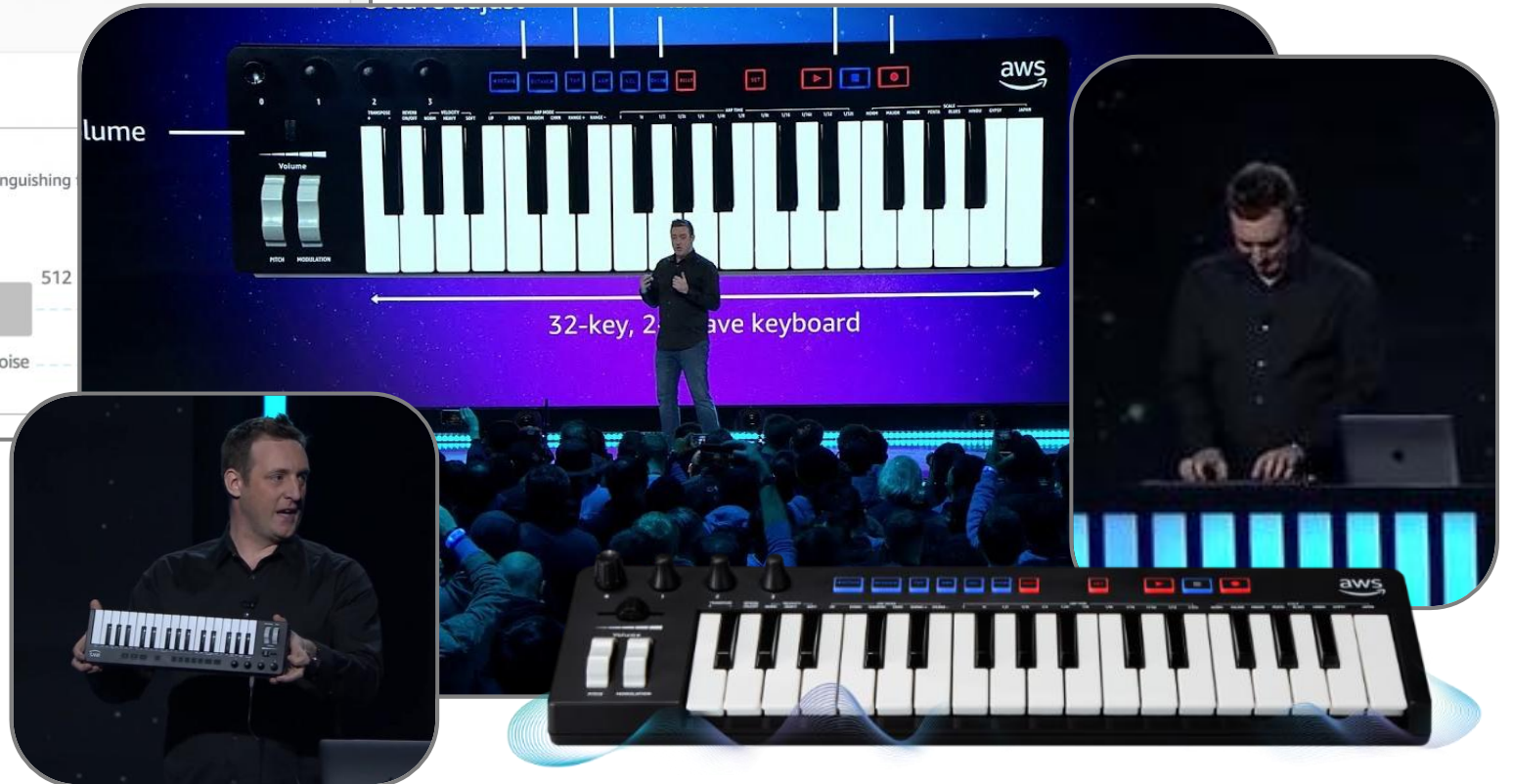
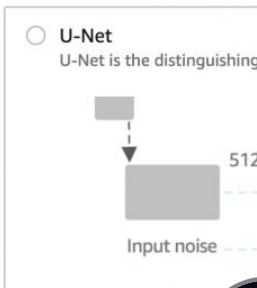
Generative algorithm [Info](#)

Choose a generative algorithm to train a model

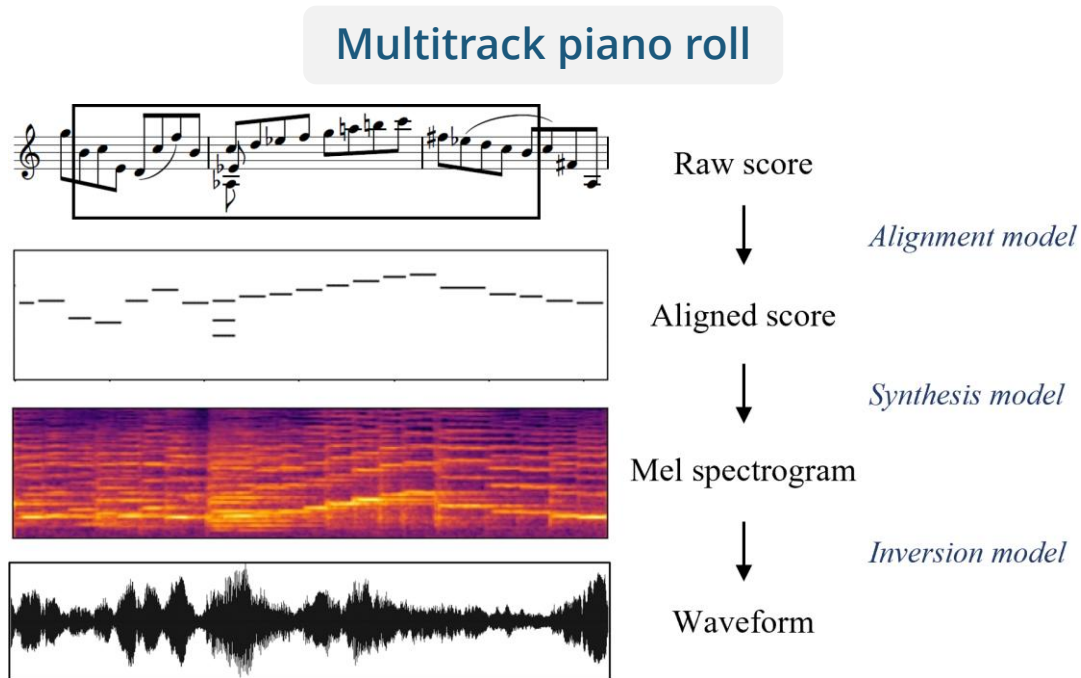
MuseGAN
GAN algorithm often used for complex music structures



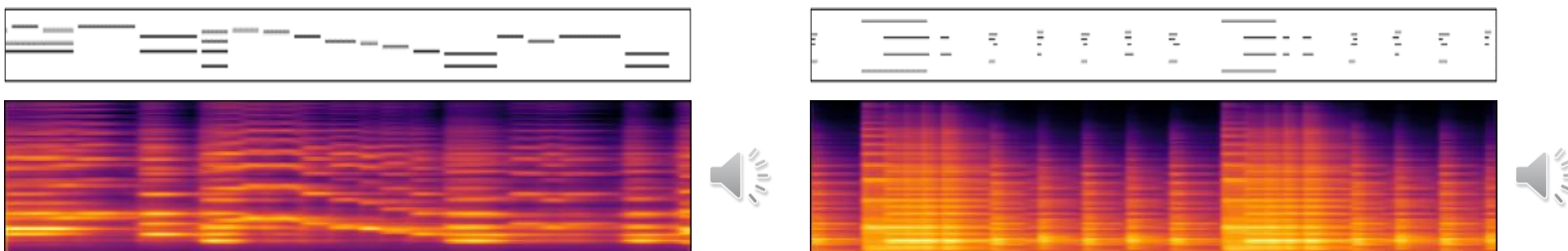
U-Net
U-Net is the distinguishing



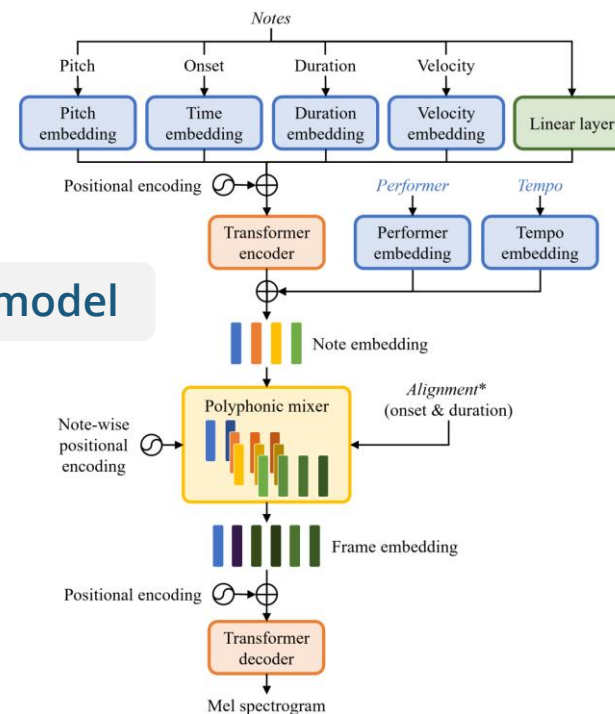
Synthesizing Natural Music Performance (ICASSP 2022)



Example results



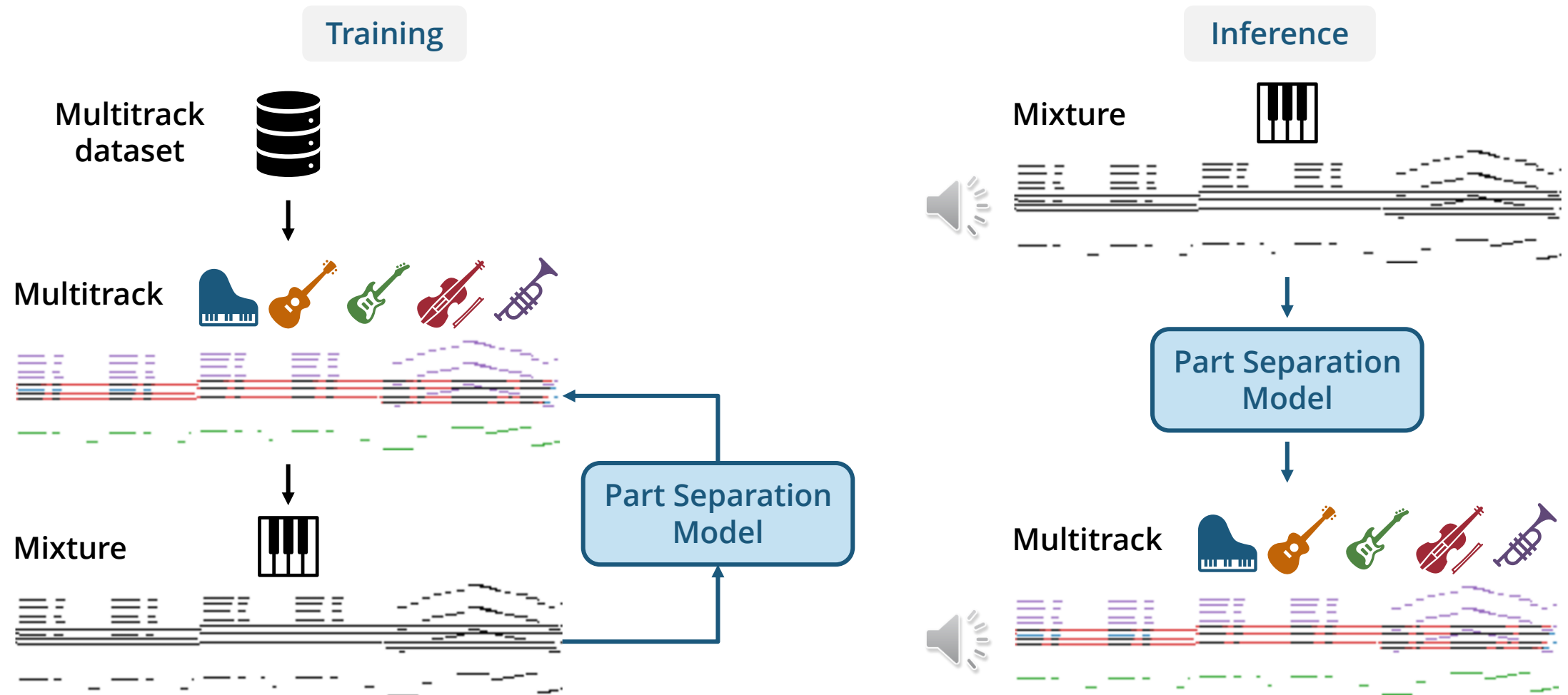
TTS-based model



Automatic Instrumentation (ISMIR 2021)



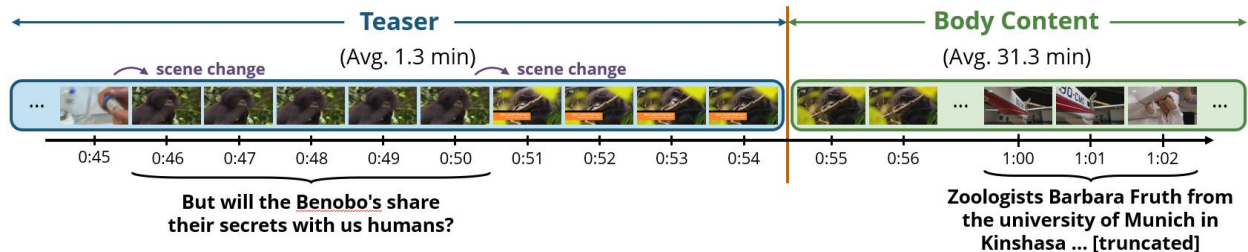
Stanford UC San Diego



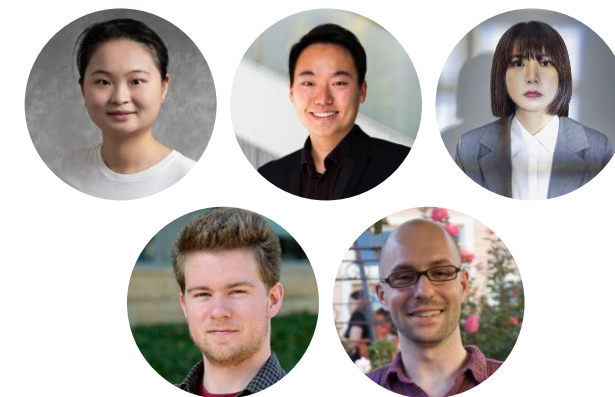
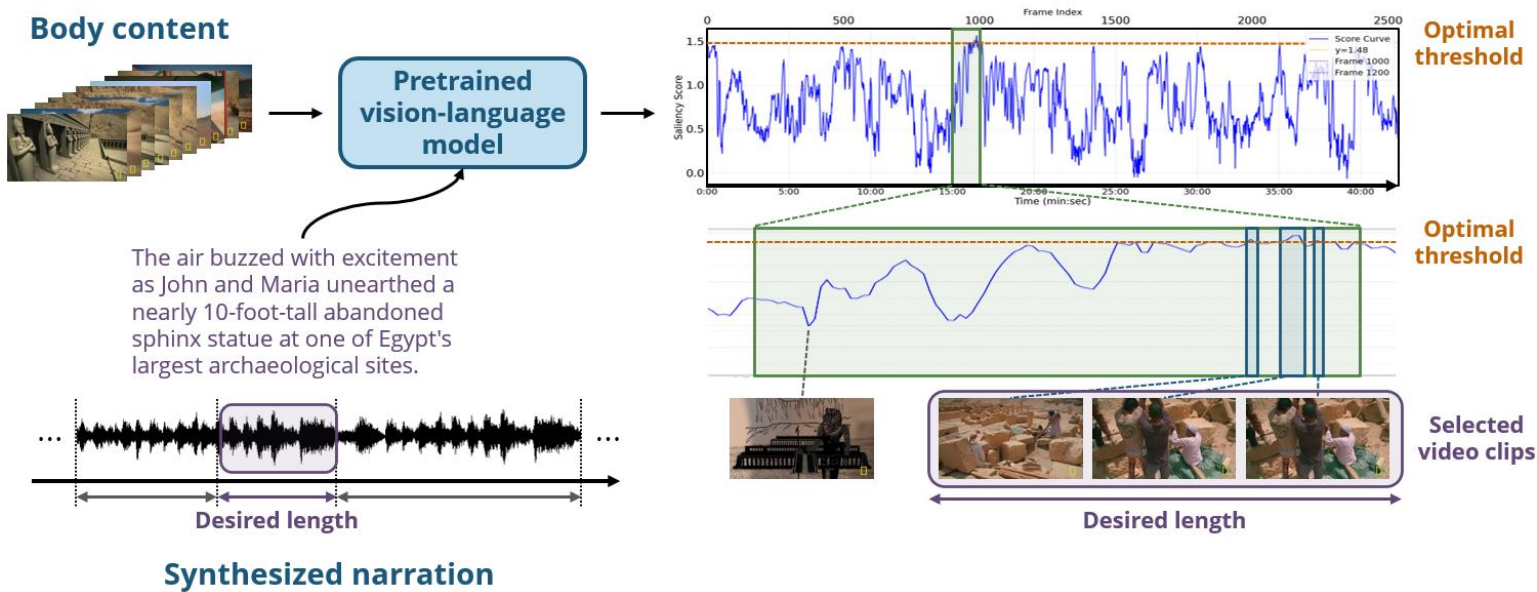
Teaser Generation (arXiv 2024)



Data



Narration-video matching model



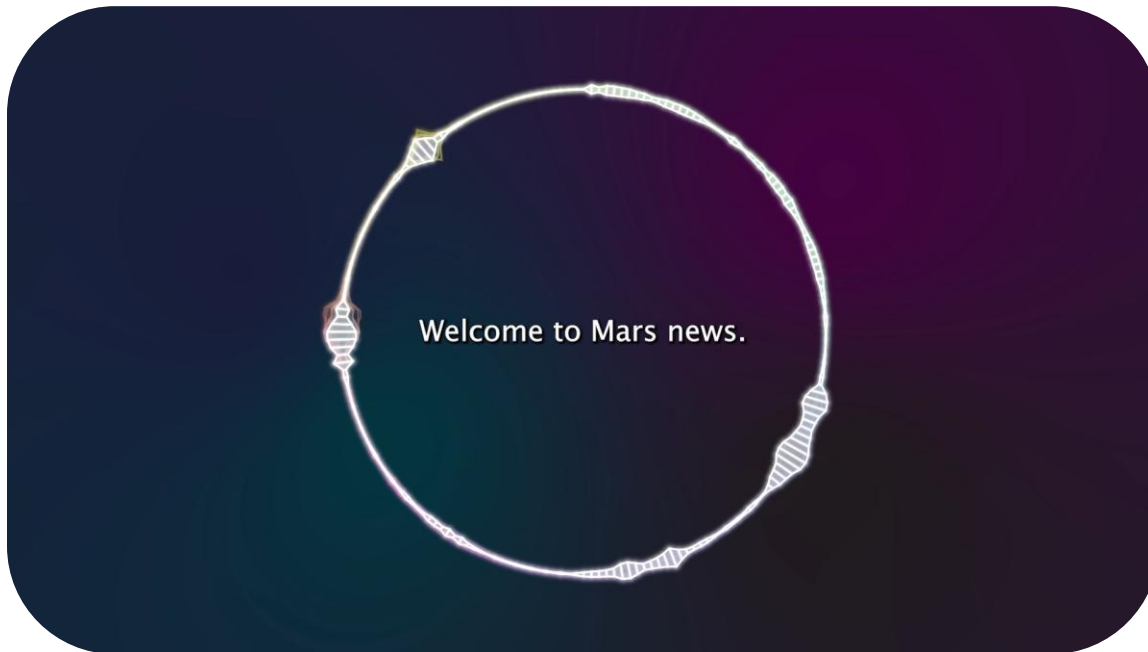
Future Directions

Human-Centered Generative AI for Content Creation

Augmenting human creativity with machine learning

- **Multimodal generative AI** for content creation
- **Human-AI co-creative tools** for music, audio and video creation
- **Human-like machine learning algorithms** for music, movies and arts

Structural Multimodal Generative AI



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**

Controllable Generative AI

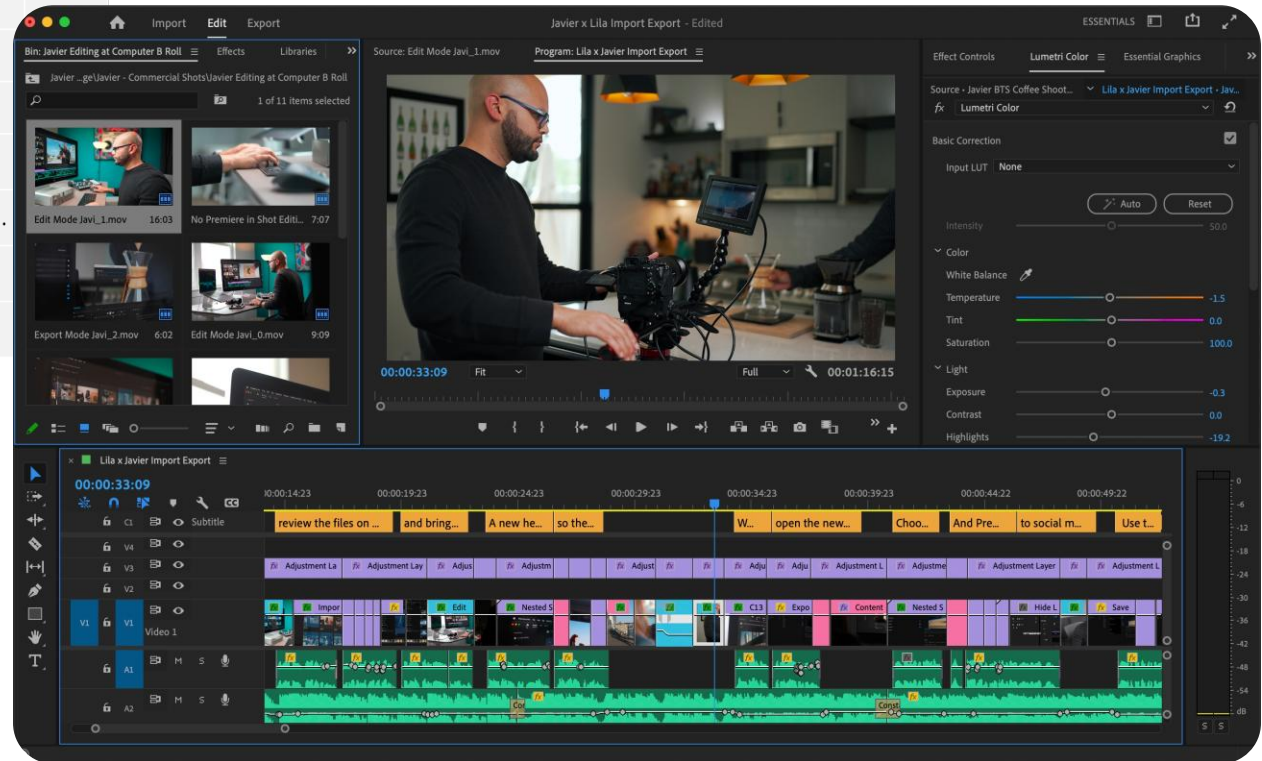


Audio Type	Layout	ID	Character	Volume	Action	Content Description	Duration
Music	Background	1	N/A	-30	Begin	Dramatic orchestral news theme.	Auto
Speech	Foreground	N/A	Host	-15	N/A	Welcome to Mars News ...	Auto
Music	Background	1	N/A	N/A	End	N/A	Auto
Speech	Foreground	N/A	Host	-15	N/A	Now let's connect 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

**Interactable
intermediate outputs**

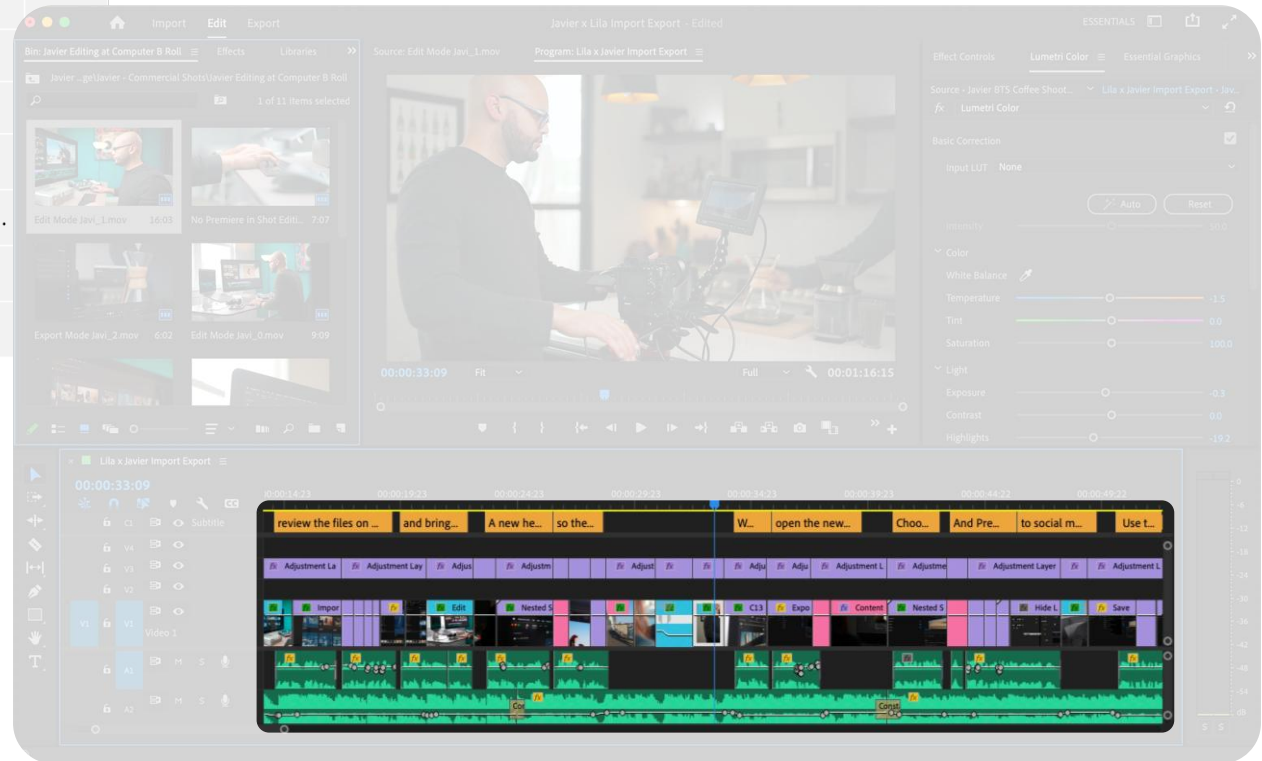
Controllable Generative AI

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 ...	



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

Human-Centered Generative AI for Content Creation

Augmenting human creativity with machine learning

- **Multimodal generative AI** for content creation
- **Human-AI co-creative tools** for music, audio and video creation
- **Human-like machine learning algorithms** for music, movies and arts

Human-Centered Generative AI for Content Creation

Augmenting human creativity with machine learning



UC San Diego



SONY

