PAT 498/598 (Fall 2024)

# Special Topics:
# Generative AI for Music and Audio Creation

**Lecture 20: Review & Discussions**

Instructor: Hao-Wen Dong

# Final Project

- Milestones (all due at the specified date at **11:59 PM ET**)

  - **Pitch**          November 6          Topic & high-level plans
  - **Proposal**      November 22       Survey & plans (1 page)
  - **Presentation**    December 9        Showcase & report
  - **Final report**     December 15      Full report (3-5 pages)

- Instructions will be released on Gradescope

- Late submissions: **NOT accepted**

# Final Project: Rubrics

- **Proposal**         **10pt**

- **Presentation**     **20pt**

- **Final report**     **30pt**
  - Implementation                                    10pt
  - Code documentation                              5pt
  - Explanation of design and implementation     5pt
  - Results, analysis and discussions            10pt

# Final Project: Presentation

- **Introduction & motivation**
  - **Why** are you interested in this topic?
  - **Who** might want to use your work?

- **Design & implementation**
  - How did you **formulate the problem**?
  - How did you **implement your idea**?

- **Results, analysis & discussions**
  - **What have you found** through your experiments?
  - What are the **implications of your results and analysis**?
  - What are the **limitations** and **future directions**?

# Review – Music & AI

# The Early Days

### Musical Dice Game
### (1792)

### ILLIAC Suite
### (1957)

### Emily Howell
### (2003)



(Source: gbrachetta)

gbrachetta.github.io/Musical-Dice/



(Source: Illinois Distributed Museum)



(Source: The Guardian)
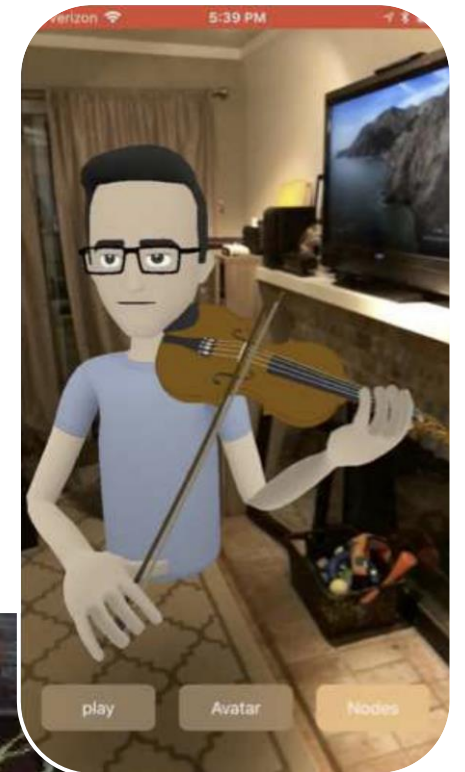
# Music & Technology

# Music & AI


(Source: Yamaha)


(Source: Sankei Shimbun)


(Shlizerman et al., 2019)


(Source: Robot Gizmos)


(Source: NBC DFW)

Shlizerman et al., "Audio to Body Dynamics," *Proc. CVPR*, 2018.
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

# Use Cases of Generative AI for Music & Audio



(Source: UploadVR)

(Source: The Denver Post)

(Source: Descript)

**Films**

**Gaming**

**Podcasts**

**Education**

**Dance**

**Theater**

**Short videos**

**Therapy**

(Source: Daily Bruin)

(Source: Wikimedia Commons)

# Review – AI/ML/DL Basics

# What is Artificial Intelligence?

AI is the study of how to make computers **do things at which, at the moment, people are better**.

– Elaine Rich and Kevin Knight, 1991

**1997**



(Source: Britannica)
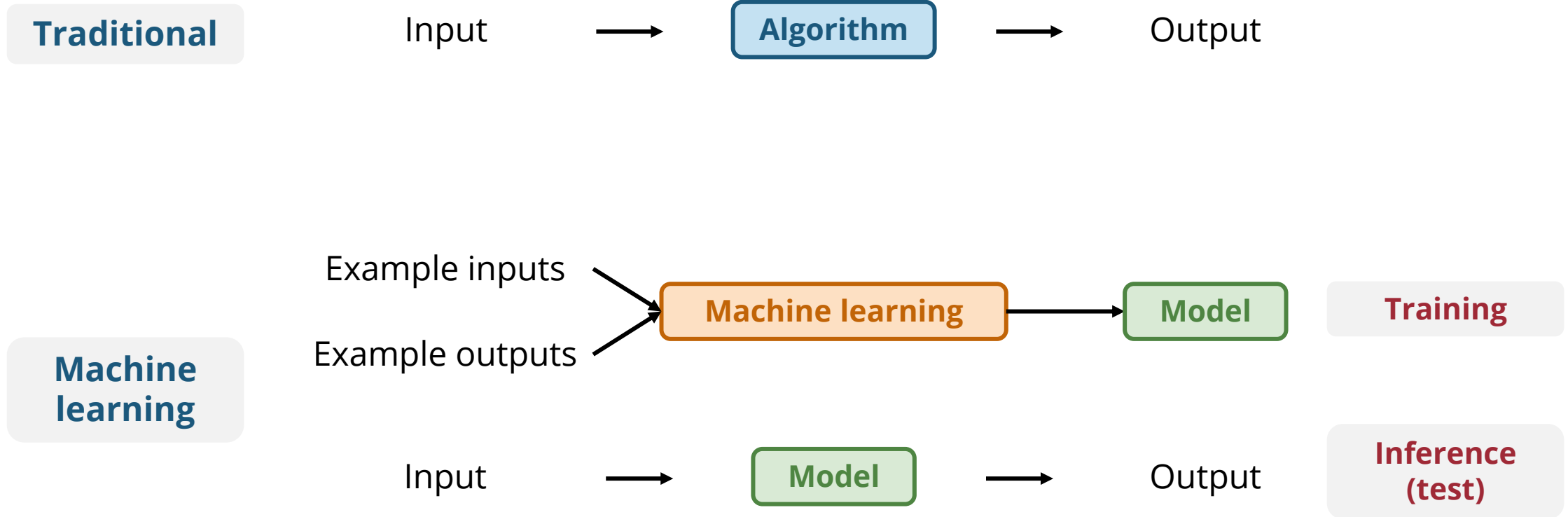
**2016**



(Source: The Guardian)

**20??**



(Source: SC2HL)

Elaine Rich and Kevin Knight, *Artificial Intelligence.* United Kingdom: McGraw-Hill, 1991.
https://www.britannica.com/topic/Deep-Blue
https://www.theguardian.com/technology/2016/mar/15/alphago-what-does-google-advanced-software-go-next
https://www.youtube.com/watch?v=PFMRDm_H9Sg

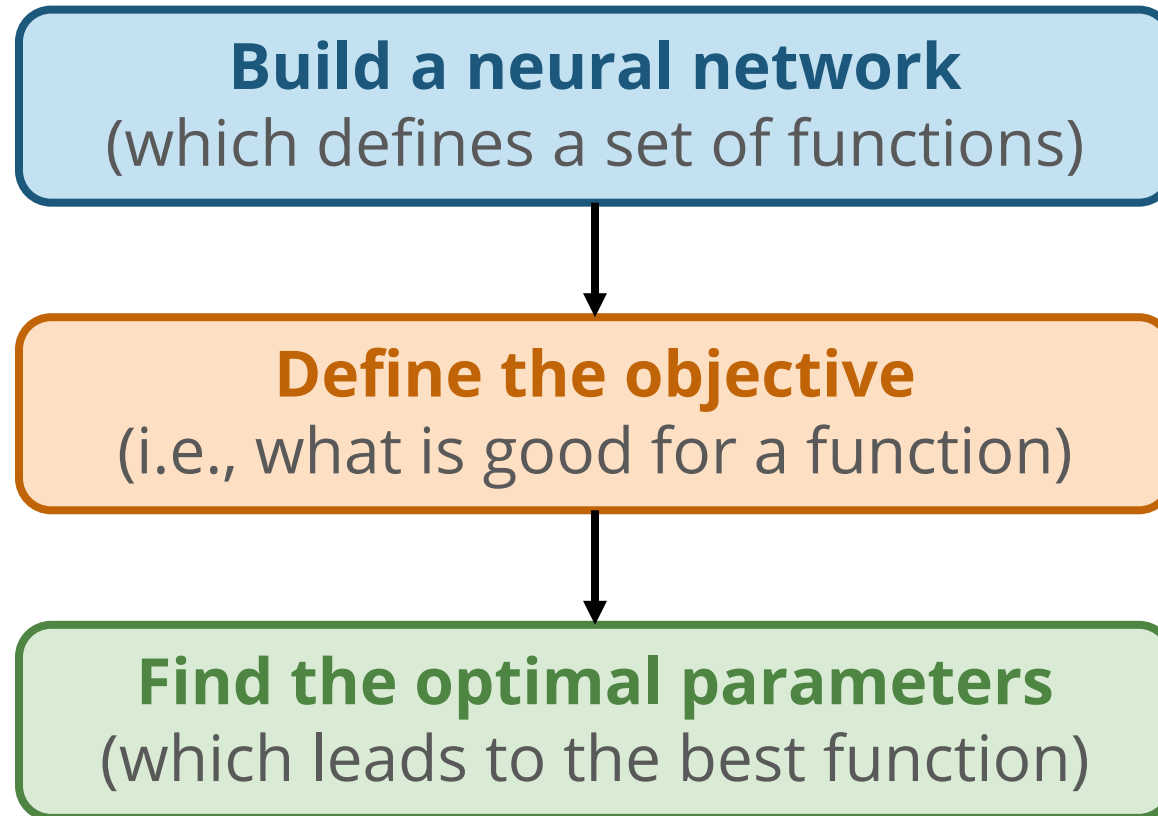# Machine Learning

**Traditional**

Input $\longrightarrow$ **Algorithm** $\longrightarrow$ Output

**Machine learning**

Example inputs
Example outputs
$\searrow$ $\nearrow$ **Machine learning** $\longrightarrow$ **Model**

**Training**

Input $\longrightarrow$ **Model** $\longrightarrow$ Output

**Inference (test)**

# Neural Networks are Parameterized Functions

- A neural network represents **a set of functions**



**Find the optimal parameters**

$$x$$

$$f_\theta(\mathbf{X})$$

**All the parameters**

$$W_1, \ldots, W_L, \mathbf{b_1}, \ldots, \mathbf{b_L}$$

$$\hat{y} \qquad y$$

**Good or bad?**

**Objective**

# Training a Neural Network

**Build a neural network**
(which defines a set of functions)

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x})$$

**Define the objective**
(i.e., what is good for a function)

$$Loss(\boldsymbol{\theta}) = \sum_{k}^{N} L(\hat{\mathbf{y}}_k, \mathbf{y}_k)$$

**Find the optimal parameters**
(which leads to the best function)

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

# Review – Training a Neural Network
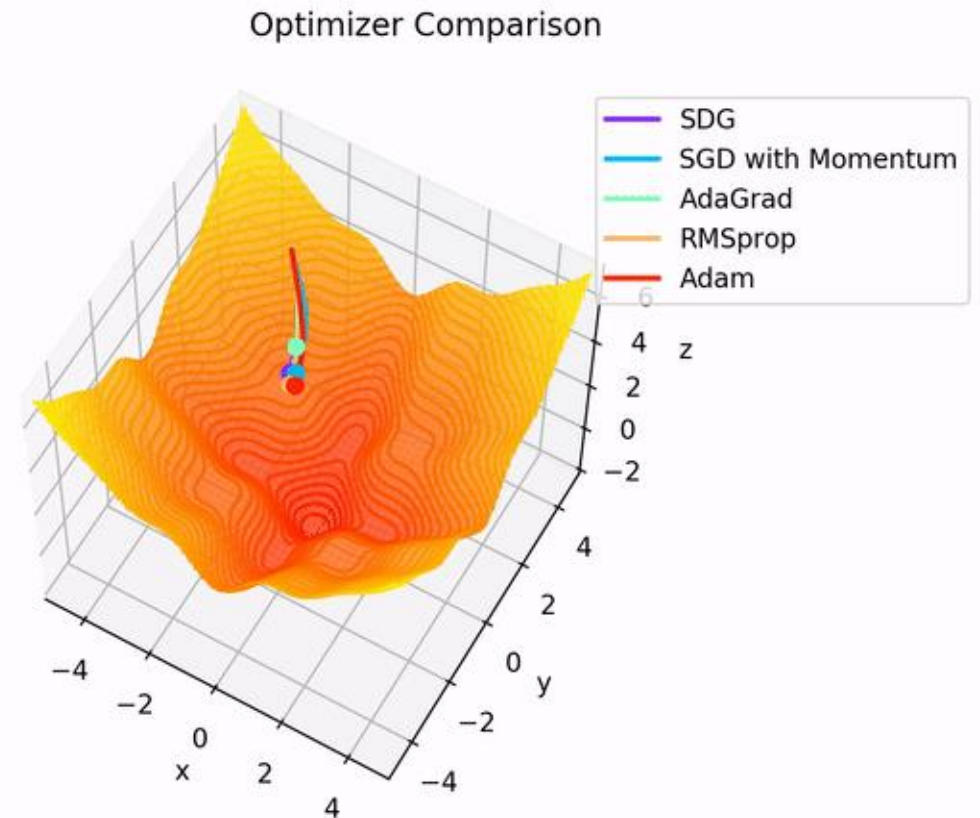
# Gradient Descent – 3D Case

# Comparison of Optimizers

- **Momentum**
  - Gets you out of spurious local minima
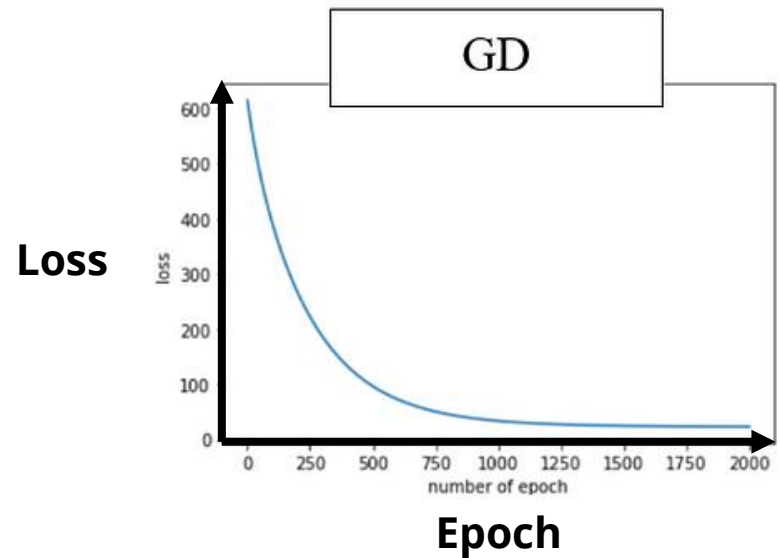  - Allows the model to explore around

- **Gradient-based adaption**
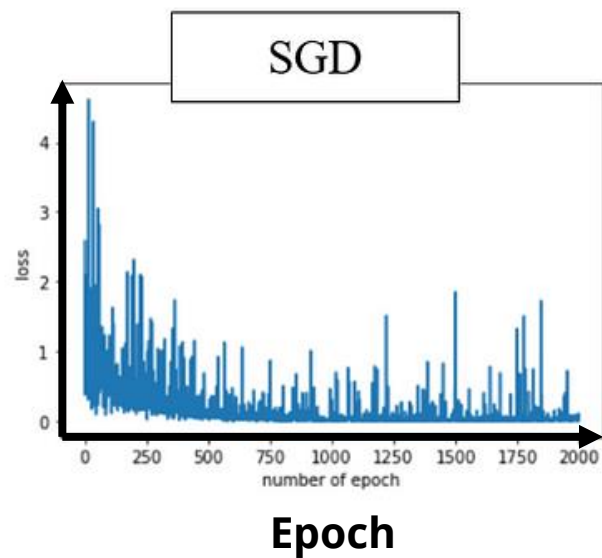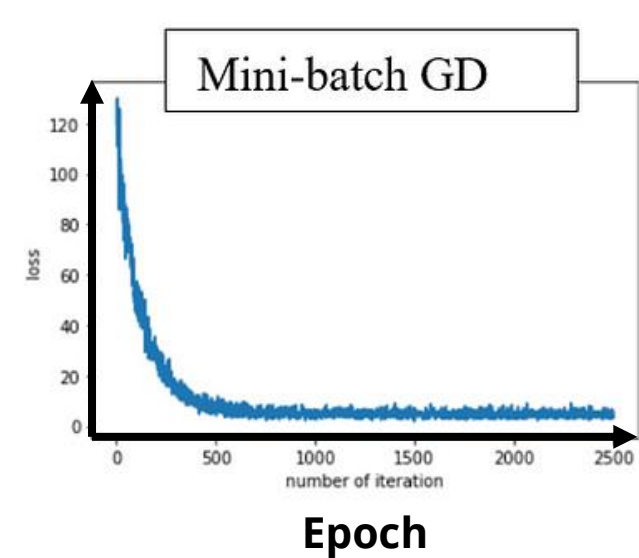  - Maintains steady improvement
  - Allows faster convergence



Optimizer Comparison

# Mini-batch Gradient Descent

- **Intuition:** **Estimate** the gradient using **several random training samples**



**Loss**

| GD | SGD | Mini-batch GD |
|---|---|---|
| Epoch | Epoch | Epoch |
| batch size = $N$ | batch size = $1$ | $1 <$ batch size $< N$ |

# Training–Validation–Test Pipeline

**Training**

**Validation**

**Test**

**Optimize**

**Select**

# Training vs Validation Losses

# Review – Neural Networks

# Network Architectures vs Training Frameworks

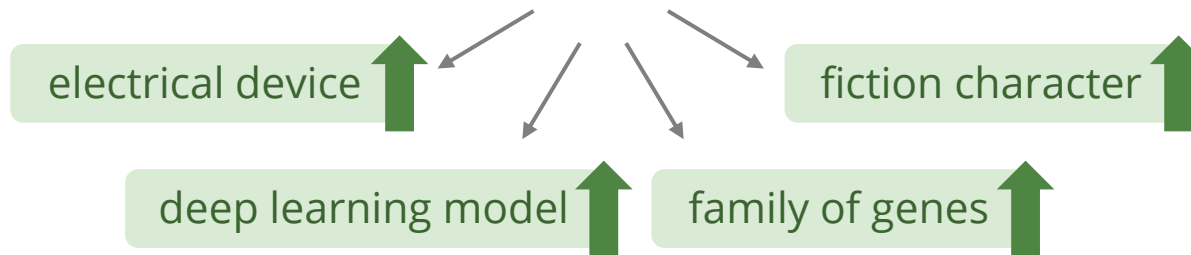| Network architectures | Training frameworks |
|:---:|:---:|
| Multilayer perceptron (MLP) | Autoregressive |
| Convolutional neural networks (CNNs) | Autoencoders |
| Recurrent neural networks (RNNs) | Variational autoencoders (VAEs) |
| Transformers | Generative adversarial networks (GANs) |
| ResNets | Diffusion models |
| U-Nets | Consistency models |
| ⋮ | ⋮ |

# Demystifying Transformers

A transformer is a _____

electrical device

deep learning model     family of genes

fiction character

**Uniform attention**     A     transformer     is     a     ?

**Variable attention**     A     transformer     is     a     ?
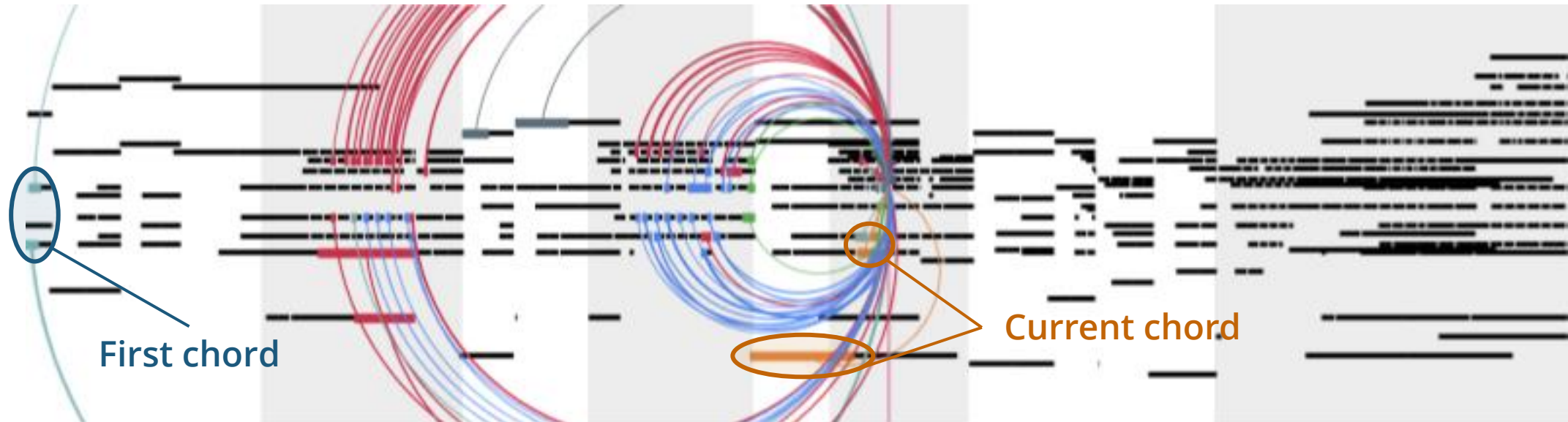
**Transformers learn what to attend to from big data!**

# What does a Transformer Learn?

(Each color represents an attention head)



**First chord**

**Current chord**

(Source: Huang et al., 2018)

Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M. Dai, Matthew D. Hoffman, Monica Dinculescu, and Douglas Eck, "Music Transformer: Generating Music with Long-Term Structure," *Magenta Blog*, December 13, 2018.

# Generating Data from a Random Distribution

**Random distribution**

**Data distribution**
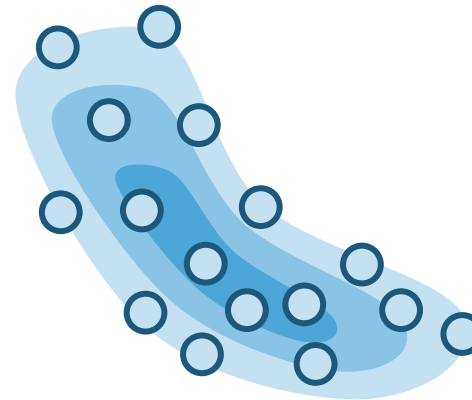
$P(z)$

$P(x)$

**If we can learn this mapping, we can easily generate new samples from the data distribution**

# Variational Autoencoders (VAEs) – Training



Reconstruction loss

KL divergence

Enc

Dec

$P(x)$

$P(\hat{x})$

$P(z)$

Diederik P Kingma and Max Welling, "Auto-Encoding Variational Bayes," *ICLR*, 2014.

# Variational Autoencoders (VAEs) – Generation



$P(z)$

Diederik P Kingma and Max Welling, "Auto-Encoding Variational Bayes," *ICLR*, 2014.

# A Loss Function for Distributions

**Random distribution**

**Data distribution**

**Loss function?**

$P(z)$

$P(\hat{x})$

$P(x)$

**Unfortunately, no easy way to measure the difference between two distributions**

**But what about another neural network!?**

# Generative Adversarial Nets (GANs) – Training



The generator aims to make the fake samples indistinguishable from the real samples for the discriminator

The discriminator aims to tell the fake samples from real samples

Random noise

$z{\sim}p_Z$

Gen

$\log(1 - Dis(Gen(z)))$

Fake samples

$G(z)$

Real samples

$x{\sim}p_X$

Dis

$\log(1 - Dis(x)) + \log(Dis(Gen(z)))$

Real/fake

1/0

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative Adversarial Networks," *NeurIPS*, 2014.

# Generative Adversarial Nets (GANs) – Generation



**Random noise**

$z \sim p_z$

$P(z)$

**Fake samples**

**Gen**

$G(z)$

$x \sim p_X$

**Real samples**

**Dis**

**Real/fake**

1/0

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative Adversarial Networks," *NeurIPS*, 2014.

# Diffusion Models

- **Intuition**: Many denoising autoencoders stacked together



**Denoising**

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

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

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

**Diffusion**

(Source: Ho et al., 2020)

Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising Diffusion Probabilistic Models," *NeurIPS*, 2020.

# Diffusion Models – Training

- **Intuition**: Many denoising autoencoders stacked together



(Source: Ho et al., 2020)

Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising Diffusion Probabilistic Models," *NeurIPS*, 2020.

# Diffusion Models

- **<u>Intuition</u>**: Many denoising autoencoders stacked together

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

**Usually, $T > 1000$**

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

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

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

(Source: Ho et al., 2020)

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

Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising Diffusion Probabilistic Models," *NeurIPS*, 2020.

# Diffusion Models – Generation

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

**Input**

**Output**



**Coarse shapes**
(low-frequency components)

**Fine details**
(high-frequency components)

(Source: Ho et al., 2020)

Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising Diffusion Probabilistic Models," *NeurIPS*, 2020.

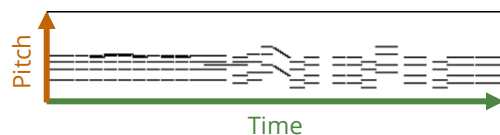# Review – Symbolic Music Generation

# Four Paradigms

**Symbolic music generation**

**Text-based**

**Image-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,
...
```

Pitch

Time

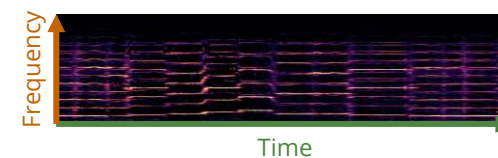MIDI

Piano roll

**Audio-domain music generation**

**Time series-based**

**Image-based**

Frequency

Time

Waveform

Spectrogram

**Today, we also have many latent-space based systems!**

# Language Models (Mathematically)

- A class of machine learning models that learn the next word probability

$$P(\ x_i\ |\ \underbrace{x_1, x_2, \dots, x_{i-1}}\ )$$

Next word          Previous words

$P(\ \text{electrical}\ |\ \text{A transformer is a}\ )$

$P(\ \text{character}\ |\ \text{A transformer is a}\ )$

$P(\quad \text{gene}\quad |\ \text{A transformer is a}\ )$

$P(\quad \text{model}\quad |\ \text{A transformer is a}\ )$

$P(\quad \text{food}\quad |\ \text{A transformer is a}\ )$

$P(\quad \text{musical}\quad |\ \text{A transformer is a}\ )$

# Representing Polyphonic Music

- We can now handle music with multi-pitch at the same time
  - In the literature, "polyphonic" & "multi-pitch" are often used interchangeably



```
Note_on_65, Note_on_68, Time_shift_eighth_note, Note_on_77, Note_on_80,
Time_shift_half_note, Note_off_77, Note_off_80, Note_on_73, Note_on_77,
Time_shift_dotted_quarter_note, Note_off_65, Note_off_68, ...
```
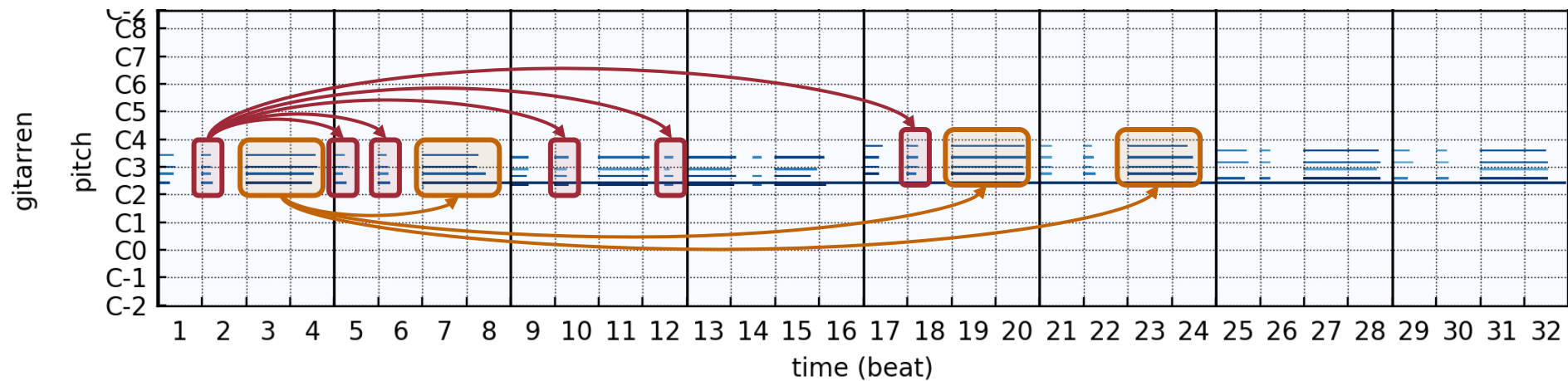
# Example: Performance RNN (Oore et al., 2020)

- **Data**
  - Yamaha e-Piano Competition dataset (MAESTRO)

- **Representation**

  **Examples of generated music**

  - 128 Note-On events
  - 128 Note-Off events
  - 125 Time-Shift events (8ms–1s)
  - 32 Set-Velocity events — Handle dynamics

- **Model**
  - LSTM

Ian Simon and Sageev Oore, "Performance RNN: Generating Music with Expressive Timing and Dynamics," *Magenta Blog*, June 29, 2017.
Sageev Oore, Ian Simon, Sander Dieleman, Douglas Eck, and Karen Simonyan, "This Time with Feeling: Learning Expressive Musical Performance", *Neural Computing and Applications*, 32, 2020.
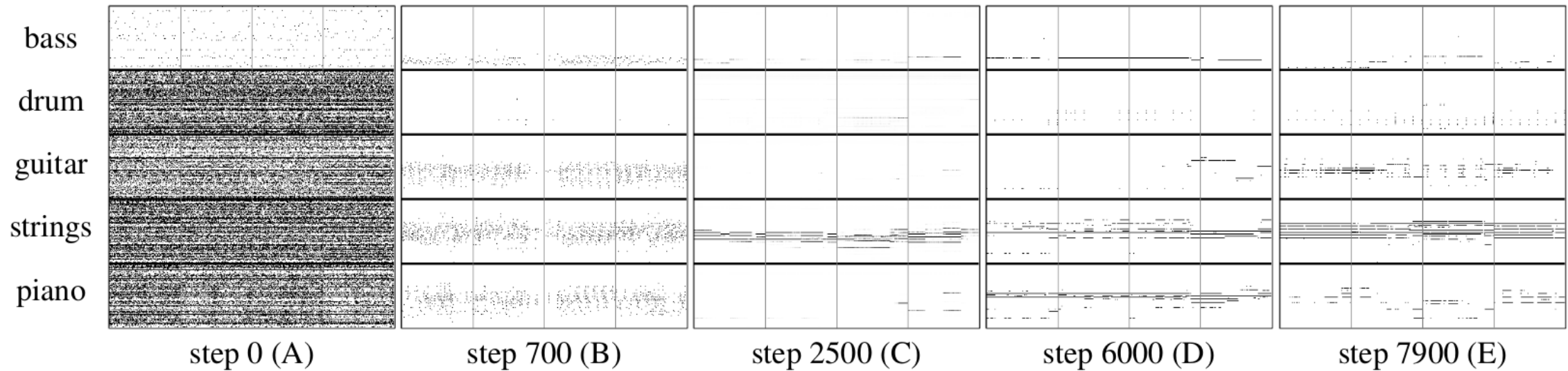
# Why Piano Rolls?



Many musical patterns like melodies, chords, scales and arpeggios
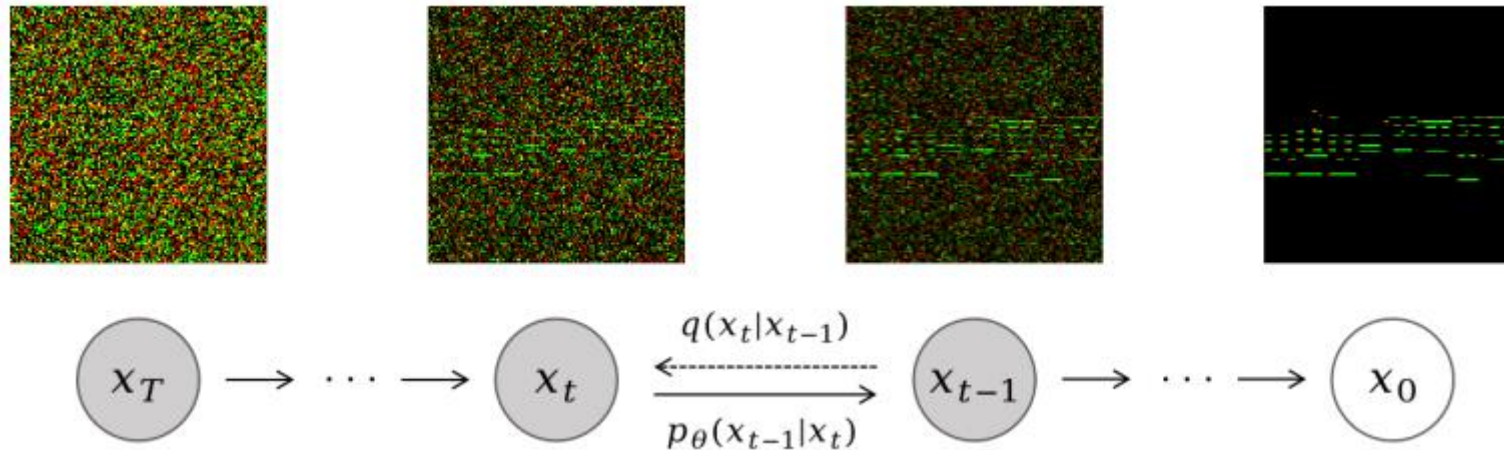are **translational invariant** in the temporal and pitch axes

# Example: MuseGAN (Dong et al., 2018)

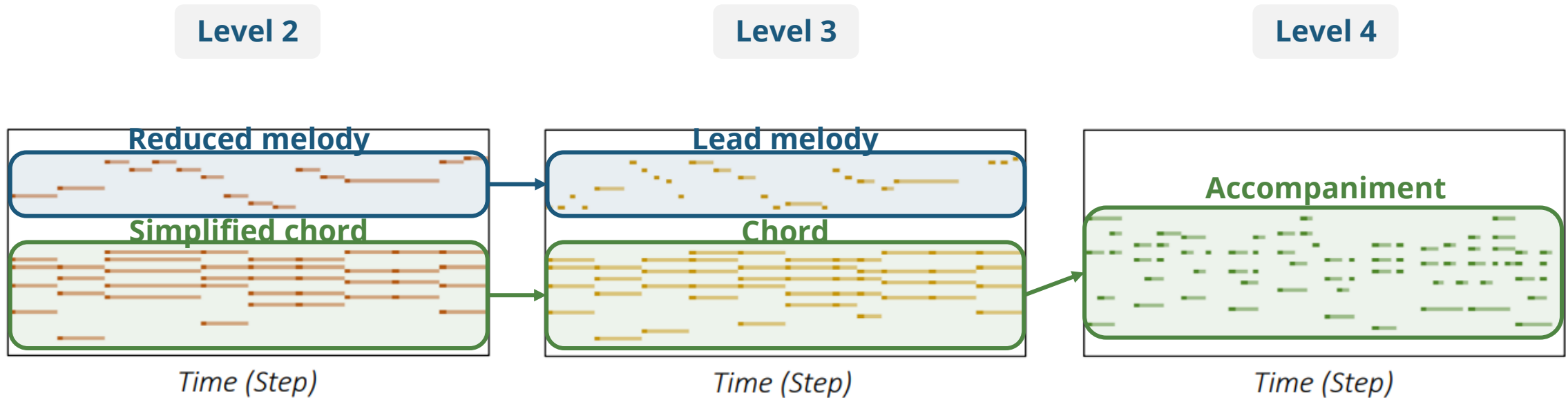**Examples of generated music**



(Source: Dong et al., 2018)

Hao-Wen Dong, Wen-Yi Hsiao, Li-Chia Yang, and Yi-Hsuan Yang, "MuseGAN: Multi-track Sequential Generative Adversarial Networks for Symbolic Music Generation and Accompaniment," *AAAI*, 2018.

# Example: Polyffusion (Min et al., 2023)



(Source: Min et al., 2023)

[polyffusion.github.io](polyffusion.github.io)

Lejun Min, Junyan Jiang, Gus Xia, and Jingwei Zhao, "Polyffusion: A Diffusion Model for Polyphonic Score Generation with Internal and External Controls," *ISMIR*, 2023.

# Example: Cascaded Diffusion Models (Wang et al., 2024)



**Level 2**   **Level 3**   **Level 4**

Reduced melody · Simplified chord → Lead melody · Chord → Accompaniment

(Source: Wang et al., 2024)

[wholesonggen.github.io](wholesonggen.github.io)

Ziyu Wang, Lejun Min, and Gus Xia, "Whole-Song Hierarchical Generation of Symbolic Music Using Cascaded Diffusion Models," *ICLR*, 2024.

# Example: Music SketchNet (Chen et al., 2020)



(Source: Chen et al., 2020)

Ke Chen, Cheng-i Wang, Taylor Berg-Kirkpatrick, and Shlomo Dubnov, "Music FaderNets: Controllable Music Generation Based On High-Level Features via Low-Level Feature Modelling," *ISMIR*, 2020.

# Review – Audio Synthesis

# Four Paradigms

**Symbolic music generation**

**Audio-domain music generation**

| Text-based | Image-based | Time series-based | Image-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,
...
```
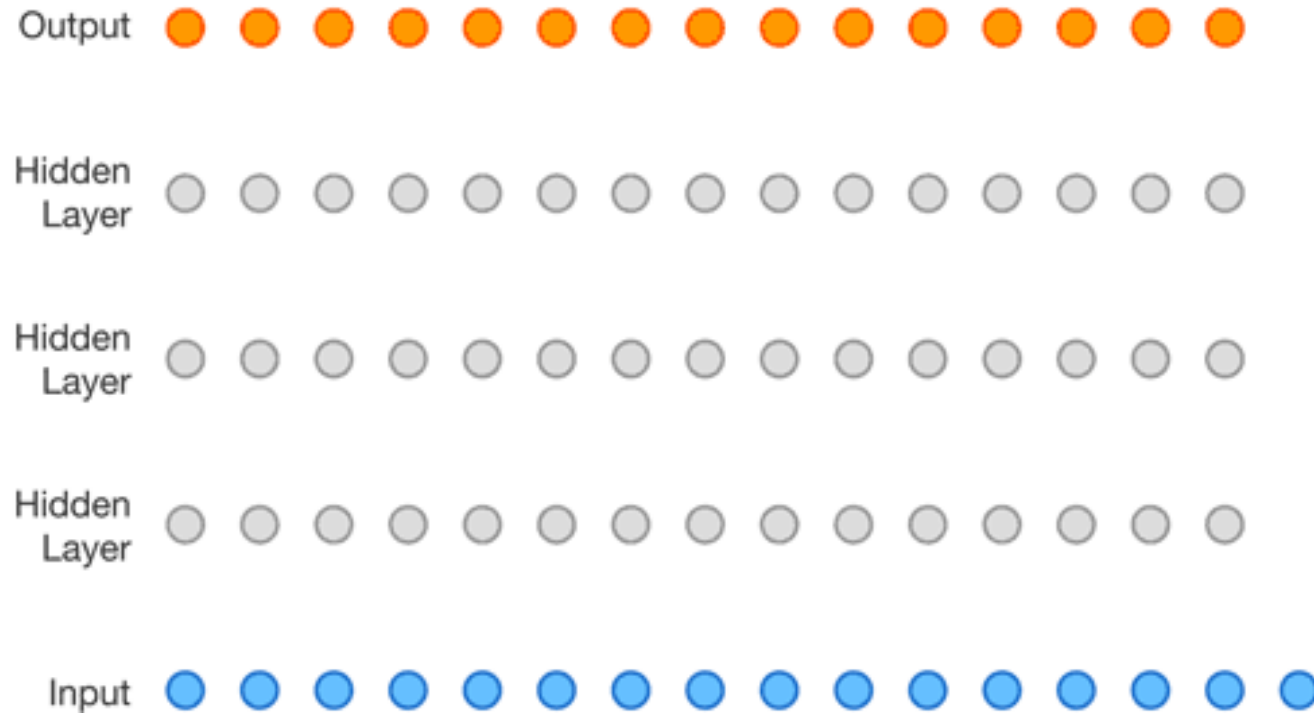
MIDI

Piano roll

Waveform

Spectrogram

**Today, we also have many latent-space based systems!**

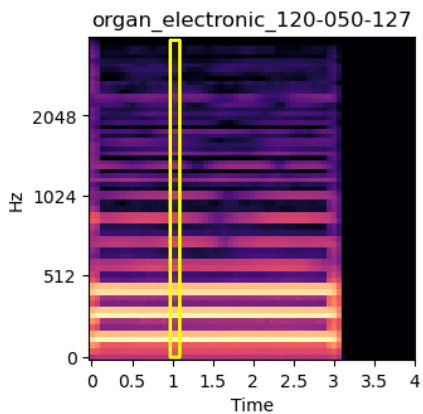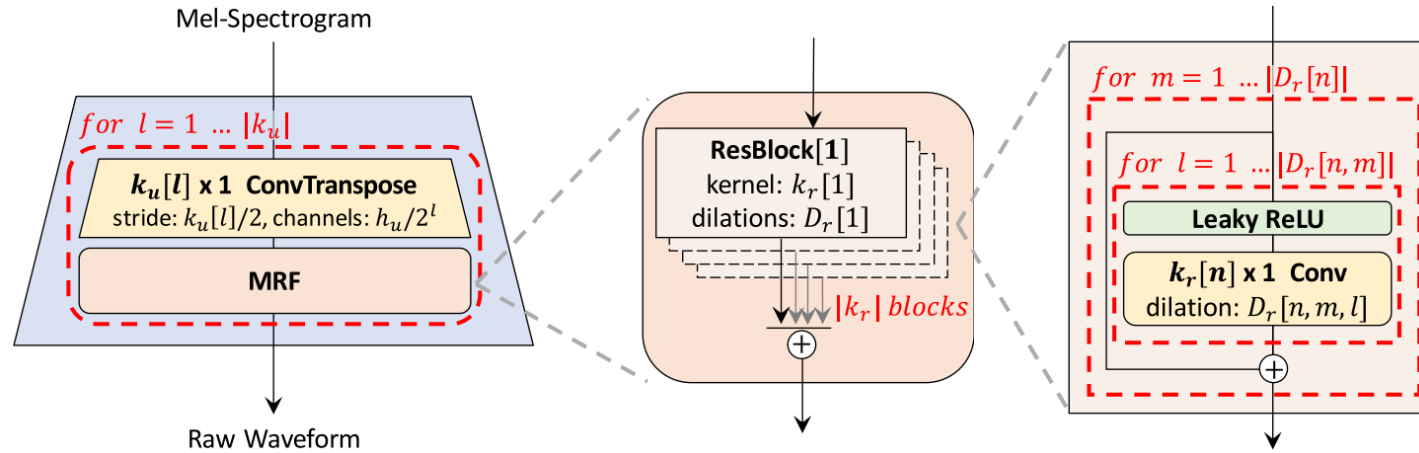# Example: WaveNet (van den Oord et al., 2016)

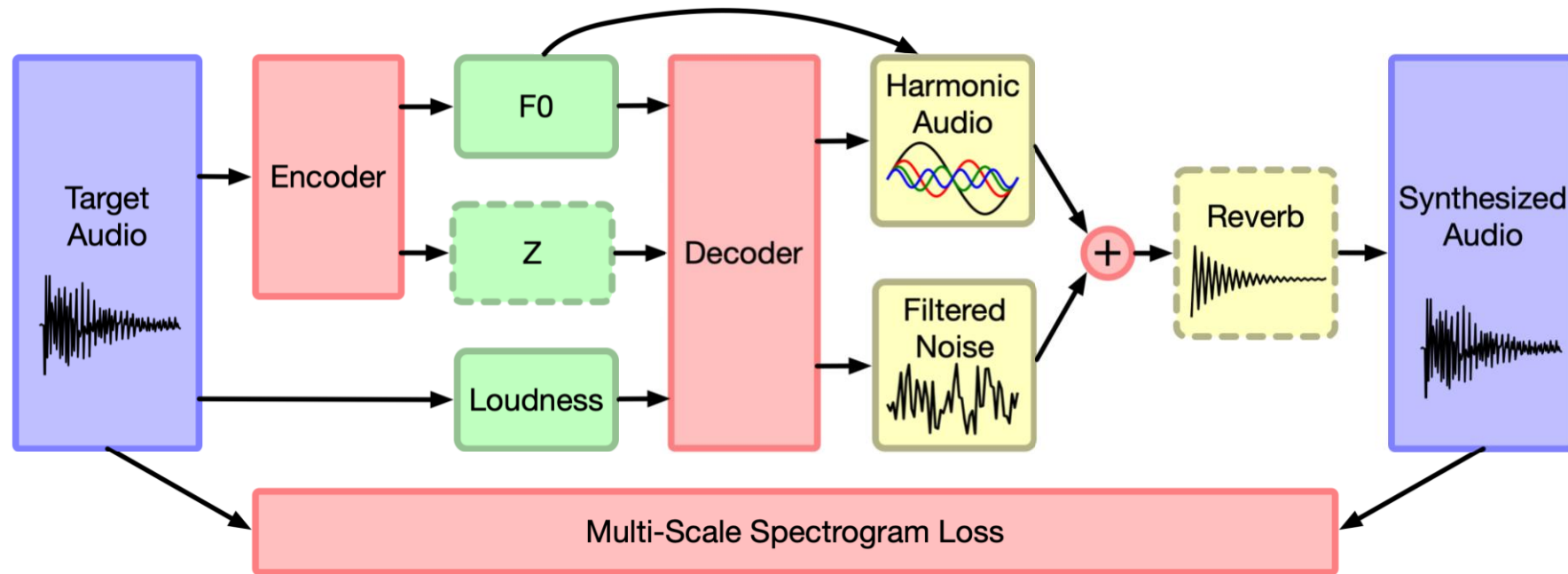

(Source: van den Oord et al., 2016)

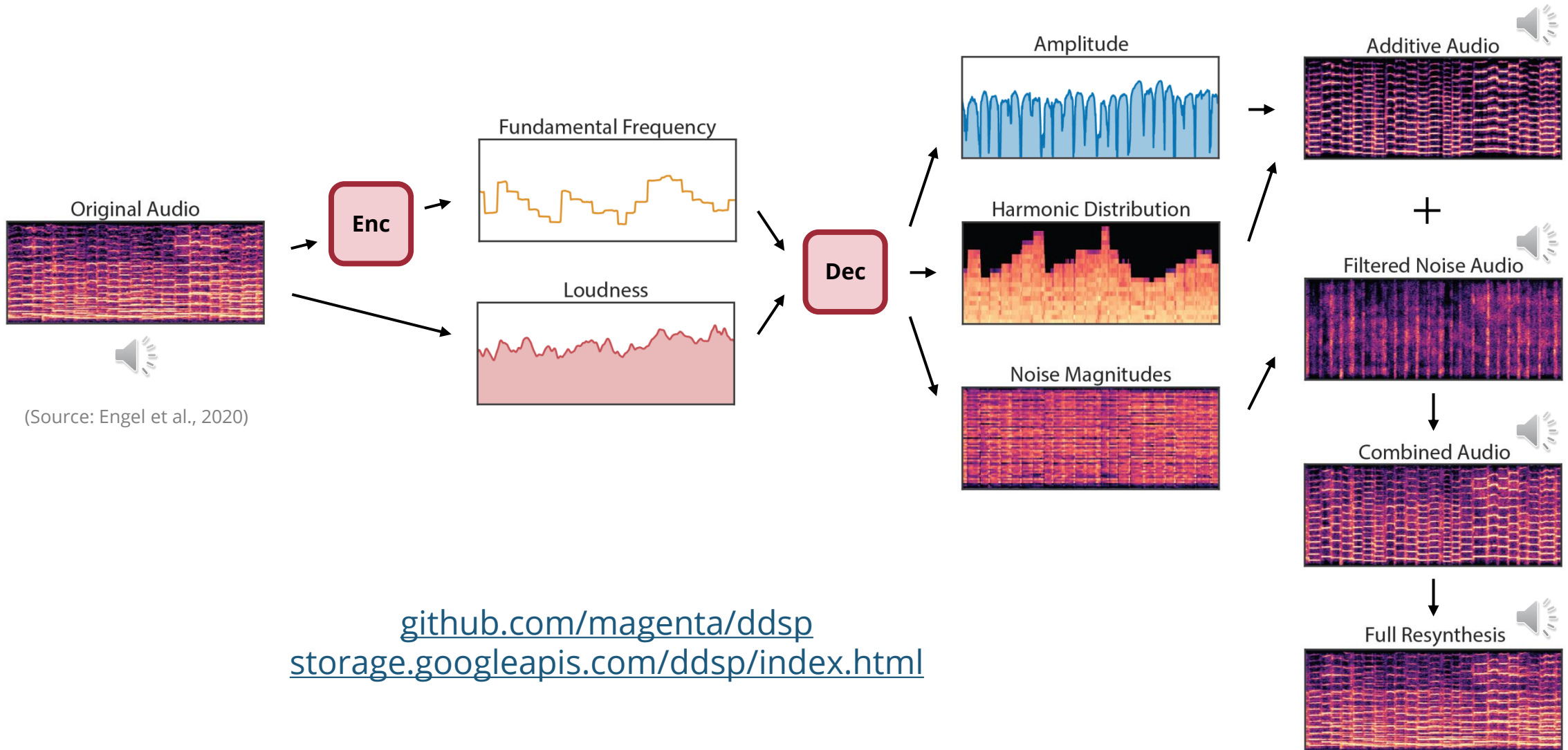**A convolutional neural network for raw waveform generation**

Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu, "WaveNet: A Generative Model for Raw Audio," *ICML*, 2016.

# Example: Hifi-GAN (Kong et al., 2020)



Mel-Spectrogram

$for\ l = 1\ ...\ |k_u|$

$k_u[l]$ x 1 ConvTranspose
stride: $k_u[l]/2$, channels: $h_u/2^l$

MRF

Raw Waveform

ResBlock[1]
kernel: $k_r[1]$
dilations: $D_r[1]$

$|k_r|\ blocks$

$for\ m = 1\ ...\ |D_r[n]|$

$for\ l = 1\ ...\ |D_r[n,m]|$

Leaky ReLU

$k_r[n]$ x 1 Conv
dilation: $D_r[n,m,l]$

organ_electronic_120-050-127

Shrinking channels

Growing resolution

[jik876.github.io/hifi-gan-demo](jik876.github.io/hifi-gan-demo)

Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae, "HiFi-GAN: Generative Adversarial Networks for Efficient and High Fidelity Speech Synthesis," *NeurIPS*, 2020.

# Example: Differentiable DSP (DDSP) (Engel et al., 2020)



(Source: Engel et al., 2020)

Jesse Engel, Lamtharn (Hanoi) Hantrakul, Chenjie Gu, Adam Roberts, "DDSP: Differentiable Digital Signal Processing," *ICLR*, 2020.

# Example: Differentiable DSP (DDSP) (Engel et al., 2020)



(Source: Engel et al., 2020)

github.com/magenta/ddsp
storage.googleapis.com/ddsp/index.html

Jesse Engel, Lamtharn (Hanoi) Hantrakul, Chenjie Gu, Adam Roberts, "DDSP: Differentiable Digital Signal Processing," *ICLR*, 2020.

# Review – Latent-based Music & Audio Synthesis

# Pipeline



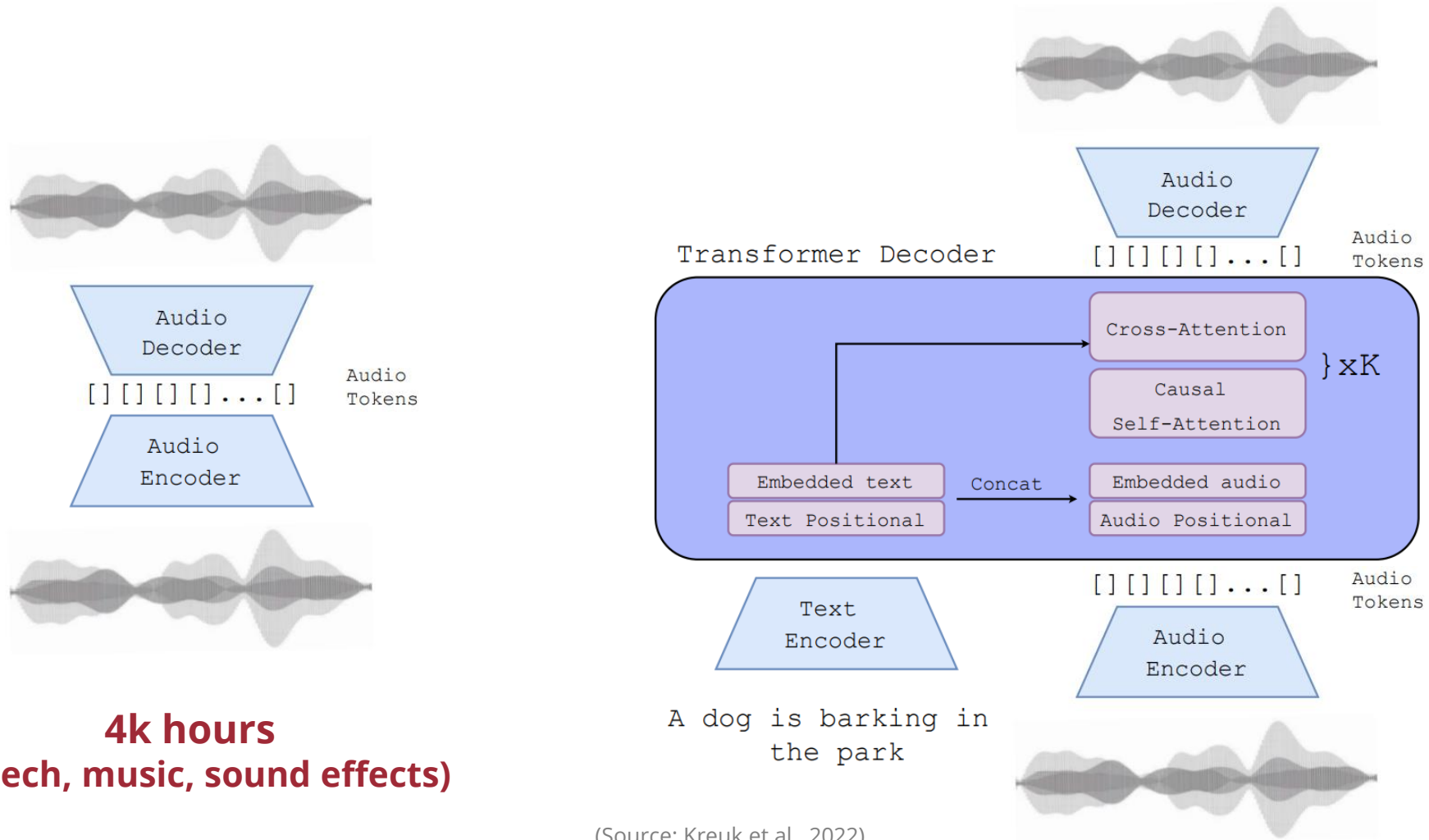**Step 1: Train an Autoencoder**

**Step 2: Compute the Latent Vectors**

**Step 3: Train a Latent Generative Model**

**Step 4: Decode the Latent Vectors**

# Example: AudioGen (Kreuk et al., 2023)
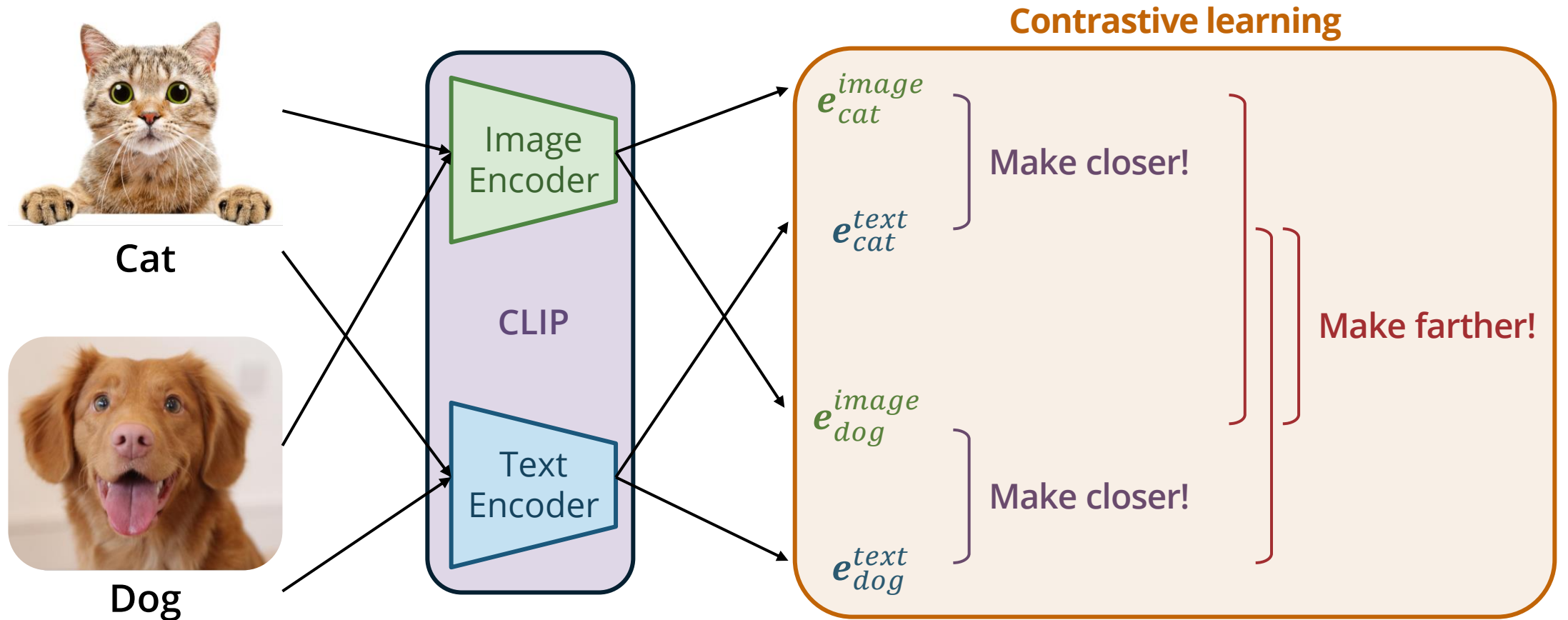


**4k hours**
**(speech, music, sound effects)**

(Source: Kreuk et al., 2022)

Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi, "AudioGen: Textually Guided Audio Generation," *ICLR*, 2023.

# Example: MusicGen (Copet et al., 2023)

- AudioGen for Music

- Use EnCodec (Défossez et al., 2022) as the autoencoder
  - instead of SoundStream for AudioGen (Kreuk et al., 2023)

- **20k hours** of licensed music
  - Internal dataset      10k      High-quality (private)
  - ShutterStock      25k      Instrument-only
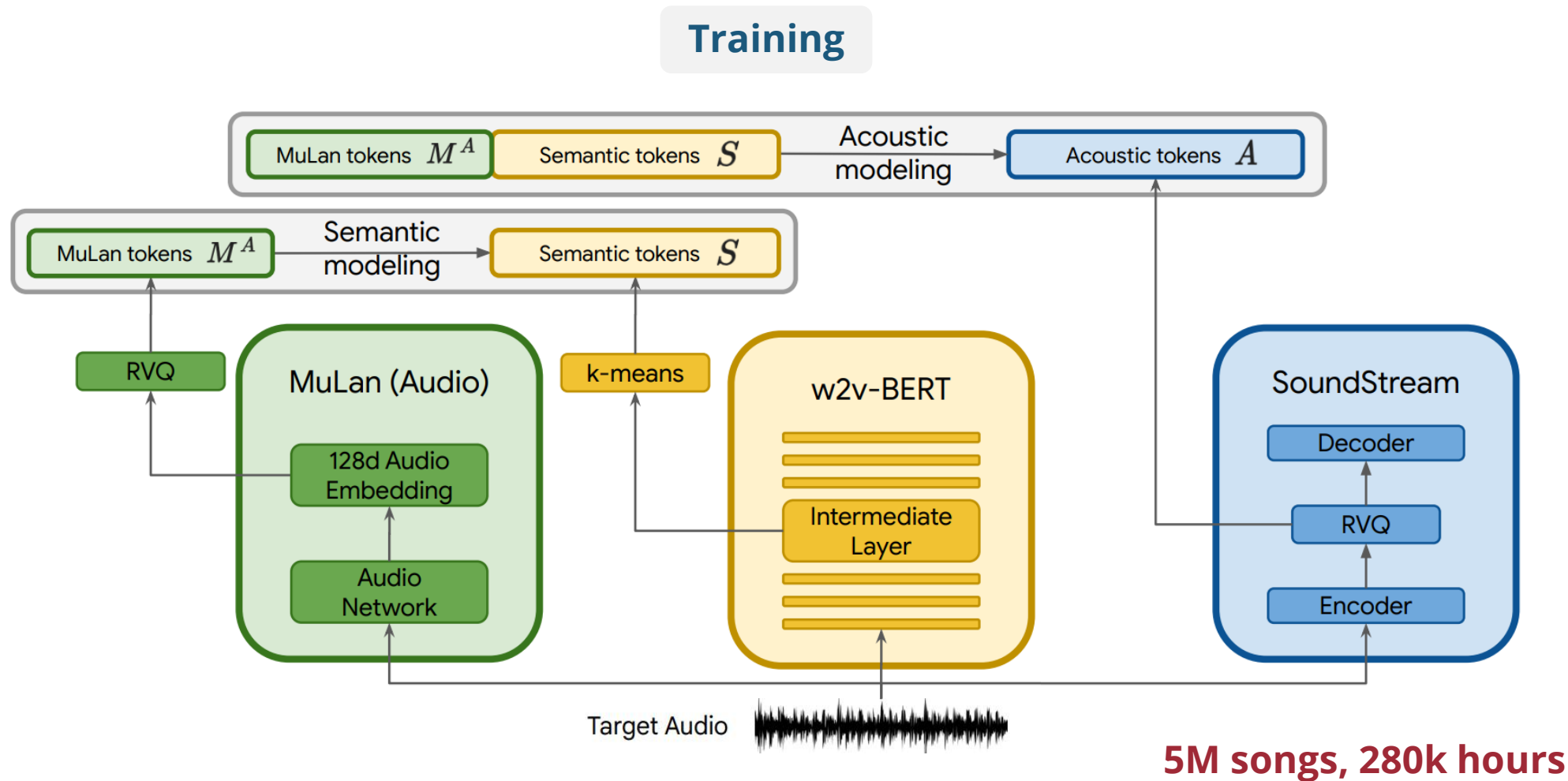  - Pond5      365k      Instrument-only

ai.honu.io/papers/musicgen/

Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez, "Simple and Controllable Music Generation," *NeurIPS*, 2023.

# Contrastive Language-Image Pretraining (CLIP)



**Contrastive learning**

$e_{cat}^{image}$

**Make closer!**

$e_{cat}^{text}$

**Make farther!**

$e_{dog}^{image}$

**Make closer!**

$e_{dog}^{text}$

**Learn a shared embedding space for images and texts**

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever, "Learning Transferable Visual Models From Natural Language Supervision," *ICML*, 2021.

# Example: MusicLM (Agostinelli et al., 2023)



(Source: Agostinelli et al., 2022)

**5M songs, 280k hours**

Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, Matt Sharifi, Neil Zeghidour, and Christian Frank, "MusicLM: Generating Music From Text," *arXiv preprint arXiv:2301.11325*, 2023.

# Example: MusicLM (Agostinelli et al., 2023)



(Source: Agostinelli et al., 2022)

google-research.github.io/seanet/musiclm/examples/

Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, Matt Sharifi, Neil Zeghidour, and Christian Frank, "MusicLM: Generating Music From Text," *arXiv preprint arXiv:2301.11325*, 2023.

# Music FX (2024)



aitestkitchen.withgoogle.com/tools/music-fx

# Music FX DJ (2024)



aitestkitchen.withgoogle.com/tools/music-fx-dj

59

# Example: MusicLDM (Chen et al., 2023)



(Source: Ke et al., 2023)

musicldm.github.io

Ke Chen, Yusong Wu, Haohe Liu, Marianna Nezhurina, Taylor Berg-Kirkpatrick, and Shlomo Dubnov, "MusicLDM: Enhancing Novelty in Text-to-Music Generation Using Beat-Synchronous Mixup Strategies," *ICASSP*, 2024.

# Example: MusicLDM (Chen et al., 2023)



youtu.be/DALv7ea6cv0

61

# Music ControlNet vs DITTO



**Music ControlNet**

**Needs some training!**

(Source: Wu et al., 2024)

**DITTO**

**No training needed!**

(Source: Novack et al., 2024)

Shih-Lun Wu, Chris Donahue, Shinji Watanabe, and Nicholas J. Bryan, "Music ControlNet: Multiple Time-varying Controls for Music Generation," *TASLP*, 2024.
Zachary Novack, Julian McAuley, Taylor Berg-Kirkpatrick, and Nicholas J. Bryan, "DITTO: Diffusion Inference-Time T-Optimization for Music Generation," *ICML*, 2024.

# Review – Neural Audio Effects

# Example: Differentiable Auto-mixing (Steinmetz et al., 2021)



(Source: Steinmetz et al., 2021)

Christian J. Steinmetz, Jordi Pons, Santiago Pascual, and Joan Serrà, "Automatic multitrack mixing with a differentiable mixing console of neural audio effects," *ICASSP*, 2021.

# Example: Differentiable Auto-mixing (Steinmetz et al., 2021)



(Source: Steinmetz et al., 2021)

(Source: Steinmetz et al., 2021)

**A differentiable (and thus trainable) mixing console!**

github.com/csteinmetz1/pymixconsole

Christian J. Steinmetz, Jordi Pons, Santiago Pascual, and Joan Serrà, "Automatic multitrack mixing with a differentiable mixing console of neural audio effects," *ICASSP*, 2021.

# Example: Audio Processing Graph (Lee et al., 2022)



**Can we predict the audio processing graph used in a reference recording?**

(Source: Lee et al., 2023)

Sungho Lee, Jaehyun Park, Seungryeol Paik, and Kyogu Lee, "Blind Estimation of Audio Processing Graph," *ICASSP*, 2023.

# Example: Audio Processing Graph (Lee et al., 2022)



**Blind estimation framework**

Sources are not given!

**Prototype decoder**

**Parameter estimator**

(Source: Lee et al., 2023)

Sungho Lee, Jaehyun Park, Seungryeol Paik, and Kyogu Lee, "Blind Estimation of Audio Processing Graph," *ICASSP*, 2023.

# Review – Interactive & Multimodal Systems

# Example: RAVE (Caillon & Esling, 2022)

**16-band decomposition, 48kHz**



| Model | CPU synthesis | GPU synthesis |
|---|---|---|
| NSynth | 18 Hz | 57 Hz |
| SING | 304 kHz | 9.8 MHz |
| RAVE (Ours) w/o multiband | 38 kHz | 3.7 MHz |
| **RAVE (Ours)** | **985 kHz** | **11.7 MHz** |

**Realtime capable on CPUs & GPUs**

[anonymous84654.github.io/RAVE_anonymous](anonymous84654.github.io/RAVE_anonymous)

# Example: A.I. Duet (Mann et al, 2016)



youtu.be/0ZE1bfPtvZo
experiments.withgoogle.com/ai/ai-duet/view

# Example: Piano Genie (Donahue et al., 2018)



piano-genie.glitch.me/



youtu.be/YRb0XAnUpIk & magenta.tensorflow.org/pianogenie

Chris Donahue, Ian Simon, and Sander Dieleman, "Piano Genie," IUI, 2019.

# Example: Piano Genie (Donahue et al., 2018)

**Input melody**

**Baseline**

**Proposed**

Encoder

Decoder

(Source: Donahue et al., 2019)

Chris Donahue, Ian Simon, and Sander Dieleman, "Piano Genie," *IUI*, 2019.

# Example: Dance-to-music Generation (Li et al., 2024)



(Source: Li et al., 2024)

Sifei Li, Weiming Dong, Yuxin Zhang, Fan Tang, Chongyang Ma, Oliver Deussen, Tong-Yee Lee, and Changsheng Xu, "Dance-to-Music Generation with Encoder-based Textual Inversion," *SIGGRAPH ASIA*, 2024.

# Example: MovieGen (2024)



(Source: Movie Gen Team, 2024)

ai.meta.com/research/movie-gen/

Movie Gen Team, "Movie Gen: A Cast of Media Foundation Models," *arXiv preprint arXiv:2410.13720*, 2024.

# Example: Brain2Music (Denk et al., 2023)

**Brain as the encoder!**



(Source: Denk et al., 2023)

**Can we decode human brain-encoded music?**

Timo I. Denk, Yu Takagi, Takuya Matsuyama, Andrea Agostinelli, Tomoya Nakai, Christian Frank, and Shinji Nishimoto, "Brain2Music: Reconstructing Music from Human Brain Activity," *arXiv preprint arXiv:2307.11078*, 2023.

# Music & AI

# Music & Technology

# Building Blocks of Modern AI Systems



**Data** × **Model** × **Use Case**

**Analysis**

**Retrieval**

**Creation**

# Final Thoughts