

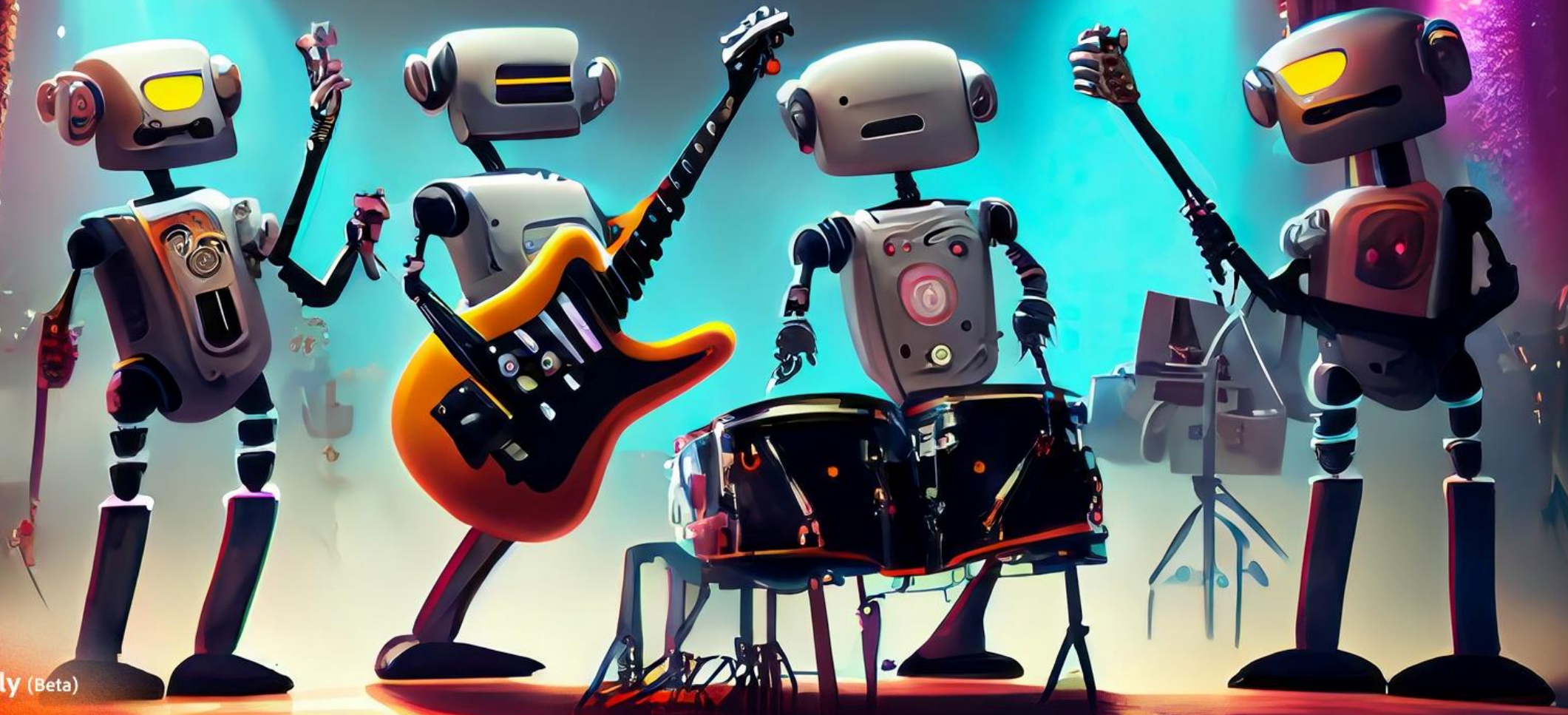
Generative AI for Music and Audio

Hao-Wen (Herman) Dong

董皓文

UC San Diego

Generative AI for Music and Audio



About me



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



B.S. in Electrical Engineering

2013 – 2017



Research Assistant

2017 – 2019

UC San Diego

M.S. in Computer Science

2019 – 2021

UC San Diego

Ph.D. in Computer Science (expected)

2019 – present

Summer 2019



Research Intern

Summer 2021



Deep Learning Audio Intern

Summer 2022



Student Intern

Fall 2022



Applied Scientist Intern

Winter 2023



Speech/Audio Deep Learning Intern

Summer 2023



Research Scientist/Engineer Intern

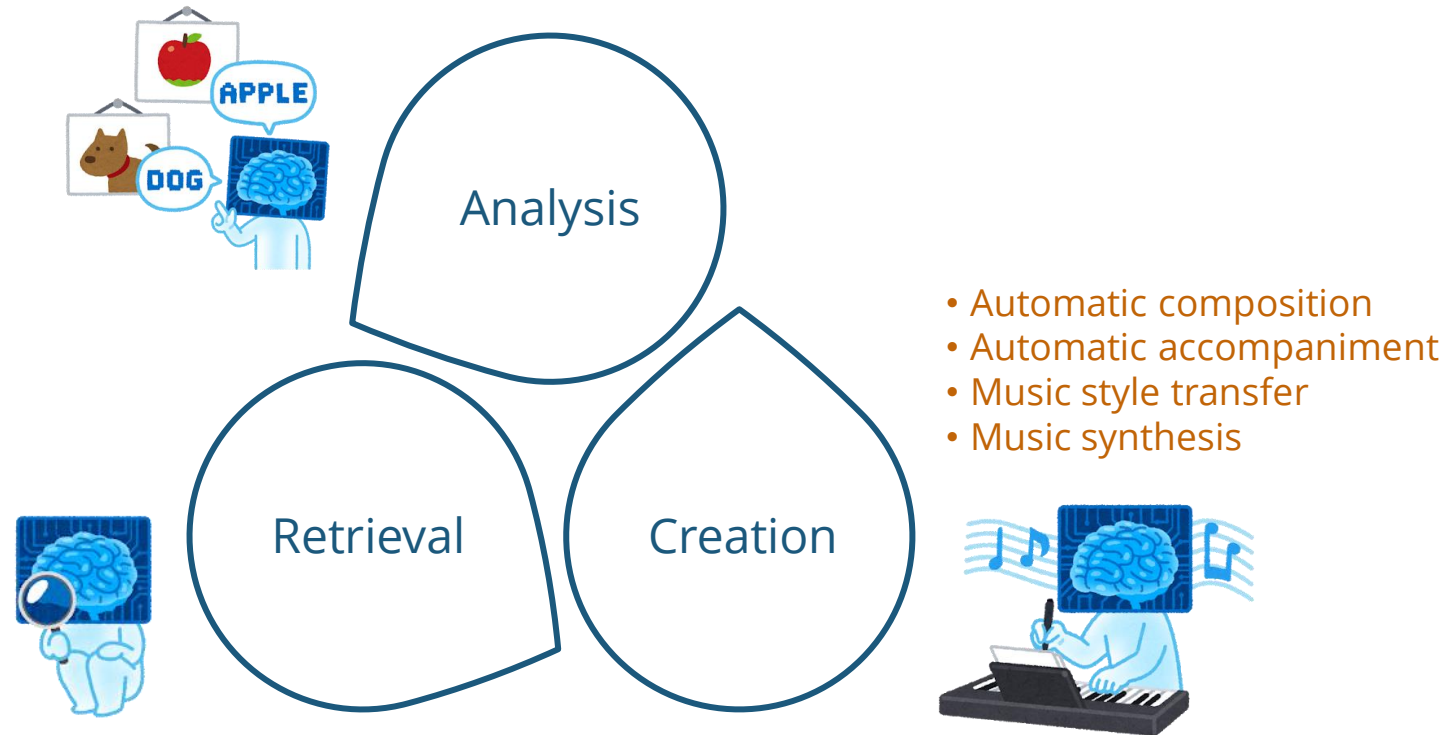
Fall 2023



Research Intern

Music Information Research (MIR)

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



My Research

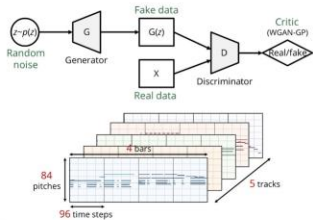


Multitrack Music Generation

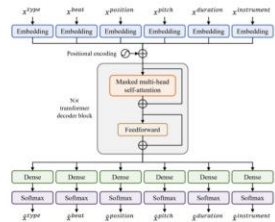
Generating new music contents automatically



MuseGAN (AAAI 2018)



Multitrack Music Transformer (ICASSP 2023)

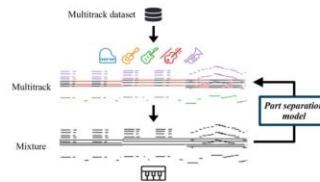


Assistive Music Creation Tools

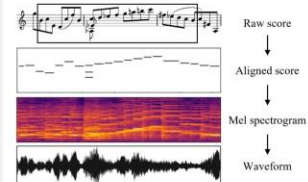
Assisting humans to create and perform music



Arranger (ISMIR 2021)



Deep Performer (ICASSP 2022)

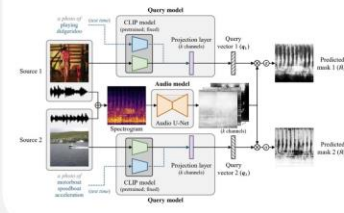


Multimodal Learning for Audio & Music

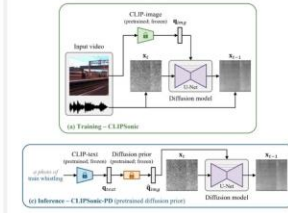
Learning sound separation and synthesis from videos



CLIPSep (ICLR 2023)



CLIPsonic (WASPAA 2023)



About me

EE



a female cat engineer making an electric chip in a classroom

Music



a cat playing heavy metal

CS



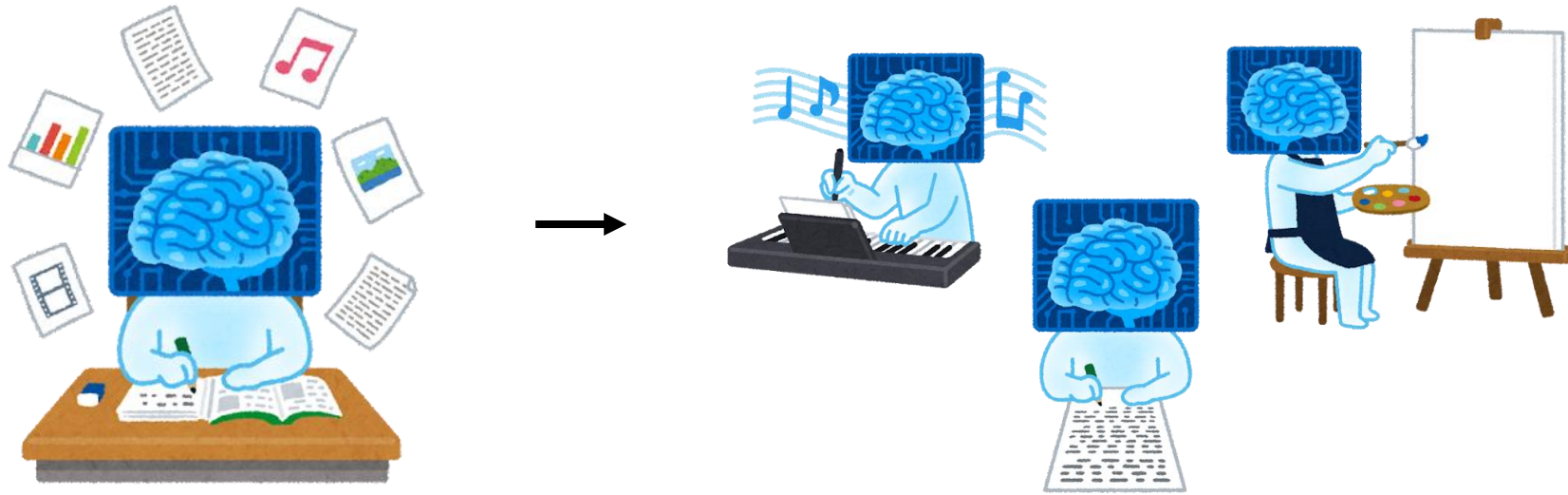
a cat engineer debugging on laptop



Introduction

What is Generative AI?

- Generative AI is **AI capable of generating text, images, or other media.**
 - Learns the patterns and structure of their input training data
 - Generates new data that has similar characteristics



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.

Landscape of Generative AI



Where is all the money going in generative AI?

Distribution of generative AI funding, Q3'22 – Q2'23

Generative interfaces

\$2,690M | 23 deals



Text
\$639M | 24 deals



Visual media
\$387M | 33 deals*



Code
\$178M | 17 deals



Speech & audio
\$121M | 15 deals

Source: CB Insights. Based on an analysis of 210+ generative AI companies building cross-industry enterprise solutions; excludes deals to industry-specific companies and model developers such as OpenAI.
*Includes 1 deal in motion capture animation and 1 deal in synthetic anonymization with undisclosed funding.



Landscape of Generative AI

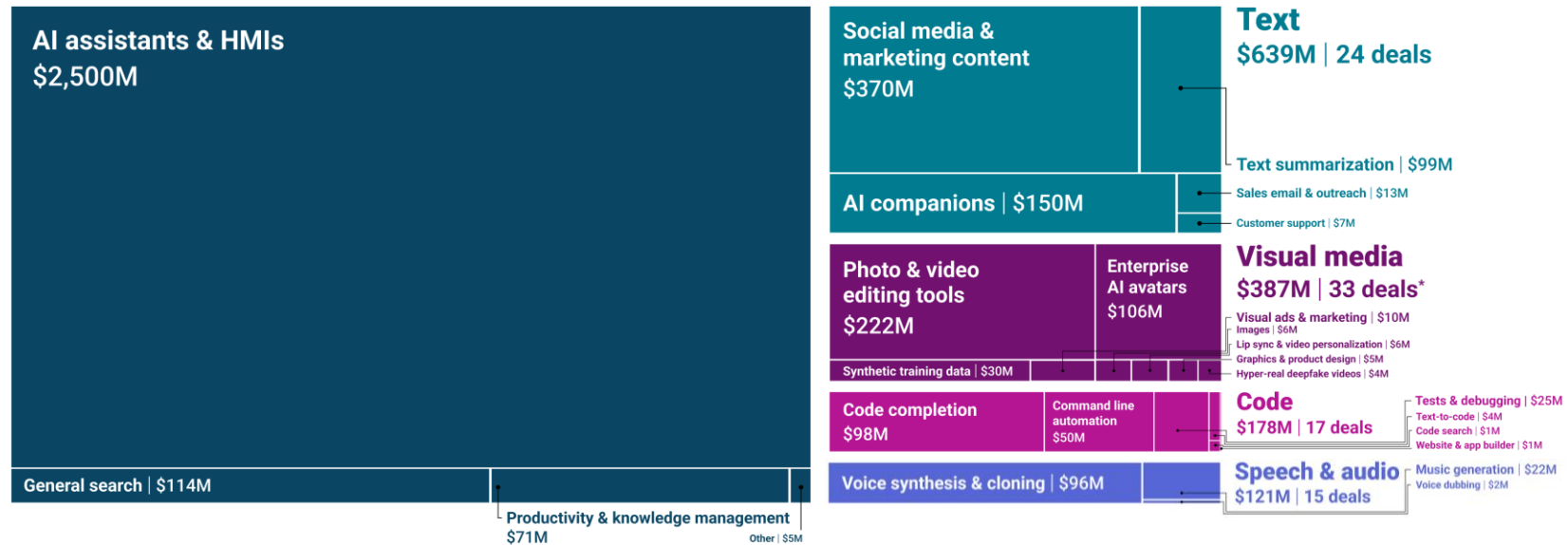


Where is all the money going in generative AI?

Distribution of generative AI funding, Q3'22 – Q2'23

Generative interfaces

\$2,690M | 23 deals



Source: CB Insights. Based on an analysis of 210+ generative AI companies building cross-industry enterprise solutions; excludes deals to industry-specific companies and model developers such as OpenAI.
*Includes 1 deal in motion capture animation and 1 deal in synthetic anonymization with undisclosed funding.



Types of Audio



Speech



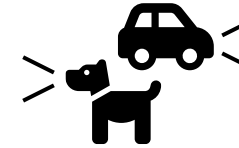
(Source: Wikimedia Commons)

Music



(Source: Wikimedia Commons)

Sound effects

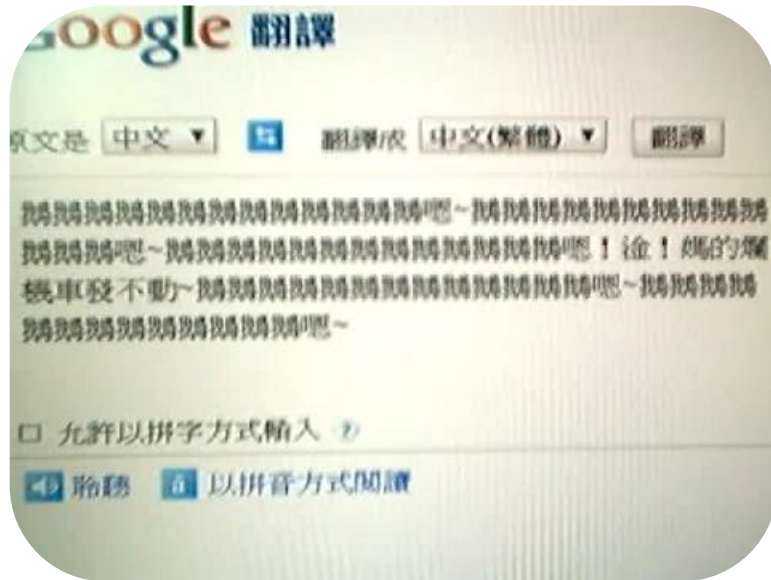


(Source: Wikimedia Commons)

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.

Generative AI for Speech

Text-to-Speech



Voice Cloning



凌小海, "超爆笑的google翻譯~我一直狂笑!" YouTube, April 27, 2011.

Eric, Twitter, <https://twitter.com/VyacheslavAI/status/1692144315055677818>, August 17, 2023.

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

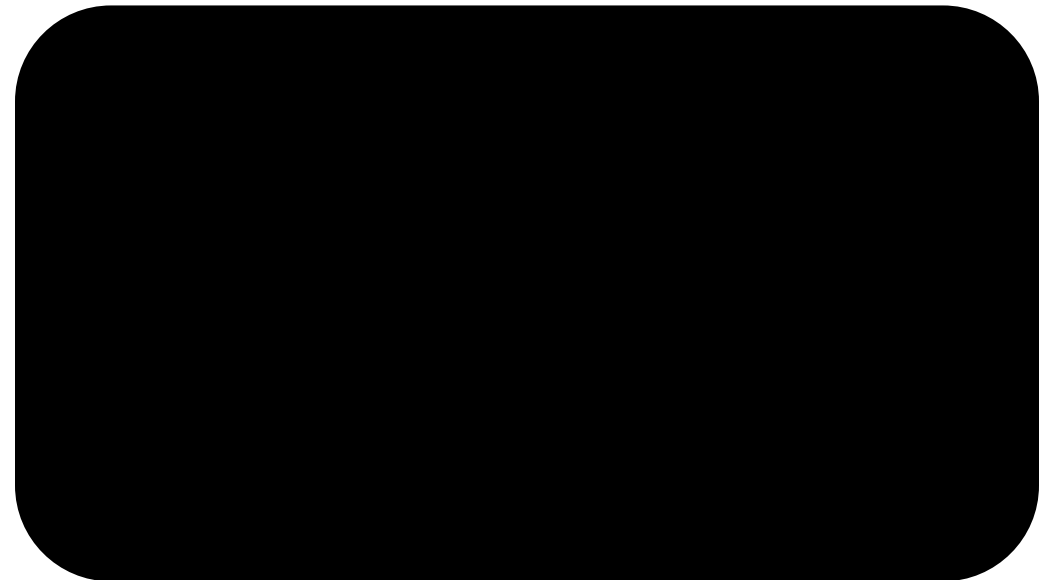


Generative AI for Sound Effects

Text-to-audio Synthesis



Image-to-audio Synthesis



My Research

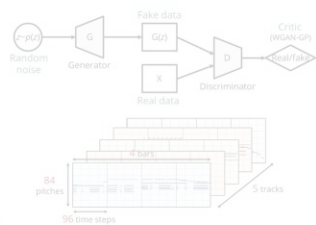


Multitrack Music Generation

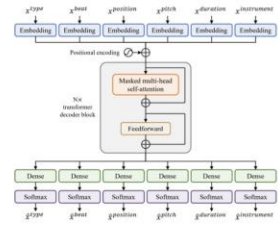
Generating new music contents automatically



MuseGAN (AAAI 2018)



Multitrack Music Transformer (ICASSP 2023)



Assistive Music Creation Tools

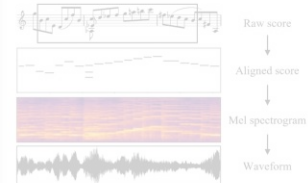
Assisting humans to create and perform music



Arranger (ISMIR 2021)



Deep Performer (ICASSP 2022)

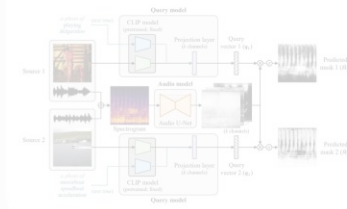


Multimodal Learning for Audio & Music

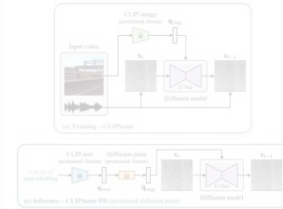
Learning sound separation and synthesis from videos



CLIPSep (ICLR 2023)



CLIPsonic (WASPAA 2023)





Multitrack Music Transformer

Hao-Wen Dong Ke Chen Shlomo Dubnov Julian McAuley Taylor Berg-Kirkpatrick
University of California San Diego



UC San Diego

Overview

Generate orchestral music

- of diverse instruments
- using a new compact representation
- with a multi-dimensional transformer



Demo



(Source: Vienna Mozart Orchestra)



Related Work (Transformers for Music Generation)

Model	Multitrack	Instrument control	Compound tokens	Generative modeling
REMI [5]				✓
MMM [10]	✓			✓
CP [6]			✓	✓
MusicBERT [15]	✓		✓	
FIGARO [11]	✓			✓
MMT (ours)	✓	✓	✓	✓

	Average sample length (sec)	Inference speed (notes per second)
MMM [10]	38.69	5.66
REMI+ [11]	28.69	3.58
MMT (ours)	100.42	11.79

→ Longer samples!
Faster inference speed!

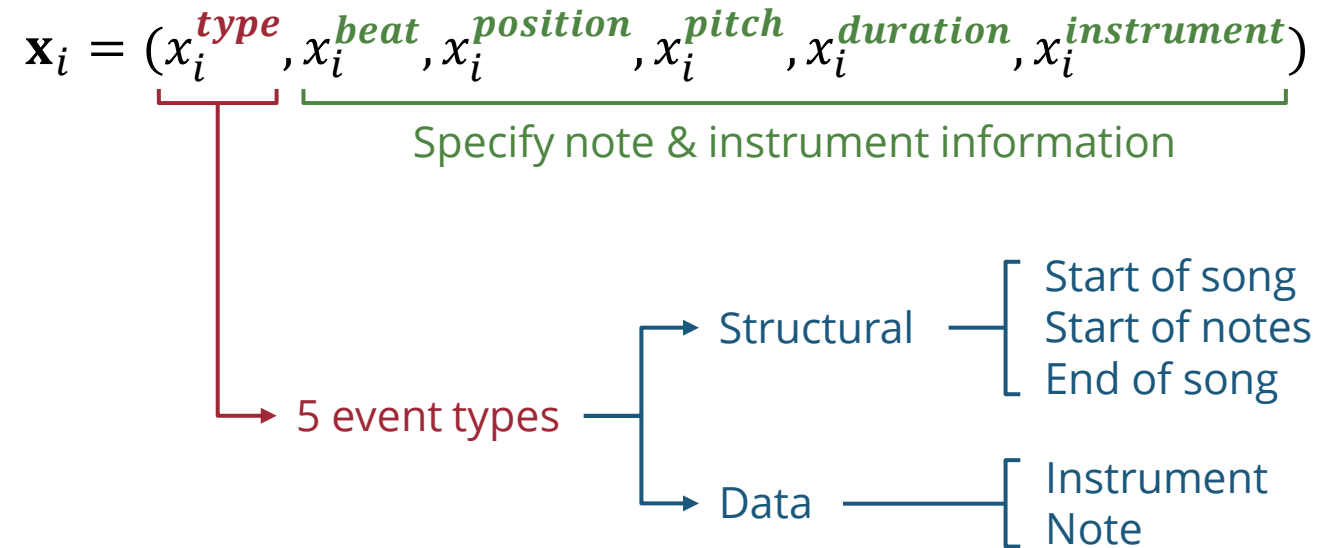
Huang and Yang, "Pop Music Transformer: Beat-based Modeling and Generation of Expressive Pop Piano Compositions," *Proc. MM*, 2020.
Ens and Pasquier, "MMM : Exploring Conditional Multi-Track Music Generation with the Transformer," *arXiv preprint arXiv:2008.06048*, 2020.
Hsiao et al., "Compound Word Transformer: Learning to Compose Full-Song Music over Dynamic Directed Hypergraphs," *Proc. AAAI*, 2023.
Zeng et al., "MusicBERT: Symbolic Music Understanding with Large-Scale Pre-Training," *Proc. Findings of ACL*, 2021.
von Rütte et al., "FIGARO: Controllable Music Generation using Learned and Expert Features," *Proc. ICLR*, 2023.

Representation

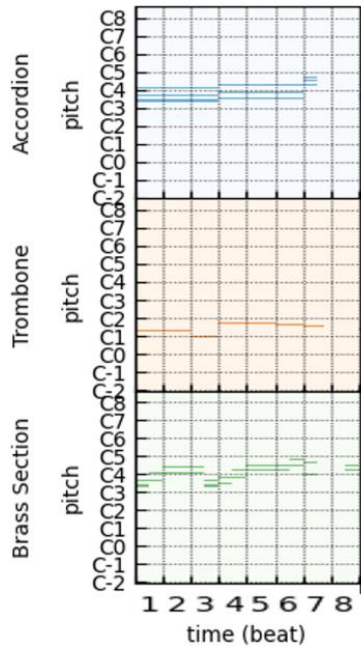
- We represent a music piece as a sequence of events

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

- Each event \mathbf{x}_i is encoded as



Representation (An Example)



Structural events

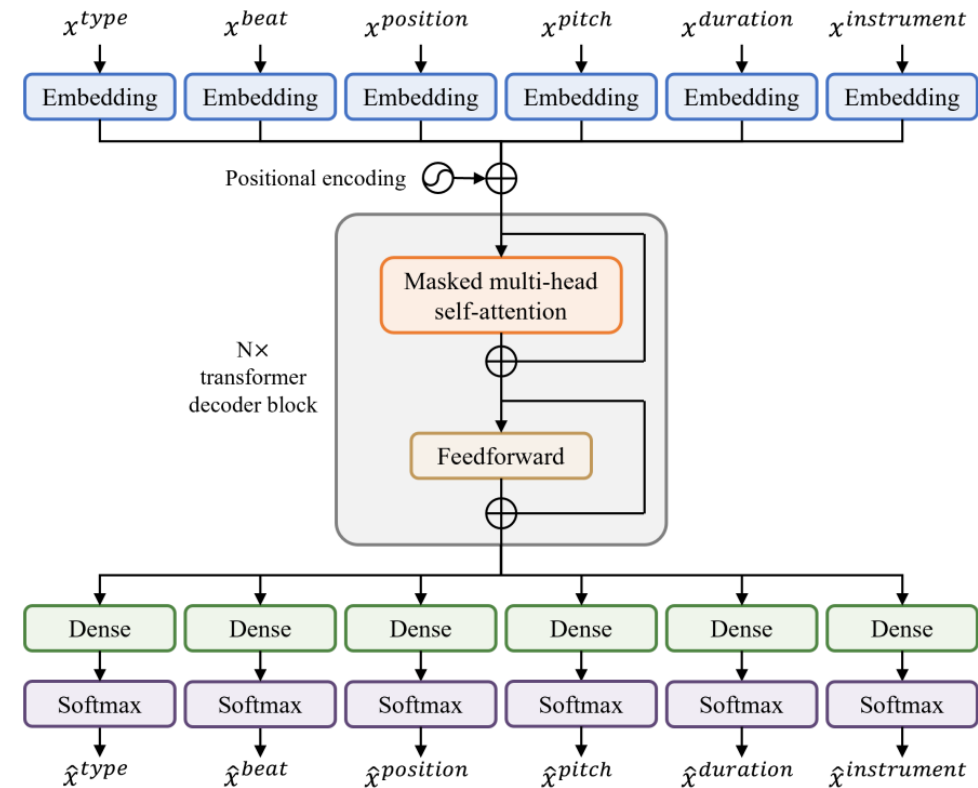
(0, 0, 0, 0, 0, 0)	Start of song
(1, 0, 0, 0, 0, 15)	Instrument: accordion
(1, 0, 0, 0, 0, 36)	Instrument: trombone
(1, 0, 0, 0, 0, 39)	Instrument: brasses
(2, 0, 0, 0, 0, 0)	Start of notes
(3, 1, 1, 41, 15, 36)	Note: beat=1, position=1, pitch=E2, duration=48, instrument=trombone
(3, 1, 1, 65, 4, 39)	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

Multitrack Music Transformer

- A multi-dimensional decoder-only transformer model
 - Predict six fields *at the same time*
- Trained autoregressively
 - Predict the next event given past events



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

Only need to train ONE model!

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

Experimental Setup

Data

- Symbolic Orchestral Database (SOD)
(Crestel et al., 2017)
 - 5,743 songs, 357 hours
- Temporal resolution: 12 time steps per quarter note
- 80% training, 10% validation, 10% test
- Data augmentation
 - Randomly shift for -5~6 semitones
 - Randomly select a starting beat

Model & Training

- 6 transformer decoder blocks
- 8 attention heads
- Model dimension: 512
- Sequence length: 1024
- Maximum number of beats: 256
- Maximum training steps: 200,000

Example Results

Unconditional
generation



Instrument-informed
generation



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

4-beat continuation



Wolfgang Amadeus Mozart's
Eine kleine Nachtmusik



More audio samples



salu133445.github.io/mmt/

Subjective Listening Test Results

	Number of parameters	Average sample length (sec)	Inference speed (notes per second)	Subjective listening test results			
				Coherence	Richness	Arrangement	Overall
MMM [10]	19.81 M	38.69	5.66	3.48 ± 0.35	3.05 ± 0.38	3.28 ± 0.37	3.17 ± 0.43
REMI+ [11]	20.72 M	28.69	3.58	3.90 ± 0.52	3.74 ± 0.21	3.74 ± 0.44	3.77 ± 0.41
MMT (ours)	19.94 M	100.42	11.79	3.55 ± 0.46	3.53 ± 0.35	3.40 ± 0.44	3.33 ± 0.47

2.6x/3.5x longer generated samples
(within the same sequence length)

2.1x/3.3x faster inference speed

Higher quality than MMM
Lower quality than REMI+

Analyzing Self-attention

- *Mean relative attention* for a field d :

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

↑ Attention weight
→ Whether the field value is of difference k

(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

$\gamma_{-8}^{(pitch)}$ (curved arrow from row 3 to row 6)
 $\gamma_{-5}^{(pitch)}$ (curved arrow from row 6 to row 9)

Analyzing Self-attention

- *Mean relative attention* for a field d :

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

Biased towards difference that occurred more frequently!

- *Mean relative attention gain* for a field d :

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

Assuming a uniform attention matrix

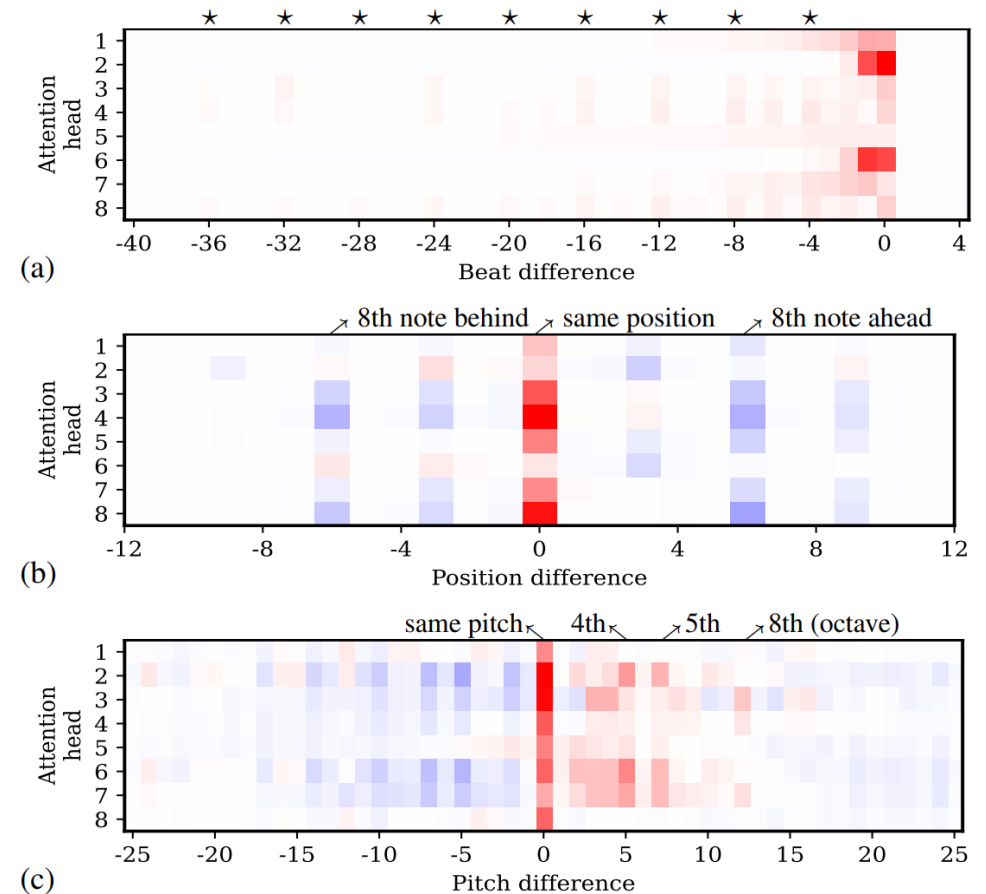
Musical Self-attention

The MMT model attends more to notes

- that are $4N$ beats away in the past
- that have the same position as the current note (A note on beat attends more to a note on beat; a note off beat attends more to a note off beat.)
- that has a pitch in an octave above which forms a consonant interval

→ MMT learns a **relative self-attention** for certain aspects of music, specifically, **beat**, **position** and **pitch**.

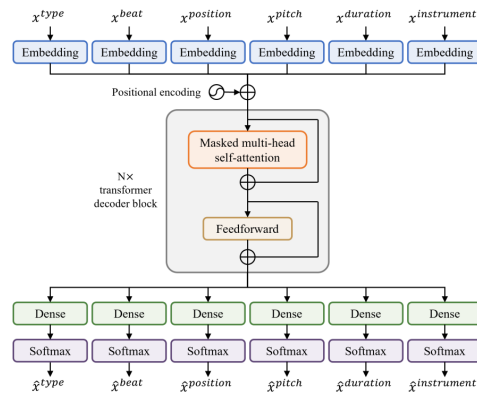
Positive and negative mean relative attention gain



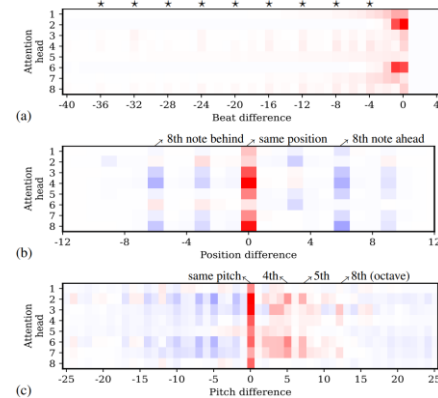
Summary

- Proposed an efficient representation and model for multitrack music generation
- Presented the first systematic analysis of musical self-attention

Multitrack Music Transformer



Musical Self-attention



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



My Research

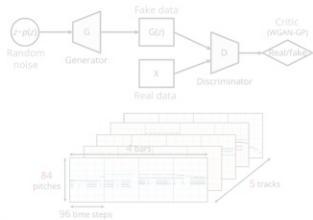


Multitrack Music Generation

Generating new music contents automatically



MuseGAN (AAAI 2018)



Multitrack Music Transformer (ICASSP 2023)



Assistive Music Creation Tools

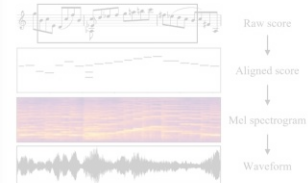
Assisting humans to create and perform music



Arranger (ISMIR 2021)



Deep Performer (ICASSP 2022)

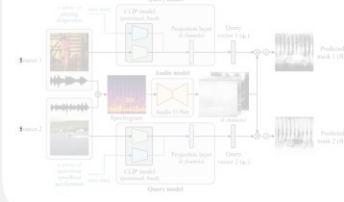


Multimodal Learning for Audio & Music

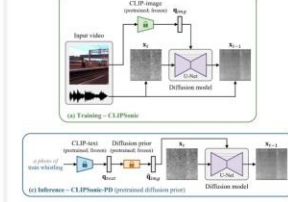
Learning sound separation and synthesis from videos



CLIPSep (ICLR 2023)



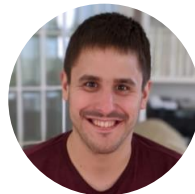
CLIPsonic (WASPAA 2023)



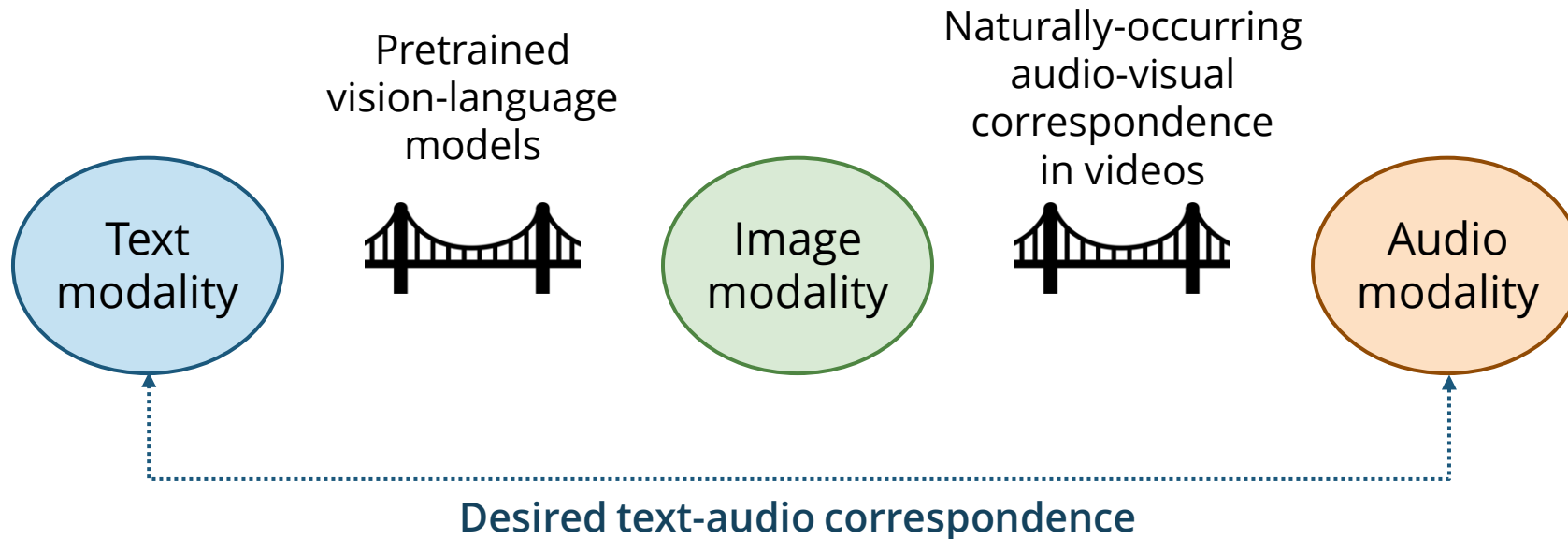
CLIPSonic: Text-to-Audio Synthesis with Unlabeled Videos and Pretrained Language-Vision Models

Hao-Wen Dong^{1,2*} Xiaoyu Liu¹ Jordi Pons¹ Gautam Bhattacharya¹
Santiago Pascual¹ Joan Serra¹ Taylor Berg-Kirkpatrick² Julian McAuley²

¹ Dolby Laboratories ² University of California San Diego
* Work done during an internship at Dolby



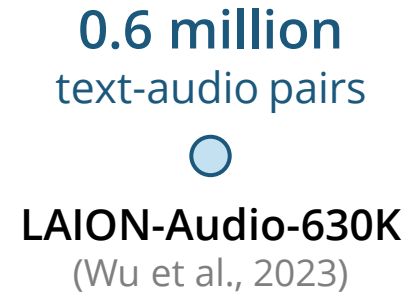
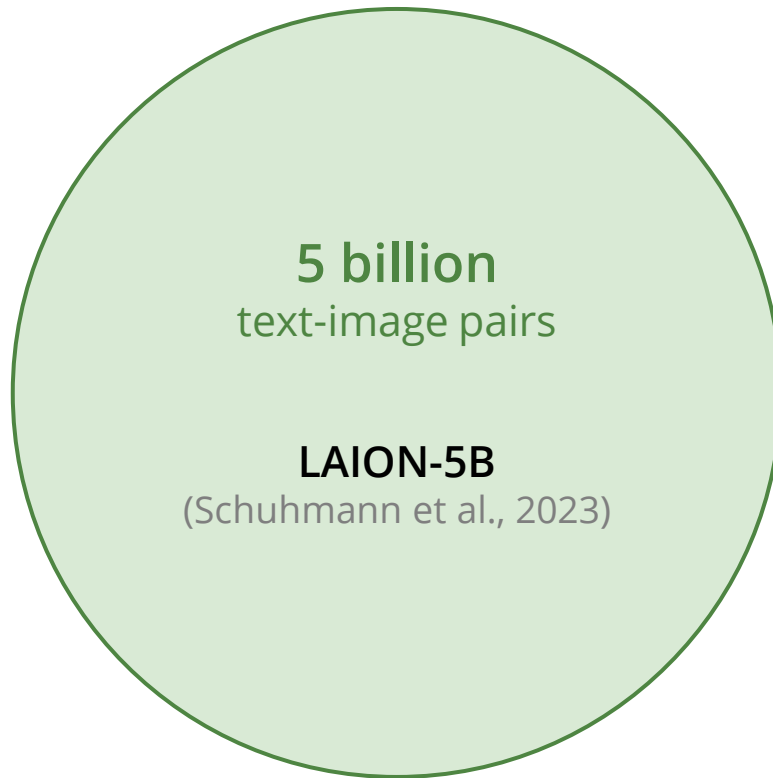
Leveraging the Visual Domain as a Bridge



No text-audio pairs required!

Scalable to large video datasets!

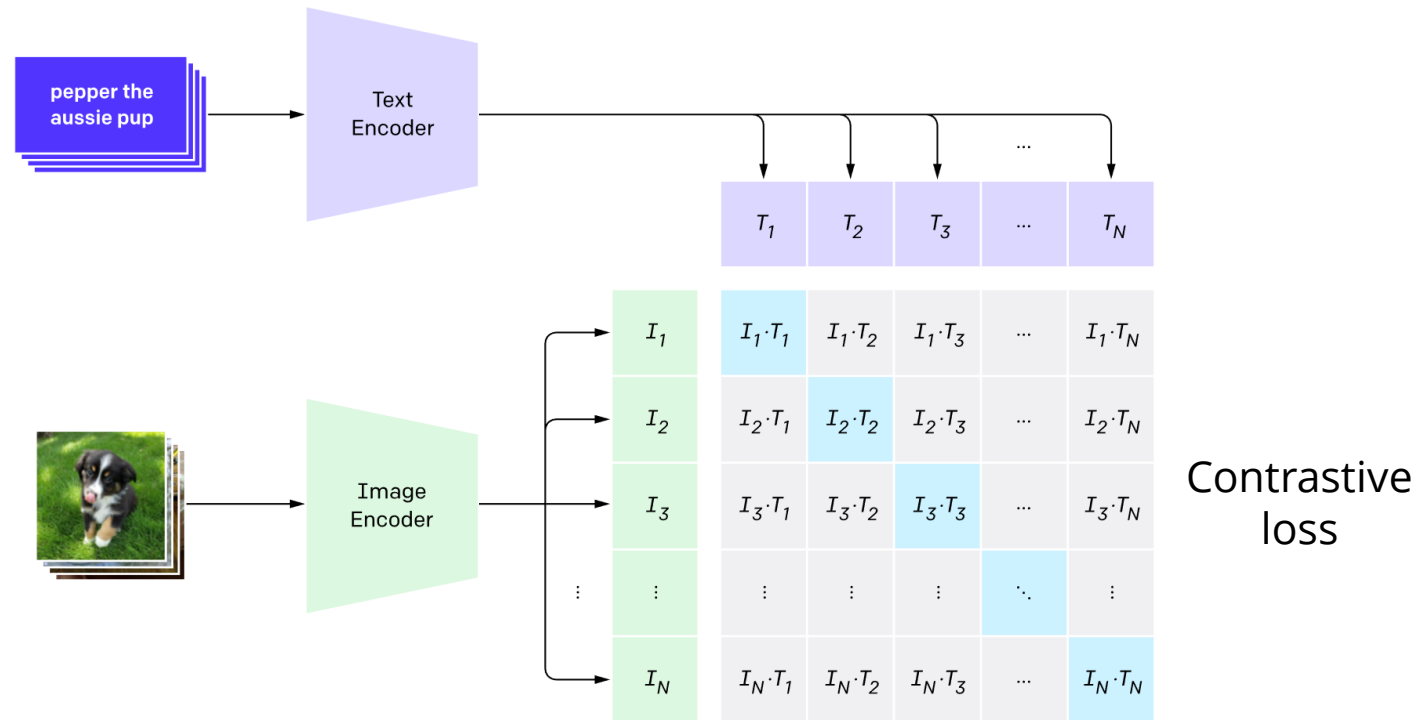
Why NOT Text-audio Pairs?



YouTube videos!
500 hours of videos
uploaded per minute

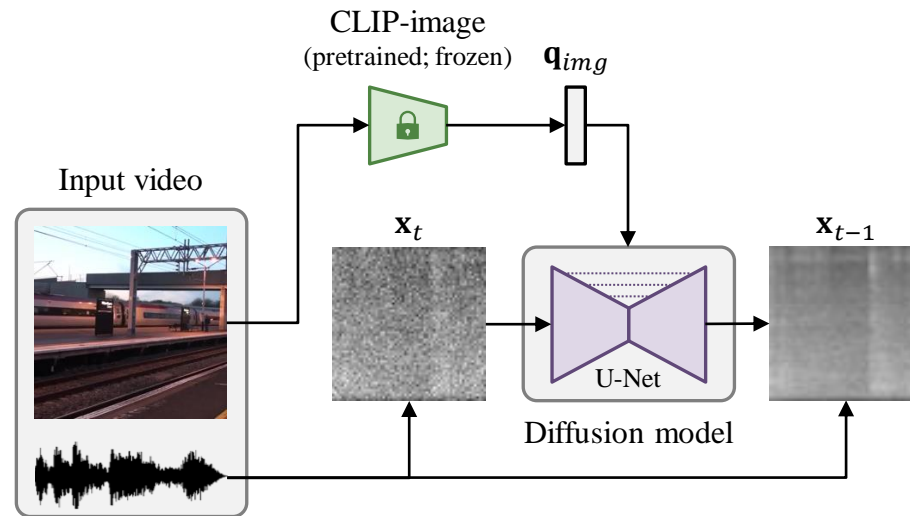
CLIP (Contrastive Language-Image Pretraining)

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



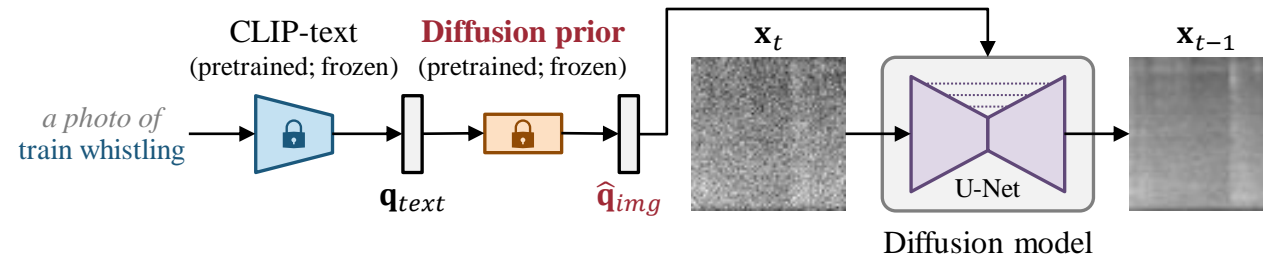
CLIP Sonic – Training

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

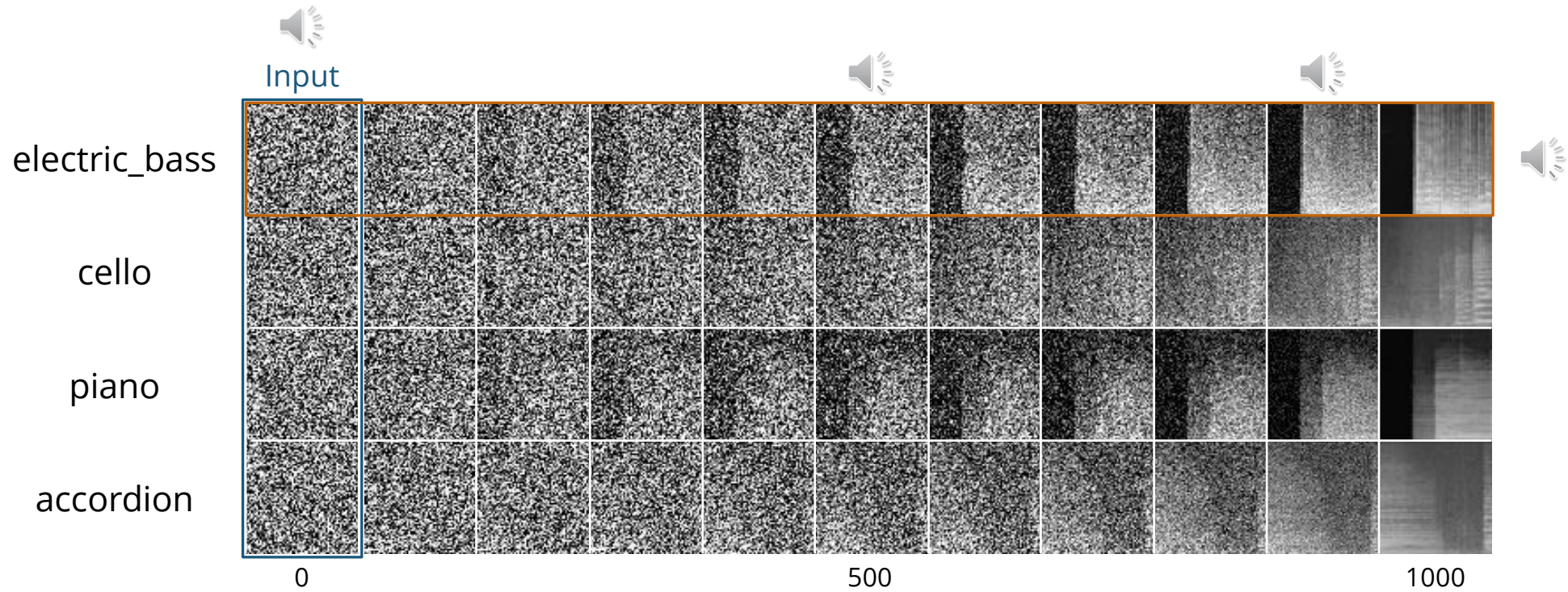


CLIPsonic – Inference

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



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

Text-to-Audio Synthesis Demo



Demo

Rapping



Sea waves



Smoke detector
beeping



Playing table
tennis



Thunder



Playing violin
fiddle



Subjective Listening Test (text-to-audio synthesis)

- CLIPSONIC-PD with a pretrained diffusion prior model performs significantly better
 - than its counterpart *without* a diffusion prior model (CLIPSONIC-ZS)
 - in terms of **fidelity** on both datasets
 - in terms of **relevance** on MUSIC

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

Significant
performance
improvement

Image-to-Audio Synthesis Demo (out-of-distribution)



Demo



CLIPsonic-IQ (ours)

Im2wav (Sheffer & Adi, 2023)

SpecVQGAN (Iashin & Rahtu, 2021)



CLIPsonic-IQ (ours)

Im2wav (Sheffer & Adi, 2023)

SpecVQGAN (Iashin & Rahtu, 2021)



Subjective Listening Test (image-to-audio synthesis)

- CLIPSonic-IQ significantly outperforms im2wav and SpecVQGAN in audio **fidelity**
- CLIPSonic-IQ significantly outperforms SpecVQGAN in text-audio **relevance**
- CLIPSonic-IQ is competitive against im2wav in text-audio **relevance**

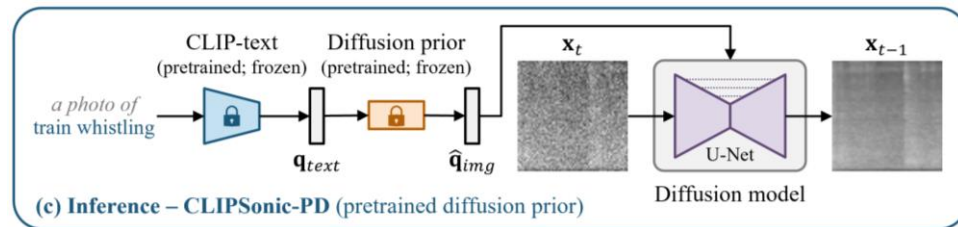
Table 4: Listening test results for image-to-audio synthesis (MOS).

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

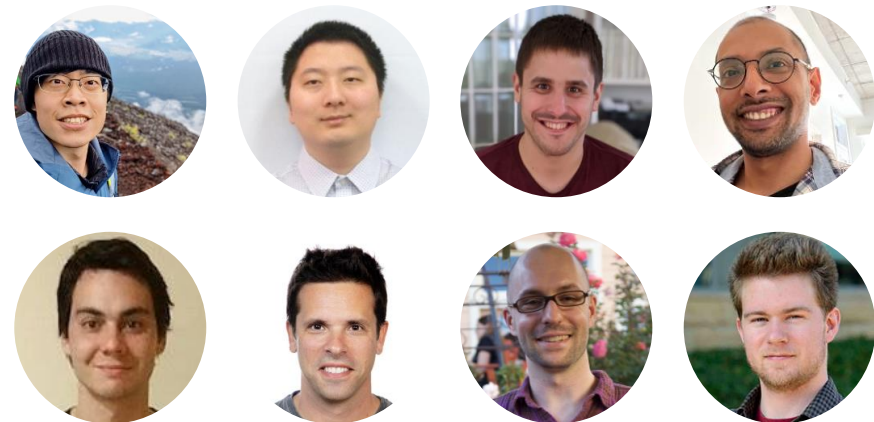
State-of-the-art
image-to-audio
performance!

Summary

- Proposed a new text-to-audio synthesis model that *requires no* text-audio pairs
- CLIPsonic-PD achieves good performance in objective and subjective evaluations
- CLIPsonic-IQ achieves state-of-the-art performance in image-to-audio synthesis



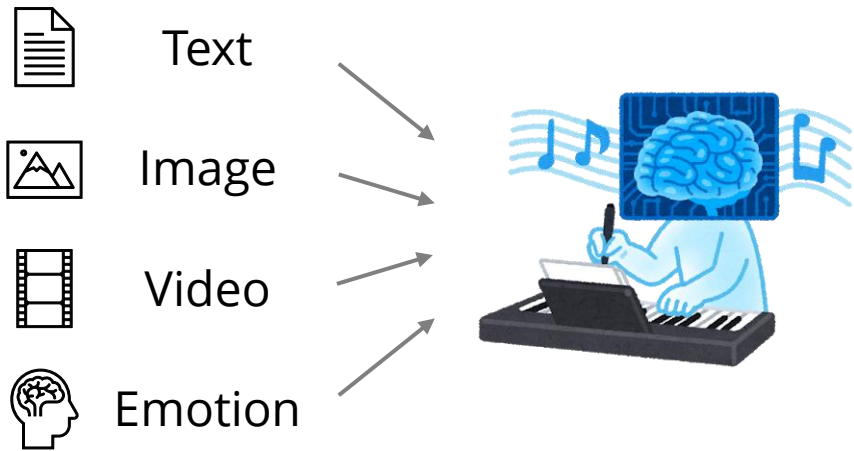
Paper: arxiv.org/abs/2306.09635
Demo: salu133445.github.io/clipsonic



Future of Generative AI

Challenges

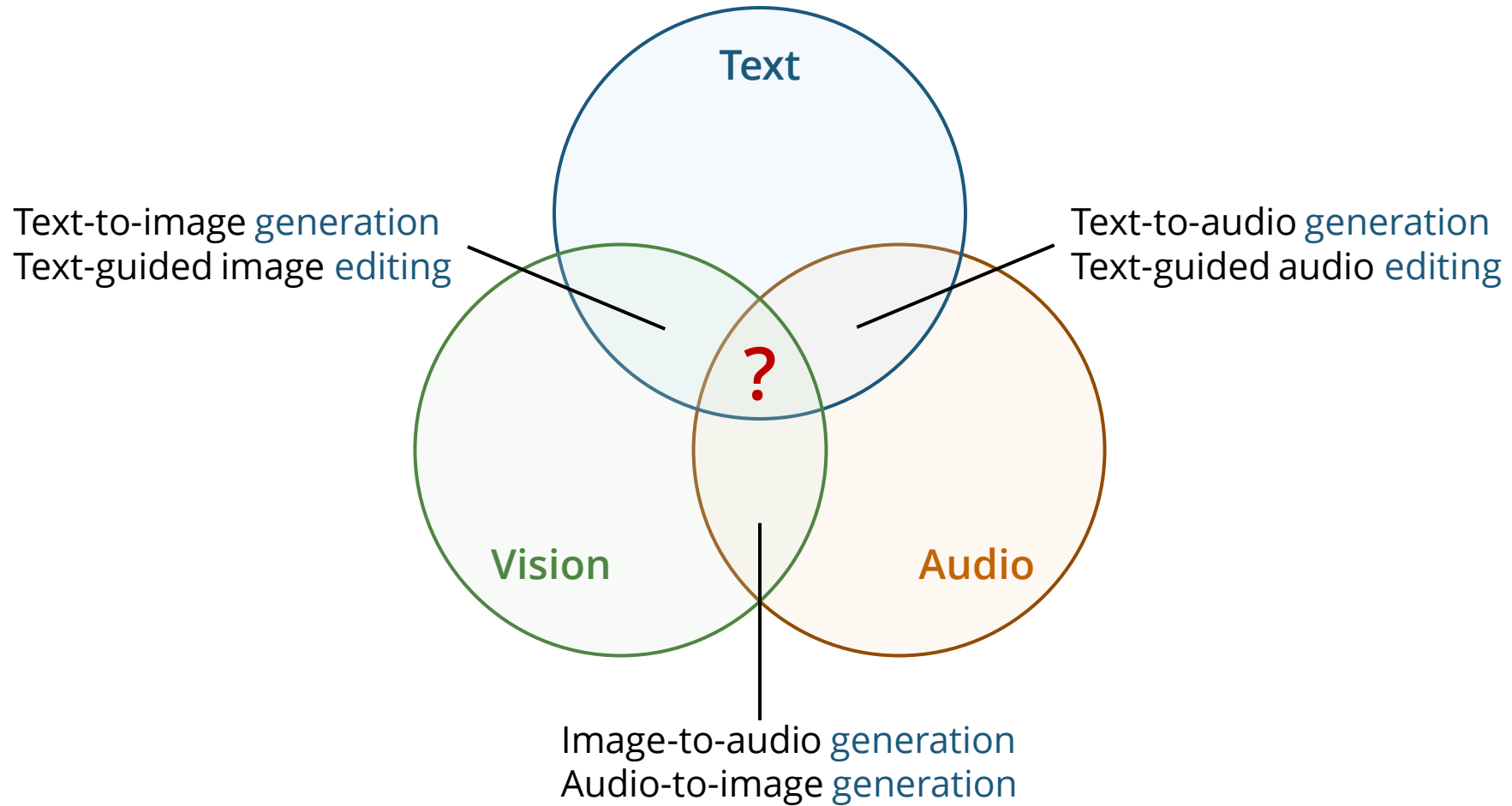
Multimodality



Usability



Multimodal Generative AI



Multimodal generative AI for Ads



Video **Runway Gen-2**
Music **MusicGen**



Multimodal generative AI for **Films**



Visuals **Midjourney**

Video **Runway**

Narration (script) **ChatGPT**

Narration (voice) **ElevenLabs**

Sound effects **Audiocraft**



Mumbai, the city of dreams.



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**

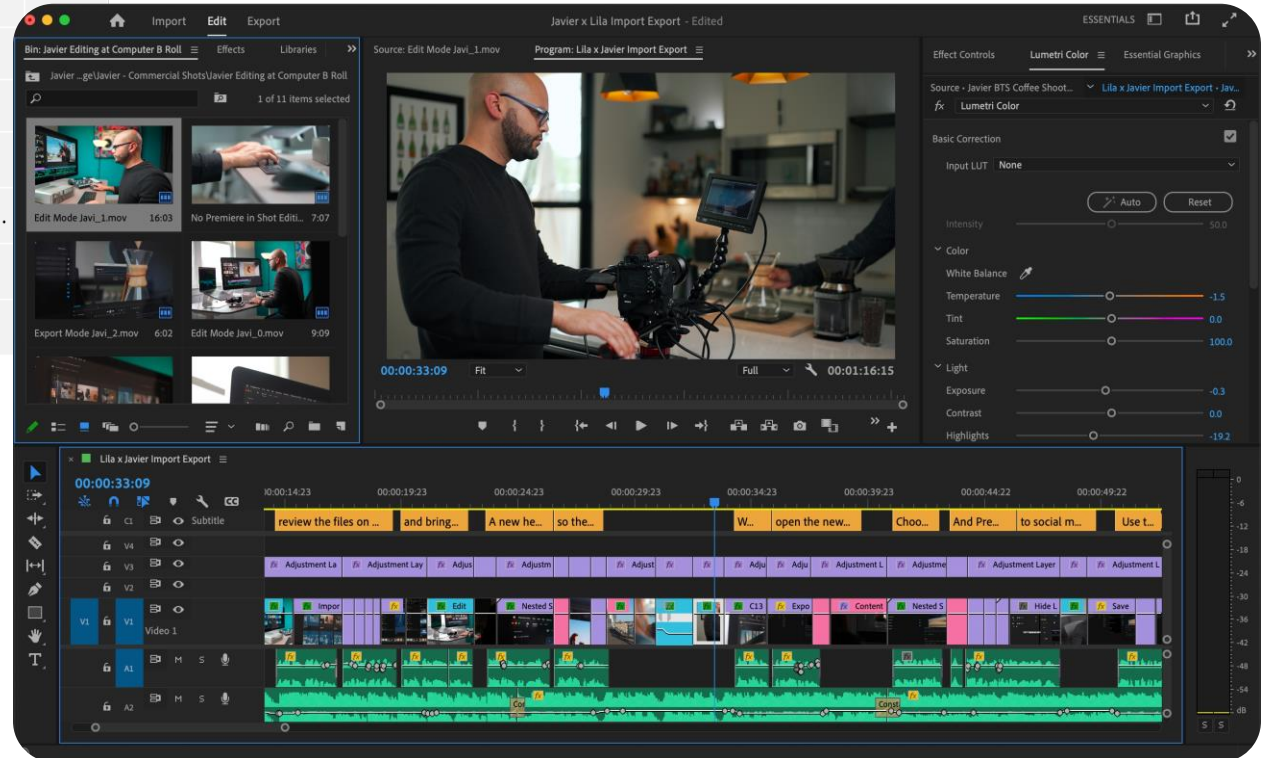
Controllable Generative AI



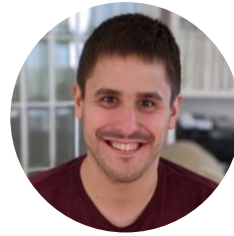
Audio Type	Layout	ID	Character	Volume	Action	Content Description	Duration
Music	Background	1	N/A	-30	Begin	Dramatic orchestral news theme.	Auto
Speech	Foreground	N/A	Host	-15	N/A	Welcome to Mars News ...	Auto
Music	Background	1	N/A	N/A	End	N/A	Auto
Speech	Foreground	N/A	Host	-15	N/A	Now let's connect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

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 ...	



Acknowledgements



Thank you!

