

PAT 498/598 (Fall 2024)

Special Topics: Generative AI for Music and Audio Creation

Lecture 14: Controllable Music Generation

Instructor: Hao-Wen Dong



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

(Recap) Four Paradigms



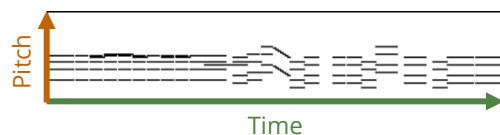
Symbolic music generation

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,  
...
```

MIDI



Piano roll



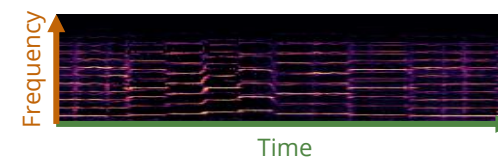
Audio-domain music generation

Time series-based

Image-based



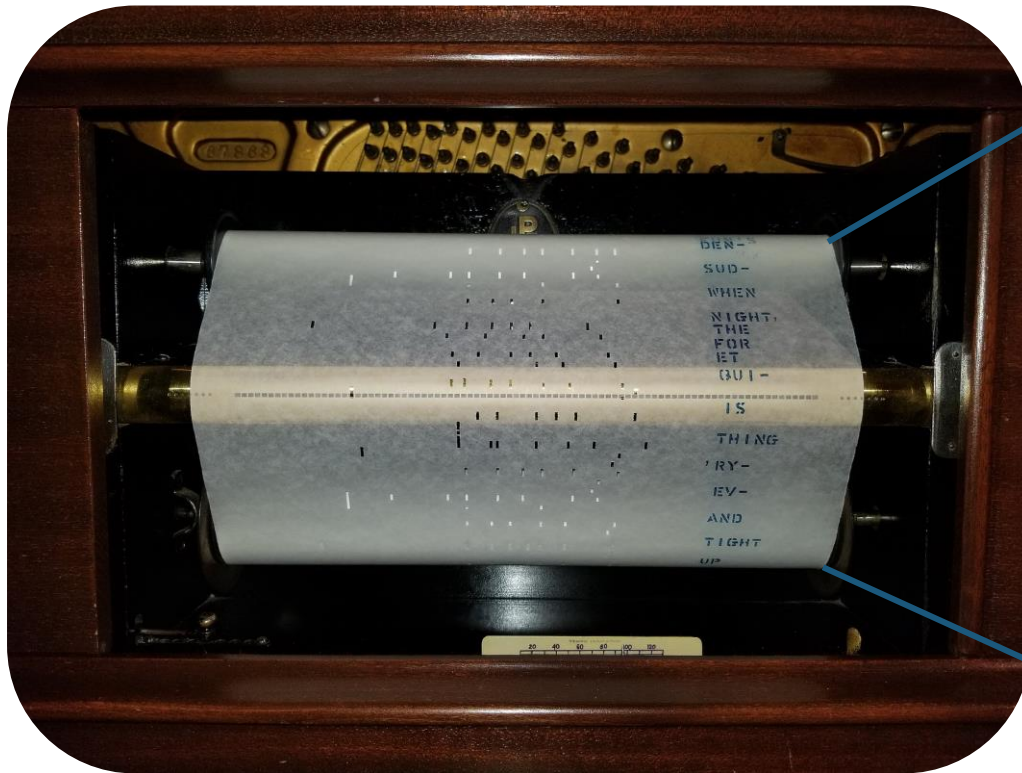
Waveform



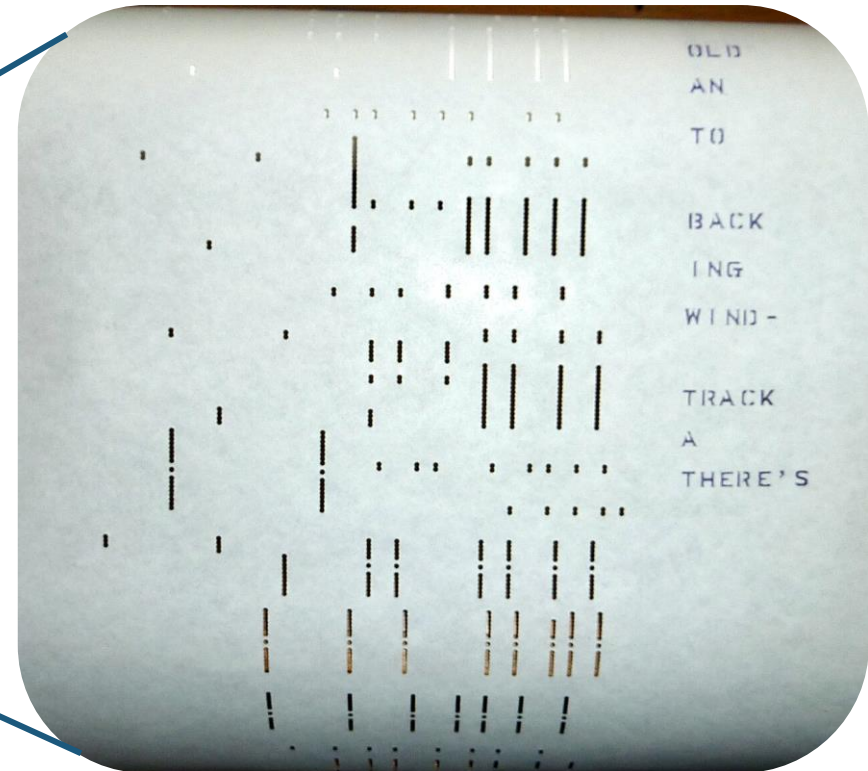
Spectrogram

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

(Recap) Piano Rolls

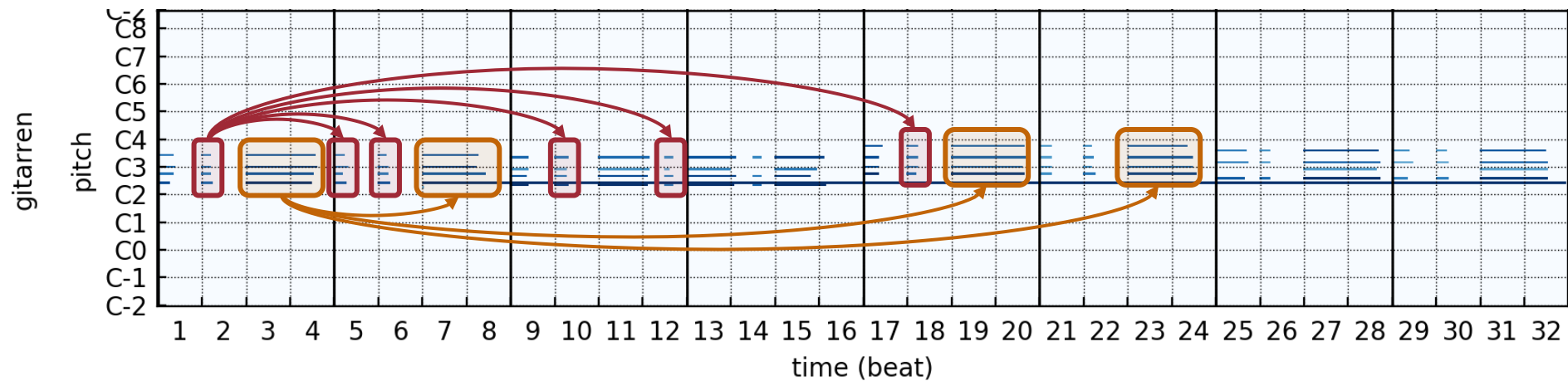


(Source: Draconichiaro)



(Source: Tangerineduel)

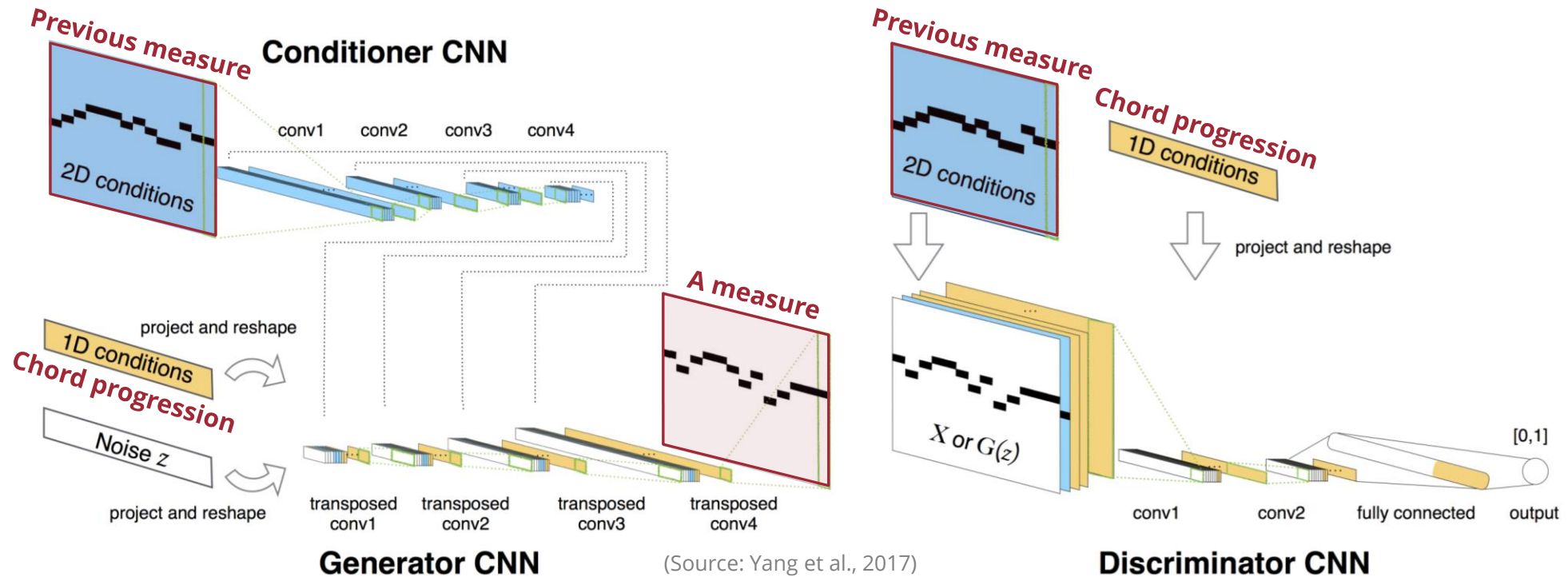
(Recap) Why Piano Rolls?



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

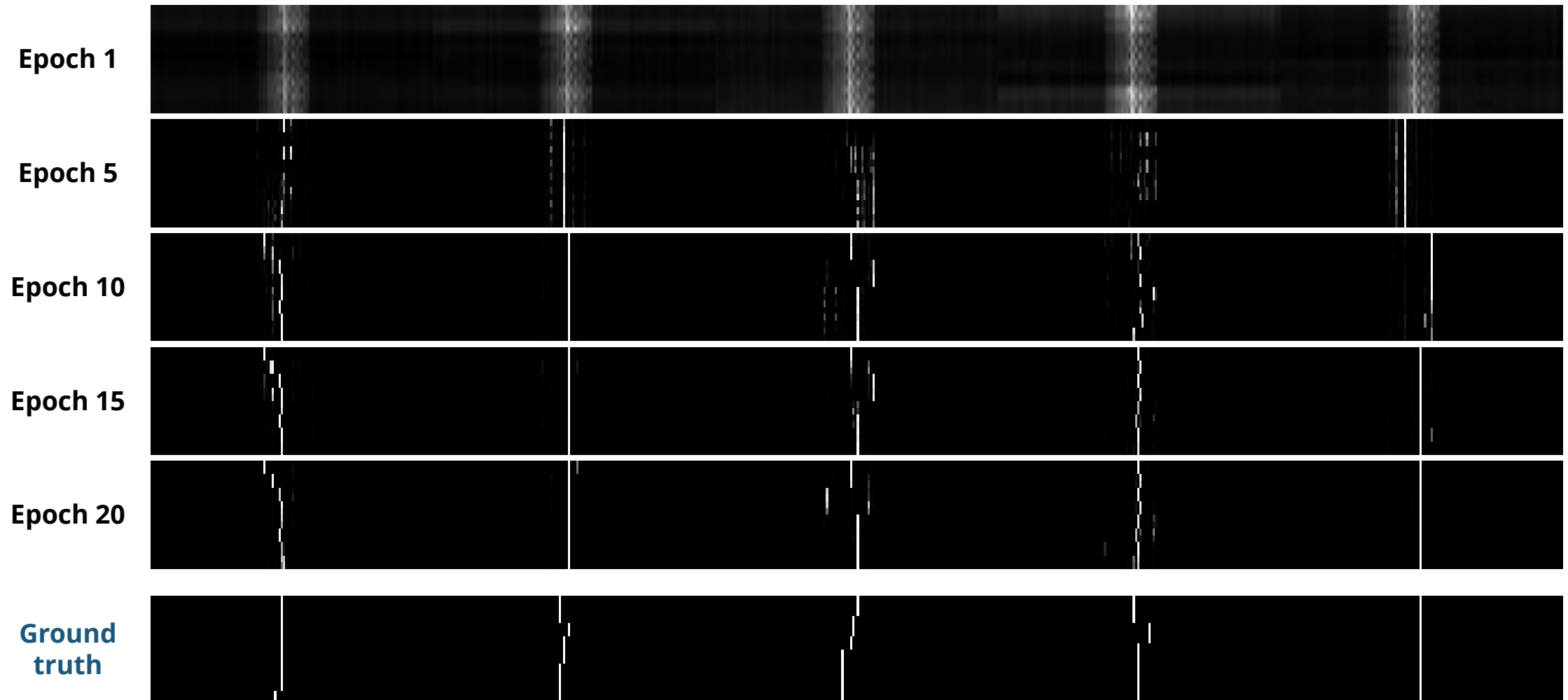
(Recap) Example: **MidiNet** (Yang et al., 2017)

Examples of generated music



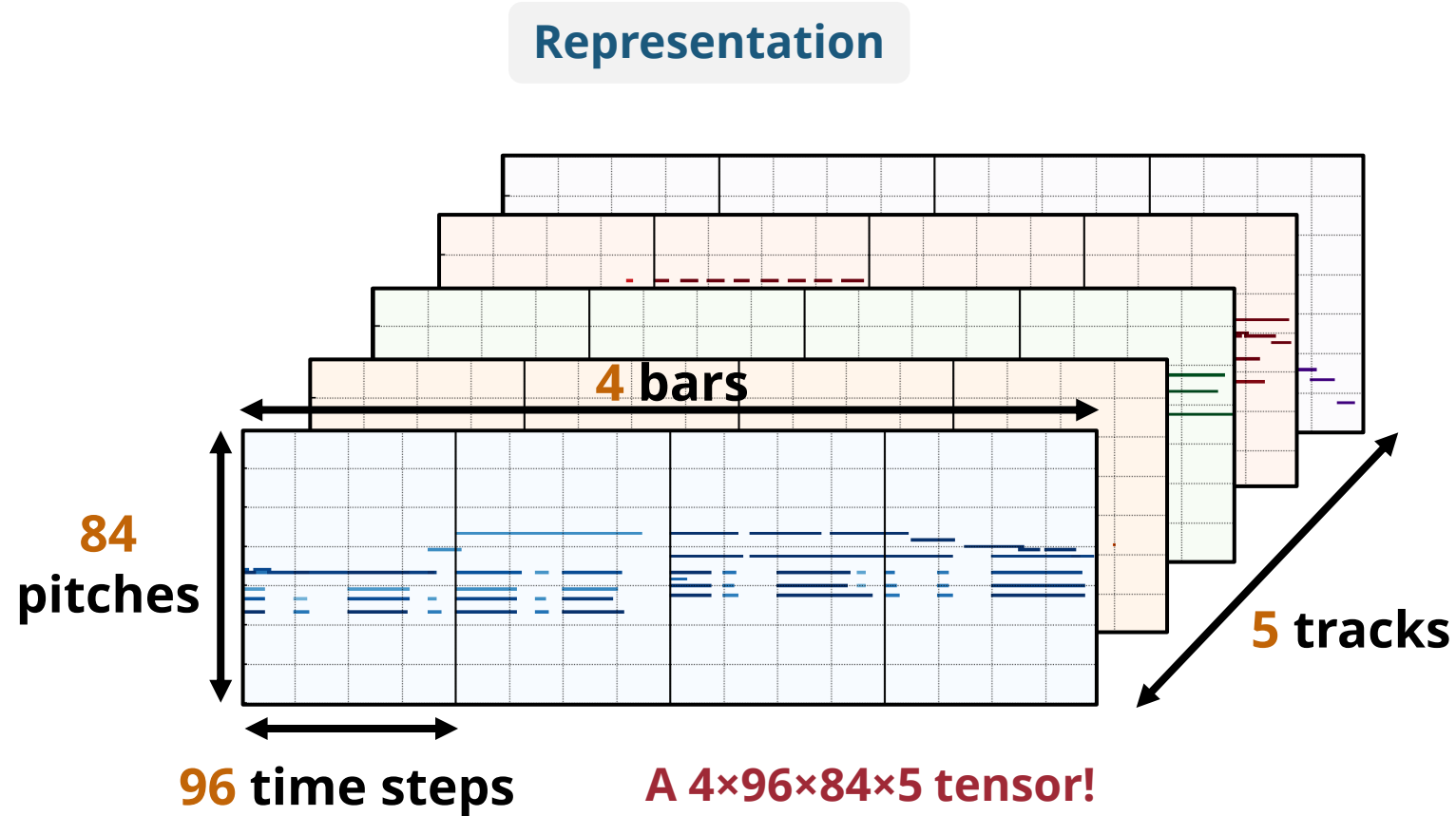
MidiNet generates music measure-by-measure by conditioning on the last measure generated

(Recap) Example: **MidiNet** (Yang et al., 2017)

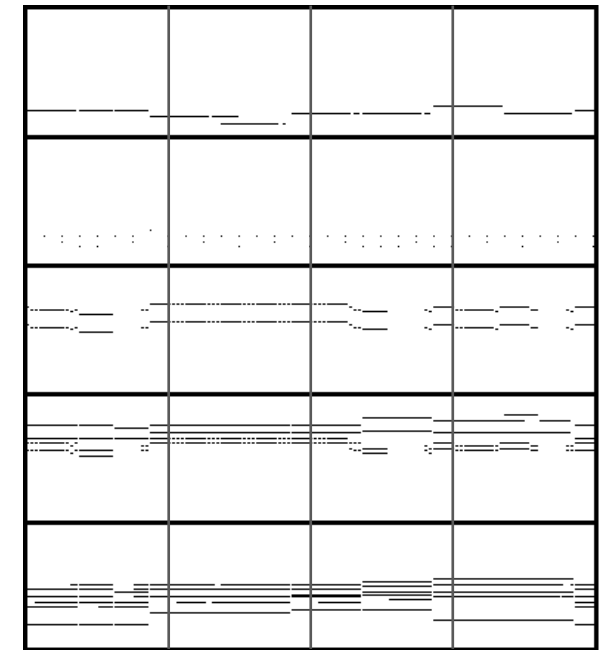


(Source: Yang et al., 2017)

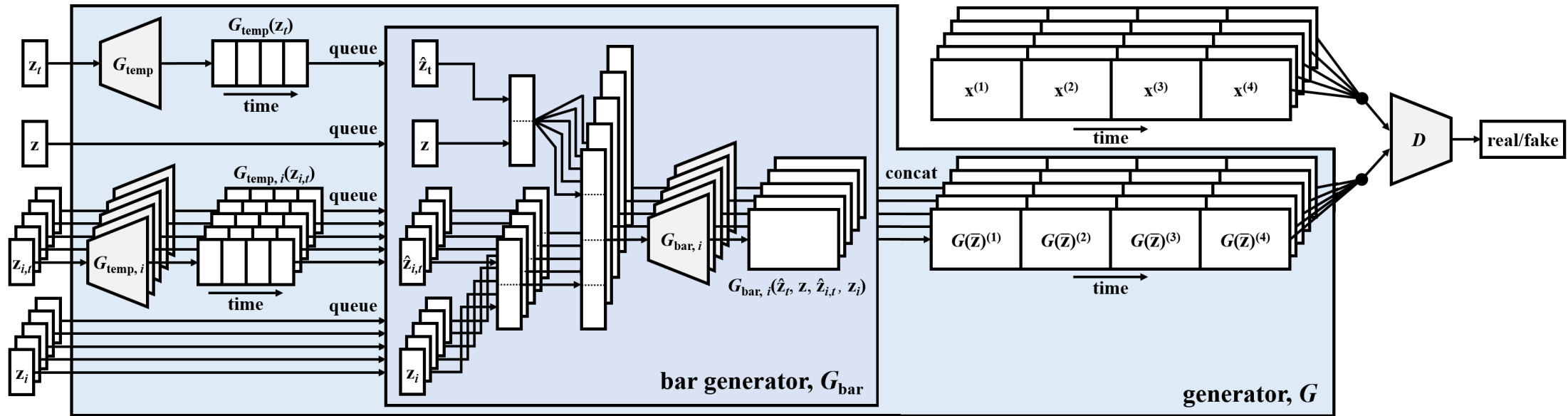
(Recap) Example: MuseGAN (Dong et al., 2018)



A training sample



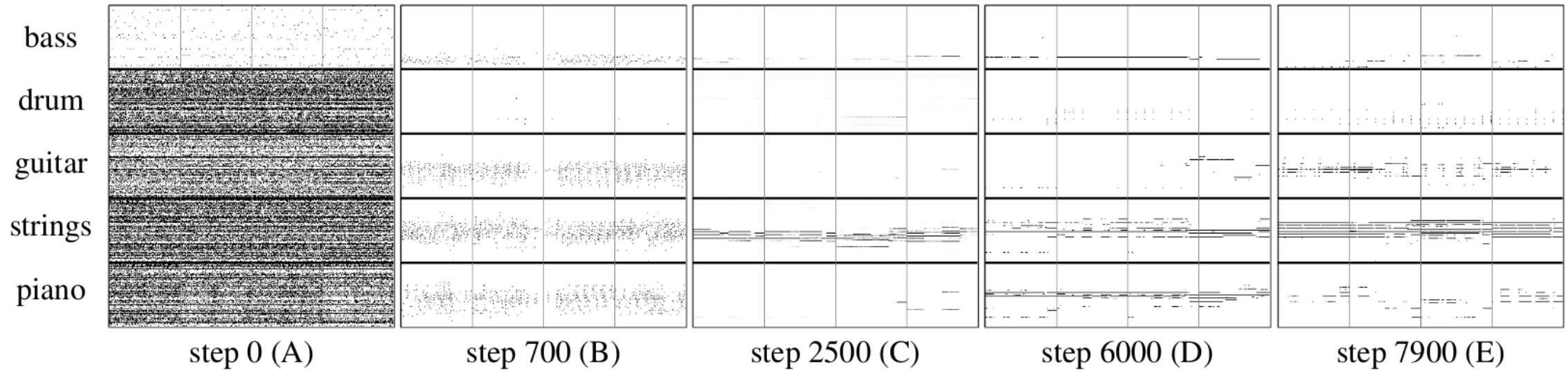
(Recap) Example: MuseGAN (Dong et al., 2018)



(Source: Dong et al., 2018)

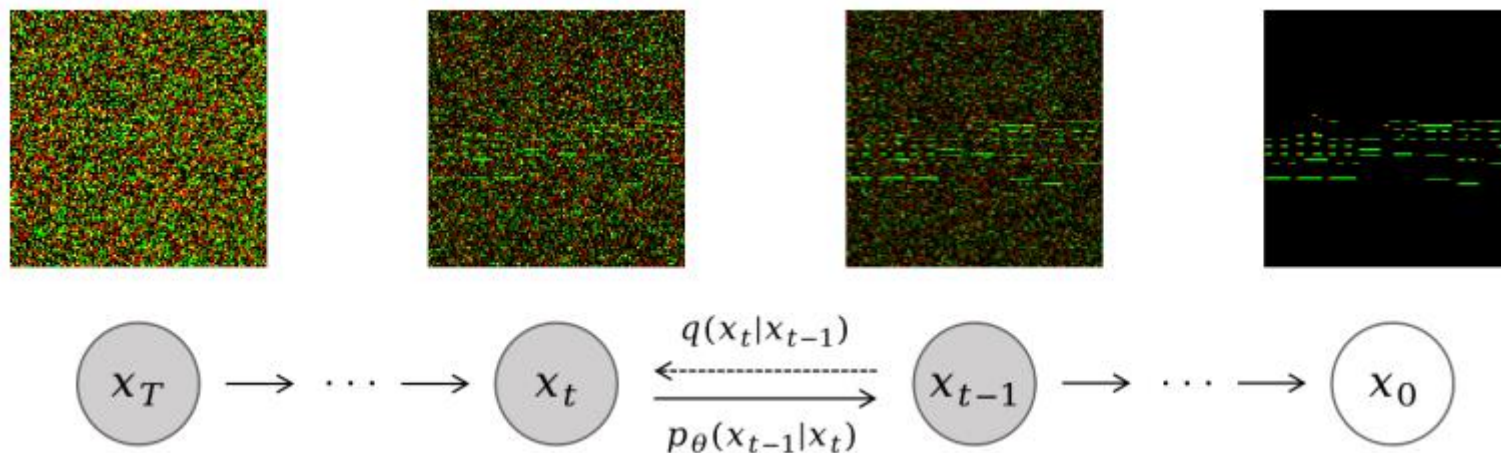
(Recap) Example: MuseGAN (Dong et al., 2018)

Examples of generated music



(Source: Dong et al., 2018)

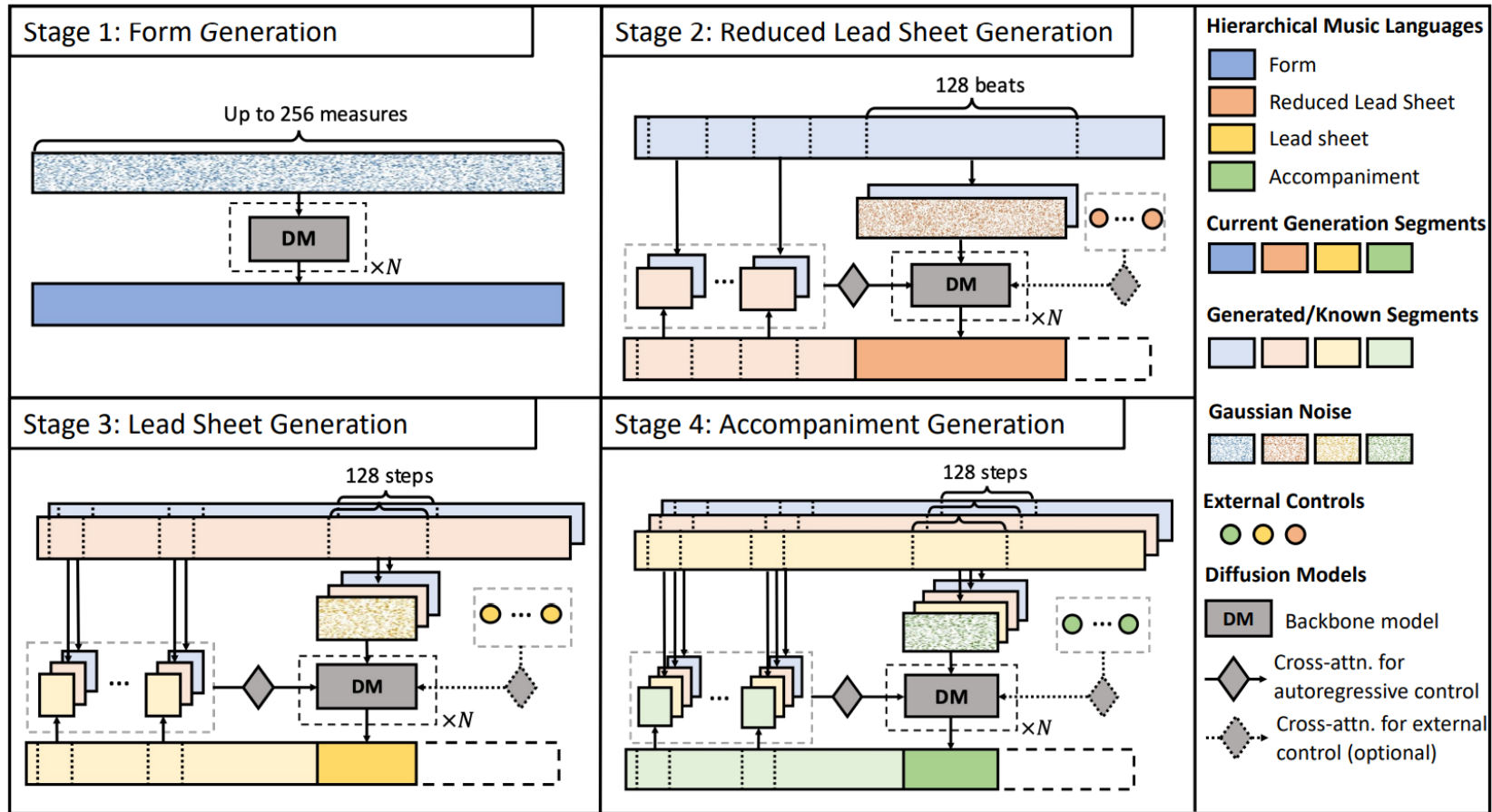
(Recap) Example: Polyffusion (Min et al., 2023)



(Source: Min et al., 2023)

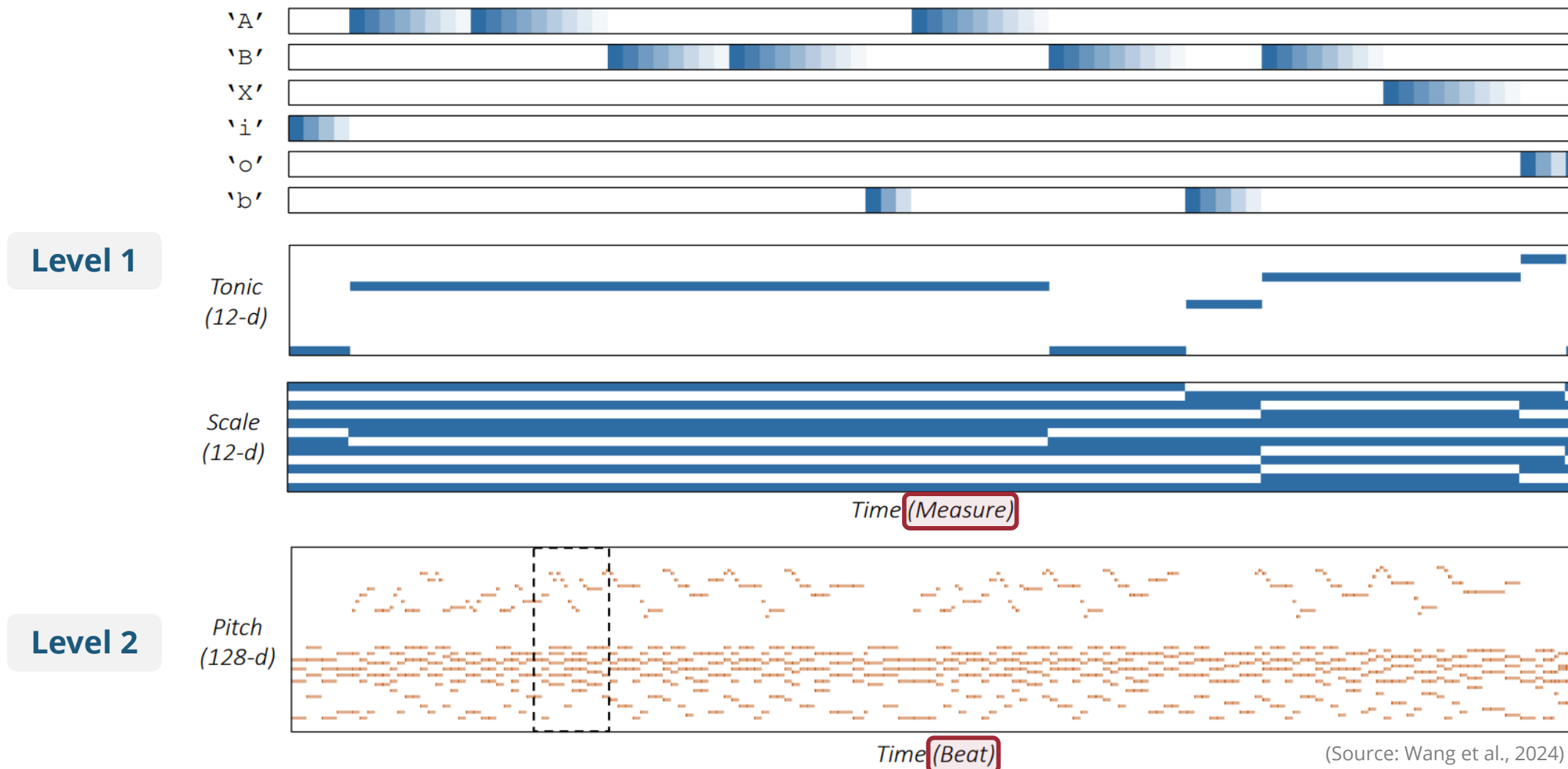
polyffusion.github.io

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



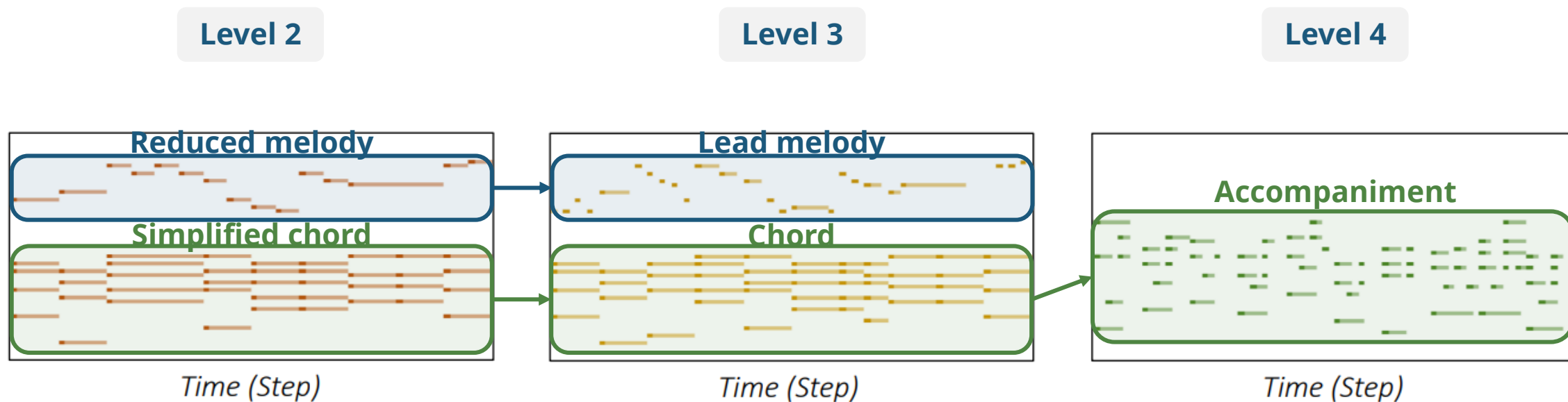
(Source: Wang et al., 2024)

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



(Source: Wang et al., 2024)

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

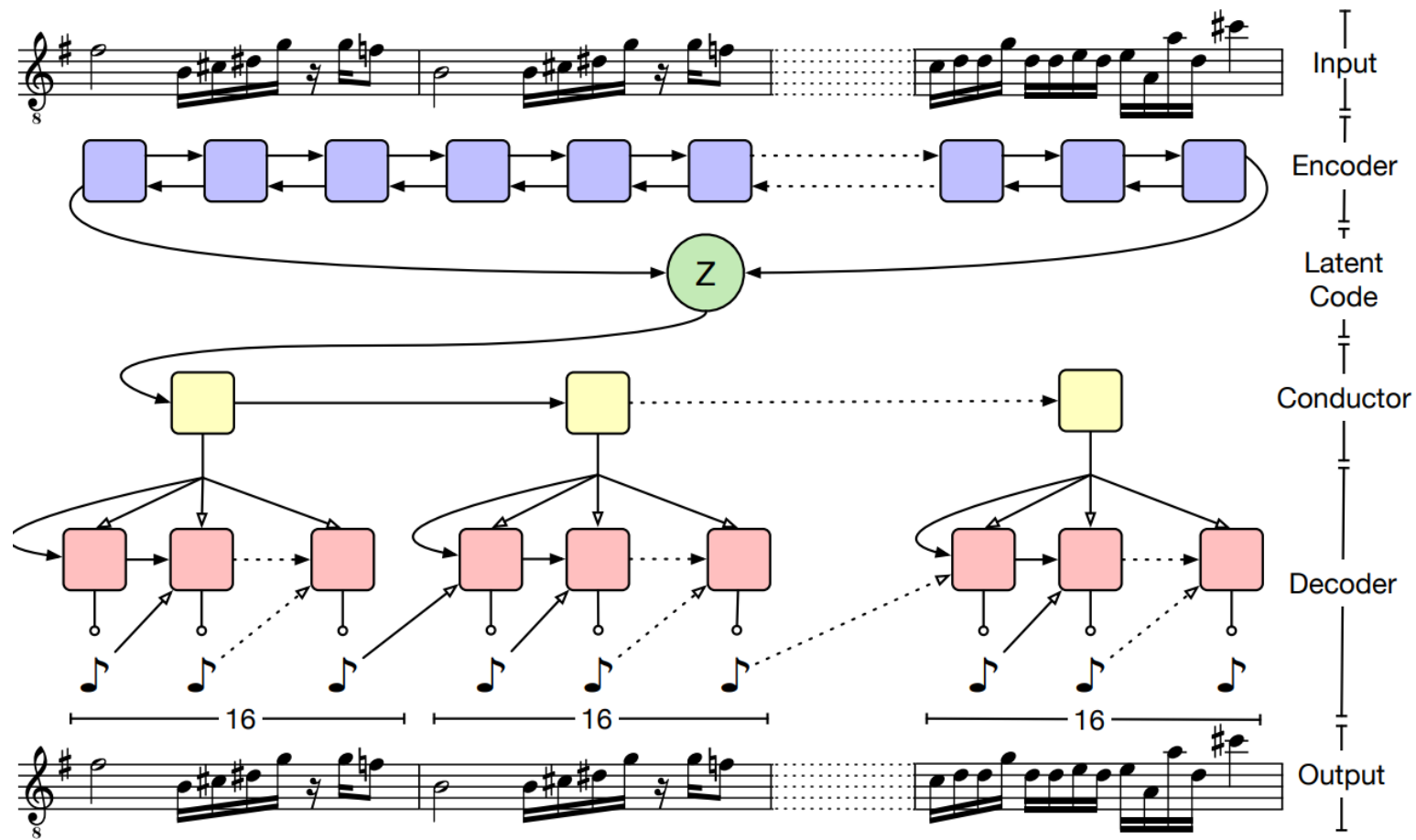


(Source: Wang et al., 2024)

wholesonggen.github.io

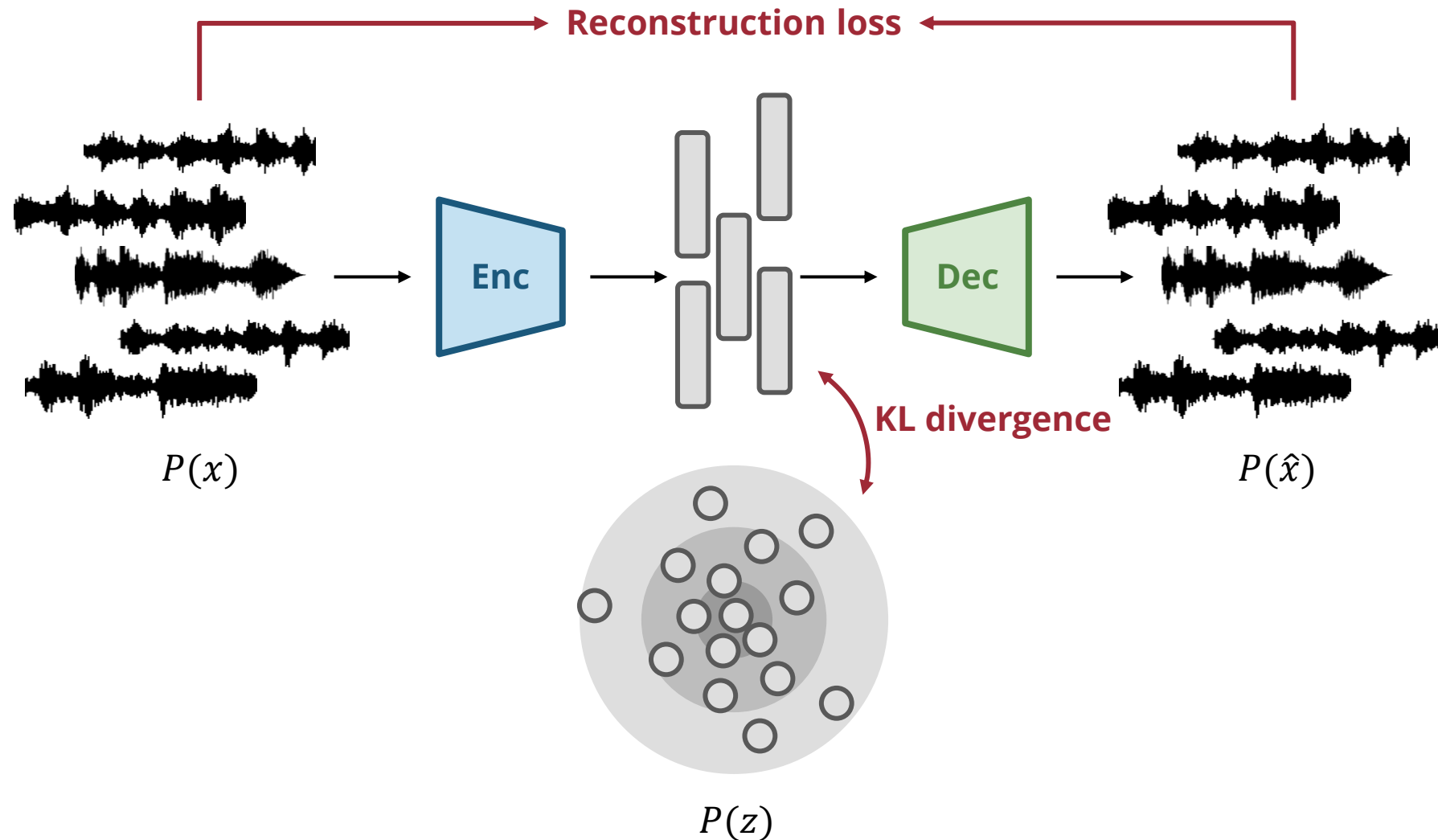
Latent Space-based Music Generation

Example: MusicVAE (Roberts et al., 2018)

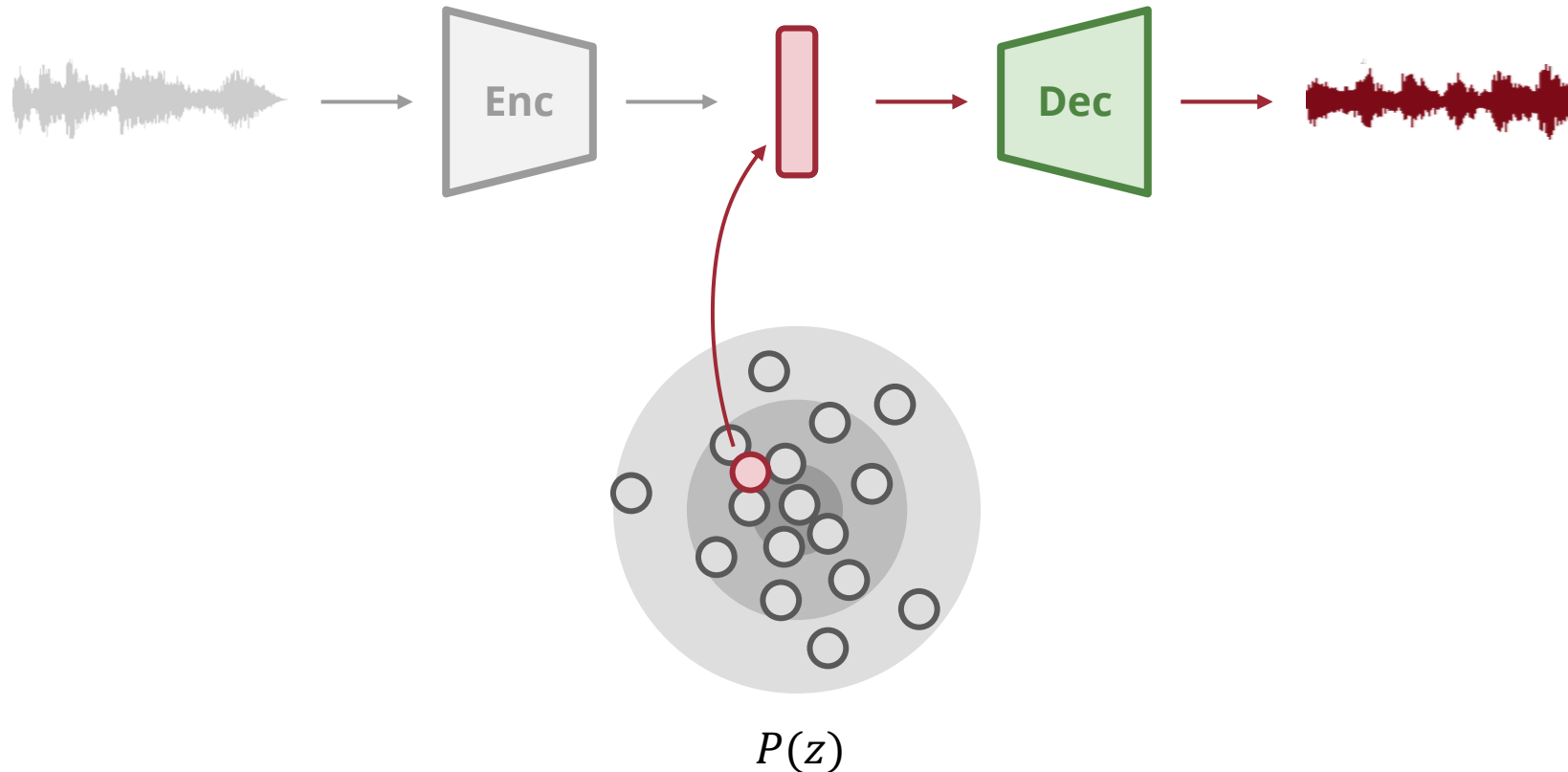


(Source: Roberts et al., 2018)

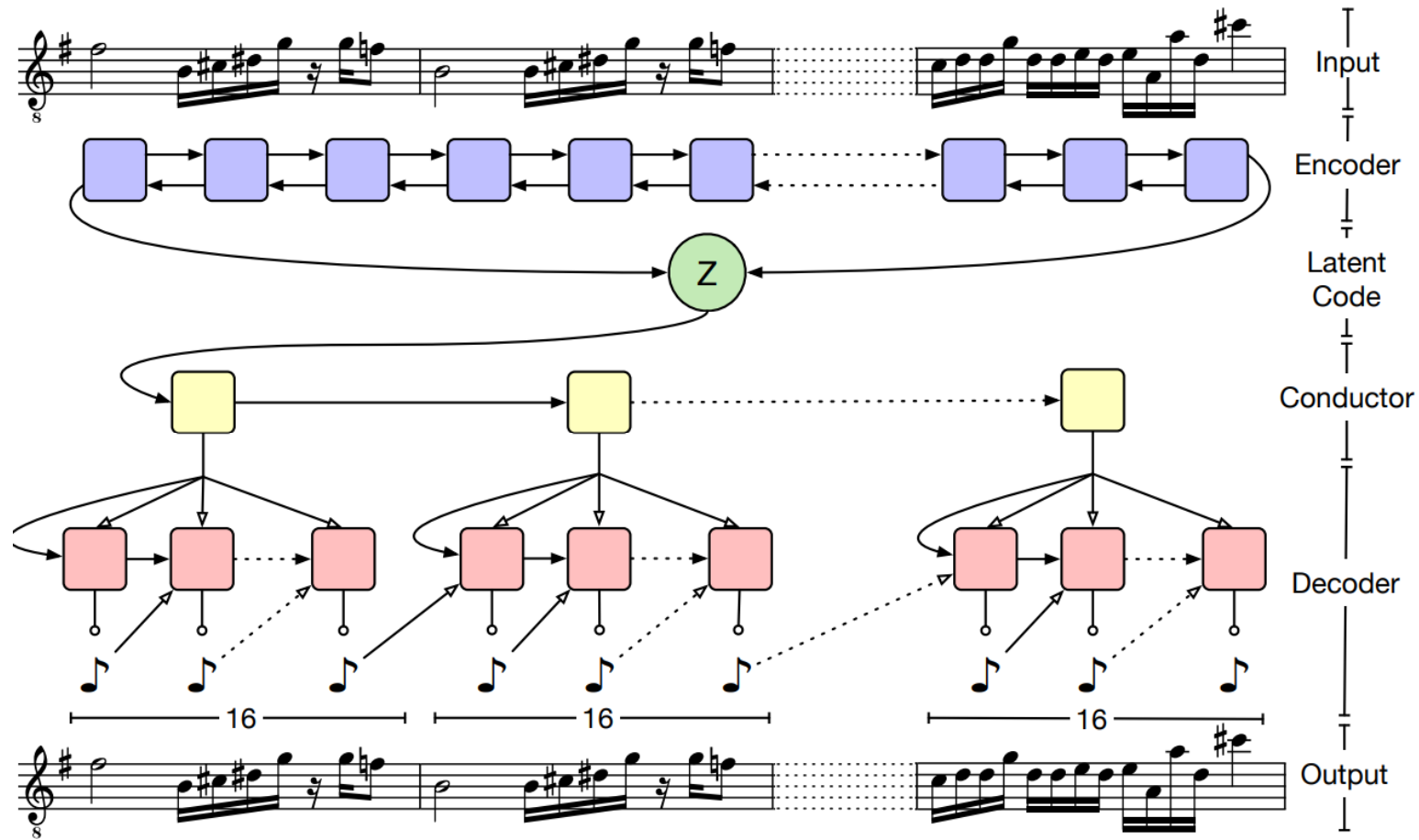
(Recap) Variational Autoencoders (VAEs) – Training



(Recap) Variational Autoencoders (VAEs) – Generation



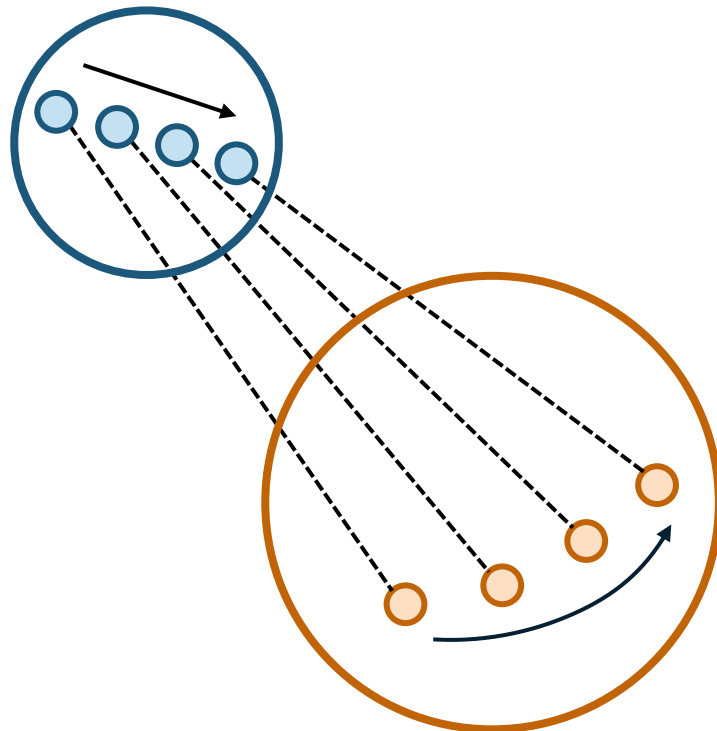
Example: MusicVAE (Roberts et al., 2018)



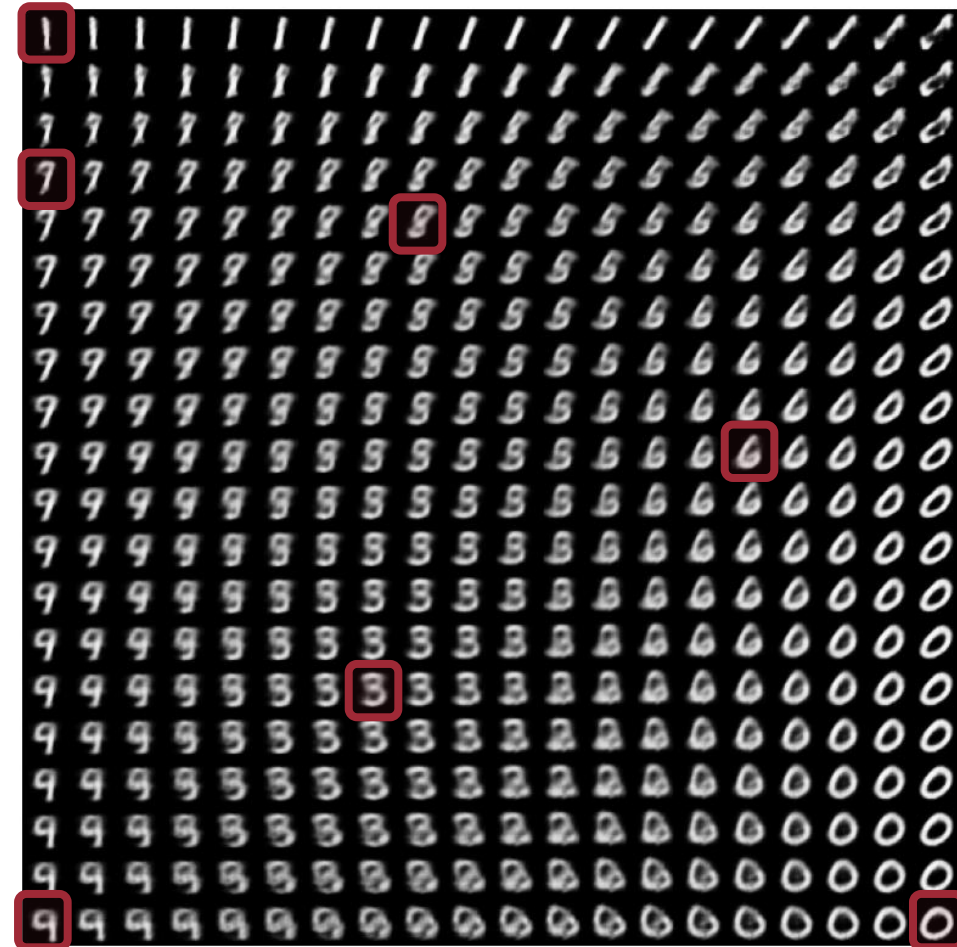
(Source: Roberts et al., 2018)

(Recap) Decoding the Latent Space of a VAE

Latent space

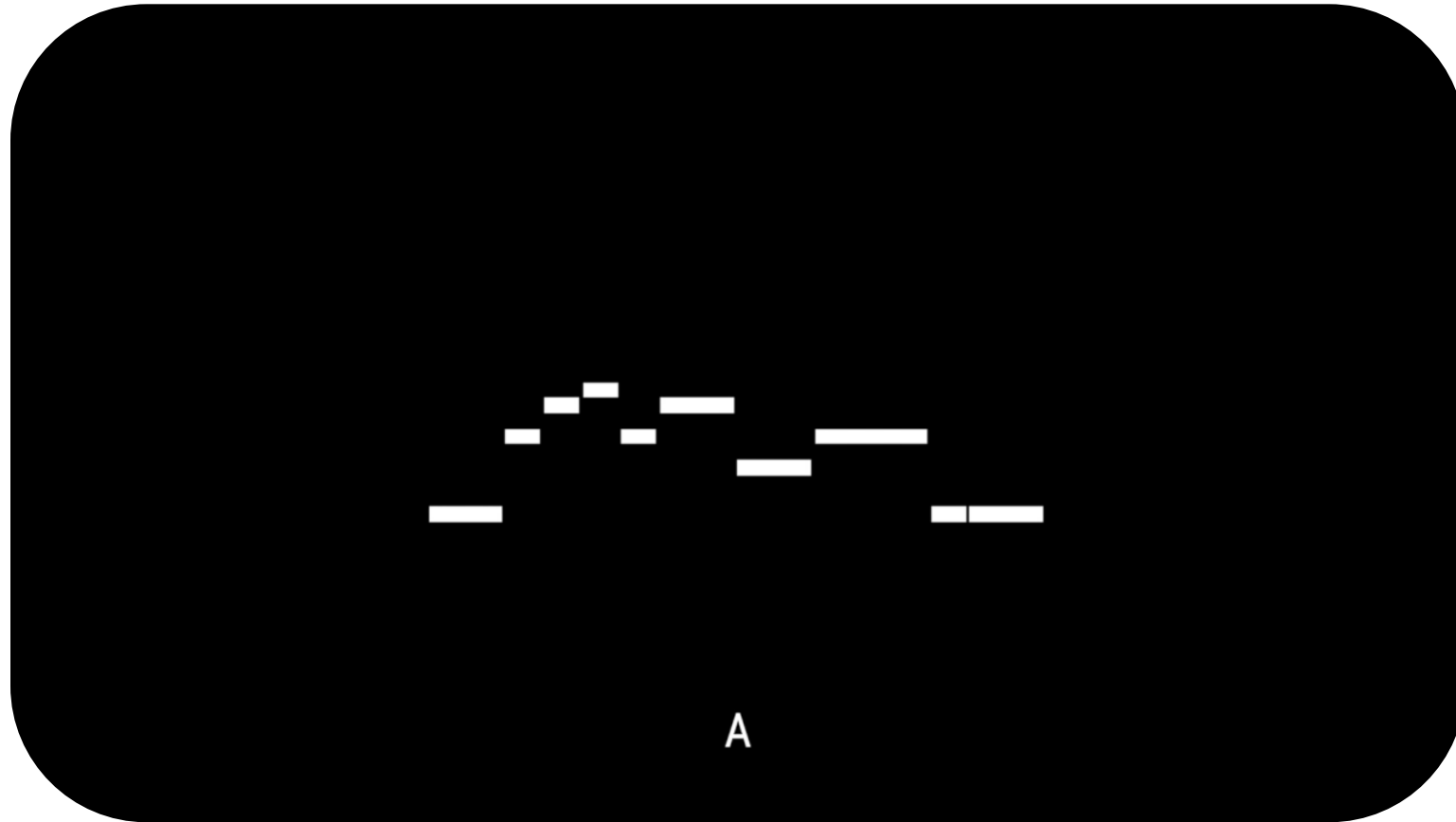


Data space



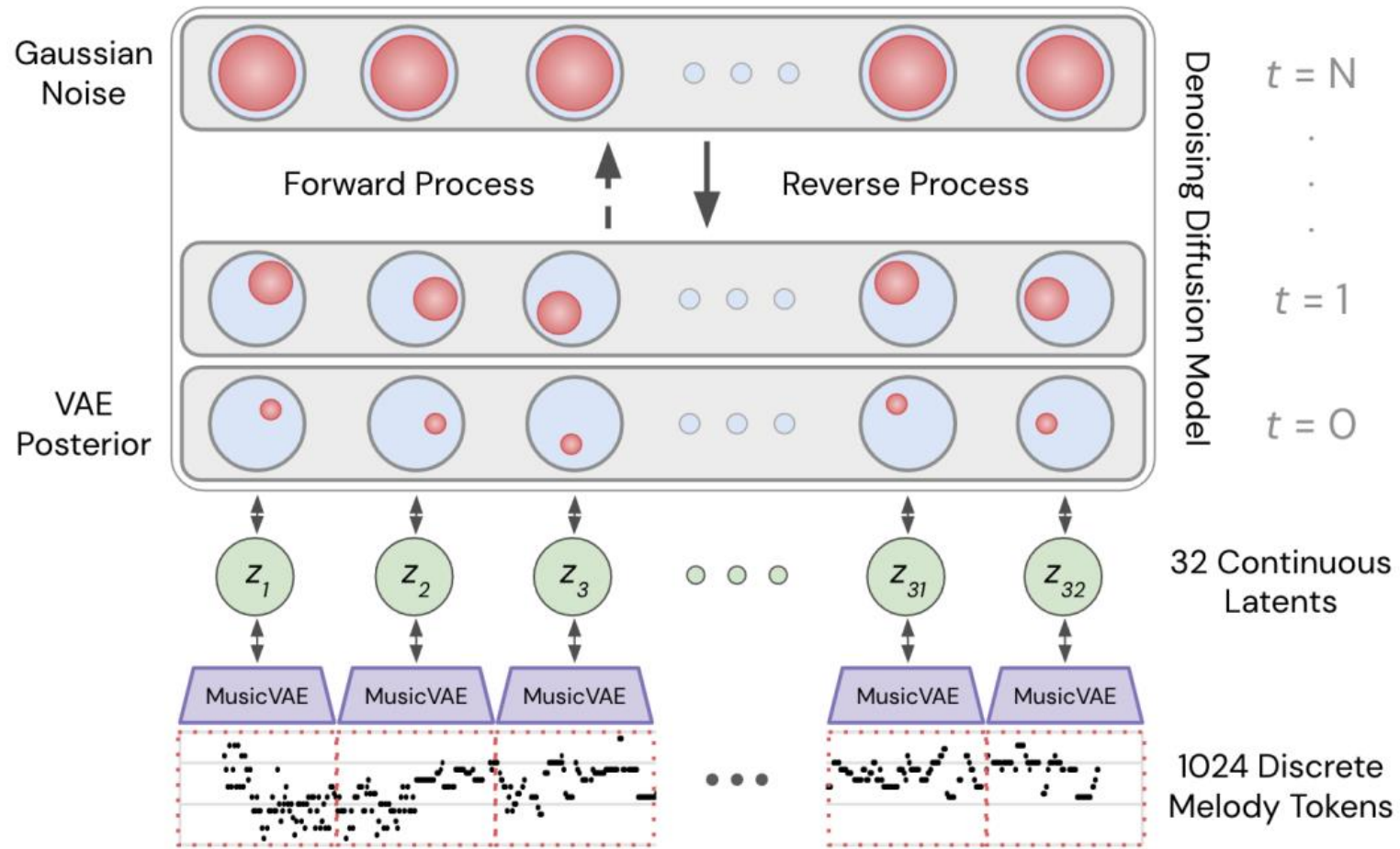
(Source: tensorflow.org)

Example: MusicVAE (Roberts et al., 2018)



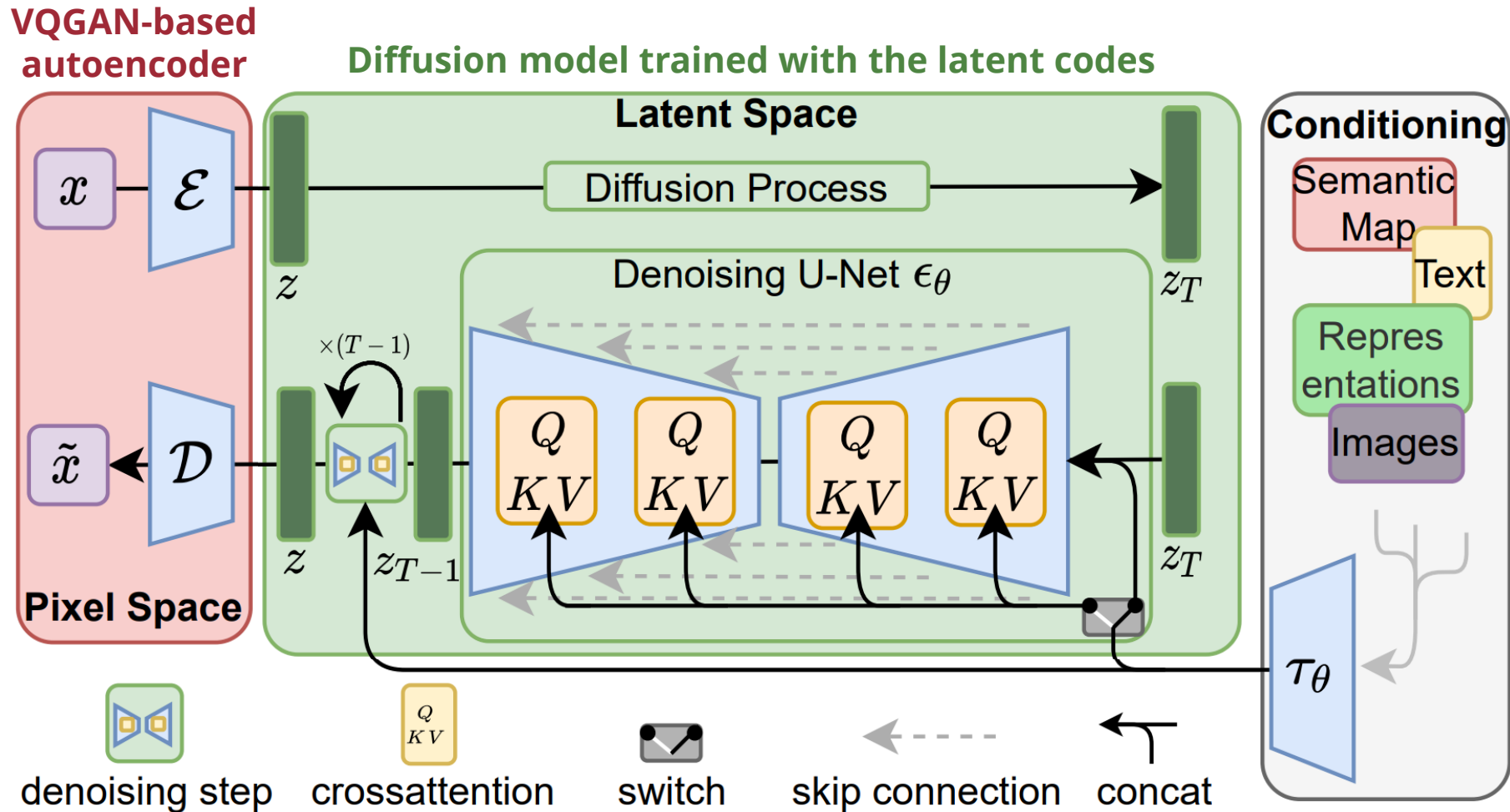
(Source: Roberts et al., 2018)

Example: Latent Diffusion (Mittal et al., 2021)



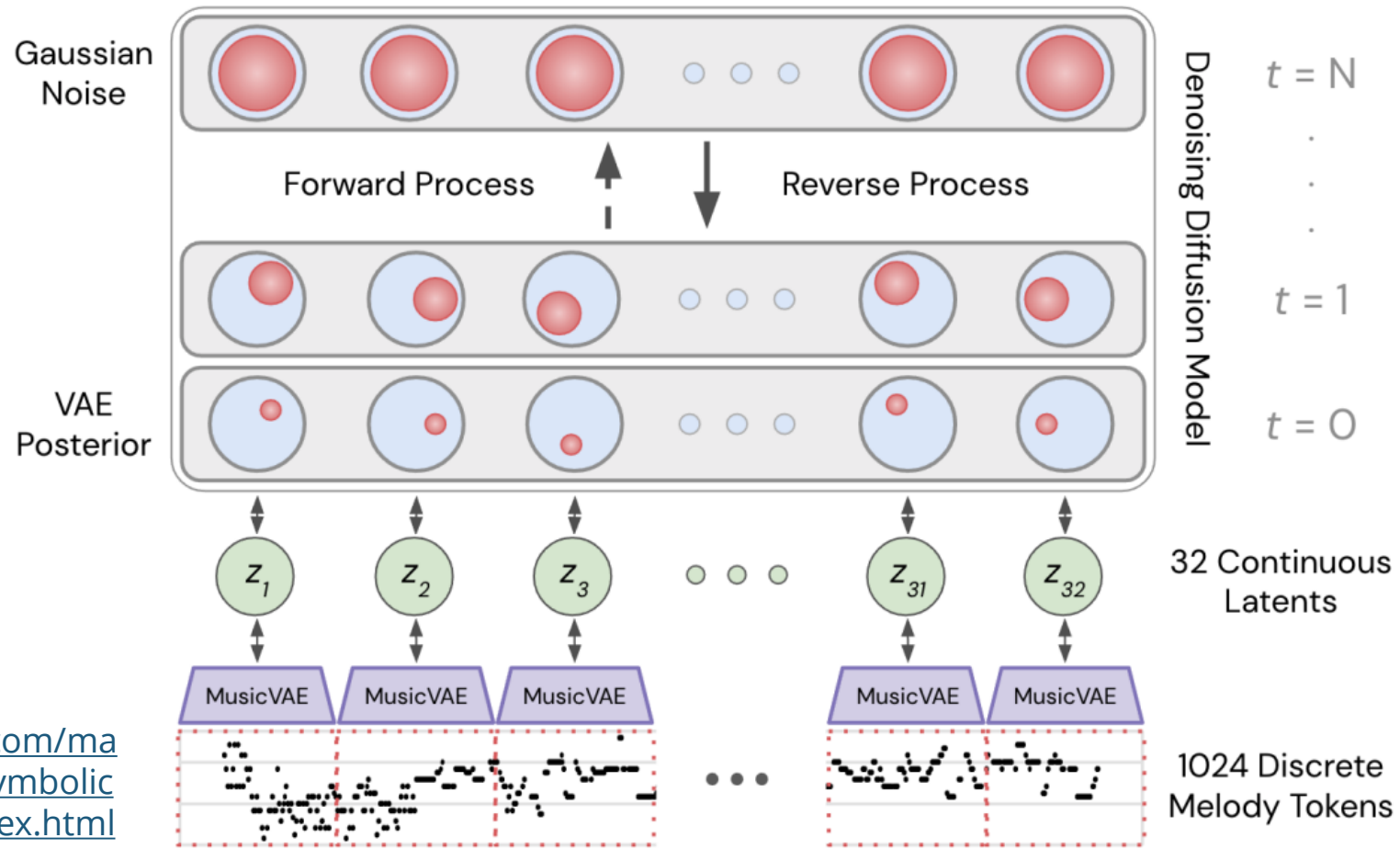
(Source: Mittal et al., 2021)

(Recap) Latent Diffusion Models (LDMs)



(Source: Rombach et al., 2022)

Example: Latent Diffusion (Mittal et al., 2021)



storage.googleapis.com/magentadata/papers/symbolic-music-diffusion/index.html

(Source: Mittal et al., 2021)

Music Infilling Models

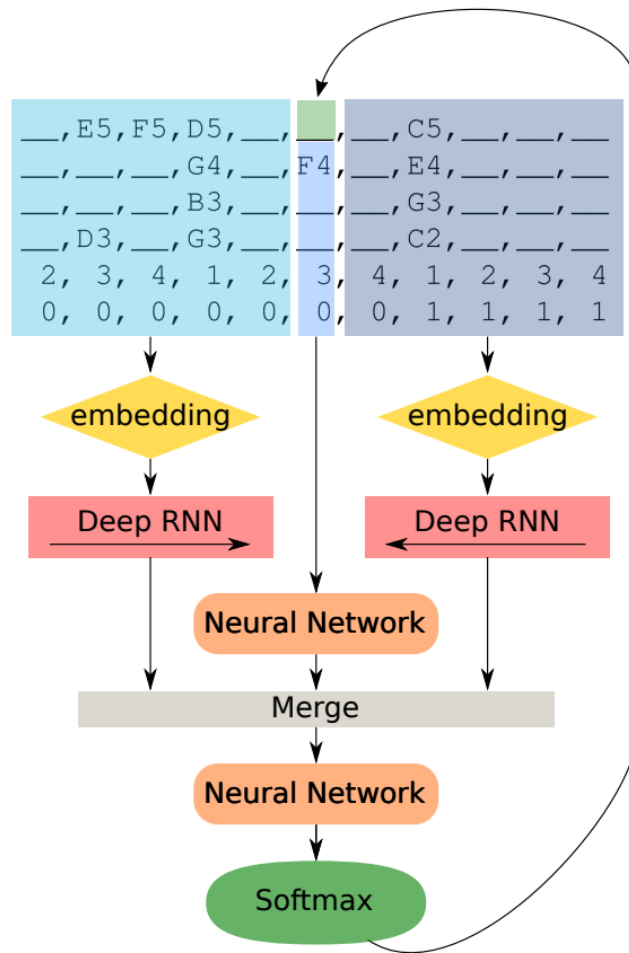
Example: DeepBach (Hadjeres et al., 2017)

The image shows a musical score for a Bach chorale fragment. The top staff is the soprano line, and the bottom staff is the bass line. A blue box highlights a group of notes in the soprano line, and an orange box highlights a note in the bass line. To the right, a MIDI representation of the notes is shown. The notes are listed in four rows, with some notes highlighted in colored boxes (blue, green, orange) corresponding to the boxes in the musical notation.

D5, __, E5, F5, D5, __, __, __, C5, __, __, __, E5
A4, __, __, __, G4, __, F4, __, E4, __, __, __, E4
C4, __, __, __, B3, __, __, __, G3, __, __, __, A3
F3, __, D3, __, G3, __, __, __, C2, __, __, __, C#2
1, 2, 3, 4, 1, 2, 3, 4, 1, 2, 3, 4, 1
0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0

(Source: Hadjeres et al., 2017)

Example: DeepBach (Hadjeres et al., 2017)



(Source: Hadjeres et al., 2017)

Algorithm 1 Pseudo-Gibbs sampling

- 1: **Input:** Chorale length L , metadata \mathcal{M} containing lists of length L , probability distributions (p_1, p_2, p_3, p_4) , maximum number of iterations M
 - 2: Create four lists $\mathcal{V} = (\mathcal{V}_1, \mathcal{V}_2, \mathcal{V}_3, \mathcal{V}_4)$ of length L
 - 3: {The lists are initialized with random notes drawn from the ranges of the corresponding voices (sampled uniformly or from the marginal distributions of the notes)}
 - 4: **for** m from 1 to M **do**
 - 5: Choose voice i uniformly between 1 and 4
 - 6: Choose time t uniformly between 1 and L
 - 7: Re-sample \mathcal{V}_i^t from $p_i(\mathcal{V}_i^t | \mathcal{V}_{\setminus i, t}, \mathcal{M}, \theta_i)$
 - 8: **end for**
 - 9: **Output:** $\mathcal{V} = (\mathcal{V}_1, \mathcal{V}_2, \mathcal{V}_3, \mathcal{V}_4)$
-

Example: DeepBach (Hadjeres et al., 2017)

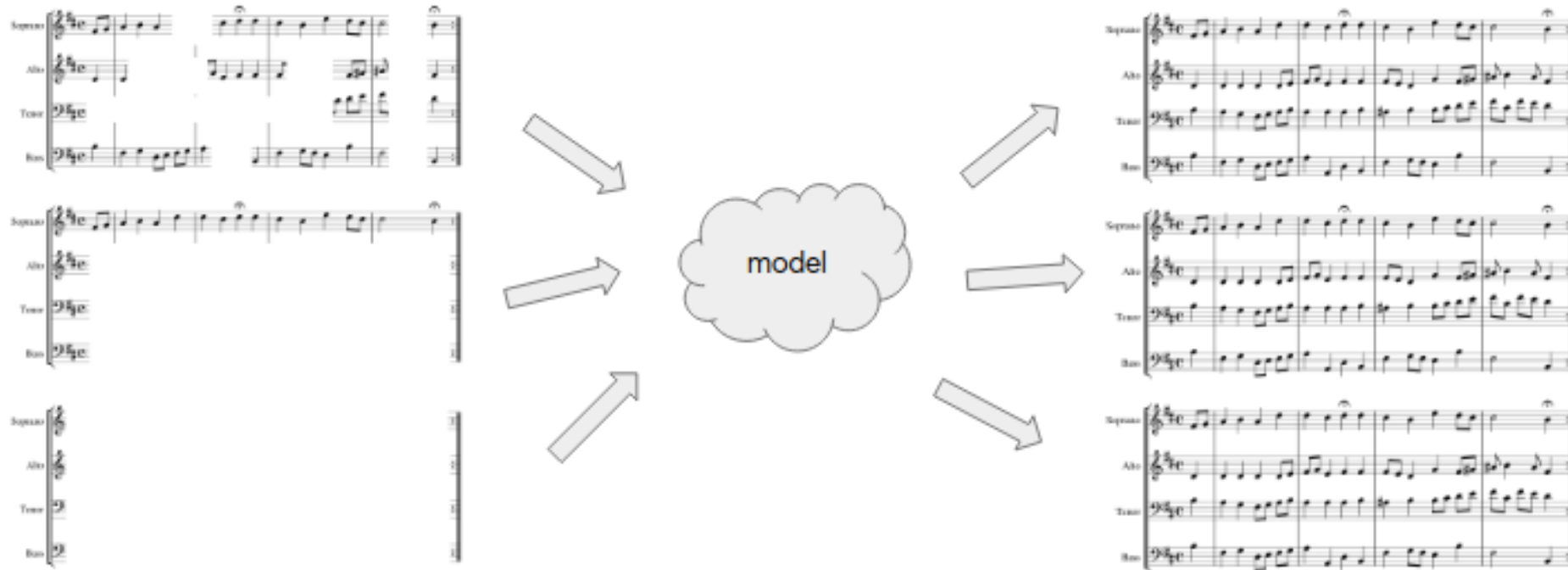
Reharmonization example



youtu.be/QiBM7-5hA6o

Example: Coconet (Huang et al., 2017)

- Based on Orderless NADE (Uribe et al., 2014)



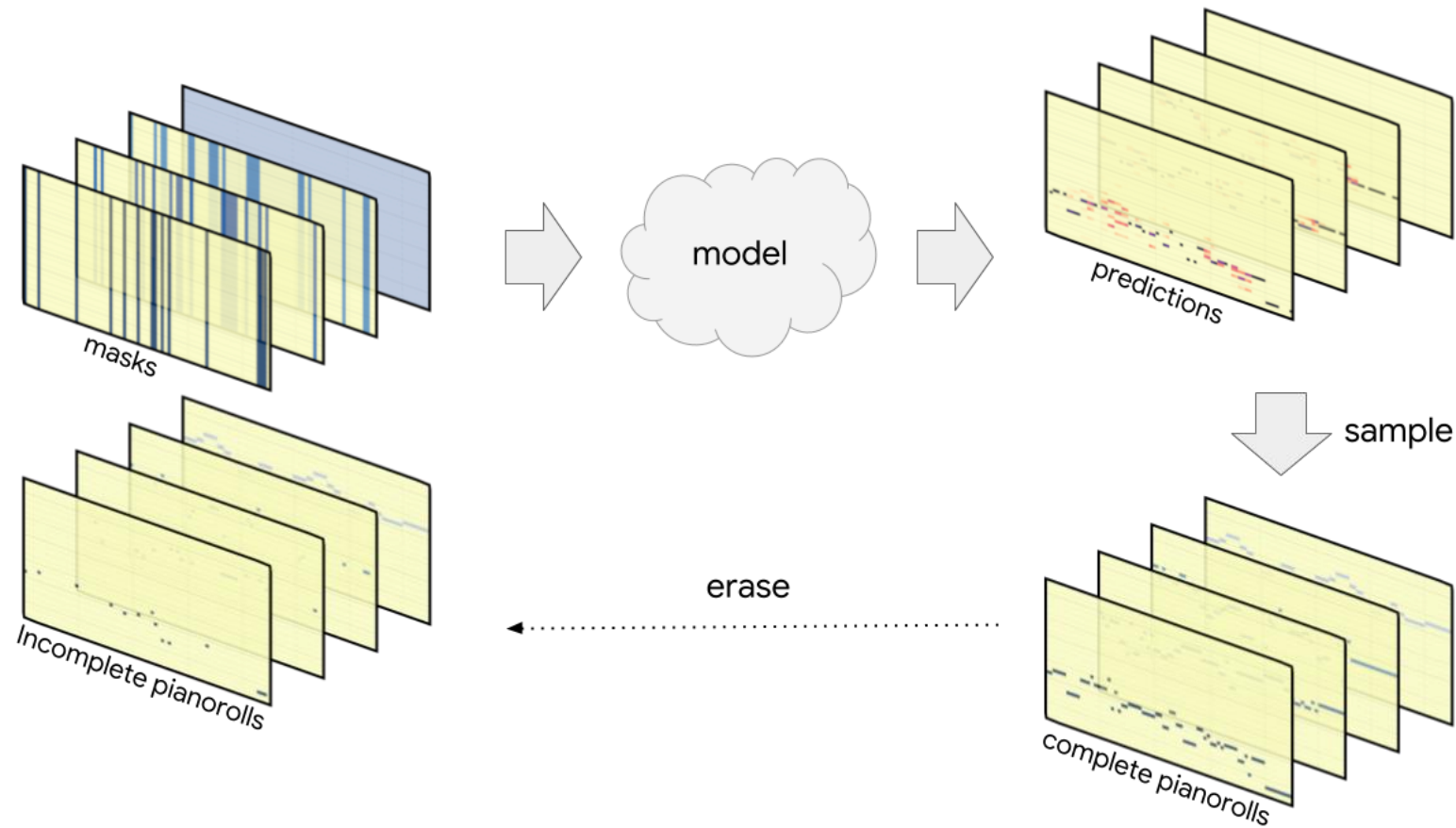
(Source: Huang et al., 2019)

Benigno Uribe, Iain Murray, and Hugo Larochelle, "A Deep and Tractable Density Estimator," *ICML*, 2014.

Cheng-Zhi Anna Huang, Tim Cooijmans, Adam Roberts, Aaron Courville, and Douglas Eck, "Counterpoint by Convolution," *ISMIR*, 2017.

Cheng-Zhi Anna Huang, Tim Cooijmans, Monica Dinulescu, Adam Roberts, and Curtis Hawthorne, "Coconet: the ML model behind today's Bach Doodle," *Magenta Blog*, 2019.

Example: Coconet (Huang et al., 2017)



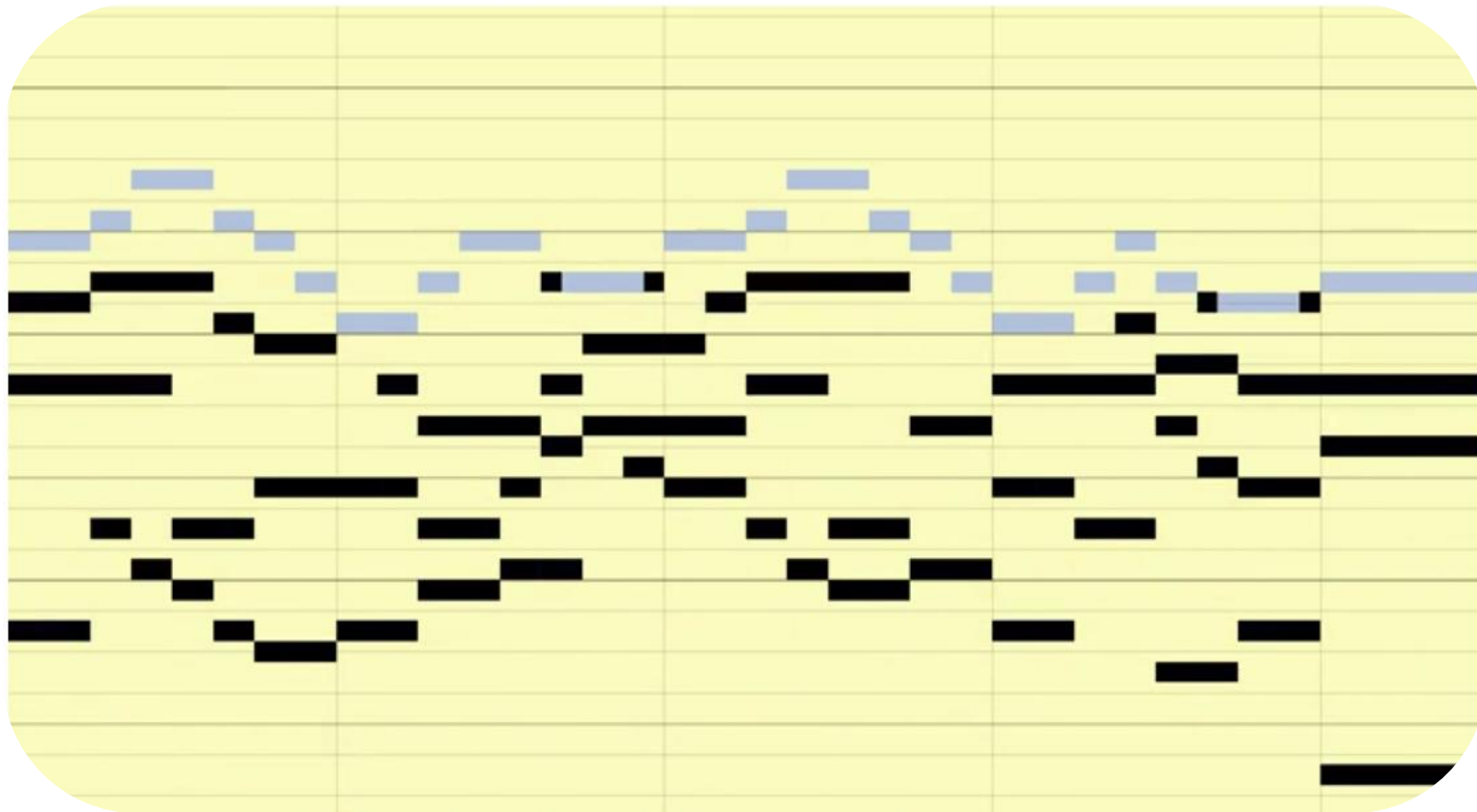
(Source: Huang et al., 2019)

Example: Coconet (Huang et al., 2017)



(Source: Huang et al., 2017)

Example: Coconet (Huang et al., 2017)



(Source: Huang et al., 2017)

Example: JS Bach Doodle (2019)



doodles.google/doodle/celebrating-johann-sebastian-bach/



youtu.be/XBfYPp6KF2g & magenta.tensorflow.org/coconet

Example: Variable-Length Infilling (VLI) (Chang et al., 2021)

$\text{♩} = 80$ C_{past} C_* C_{future}

n_0 n_1 n_2 n_3 n_4 n_5

ILM

Input:	n_0	n_1	[BLANK]	n_5	[SEP]
target:		n_2	n_3	n_4	[EOS]

FELIX

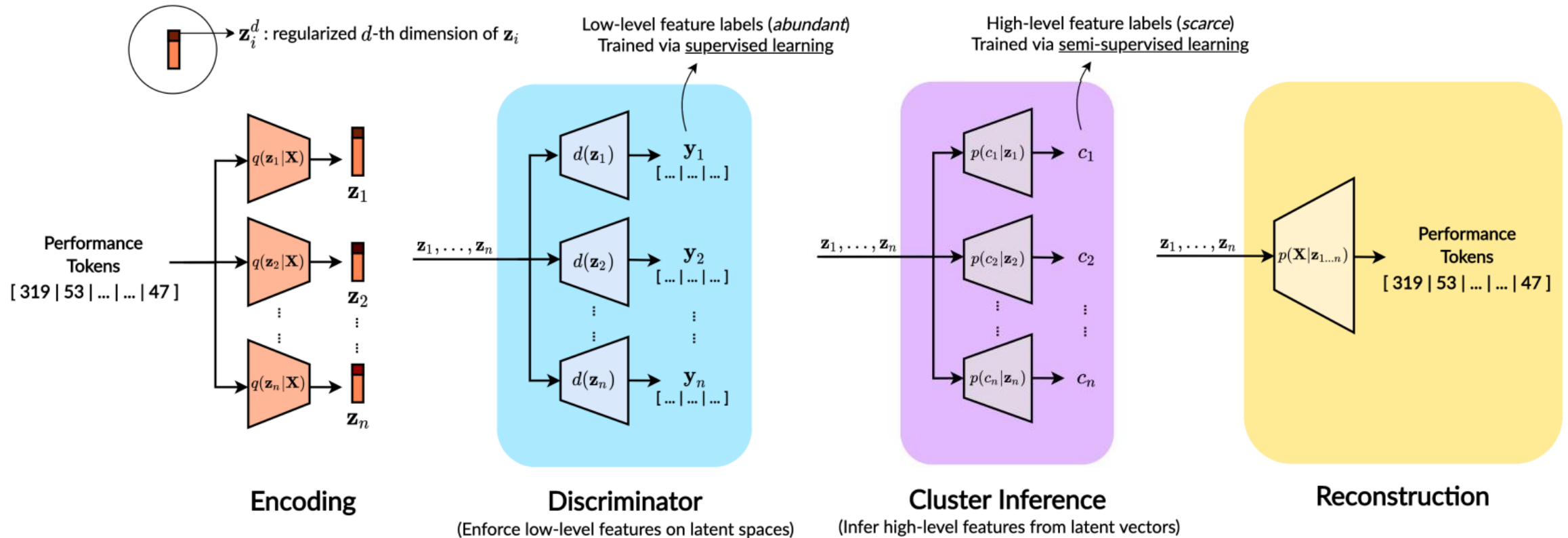
Input:	n_0	n_1	[MASK]	[MASK]	[MASK]	[MASK]	[MASK]	n_5
target:			n_2	n_3	n_4	[PAD]	[PAD]	

(Source: Chang et al., 2021)

jackyhsiung.github.io/piano-infilling-demo

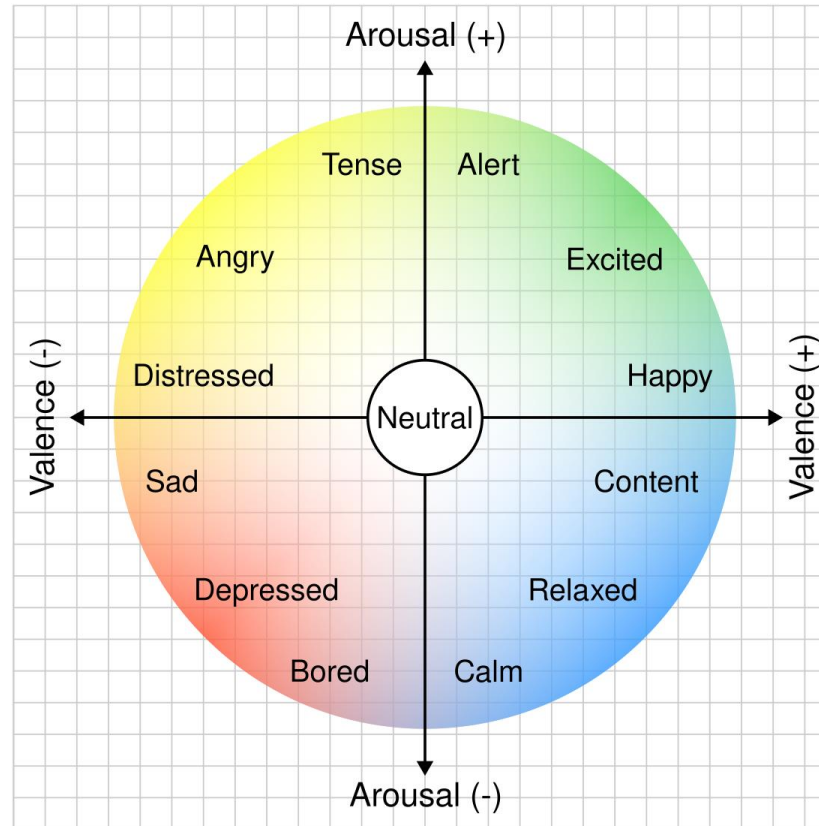
Controllable Music Generation

Example: Music FaderNet (Tan & Herremans, 2020)



(Source: Tan & Herremans, 2020)

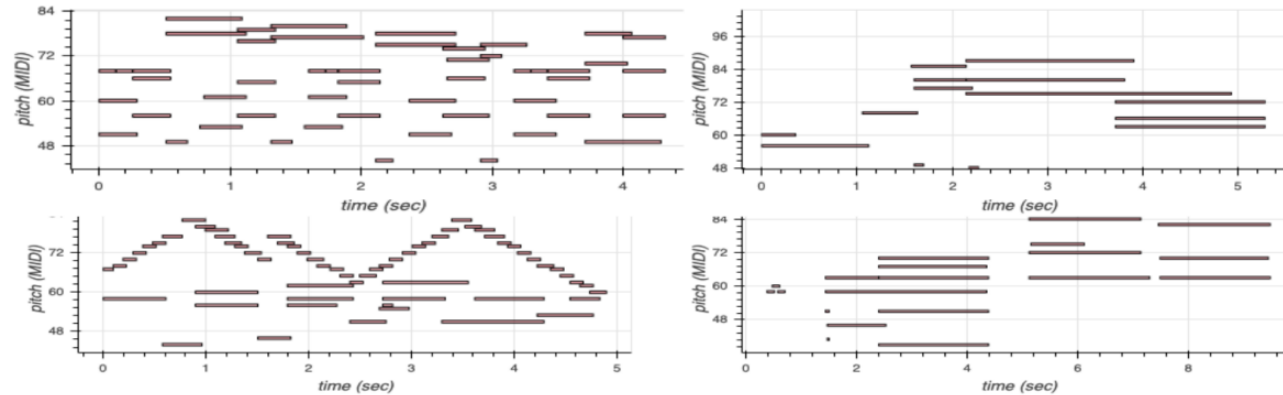
Valence-Arousal Model for Emotion



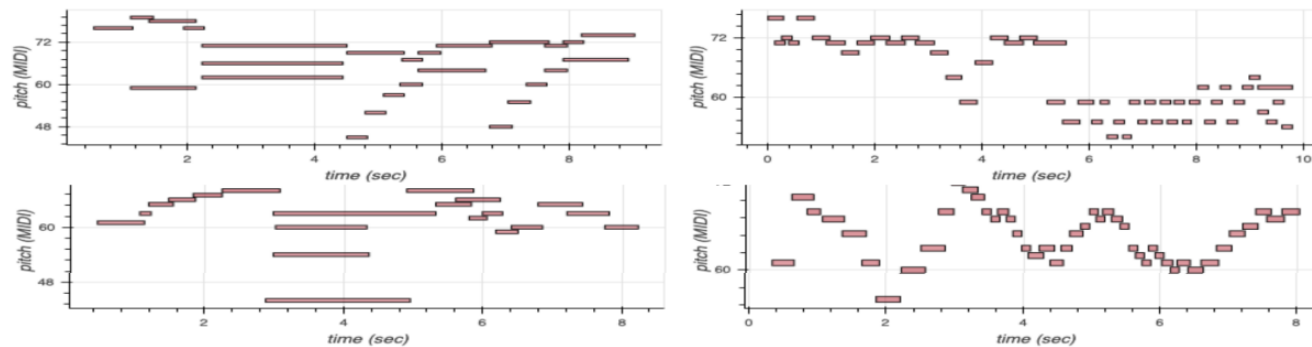
(Source: mrAnmol)

Example: Music FaderNet (Tan & Herremans, 2020)

High Arousal → Low Arousal



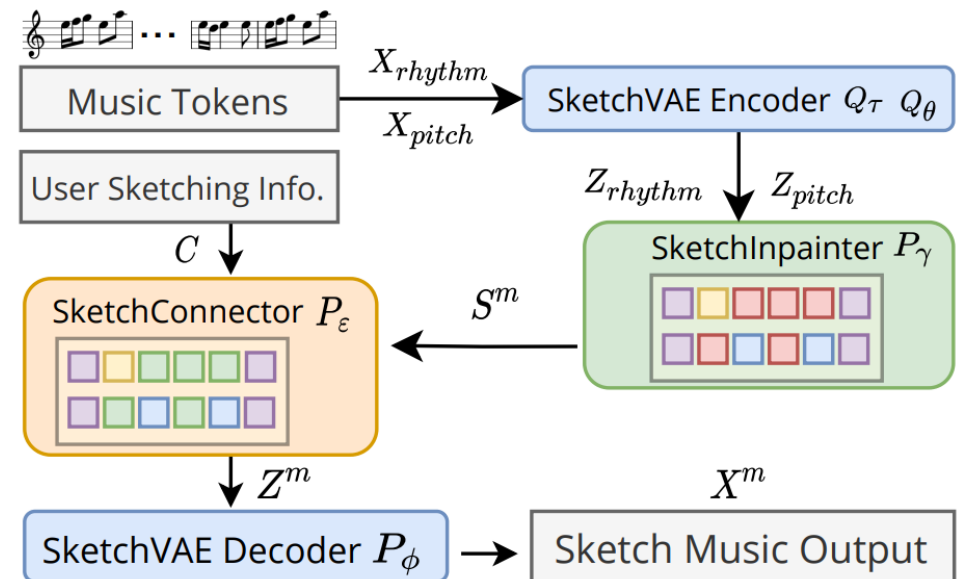
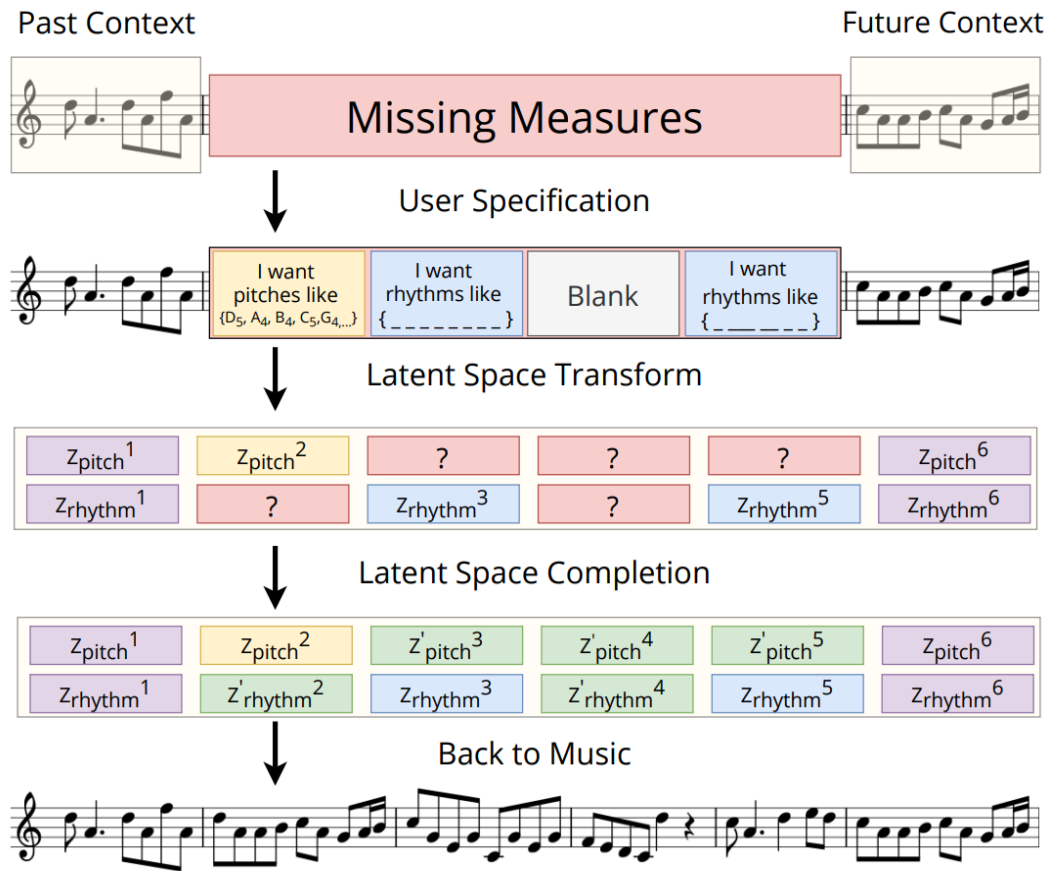
Low Arousal → High Arousal



(Source: Tan & Herremans, 2020)

music-fadernets.github.io

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



(Source: Chen et al., 2020)

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

The diagram illustrates the Music SketchNet architecture across three phases: Past Context, Generation, and Future Context. It features four staves: Original, Control Pitch, Control Rhythm, and Control Both.

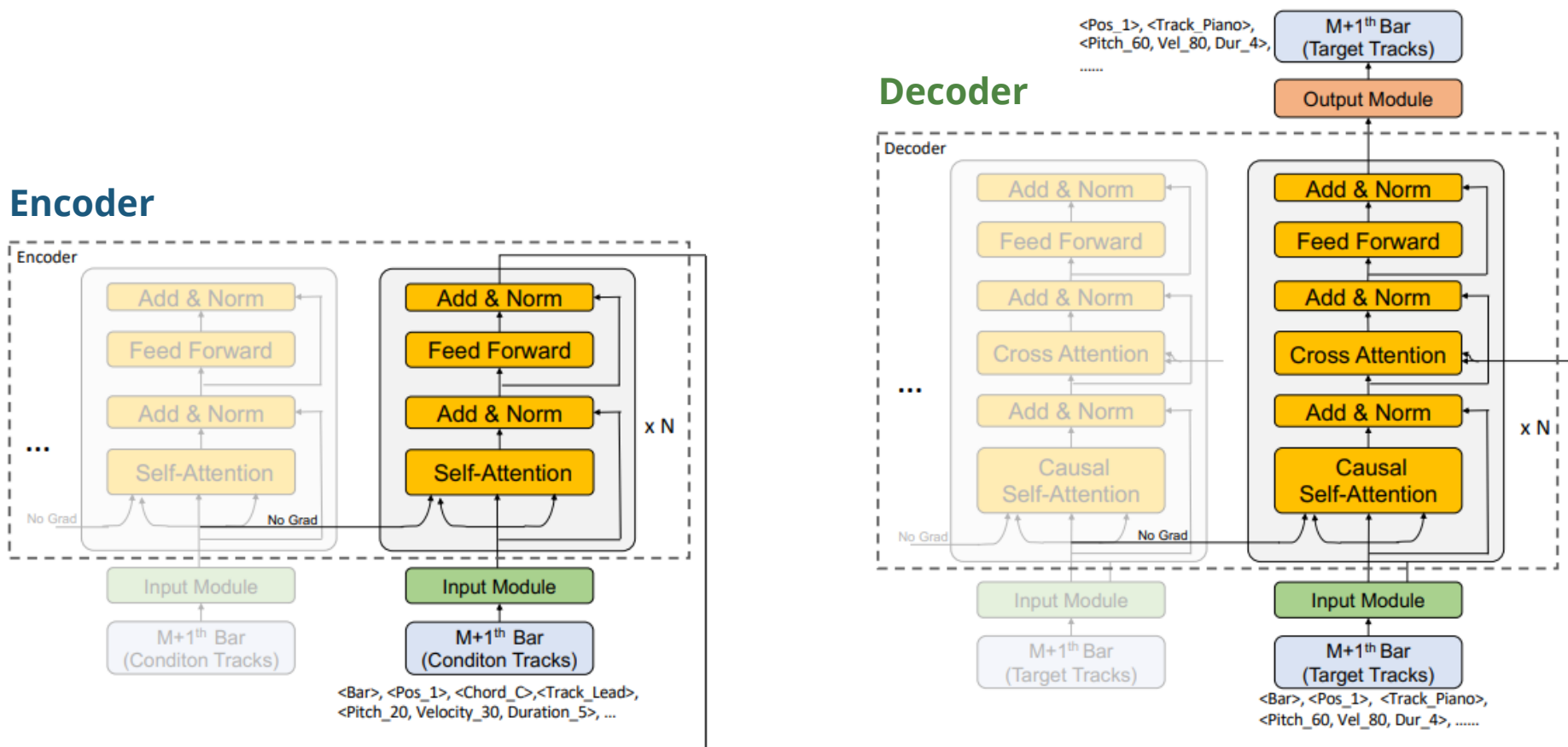
- Original:** The target musical piece, shown in 4/4 time.
- Control Pitch:** A sequence of chords and triplets that guide the pitch of the generated music. Chords are labeled as $\{Ab5, Db6, Eb6, Gb6\}$, $\{C6, Eb6, Db6, F6, Db6\}$, $\{F6, Gb6, Ab6, Ab6, F6\}$, and $\{Db6, F6, Ab6, Bb6, Db6\}$.
- Control Rhythm:** A sequence of rhythmic patterns, represented by pink bars, that guide the rhythm of the generated music.
- Control Both:** A sequence of chords and triplets that guide both pitch and rhythm. It includes the chord $\{Ab5, Db6, Eb6, Gb6\}$ and a grey bar labeled "No Sketch" indicating a period where no sketch is provided.

The timeline is divided into three sections: Past Context, Generation, and Future Context. The Generation section is further divided into four sub-sections, each corresponding to a set of control parameters.

(Source: Chen et al., 2020)

Music Accompaniment

Example: PopMAG (Ren et al., 2020)

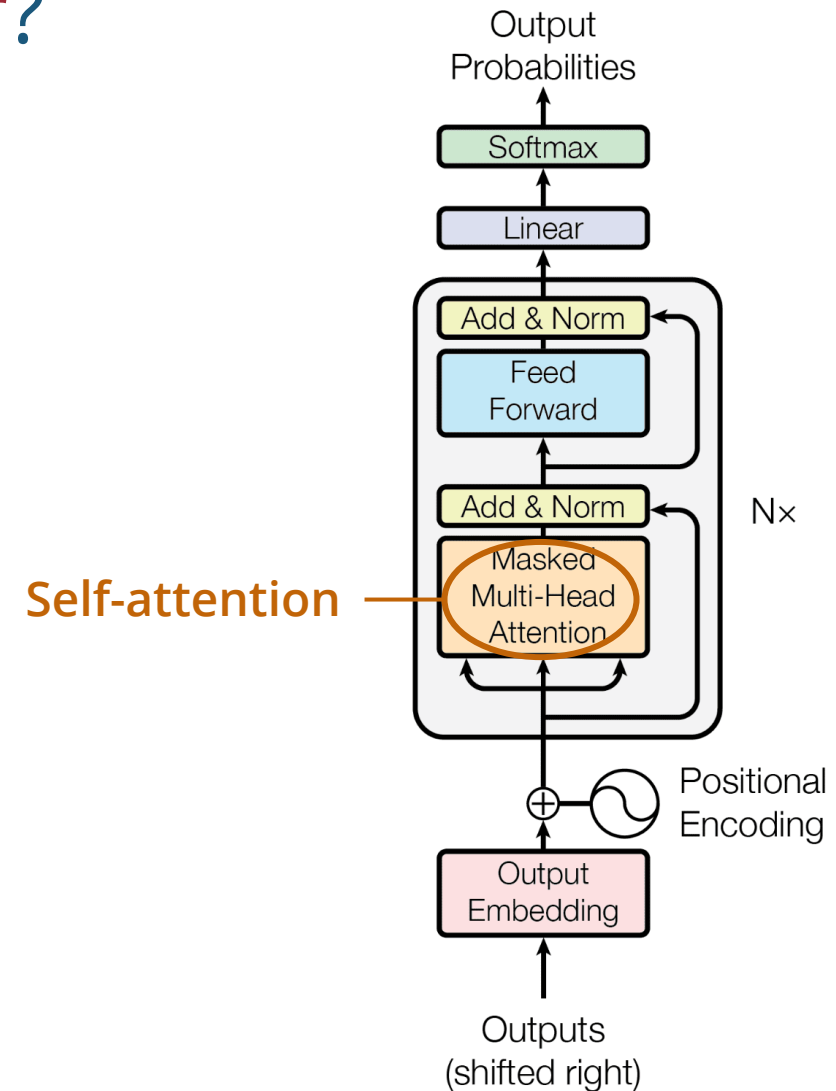


(Source: Ren et al., 2020)

ai-music.github.io/popmag

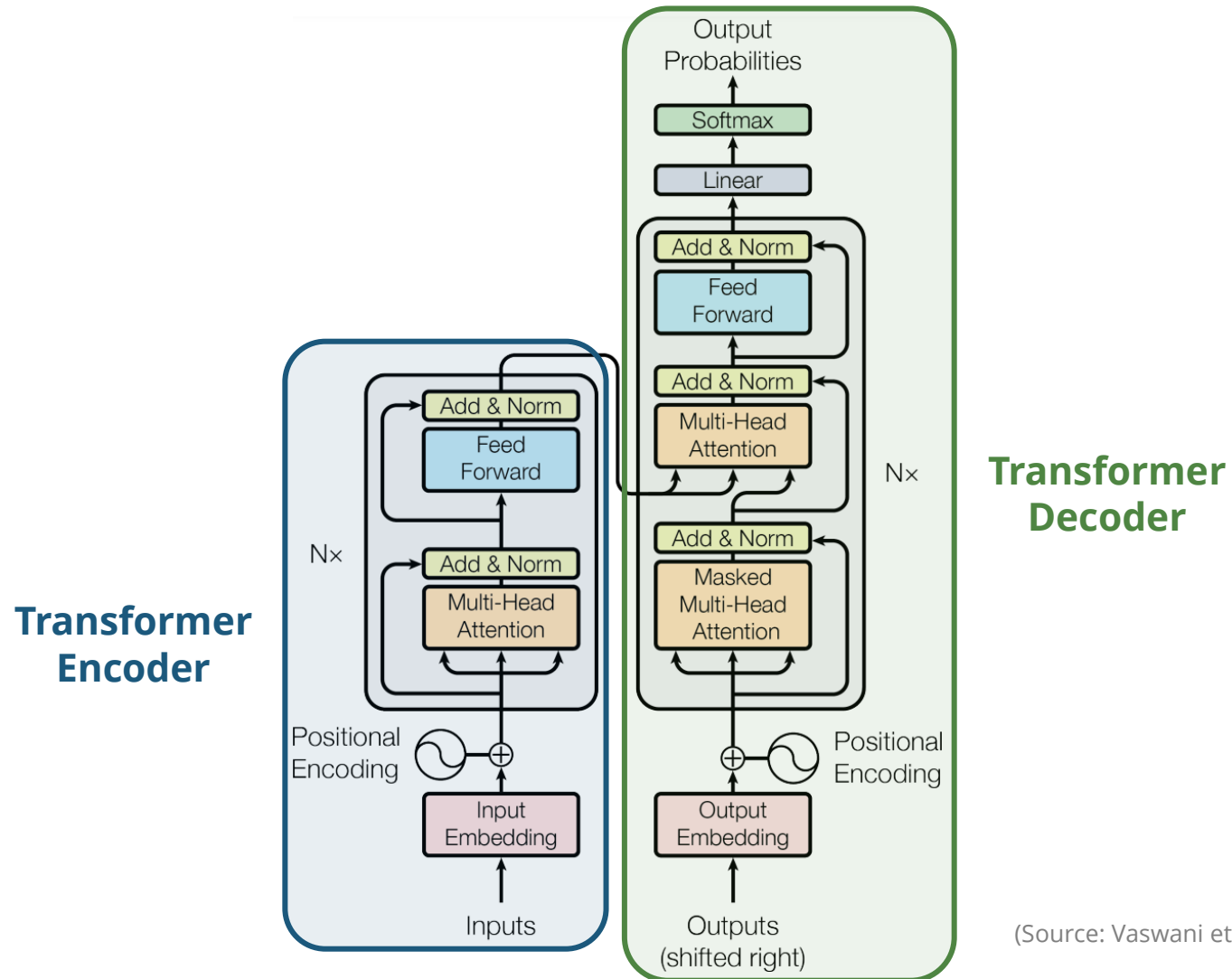
(Recap) What is a Transformer?

- A type of neural network that use the **self-attention mechanism**



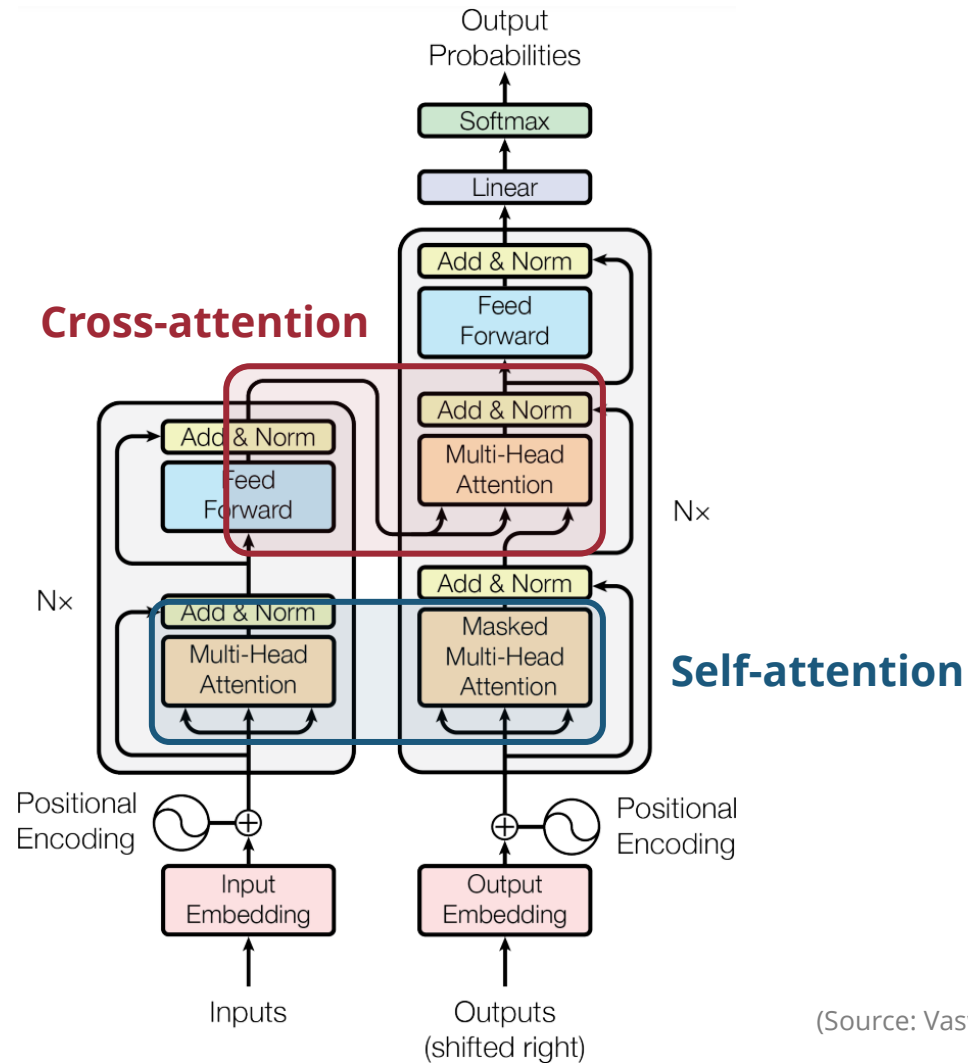
(Source: Vaswani et al., 2017; adapted)

The Original Transformer – Encoder & Decoder

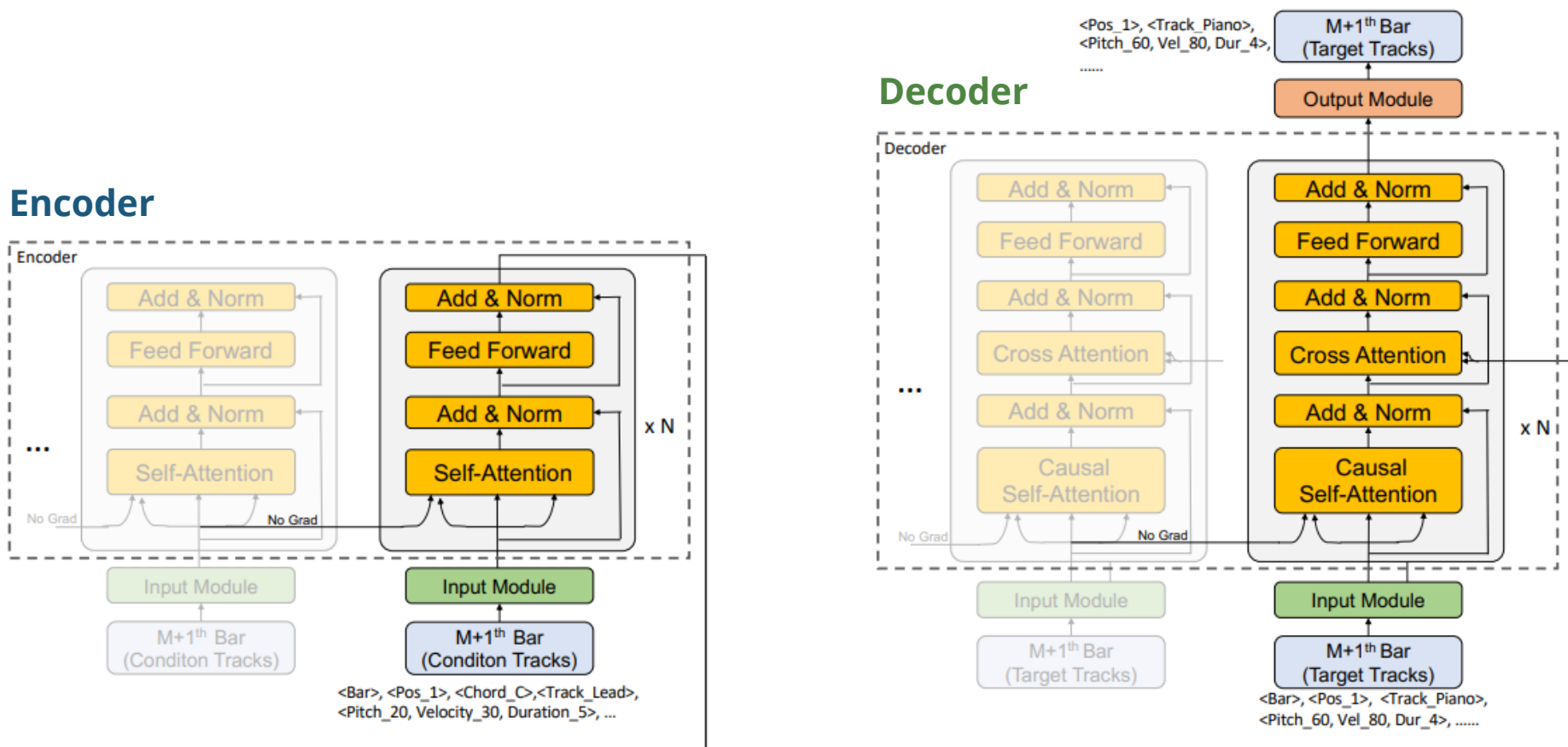


(Source: Vaswani et al., 2017)

The Original Transformer – Cross-attention



Example: PopMAG (Ren et al., 2020)

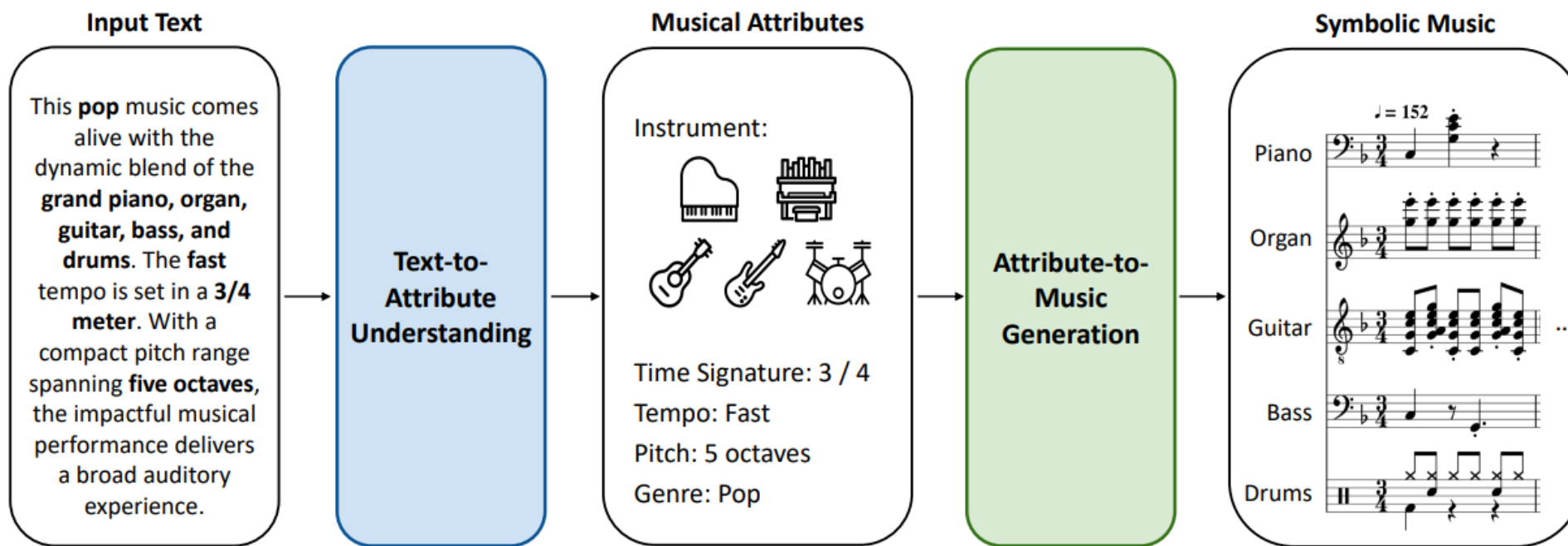


(Source: Ren et al., 2020)

ai-music.github.io/popmag

Multimodal Music Generation

Example: MuseCoco (Lu et al., 2023)



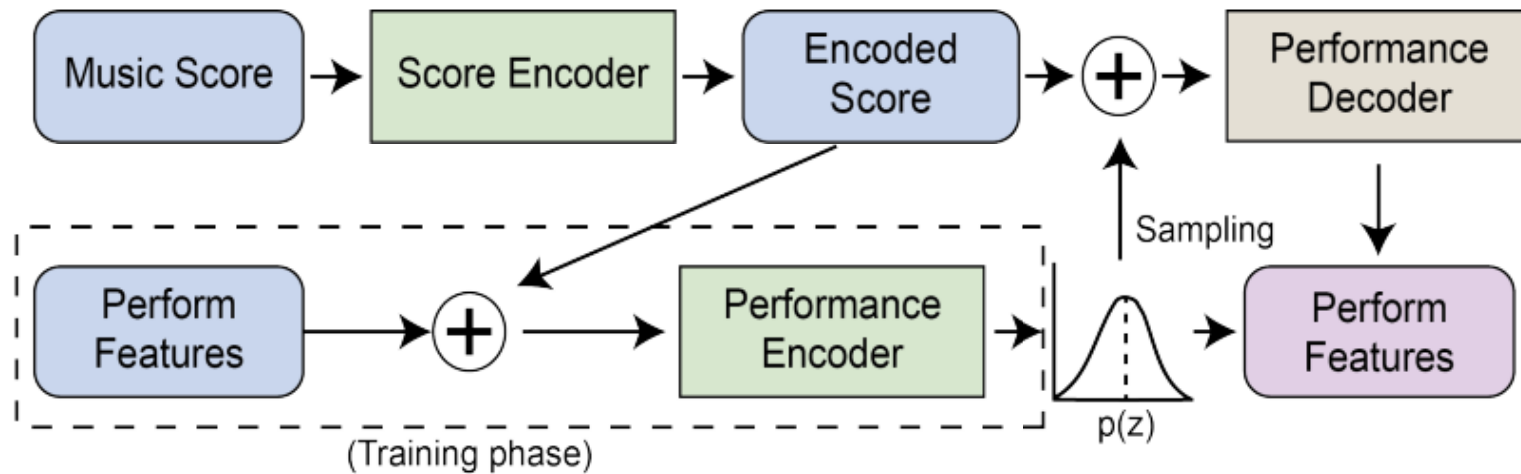
(Source: Lu et al., 2023)

ai-music.github.io/musecoco

Performance Rendering

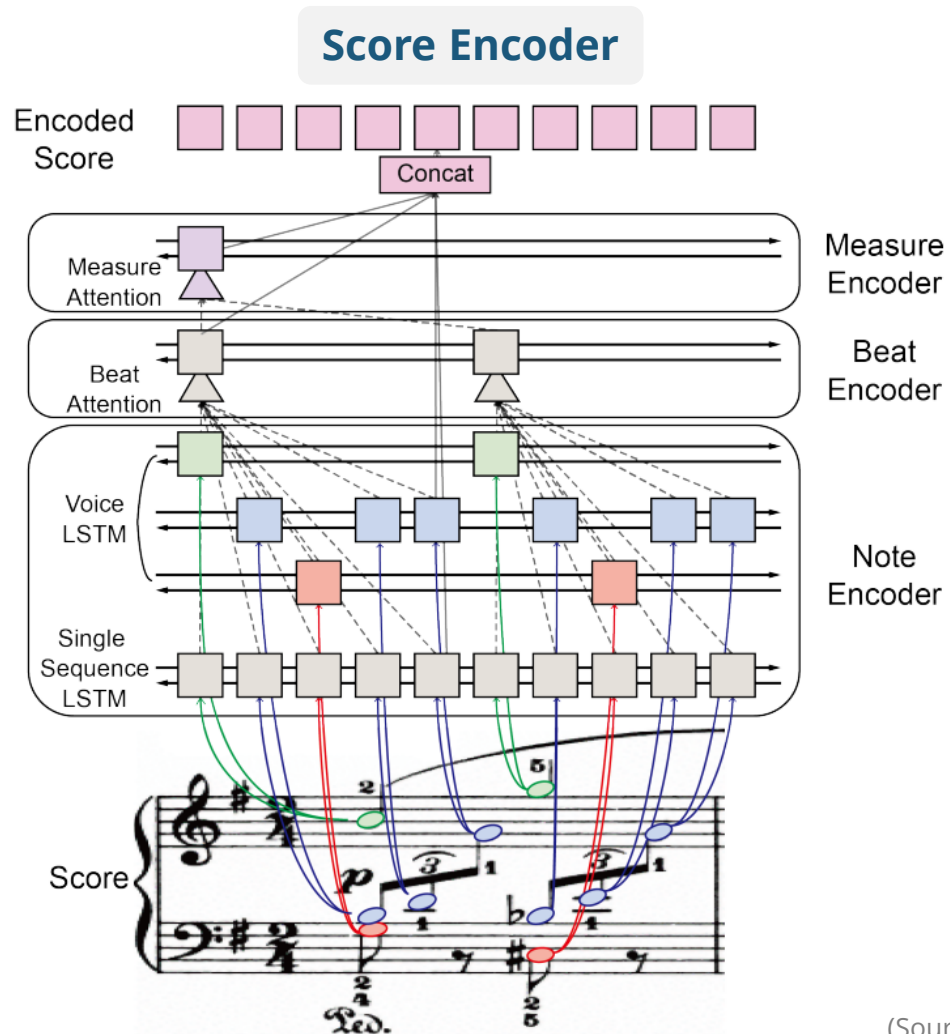
Example: VirtuosoNet (Jeong et al., 2019)

- **Input:** pitch, duration, articulation marking, slur and beam status, tempo marking, and dynamic marking, etc.
- **Output:** absolute tempo, velocity, onset deviation, articulation, pedal usages

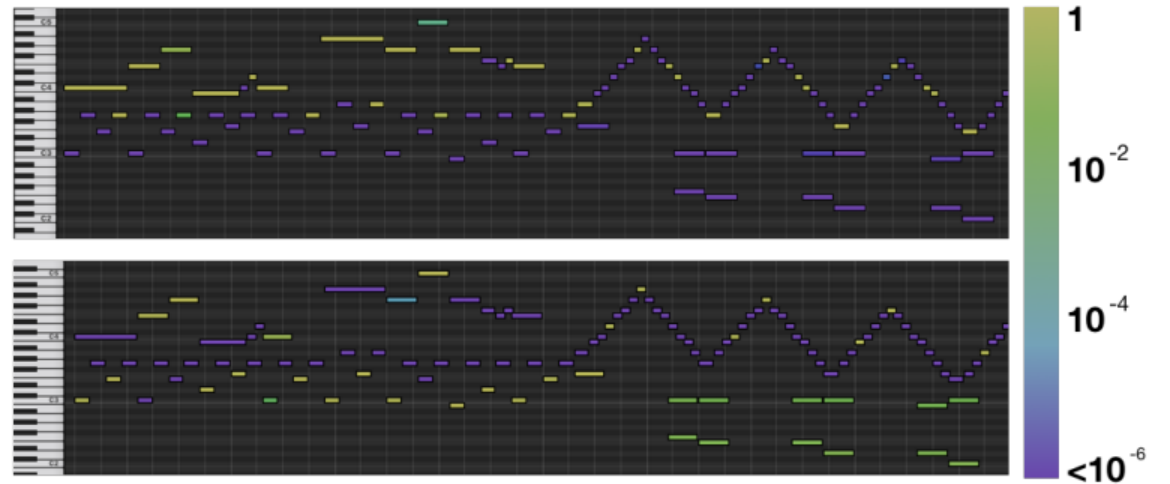


(Source: Jeong et al., 2019)

Example: VirtuosoNet (Jeong et al., 2019)



Attention visualization



(Source: Jeong et al., 2019)

Example: **VirtuosoNet** (Jeong et al., 2019)



youtu.be/6HeFJQf2h2o

Example: **VirtuosoNet** (Jeong et al., 2019)



youtu.be/BN0ZCBS9q0Y

Example: **VirtuosoNet** (Jeong et al., 2019)



youtu.be/hPBR2Rxu3-s

Open Questions

Open Questions

- How to generate **long-term structure**?
- How to enable more **intuitive controls** for music generation systems?
- How to adopt these models for **improvisation**?
- We are **running out of symbolic music data**
 - Can we learn symbolic music composition from **listening to raw audio**?
- **Is symbolic reasoning a must** for true AI music generation?
 - Can an AI play perfect music **without processing it into symbolic music internally**?
- How do humans **learn to create music**?

(Recap) A Simplified Music Production Workflow

