

PAT 498/598 (Fall 2024)

Special Topics: Generative AI for Music and Audio Creation

Lecture 15: Time-domain Audio Synthesis

Instructor: Hao-Wen Dong



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

(Recap) Four Paradigms



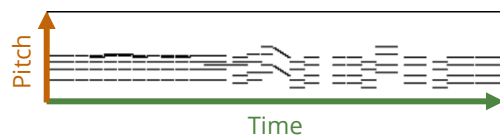
Symbolic music generation

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



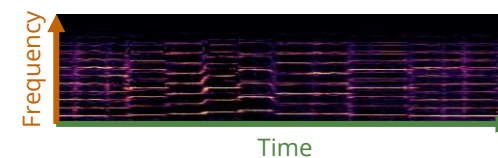
Audio-domain music generation

Time series-based

Image-based



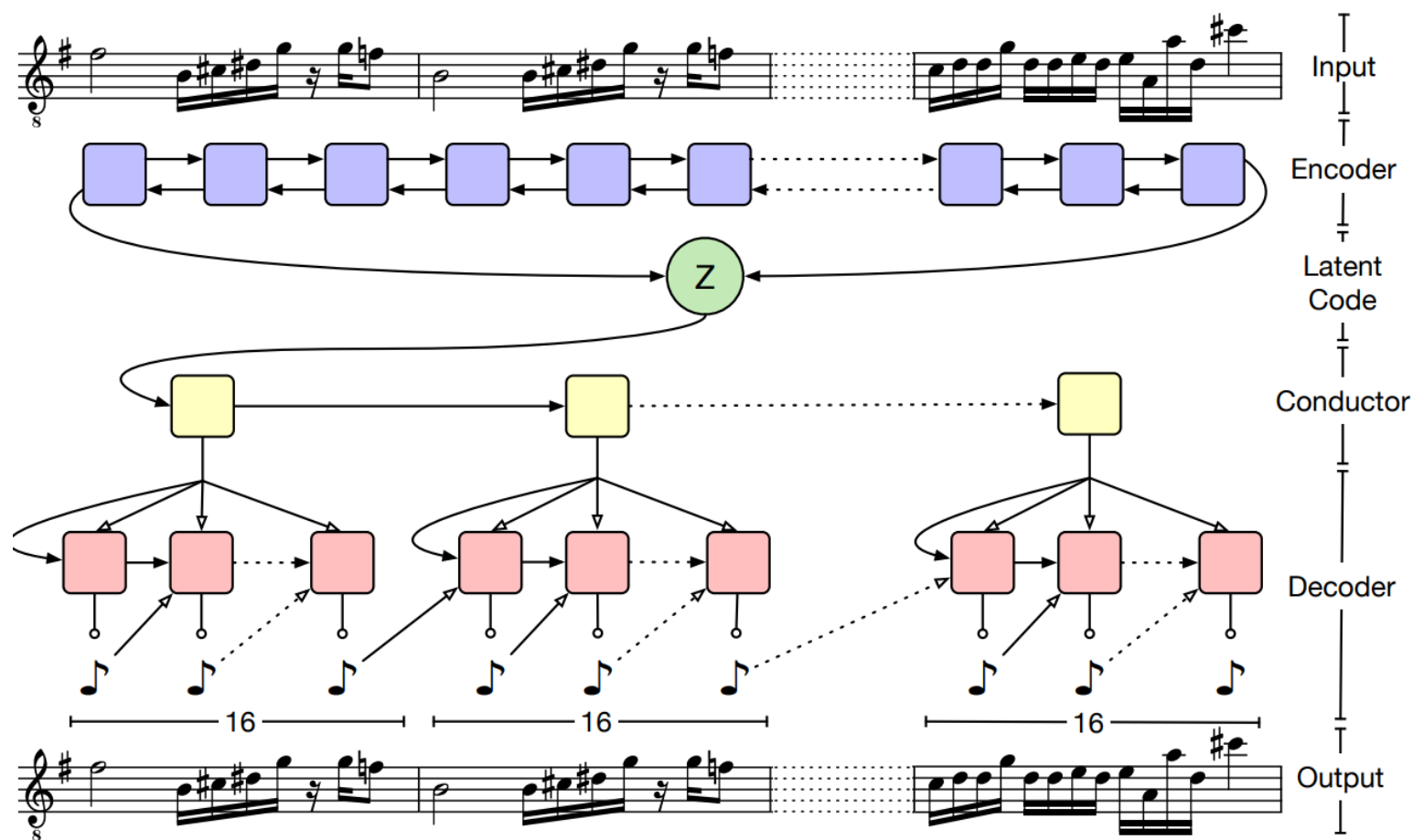
Waveform



Spectrogram

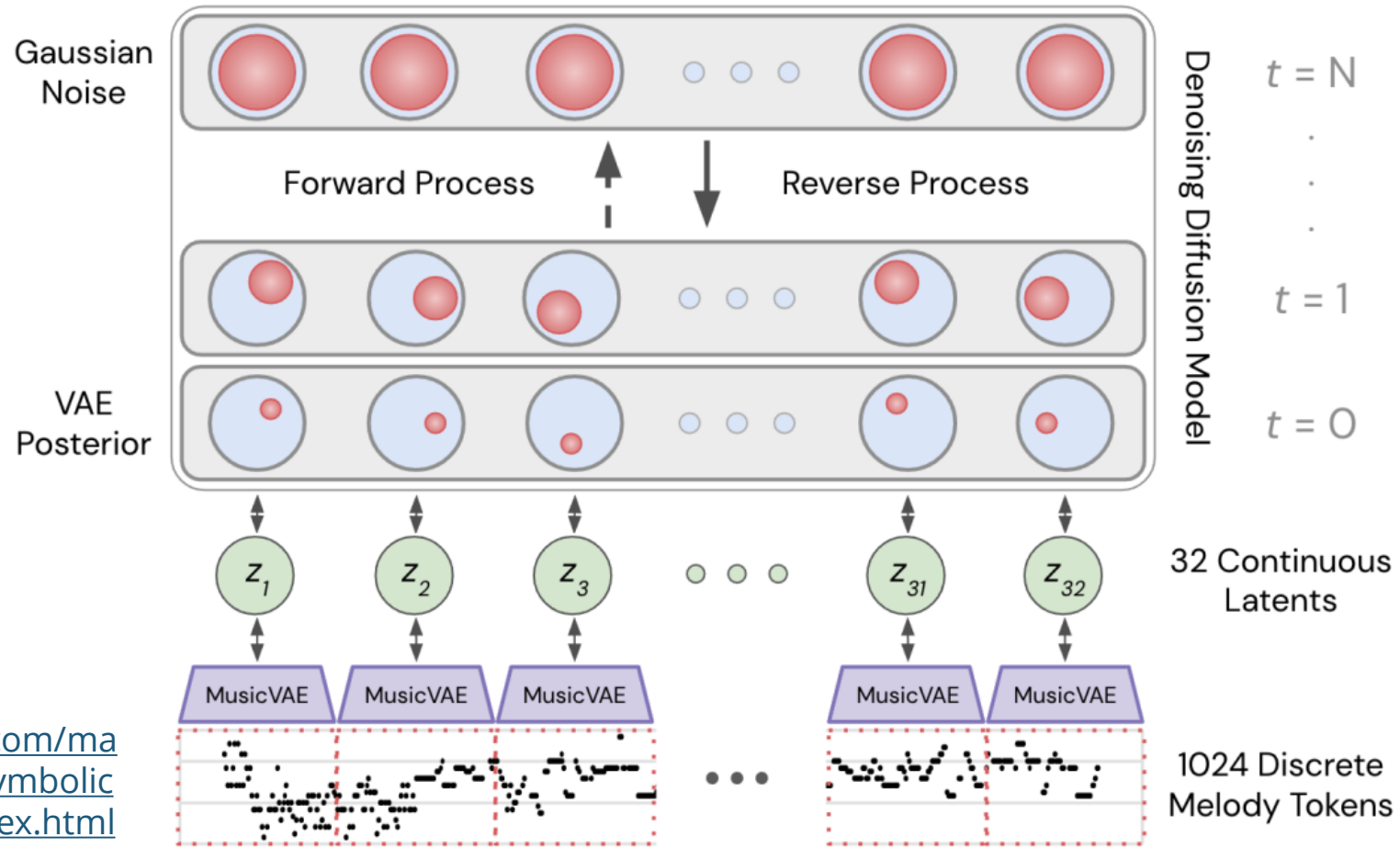
Today, we also have many **latent-space based systems!**

(Recap) Example: MusicVAE (Roberts et al., 2018)



(Source: Roberts et al., 2018)

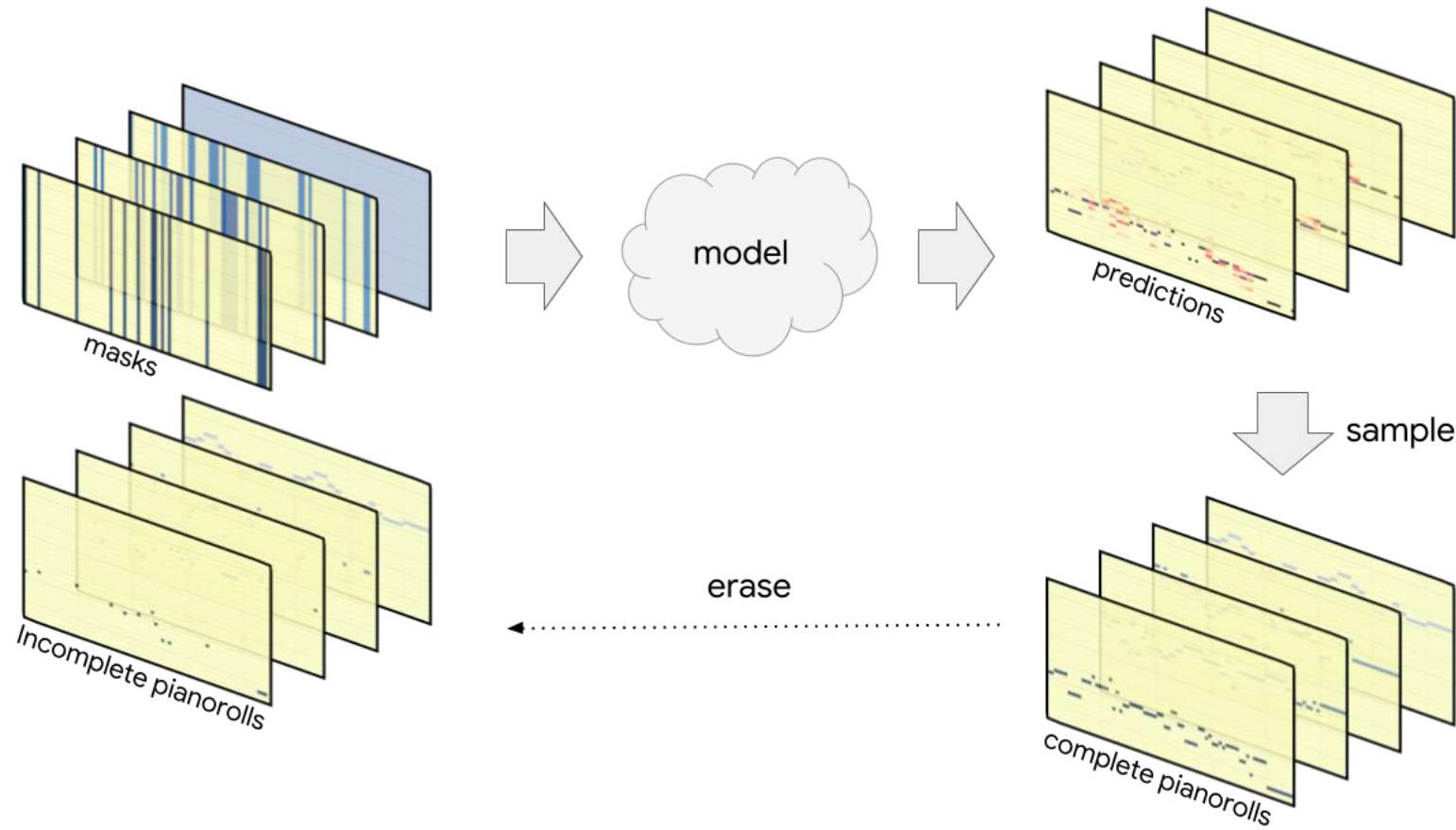
(Recap) Example: Latent Diffusion (Mittal et al., 2021)



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

(Source: Mittal et al., 2021)

(Recap) Example: **Coconet** (Huang et al., 2017)



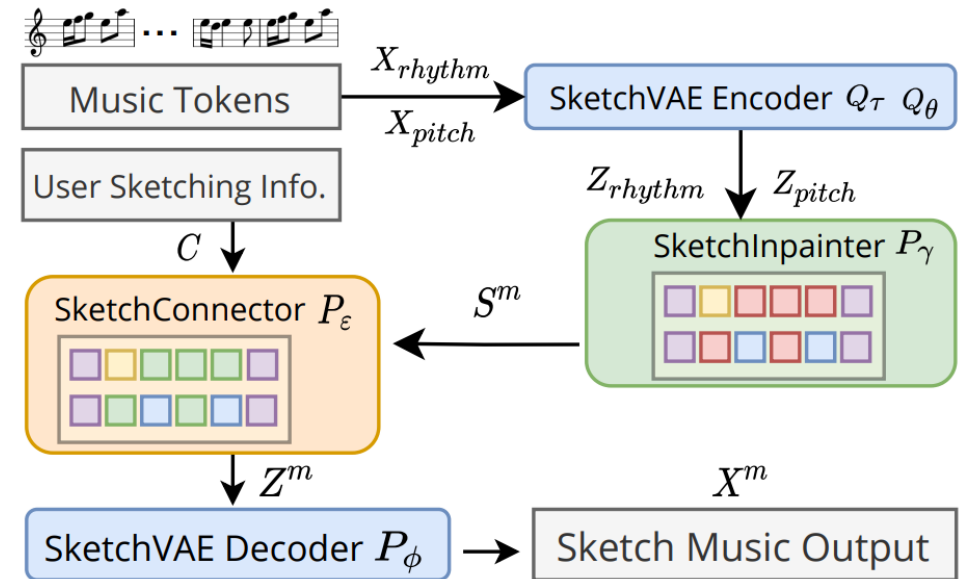
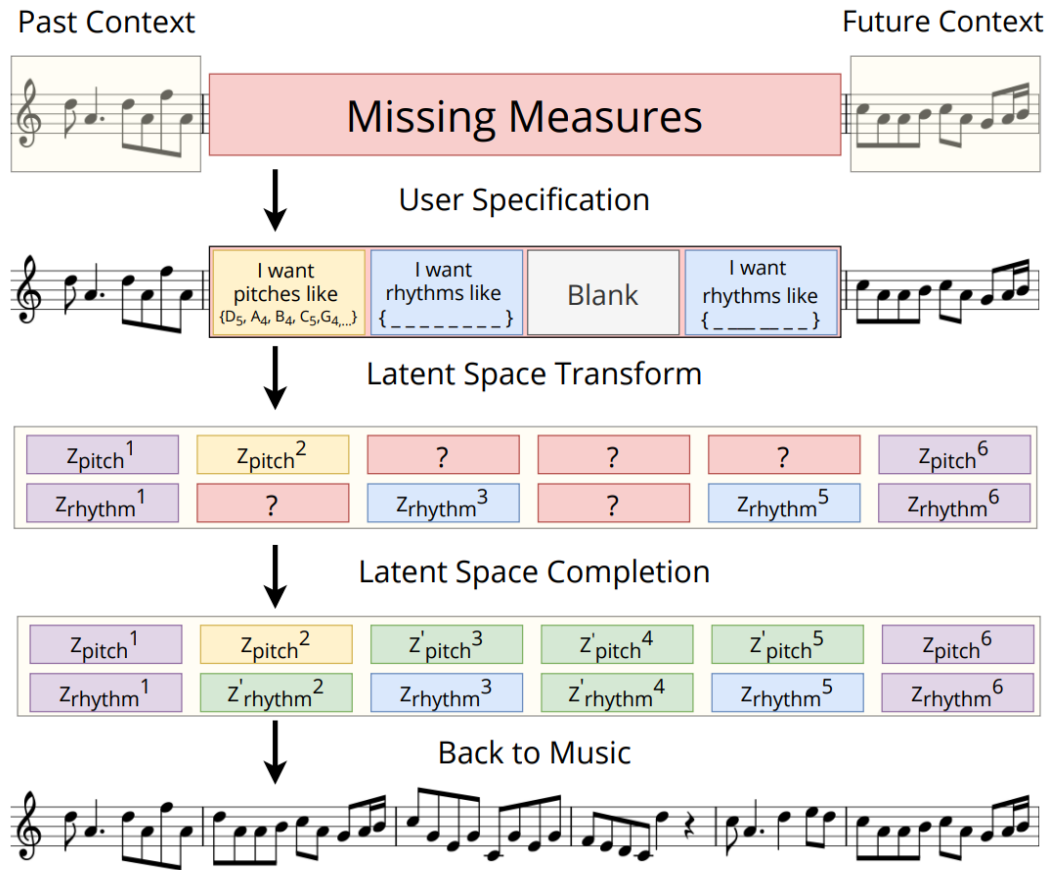
(Source: Huang et al., 2019)

(Recap) Example: Coconet (Huang et al., 2017)



(Source: Huang et al., 2017)

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



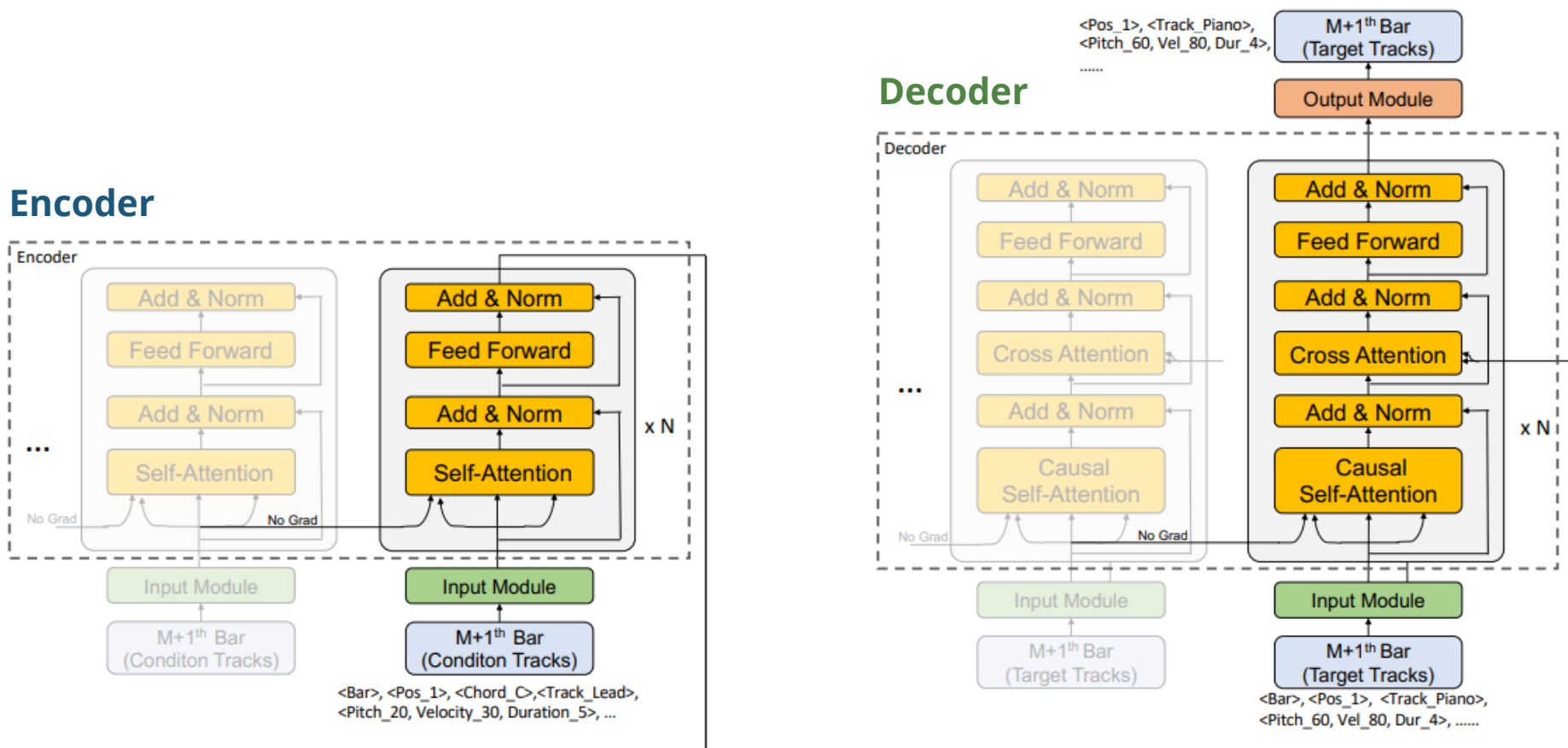
(Source: Chen et al., 2020)

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

The diagram illustrates the Music SketchNet architecture. It consists of four staves: Original, Control Pitch, Control Rhythm, and Control Both. The score is divided into three sections: Past Context, Generation, and Future Context. The Control Pitch track shows chord sets: {Ab5, Db6, Eb6, Gb6} in the Past Context; {C6, Eb6, Db6, F6, Db6} and {F6, Gb6, Ab6, Ab6, F6} in the Generation section; and {Db6, F6, Ab6, Bb6, Db6} in the Future Context. The Control Rhythm track uses pink bars to indicate rhythmic sketches, with a grey bar labeled 'No Sketch' in the Future Context. The Control Both track combines pitch and rhythm information, with a '3' indicating a triplet in the Past Context.

(Source: Chen et al., 2020)

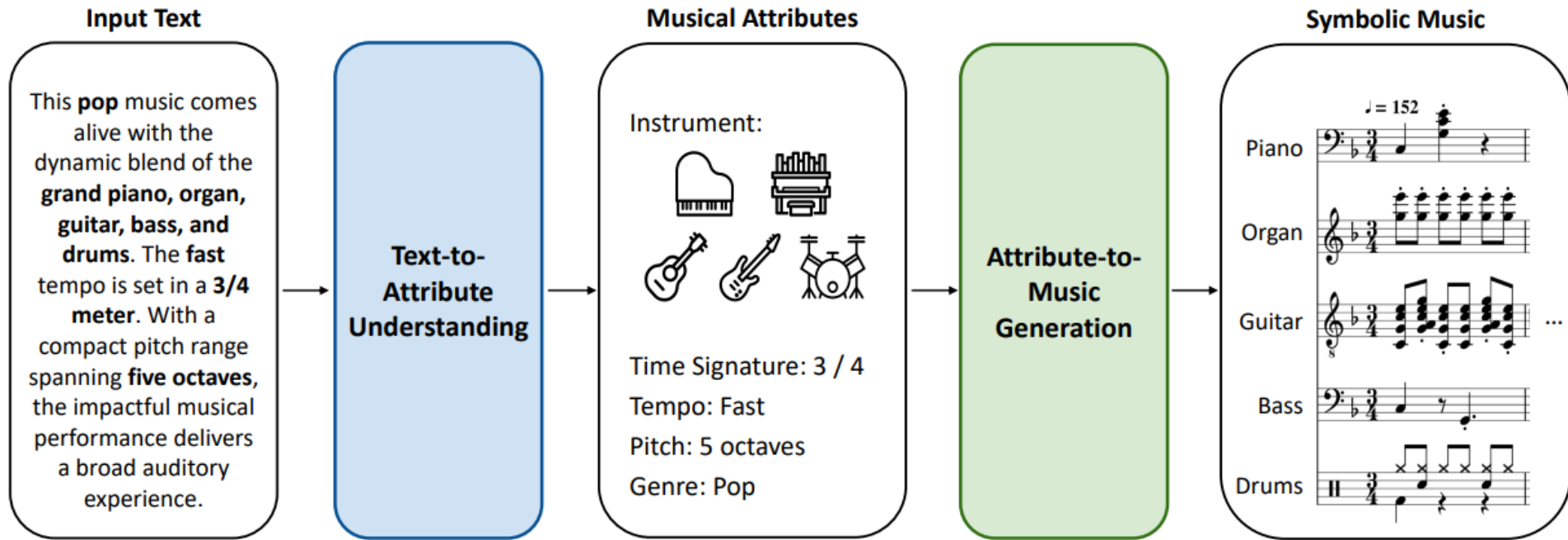
(Recap) Example: PopMAG (Ren et al., 2020)



(Source: Ren et al., 2020)

ai-music.github.io/popmag

(Recap) Example: MuseCoco (Lu et al., 2023)



(Source: Lu et al., 2023)

ai-music.github.io/musecoco

(Recap) Example: MetaScore Transformer (Xu et al., 2024)

Examples

Input: Chopin, piano, easy, all rights preserved
Output: an easy piano piece by Chopin. All rights preserved

Input: classical/traditional, Programmatic, robot
Output: A classical music piece.

Input: classical/traditional, Heaven's, His, Jesus, Lord, Son, accords, dawn, day, grace, hymn, light, peace, rest, soul, sovereign
Output: This is a classical/traditional music piece.


Input: Michael Jackson, bass, guitar, rock/metal, technical, electronic/dance. Creative Commons Copyright Waiver
Output: A mix of pop, rock and electronic music composed by Michael Jackson. The music has bass and guitar in it. This music piece is free for use by anyone for any purpose.


Input: William Marshall, adele, advanced, piano, violin, folk/country, soundtrack/stage. Attribution
Output: An advanced piano and violin music piece composed by William Marshall and Adele. This is a soundtrack with folk vibe. Users must give credit to the creator when using the music.

Input: classical/traditional, piano, violin, intermediate, Rm, 20
Output: A intermediate classical music piece. This music piece has piano and violin in it.

Query

Input: {classical/traditional, easy, alexander walker, public domain}
Output:

Prompt: This is an easy classical piano piece composed by wolfgang amadeus mozart. 

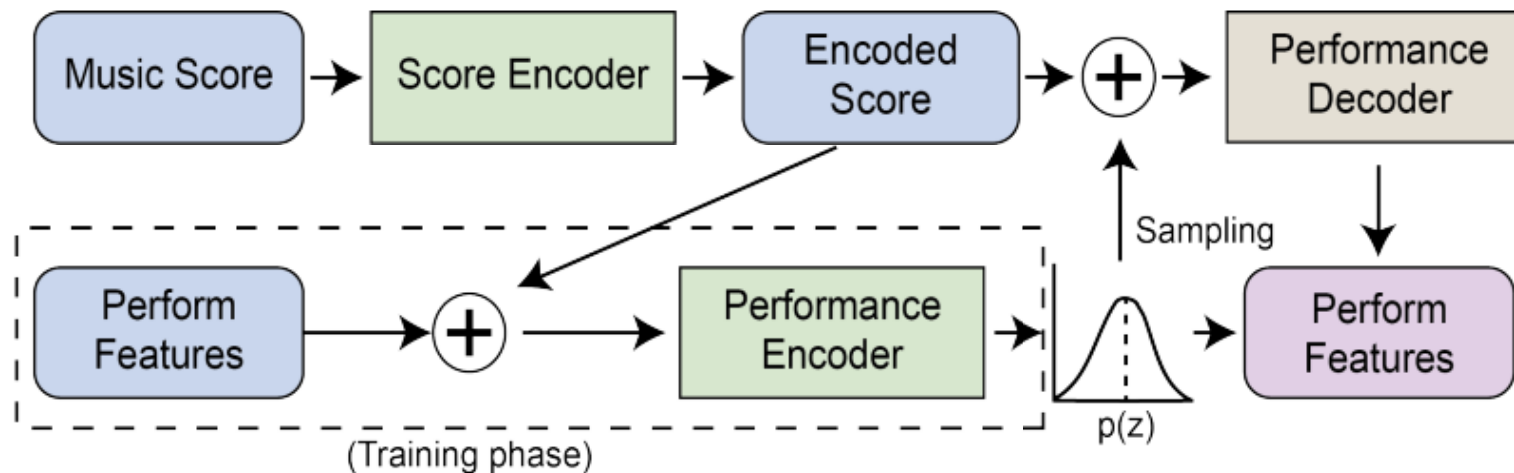
Prompt: A short and emotional music piece inspired by an anime scene. 

Prompt: A powerful orchestral music piece. 

(Source: Xu et al., 2024)

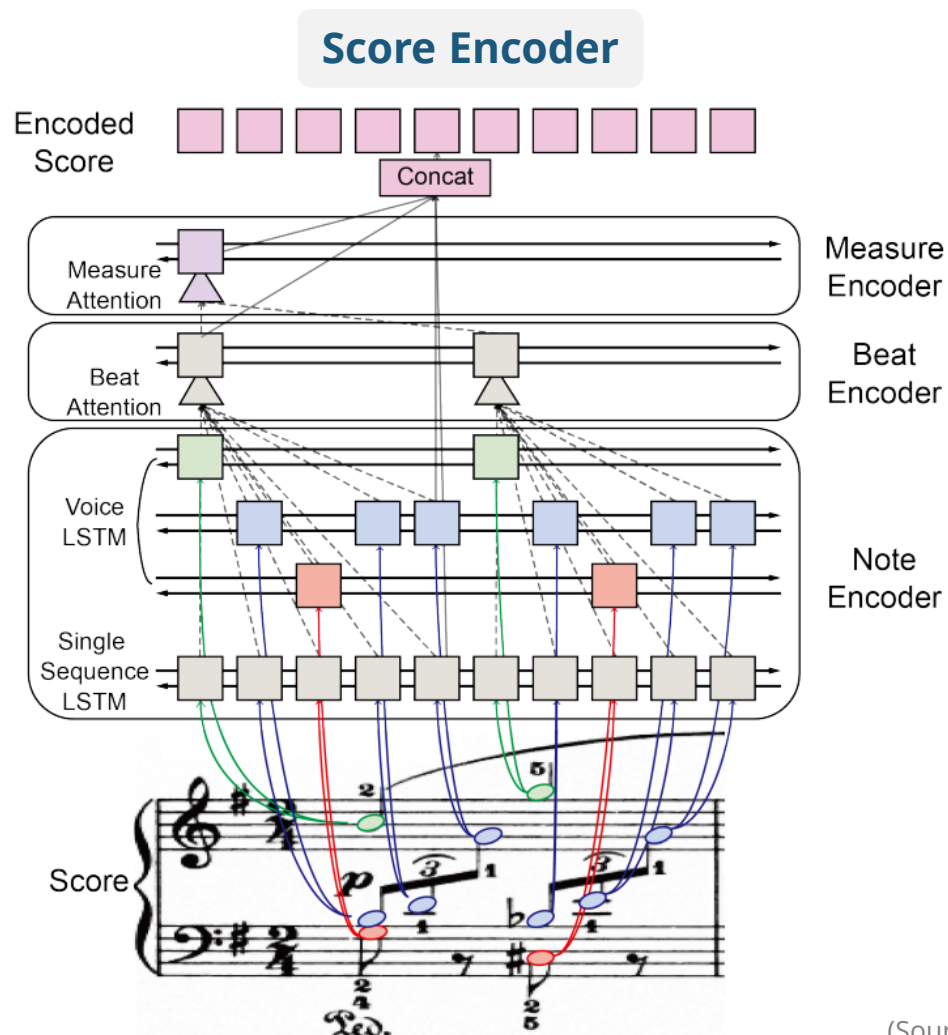
(Recap) Example: **VirtuosoNet** (Jeong et al., 2019)

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

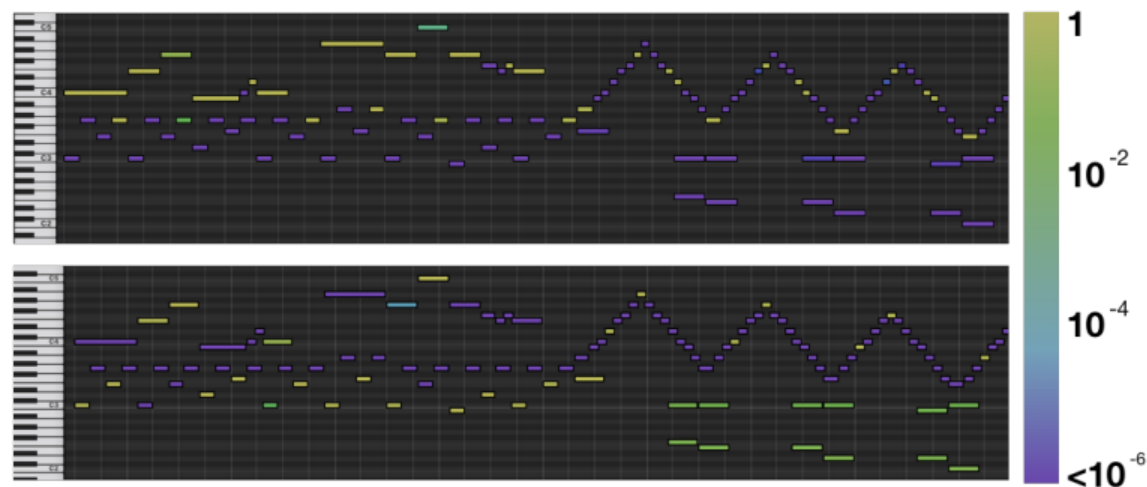


(Source: Jeong et al., 2019)

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



Attention visualization



(Source: Jeong et al., 2019)

(Recap) Four Paradigms



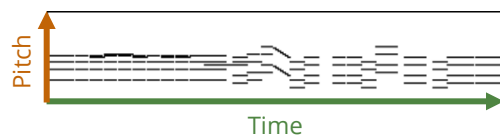
Symbolic music generation

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



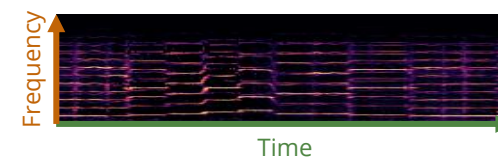
Audio-domain music generation

Time series-based

Image-based



Waveform

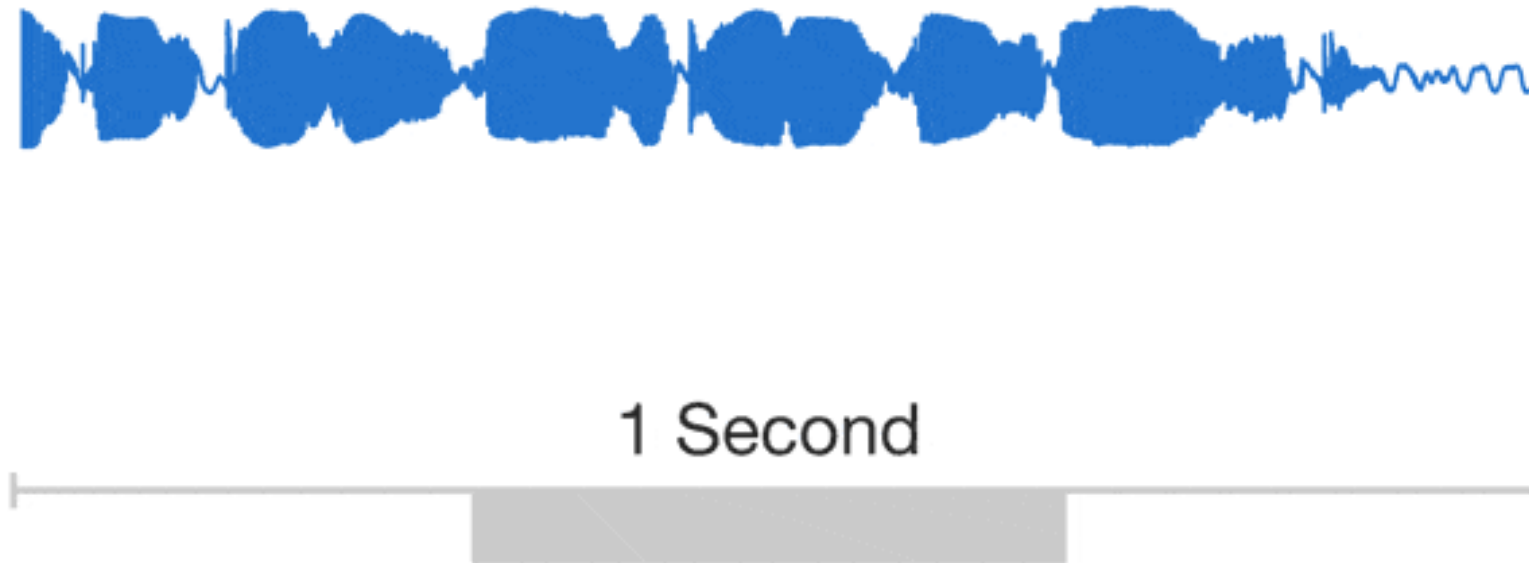


Spectrogram

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

Autoregressive Waveform Synthesis

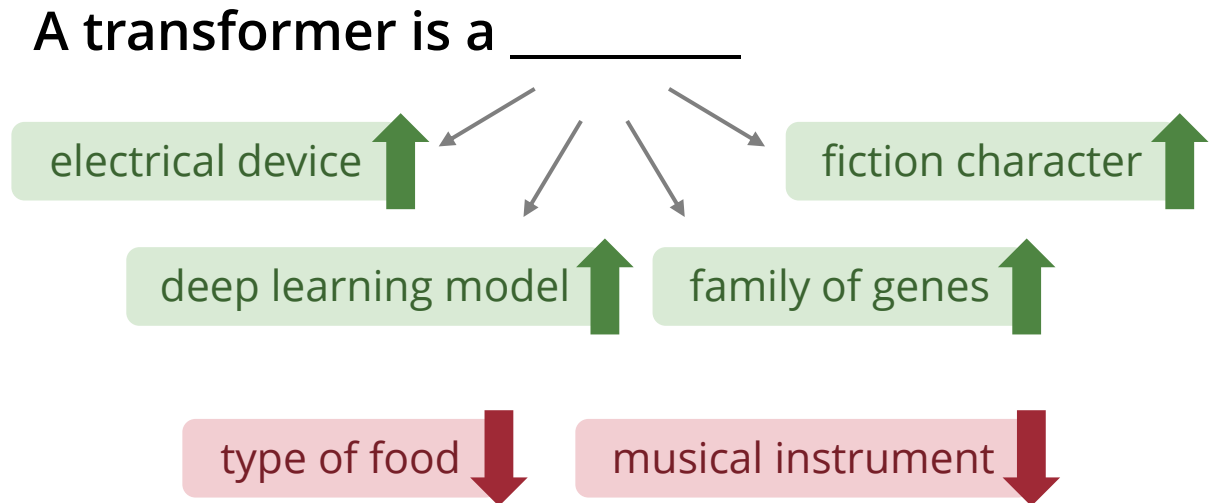
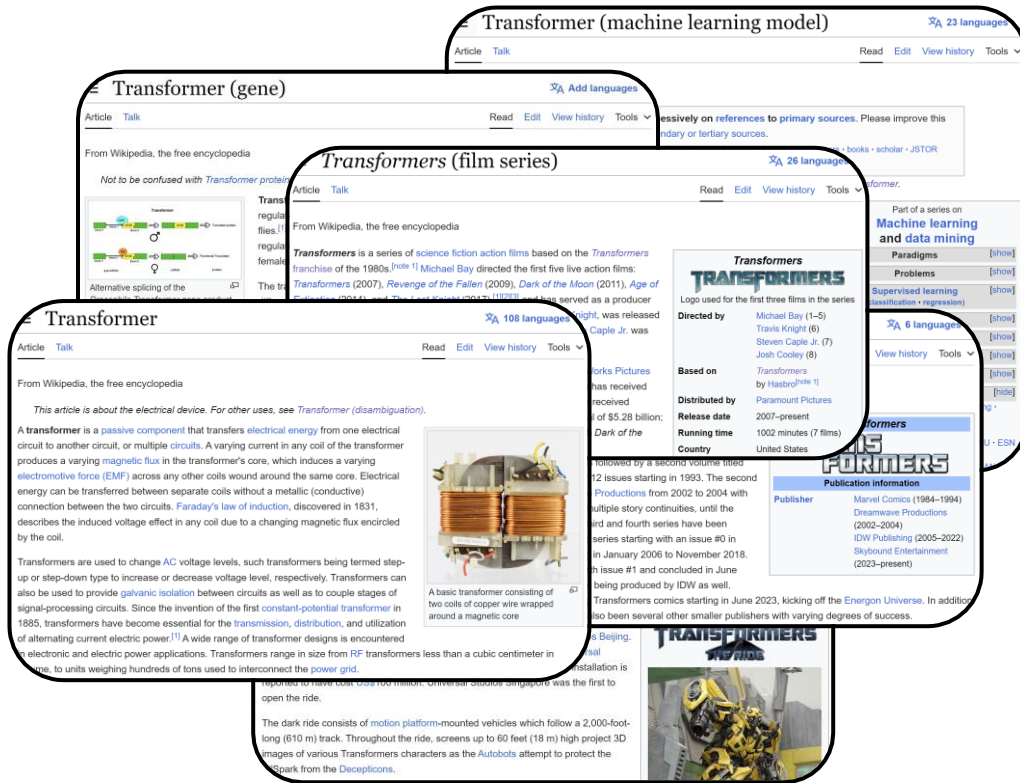
Generating Waveforms using a Neural Network



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

(Recap) Language Models

- Predicting the next word **given the past sequence of words**



(Recap) Language Models (Mathematically)

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

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

Next word Previous words

$P(\text{electrical} \mid \text{A transformer is a})$	↑
$P(\text{character} \mid \text{A transformer is a})$	↑
$P(\text{gene} \mid \text{A transformer is a})$	↑
$P(\text{model} \mid \text{A transformer is a})$	↑
$P(\text{food} \mid \text{A transformer is a})$	↓
$P(\text{musical} \mid \text{A transformer is a})$	↓

Autoregressive Models (Mathematically)

- A class of machine learning models that **learn** the probability of the next value given previous values

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

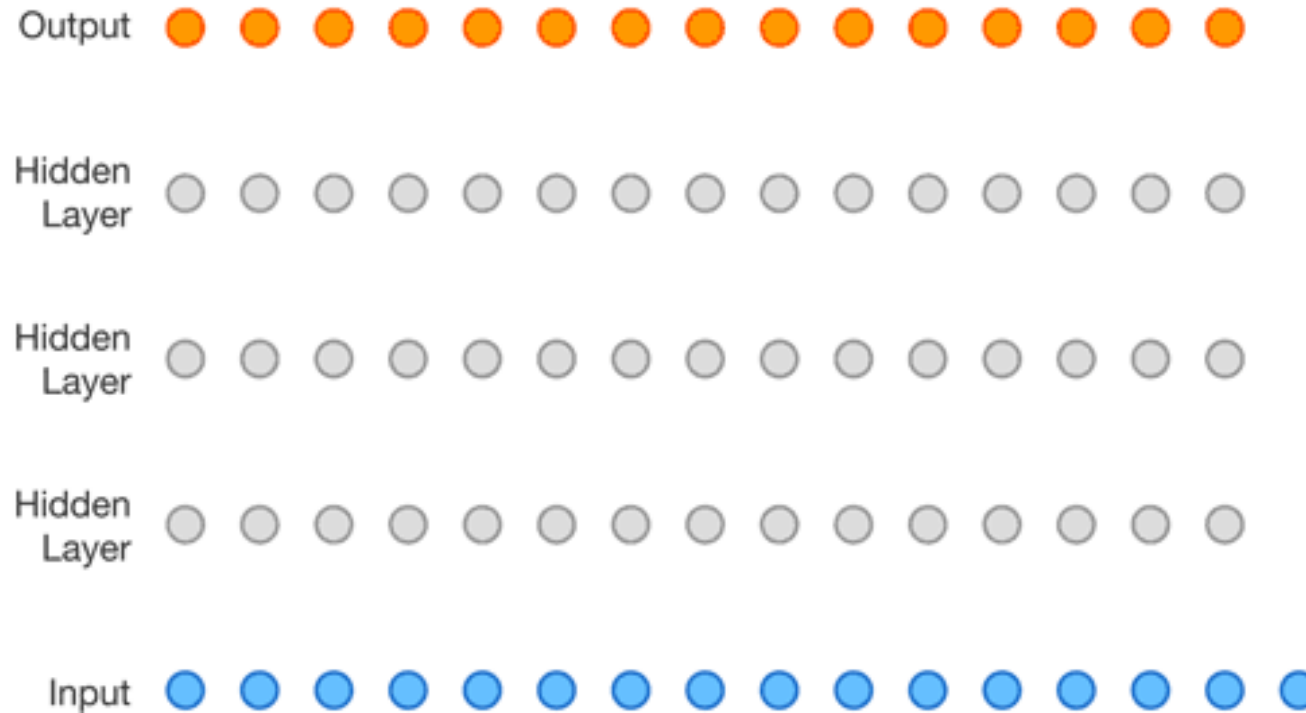
Next number Previous numbers

$$P(0.1 \mid 0.5, 0.4, 0.3, 0.2) \quad \uparrow$$
$$P(0.09 \mid 0.5, 0.4, 0.3, 0.2) \quad \uparrow$$
$$P(0.11 \mid 0.5, 0.4, 0.3, 0.2) \quad \uparrow$$
$$P(99 \mid 0.5, 0.4, 0.3, 0.2) \quad \downarrow$$
$$P(-1 \mid 0.5, 0.4, 0.3, 0.2) \quad \downarrow$$

The term “autoregressive” has different definitions in machine learning and signal processing.
In signal processing, an autoregressive model needs to be a linear model.

Unconditional Audio Synthesis using CNNs

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

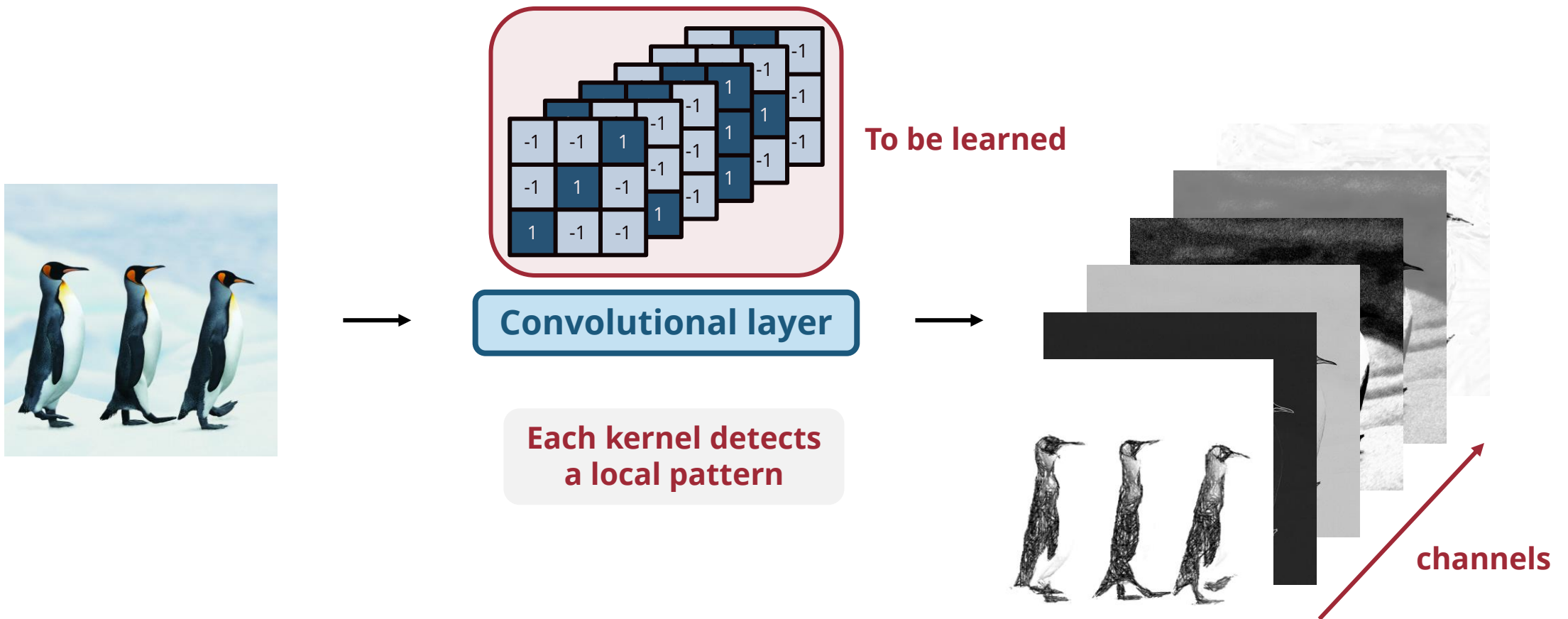


(Source: van den Oord et al., 2016)

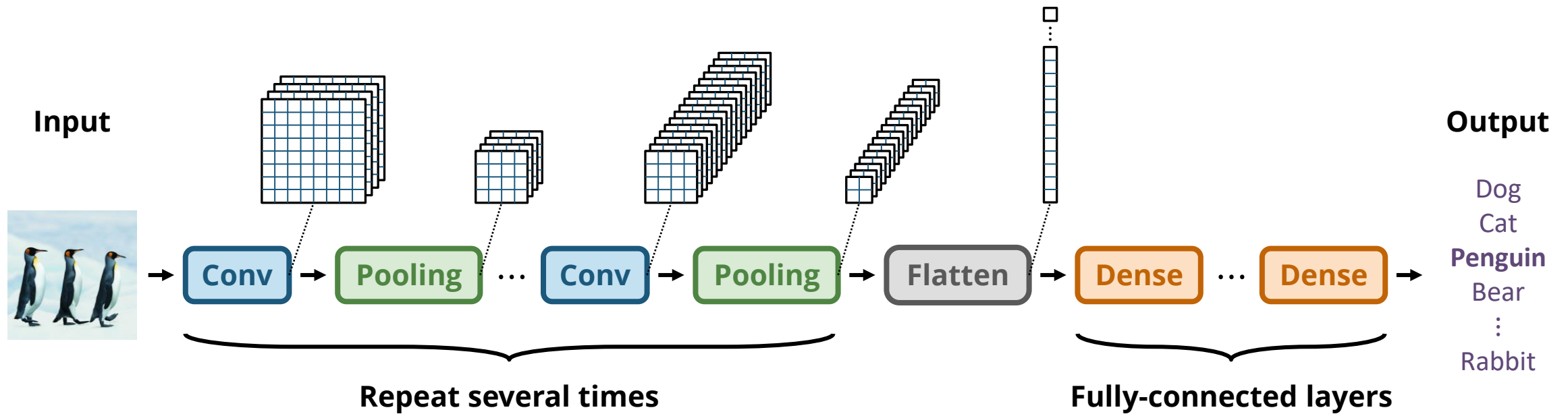
A convolutional neural network for raw waveform generation

(Recap) Convolutional Layer

- A convolutional layer consists of many **learnable kernels** (channels)

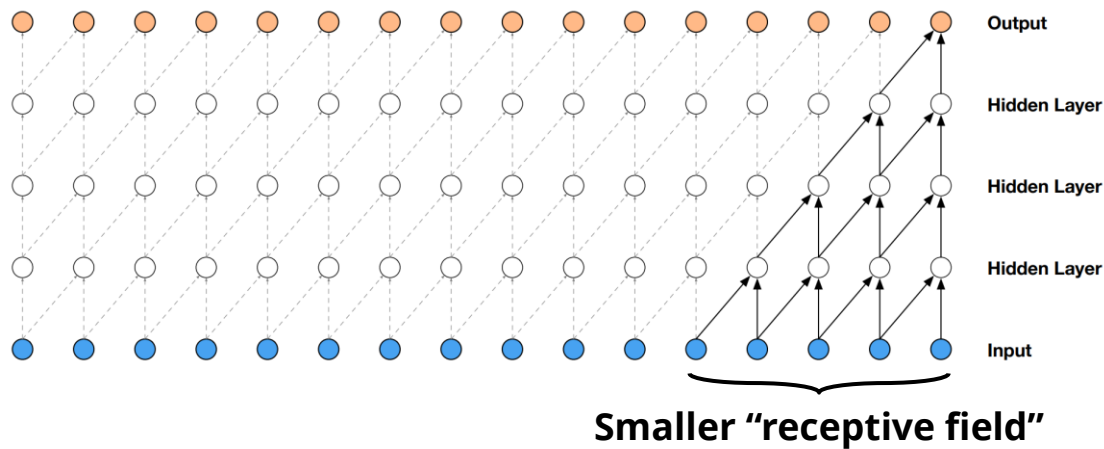


(Recap) Convolutional Neural Network (CNNs)

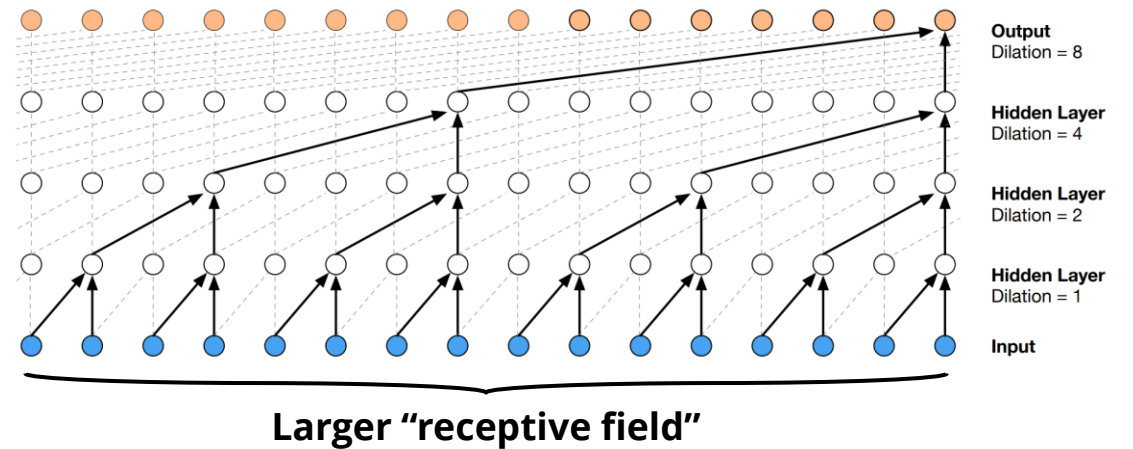


Example: WaveNet (van den Oord et al., 2016)

Standard convolution



Dilated convolution



deepmind.google/discover/blog/wavenet-a-generative-model-for-raw-audio

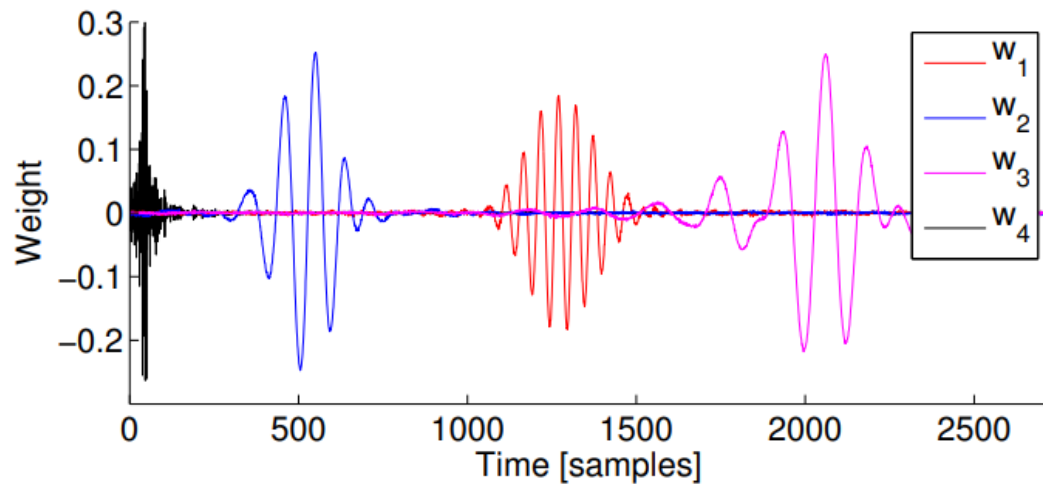
Example of generated music



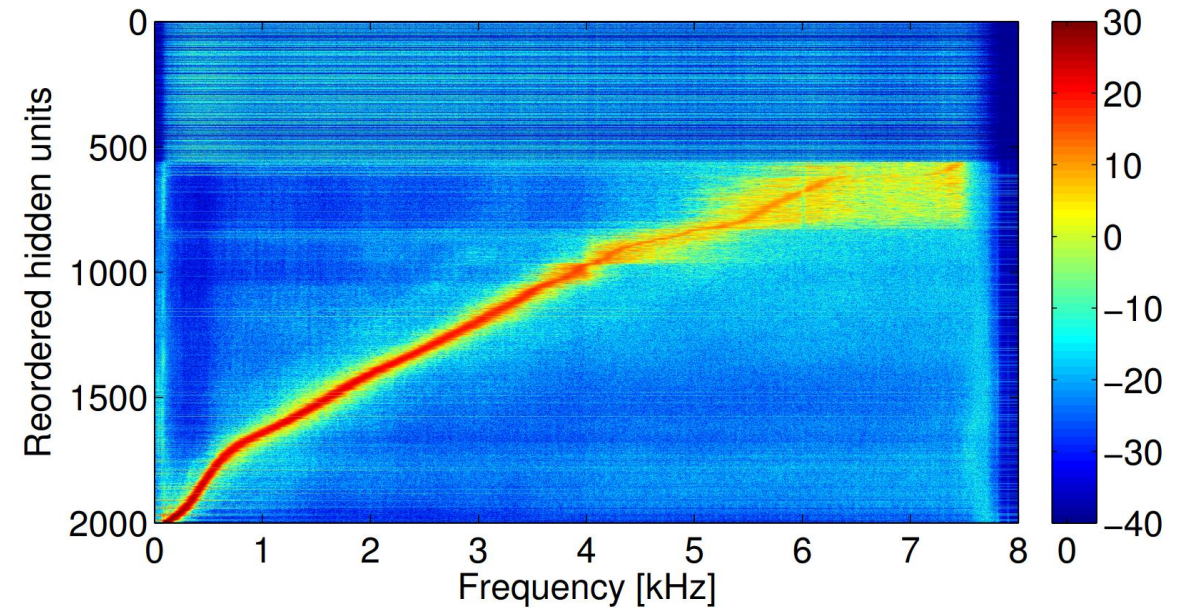
(Source: van den Oord et al., 2016)

1D CNNs & Fourier Transform

Convolution kernels learned

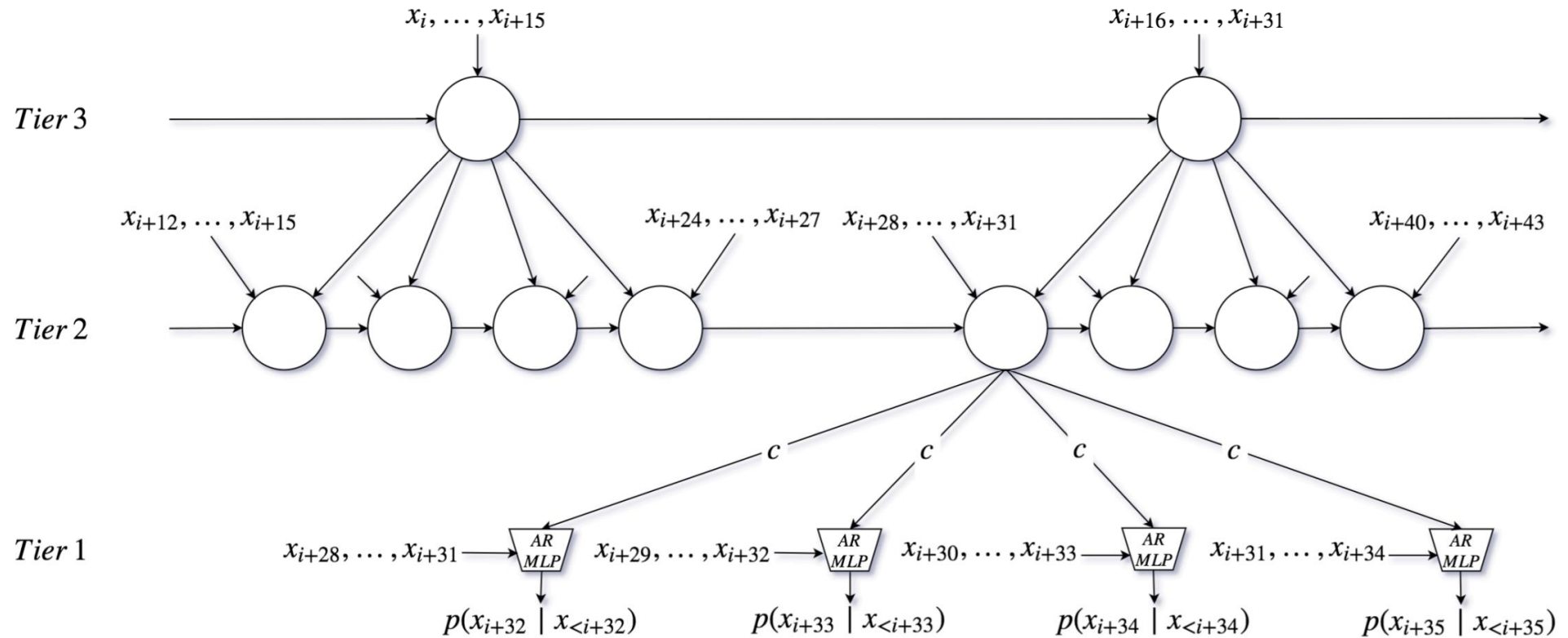


Peak frequency detected by the learned kernels



Unconditional Audio Synthesis using RNNs

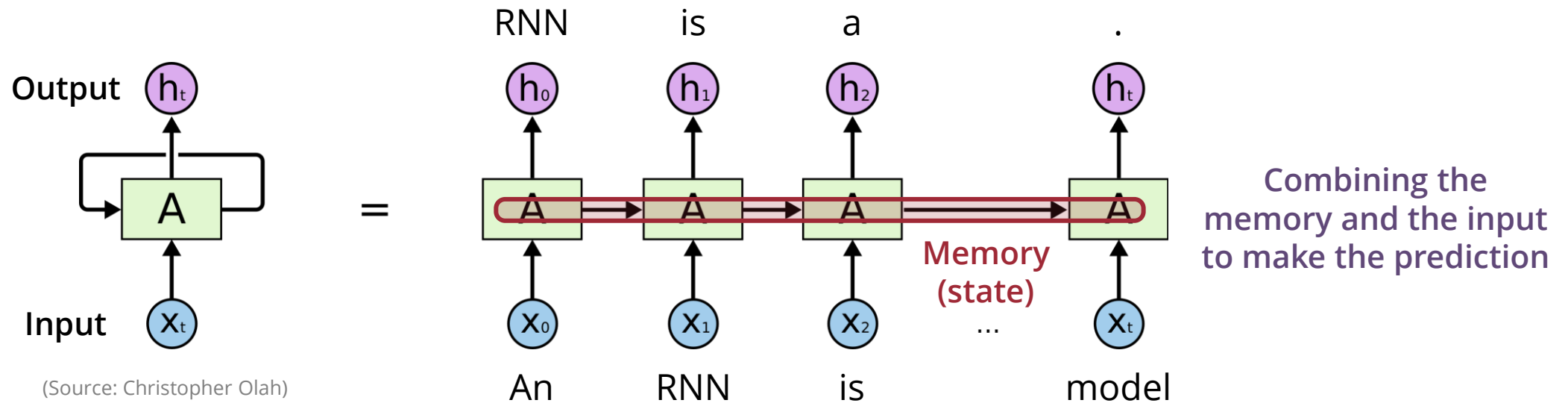
Example: SampleRNN (Mehri et al., 2017)



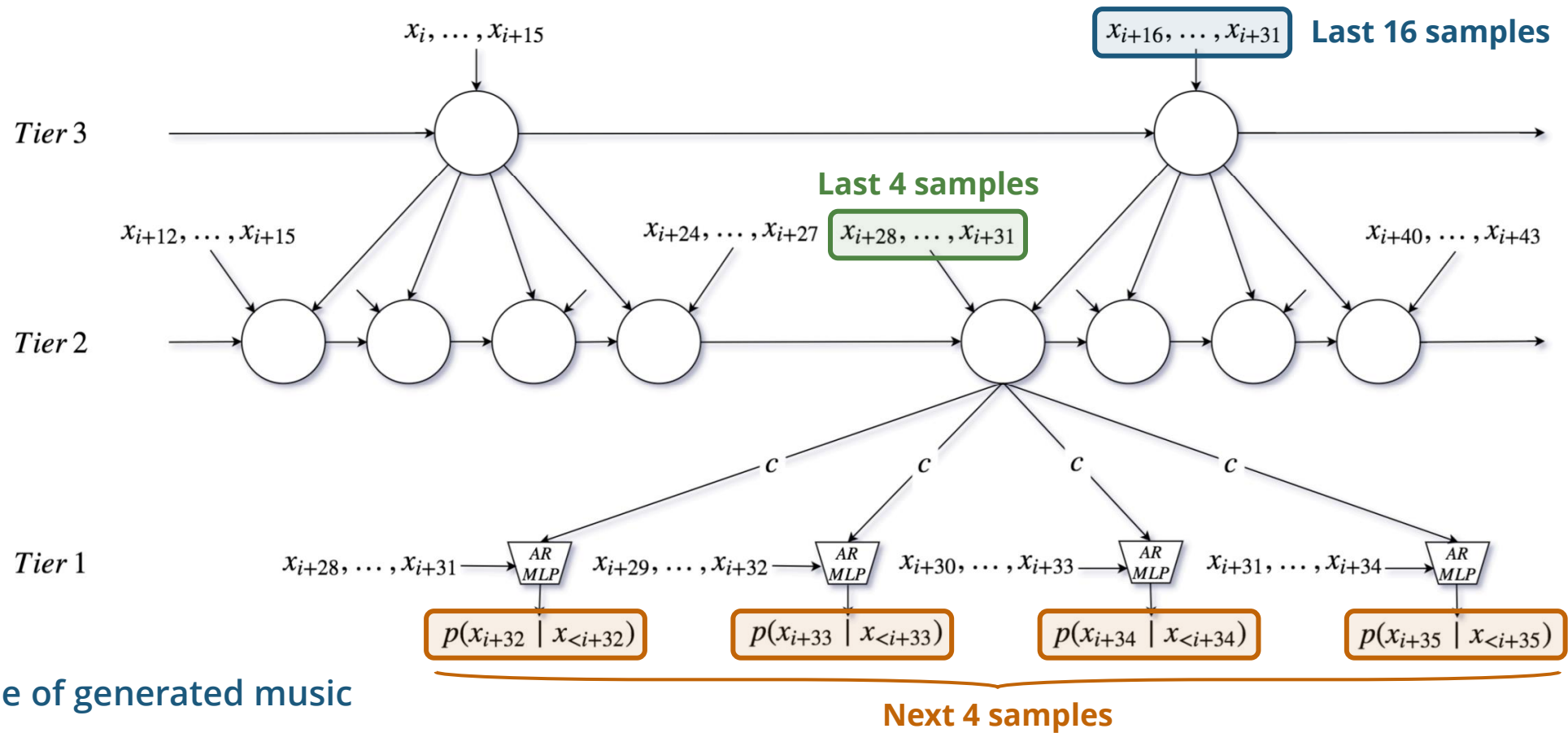
(Source: Mehri et al., 2017)

(Recap) What is an RNN (Recurrent Neural Network)?

- A type of neural networks that have **loops**
- Widely used for **modeling sequences** (e.g., in natural language processing)



Example: SampleRNN (Mehri et al., 2017)



Example of generated music

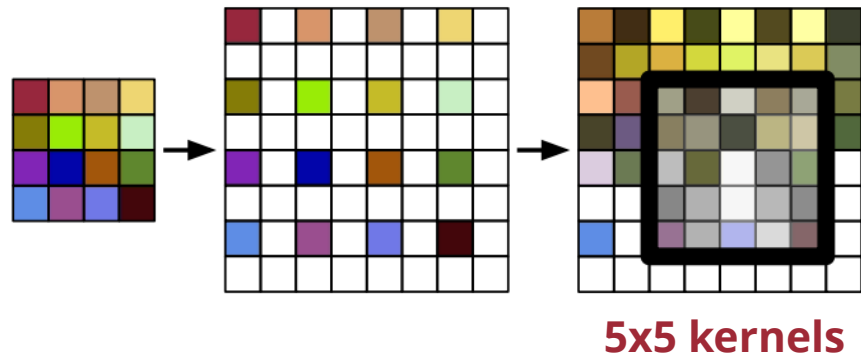


(Source: Mehri et al., 2017)

Unconditional Audio Synthesis using GANs

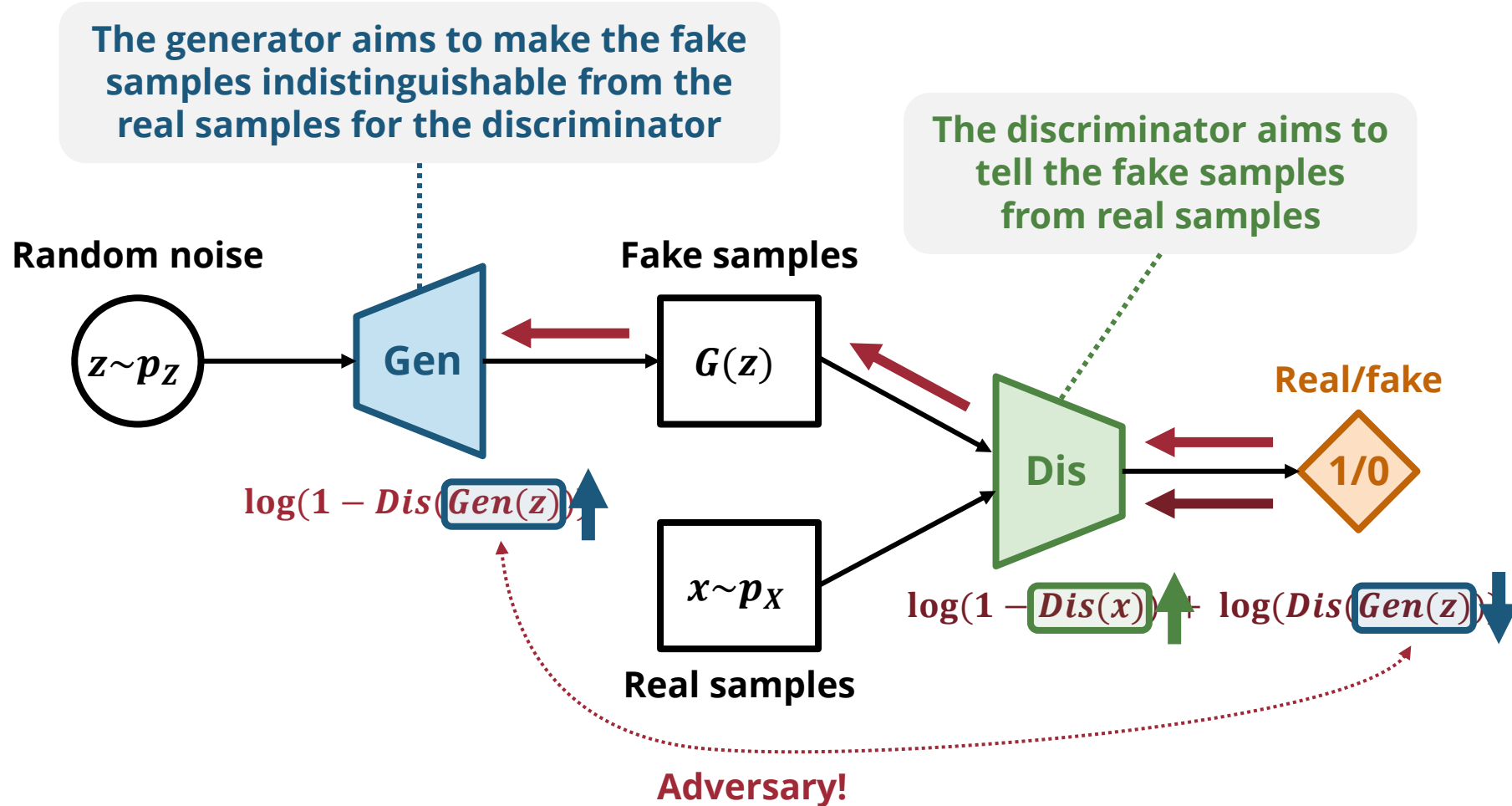
Example: WaveGAN (Donahue et al., 2019)

DCGAN for images

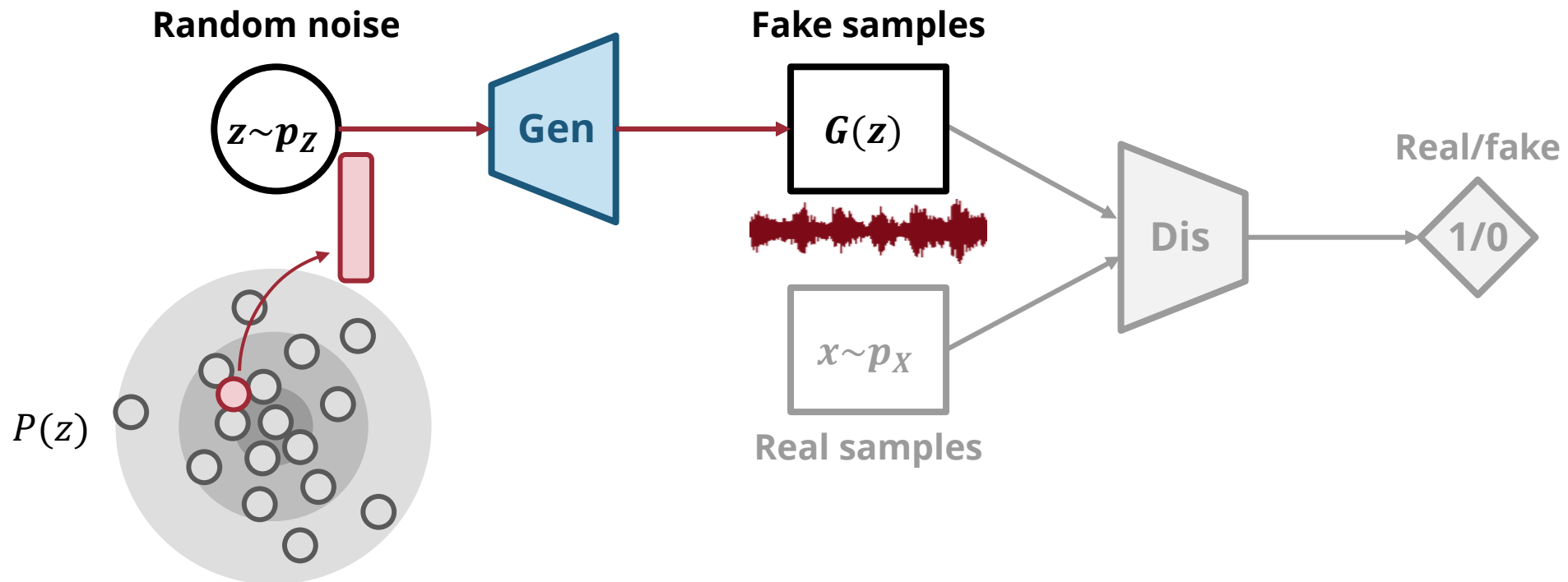


(Source: Donahue et al., 2019)

(Recap) Generative Adversarial Nets (GANs) – Training

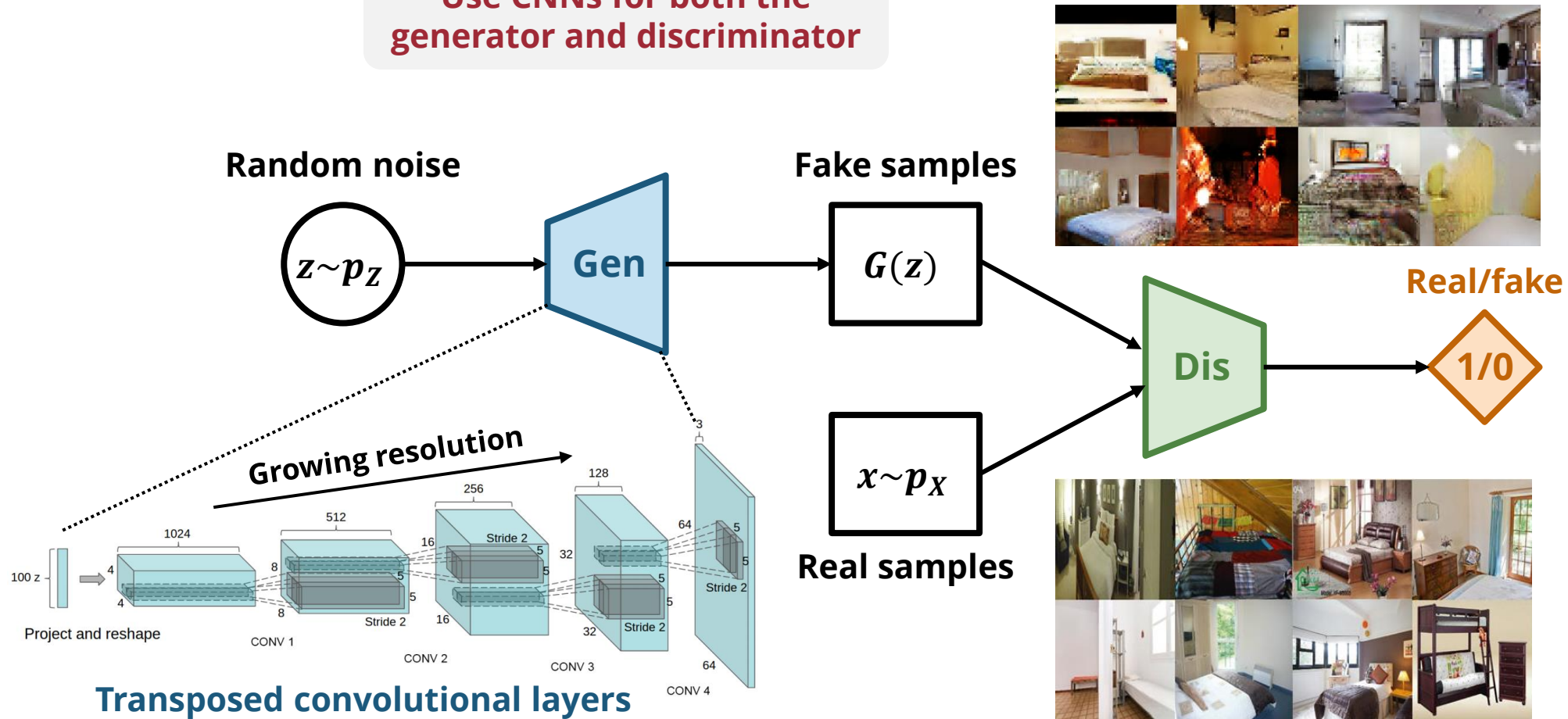


(Recap) Generative Adversarial Nets (GANs) – Generation



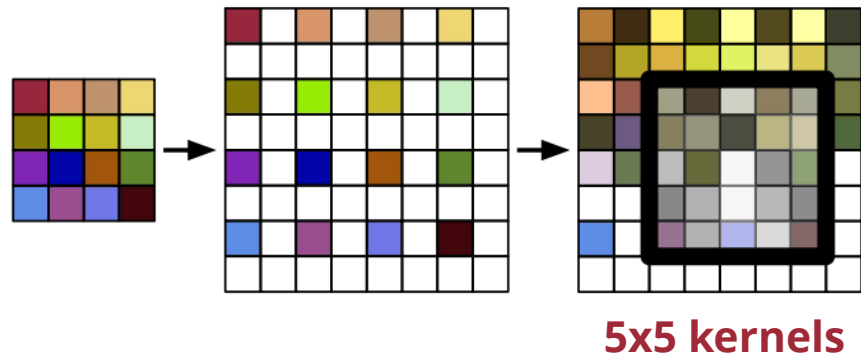
(Recap) Deep Convolutional GANs (DCGANs)

Use CNNs for both the generator and discriminator

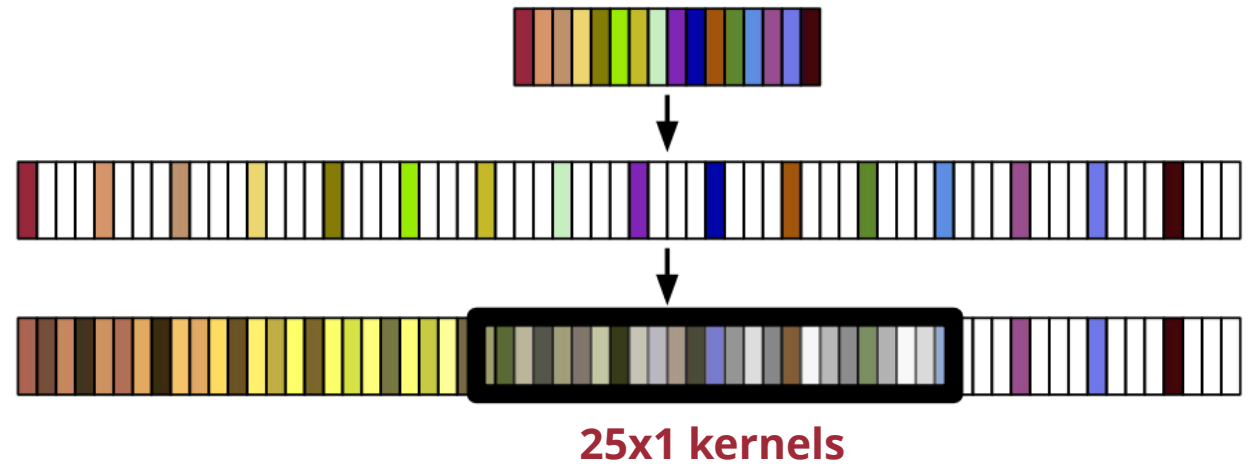


Example: WaveGAN (Donahue et al., 2019)

DCGAN for images



WaveGAN for audio



(Source: Donahue et al., 2019)

chrisdonahue.com/wavegan_examples
chrisdonahue.com/wavegan

Example of generated music

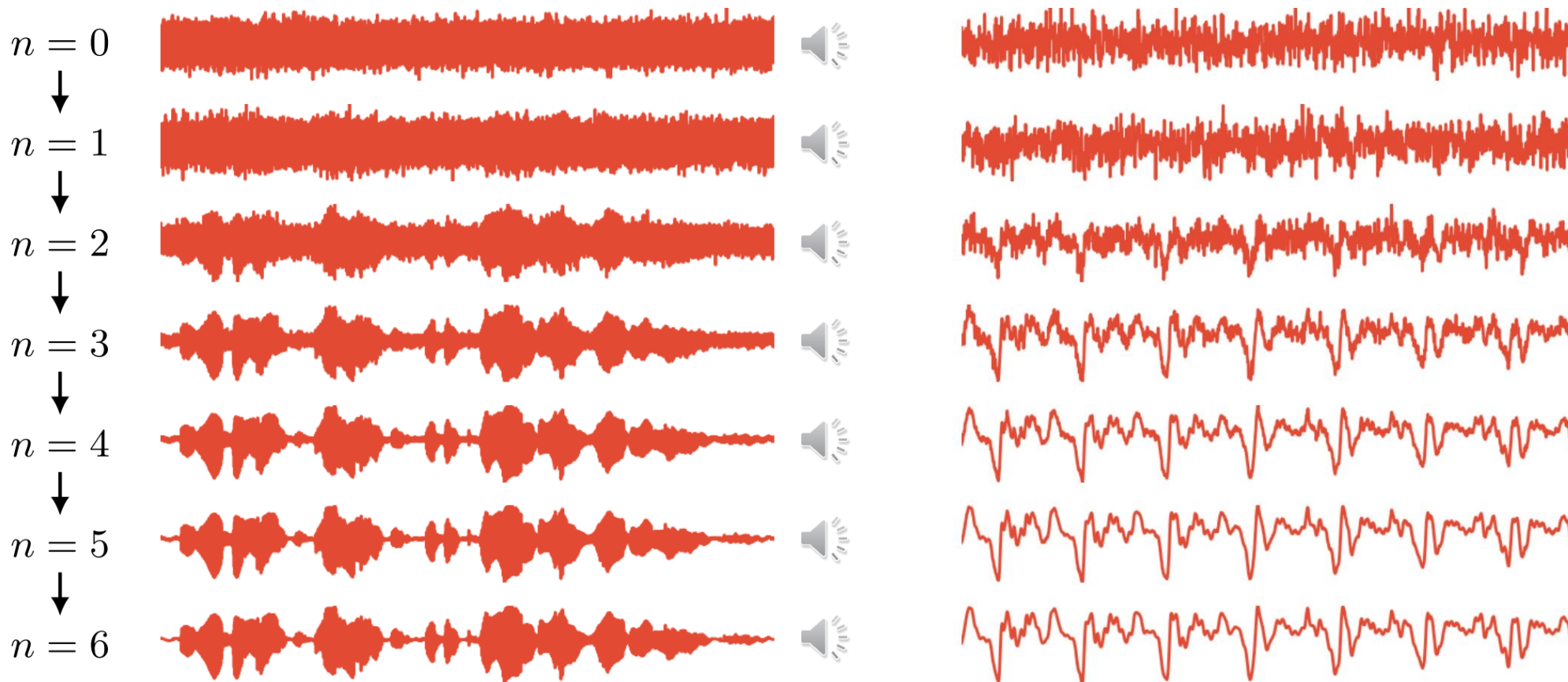


Unconditional Audio Synthesis using Diffusions

Example: WaveGrad (Chen et al., 2021)

Text: Here are the match lineups for the Colombia Haiti match.

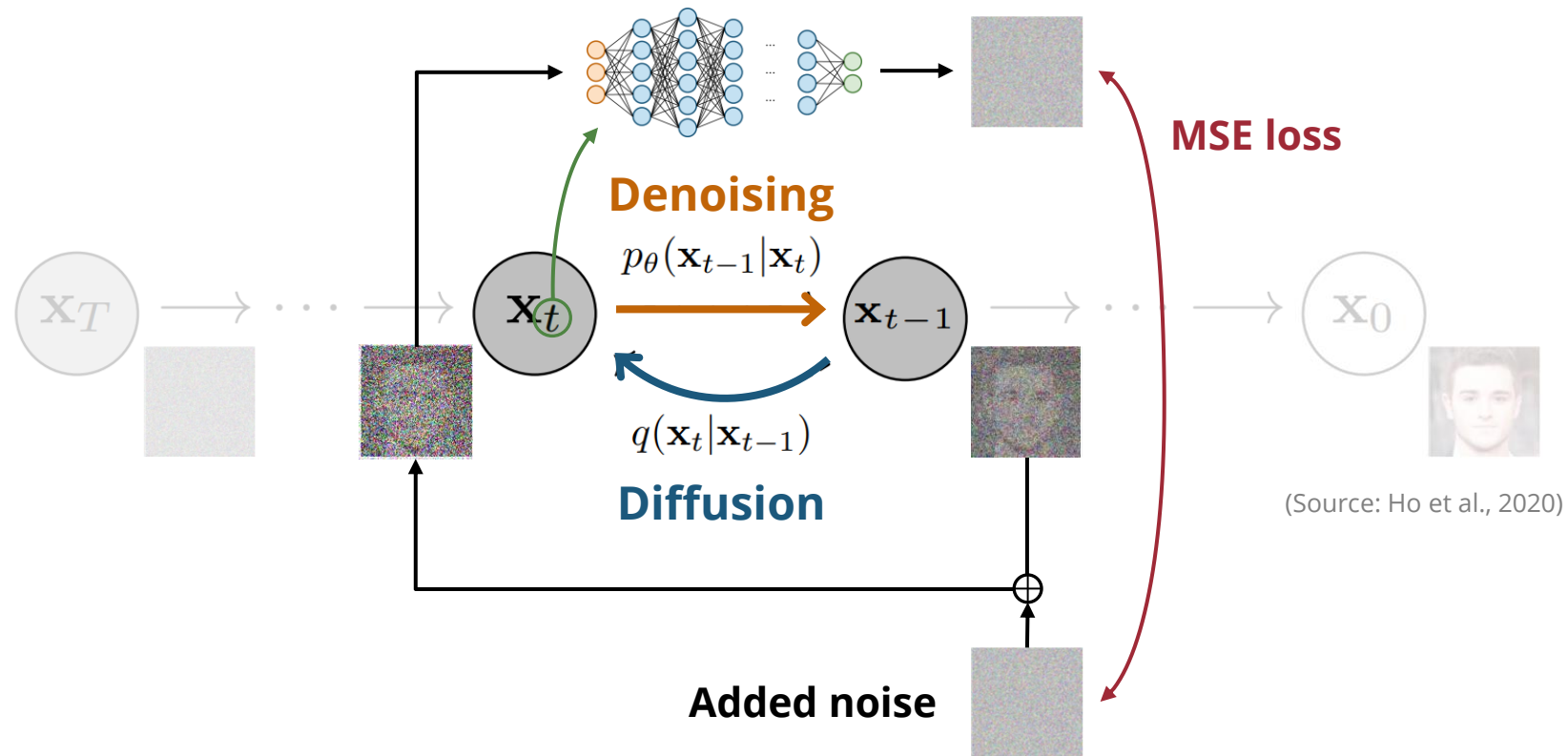
Zoom in



(Source: Chen et al., 2021)

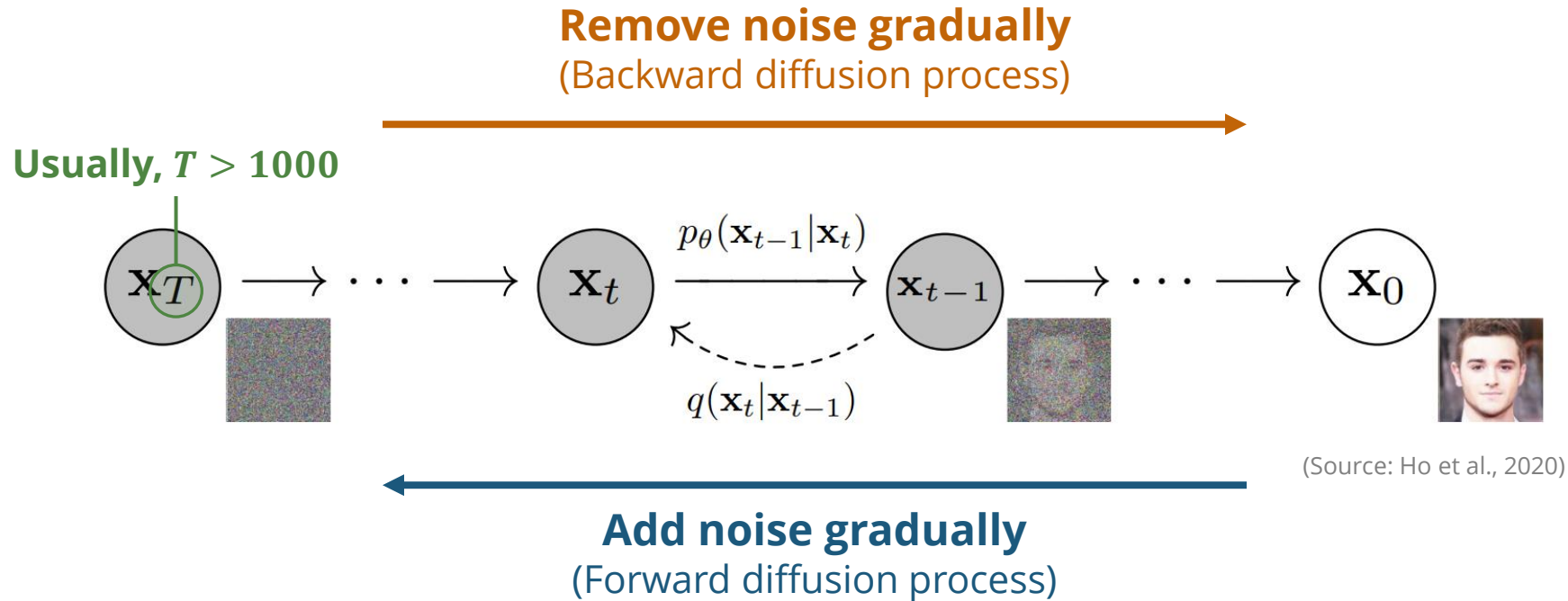
(Recap) Diffusion Models – Training

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

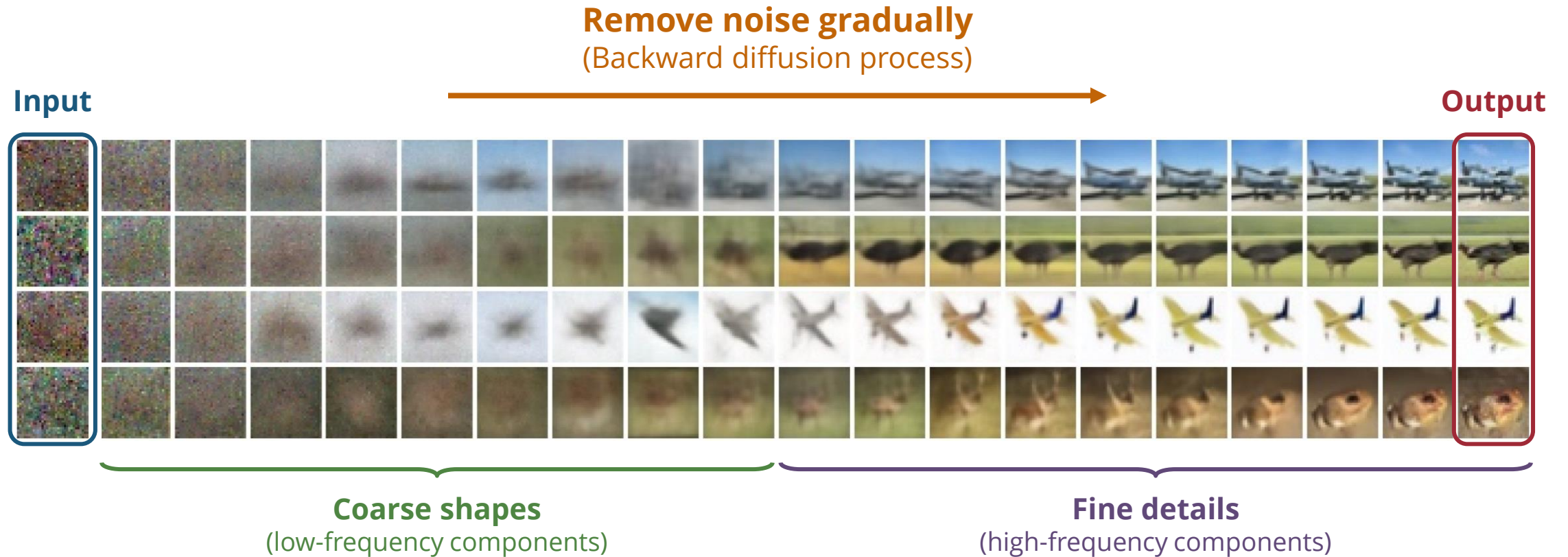


(Recap) Diffusion Models

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

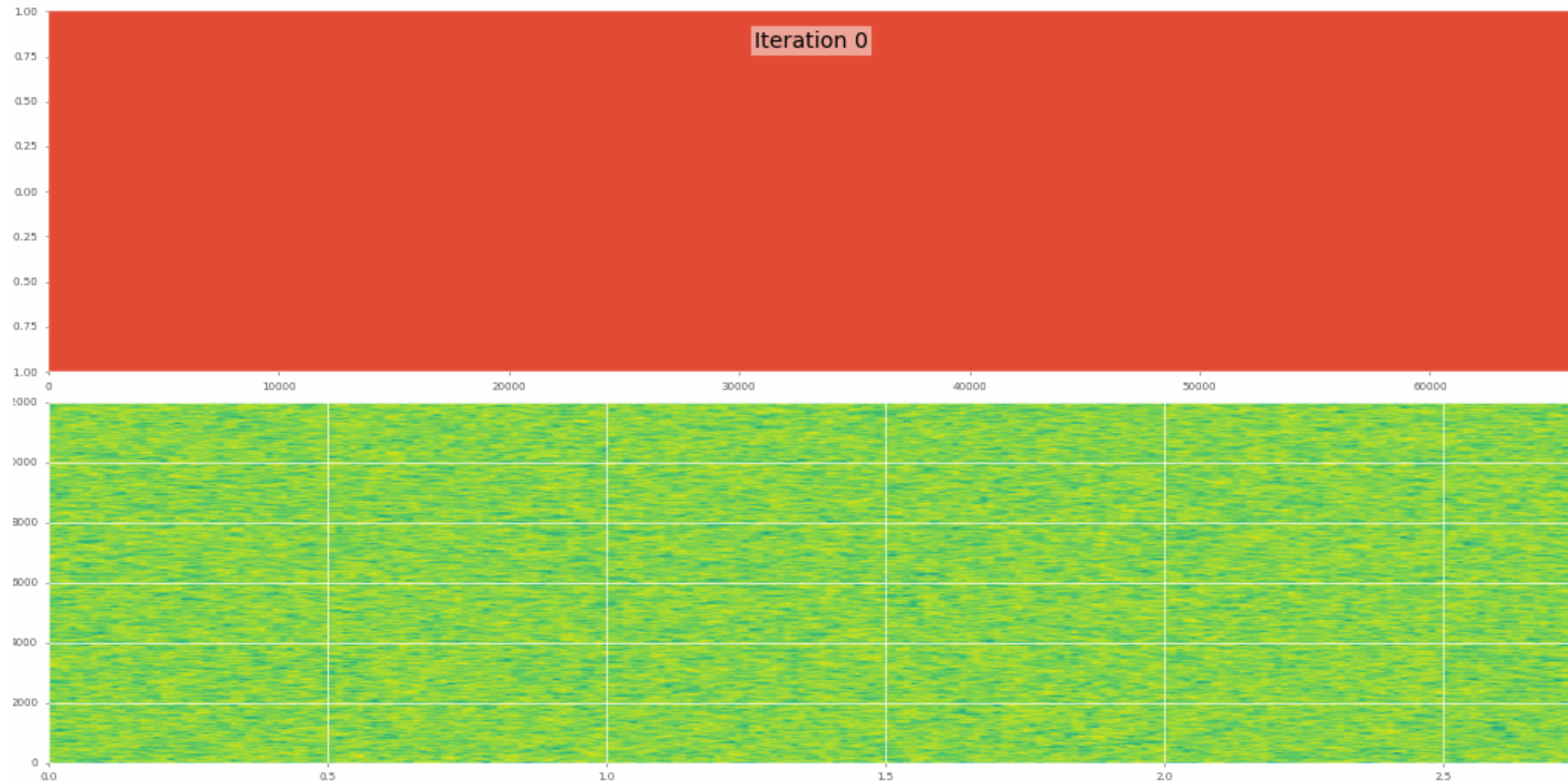


(Recap) Diffusion Models – Generation



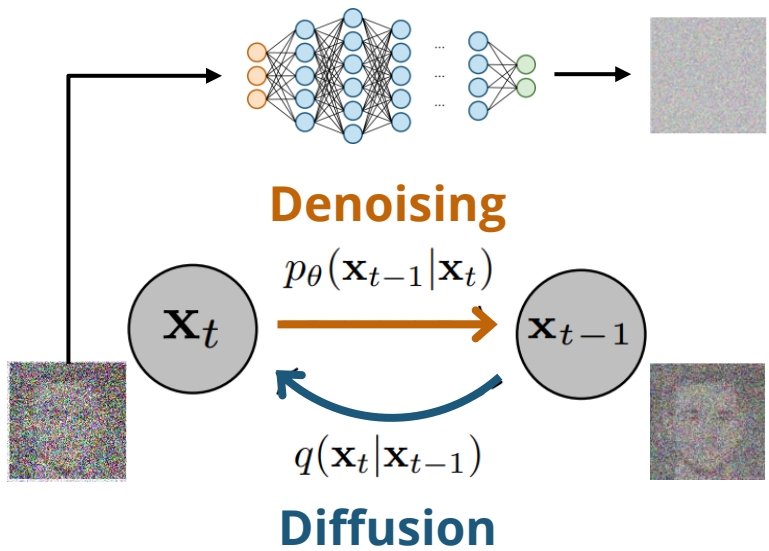
(Source: Ho et al., 2020)

Example: WaveGrad (Chen et al., 2021)

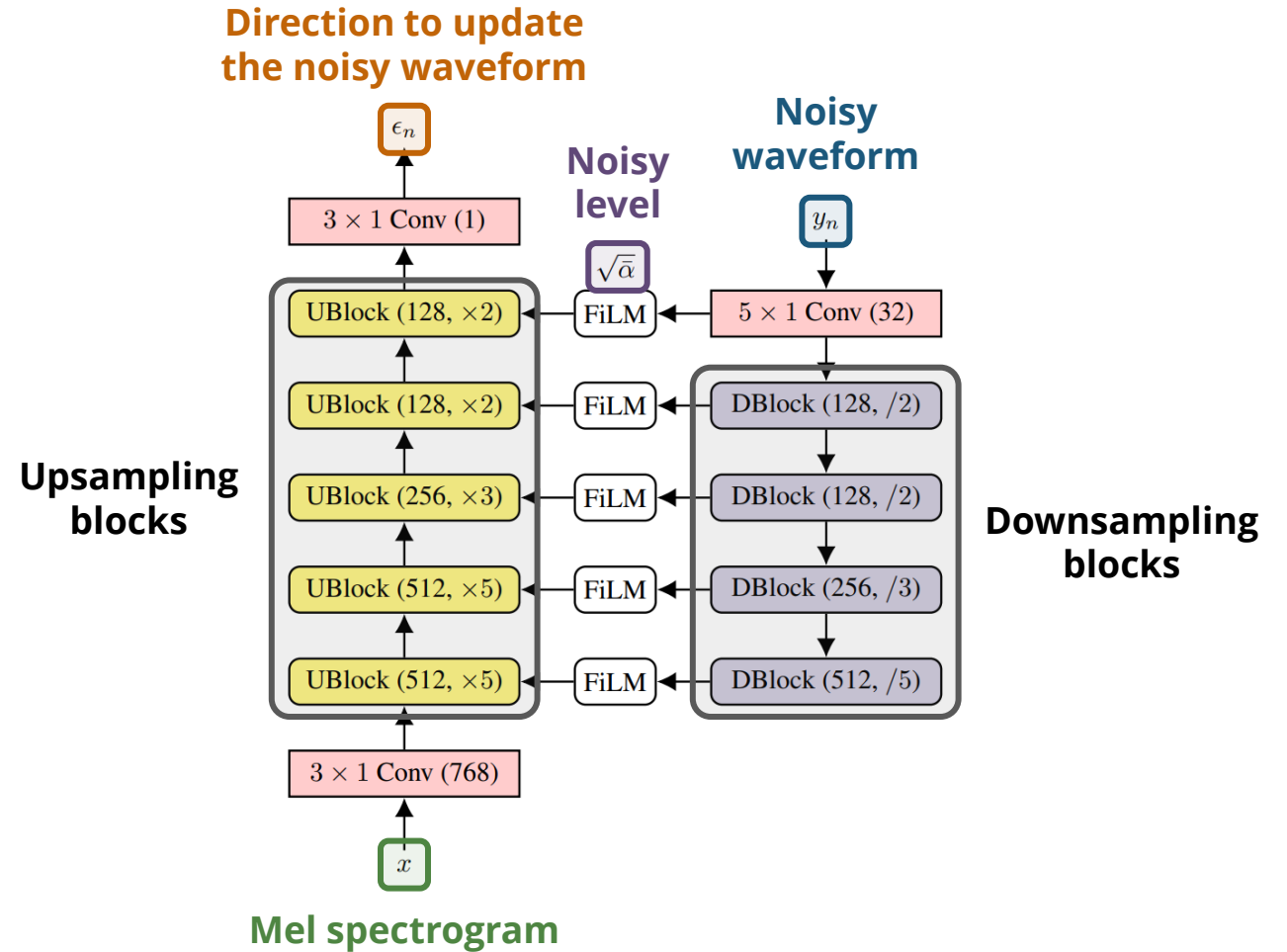


(Source: Chen et al., 2021)

Example: WaveGrad (Chen et al., 2021)

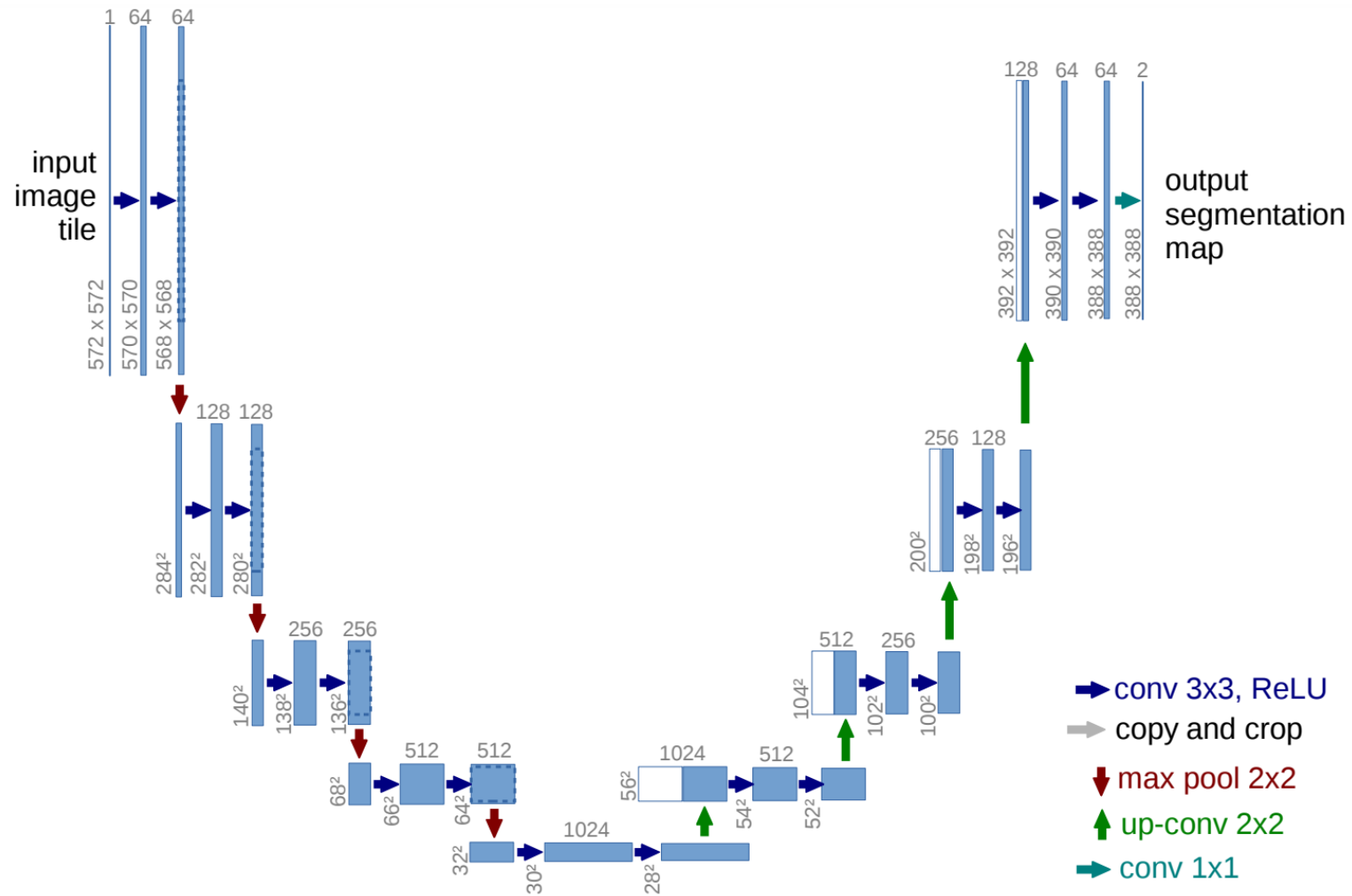


(Source: Ho et al., 2020)



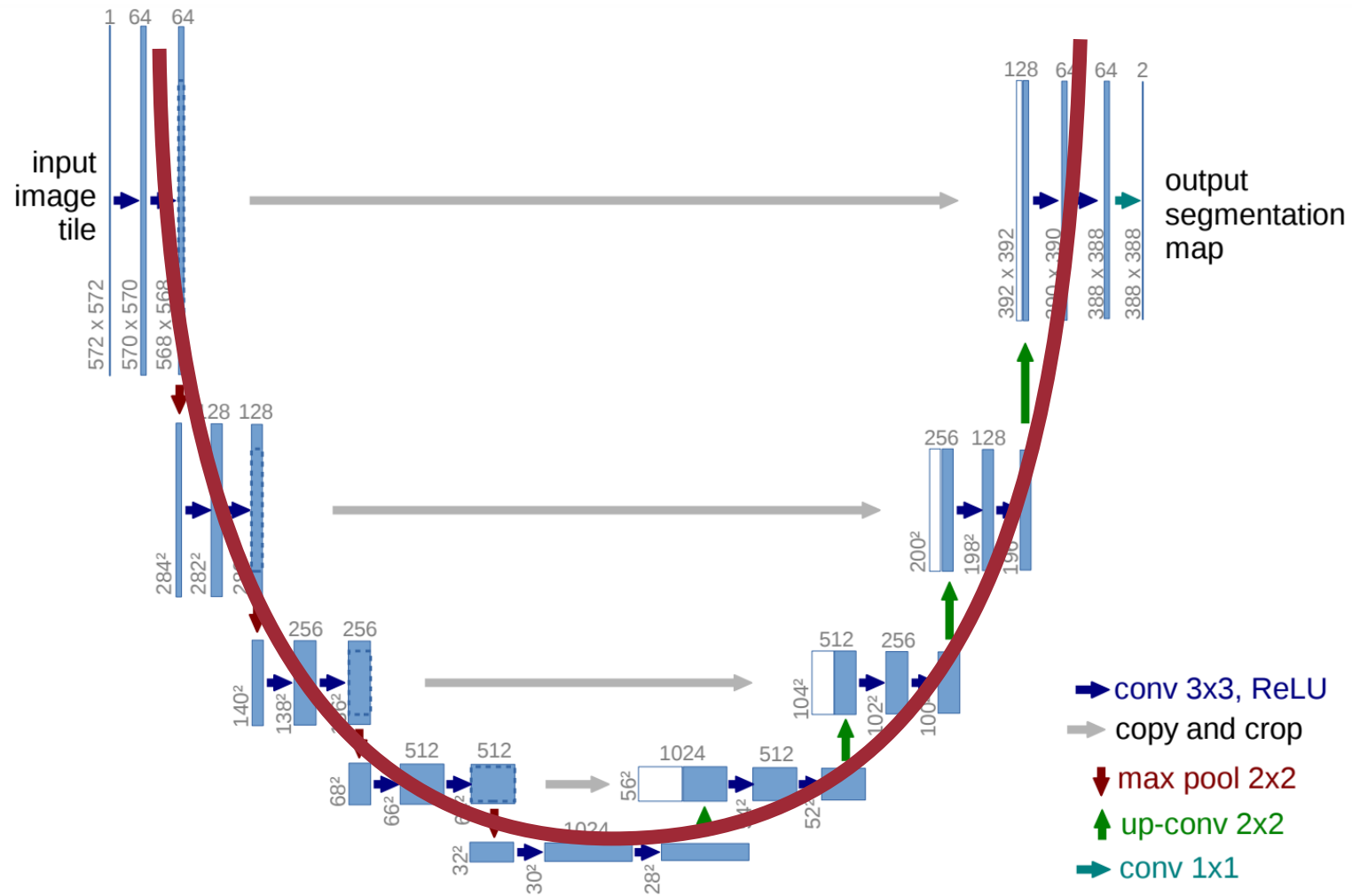
(Source: Chen et al., 2021)

U-Net (Ronneberger et al., 2015)



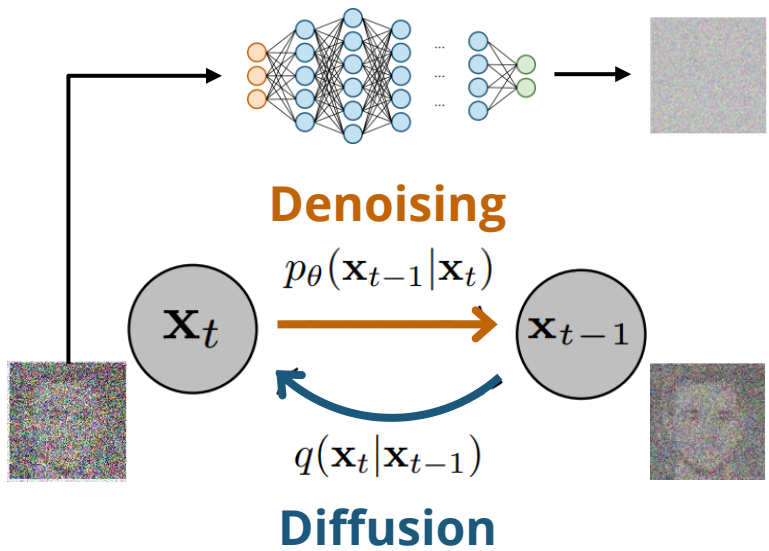
(Source: Ronneberger et al., 2015)

U-Net (Ronneberger et al., 2015)

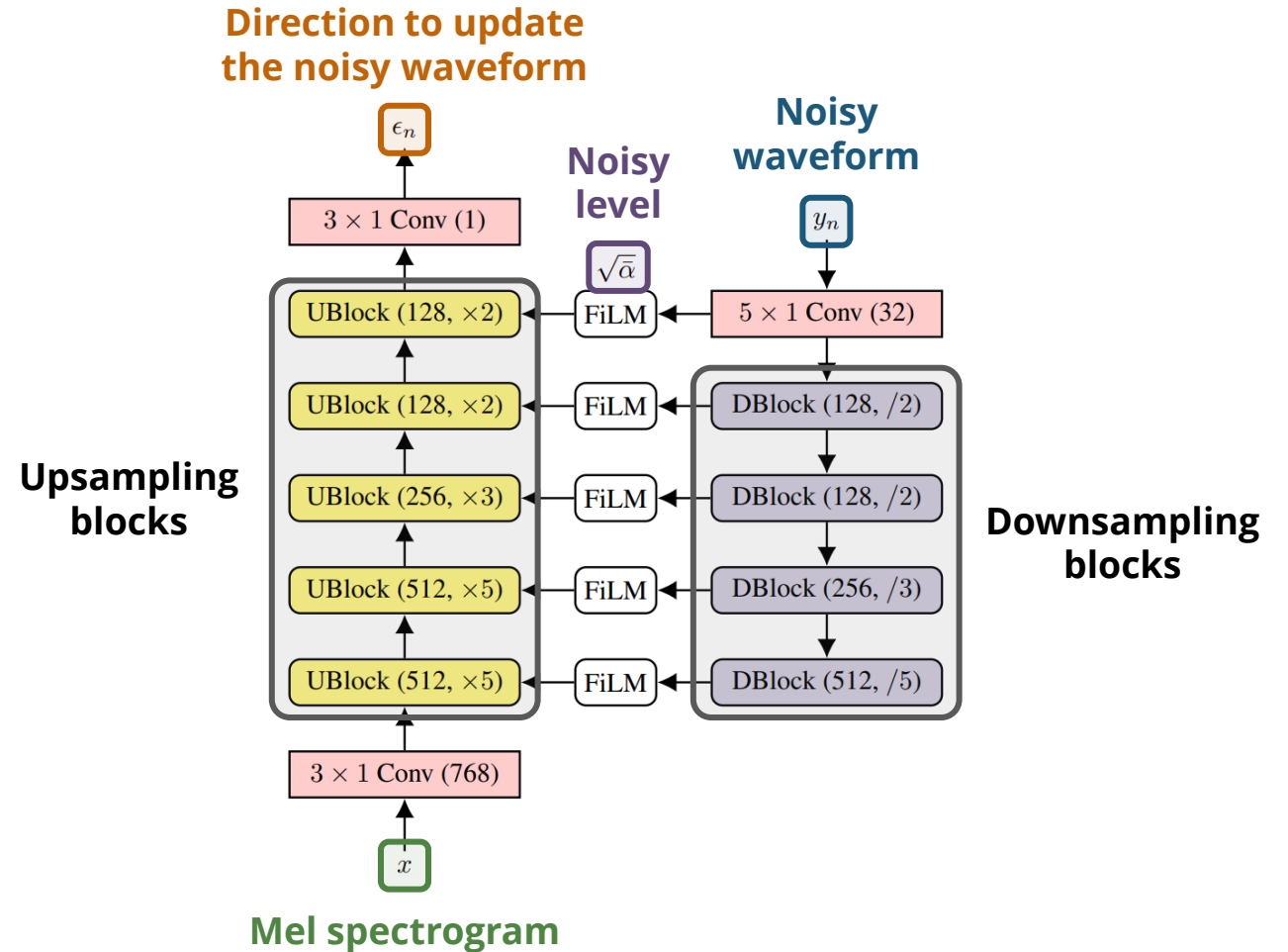


(Source: Ronneberger et al., 2015)

Example: WaveGrad (Chen et al., 2021)

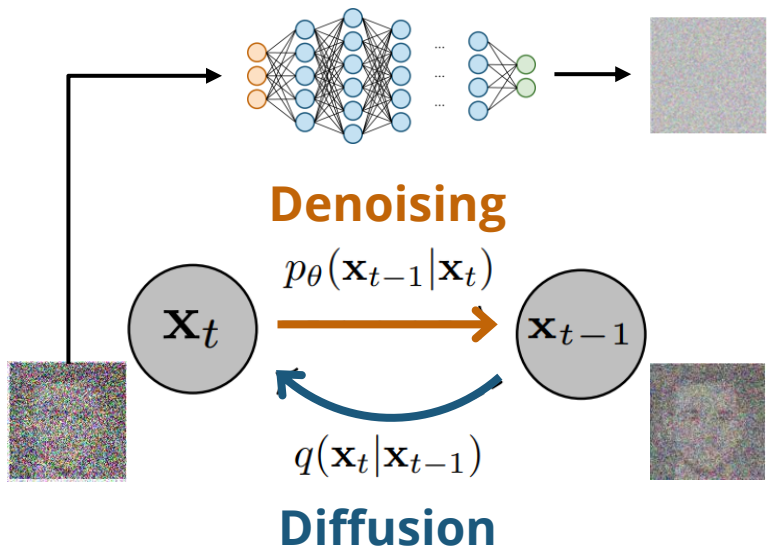


(Source: Ho et al., 2020)

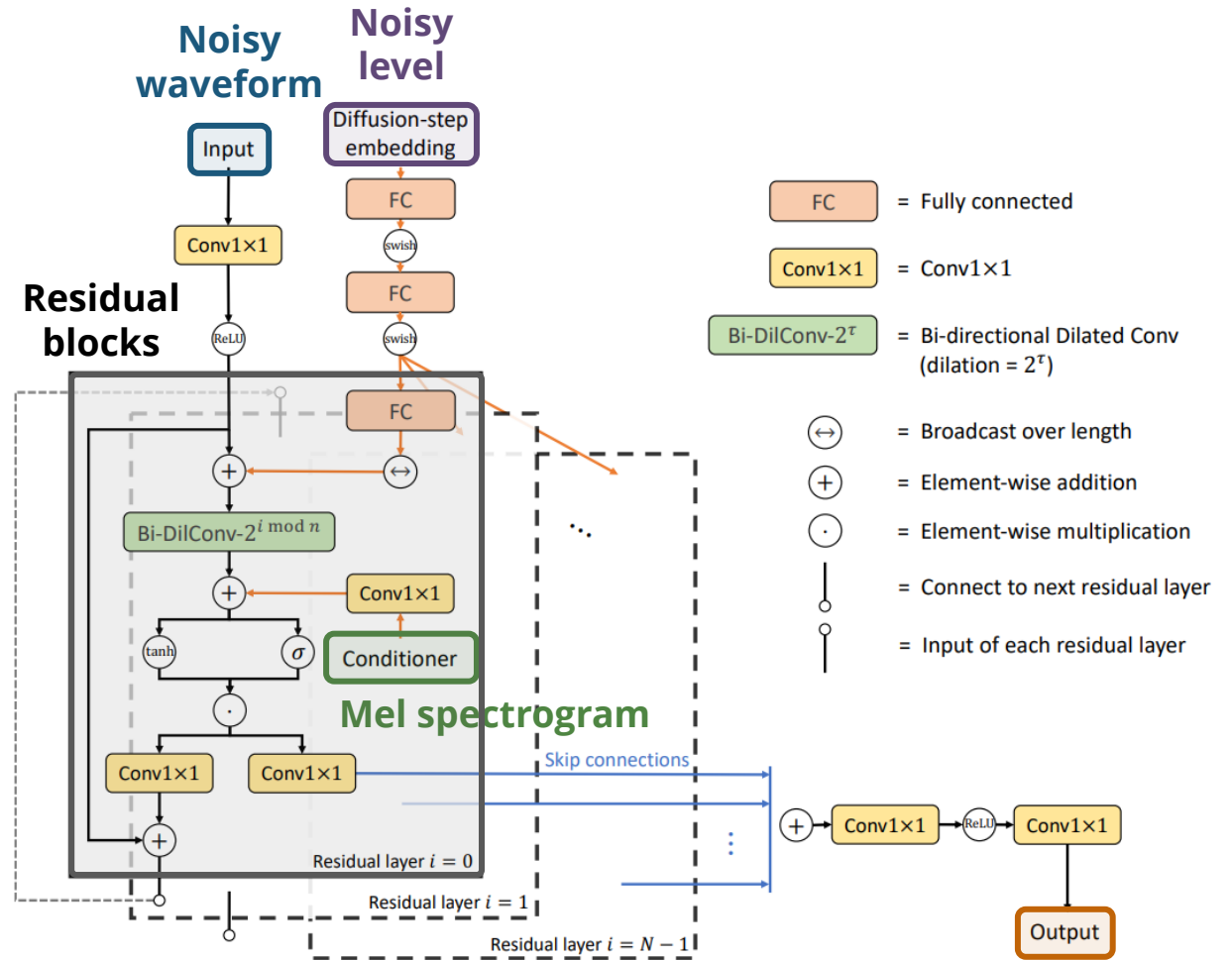


(Source: Chen et al., 2021)

Example: DiffWave (Kong et al., 2021)

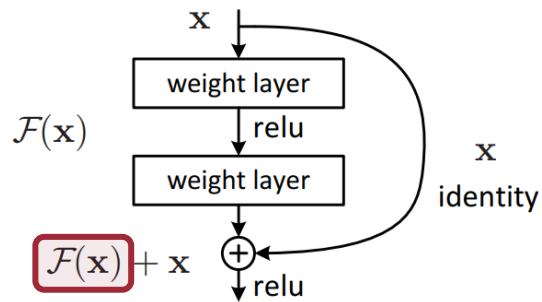


(Source: Ho et al., 2020)



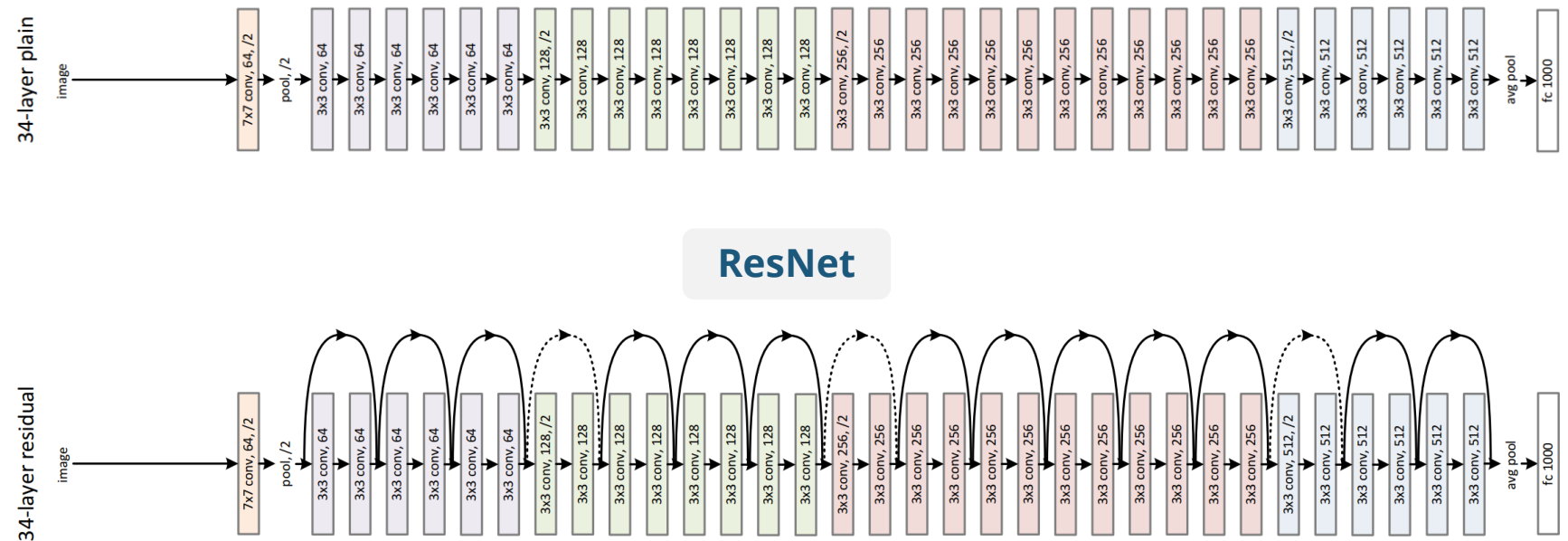
(Source: Kong et al., 2021)

Deep Residual Nets (ResNets) (He et al., 2016)



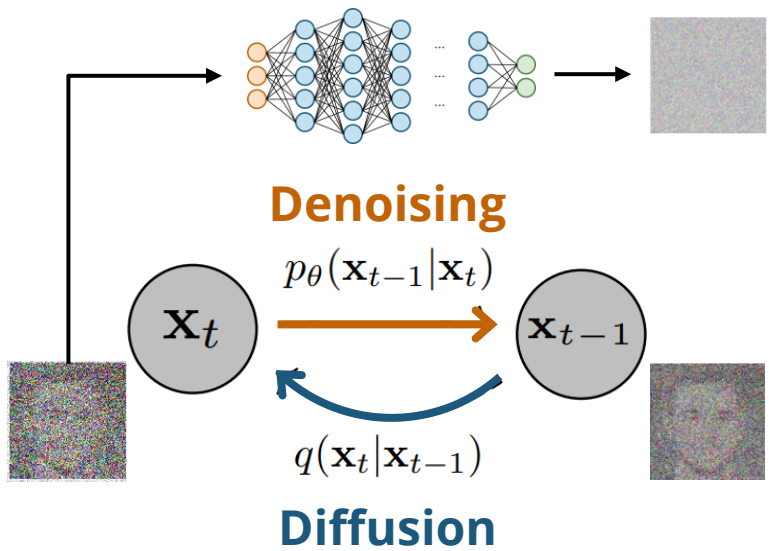
Learn the residuals (changes)

Without skip connections

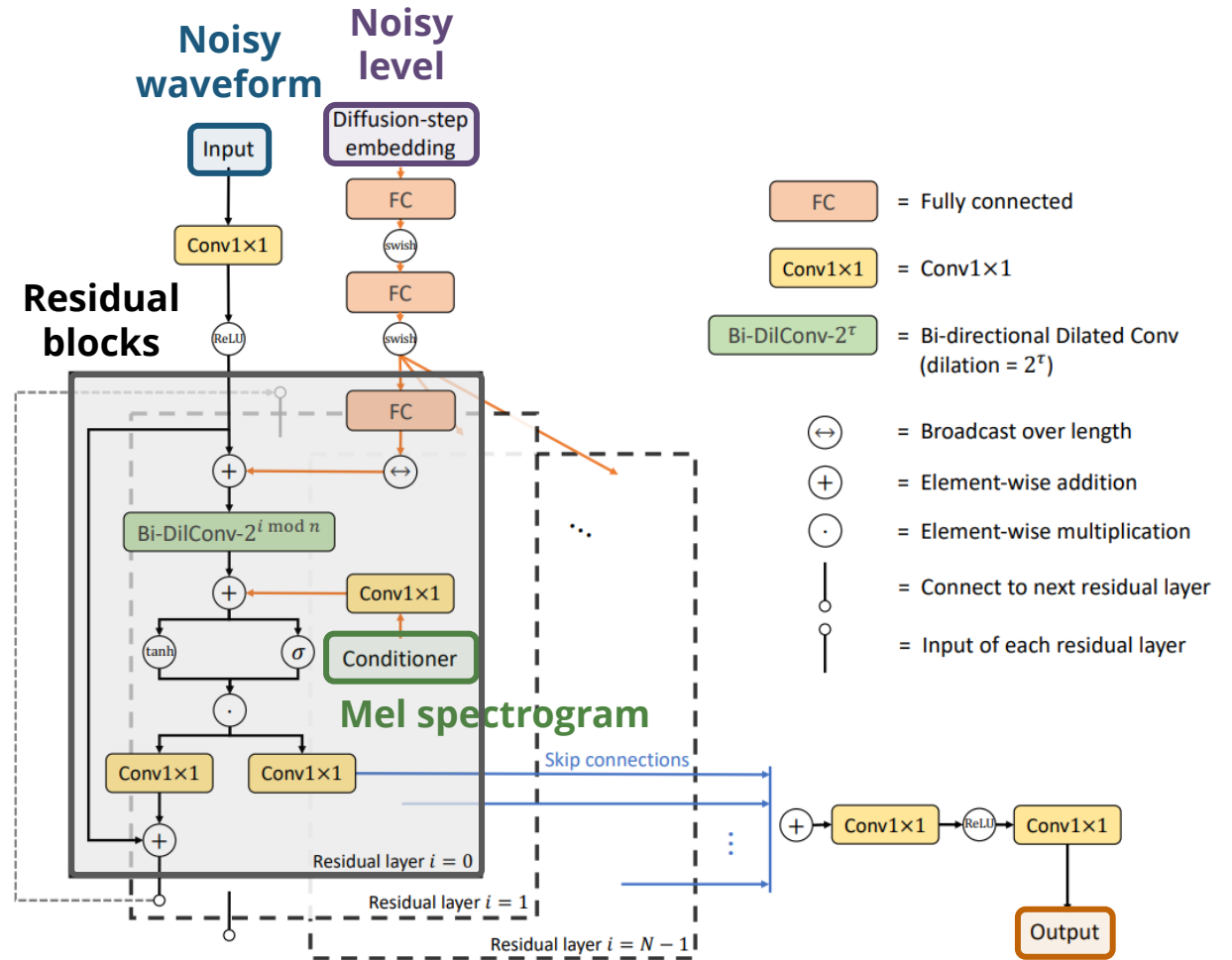


(Source: He et al., 2016)

Example: DiffWave (Kong et al., 2021)



(Source: Ho et al., 2020)

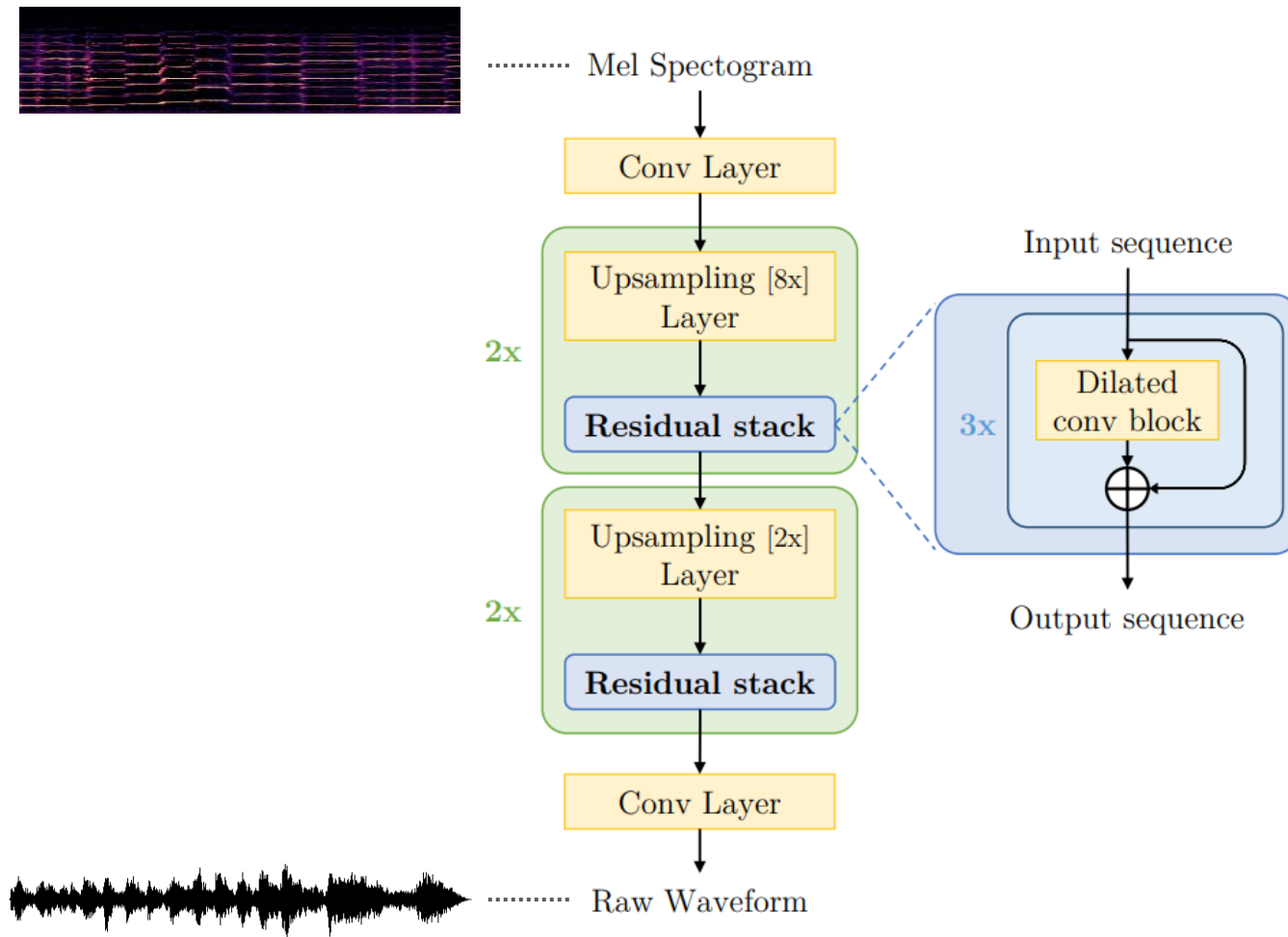


- FC = Fully connected
- Conv1x1 = Conv1x1
- Bi-DilConv-2 ^{τ} = Bi-directional Dilated Conv (dilation = 2 ^{τ})
- \leftrightarrow = Broadcast over length
- $+$ = Element-wise addition
- \cdot = Element-wise multiplication
- \circ = Connect to next residual layer
- \circ = Input of each residual layer

(Source: Kong et al., 2021)

Conditional Audio Synthesis

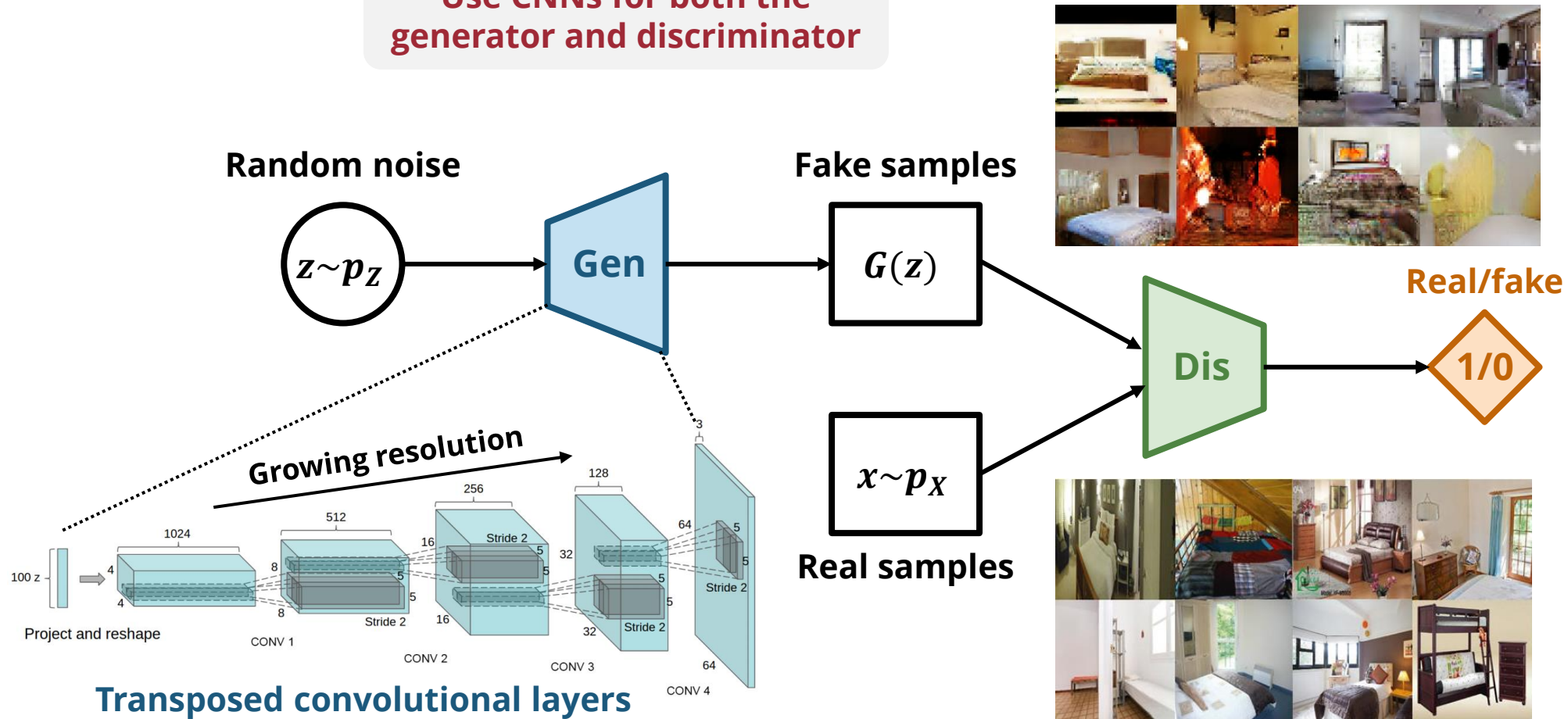
Example: MelGAN (Kumar et al., 2019)



(Source: Kumar et al., 2019)

(Recap) Deep Convolutional GANs (DCGANs)

Use CNNs for both the generator and discriminator



(Recap) Transposed Convolution

Convolution

1	-1	-1	-1
-1	1	-1	-1
-1	-1	1	-1
-1	-1	-1	1

*

1	-1	-1
-1	1	-1
-1	-1	1

=

9	-1
-1	9

Transposed convolution

1	-1
-1	1

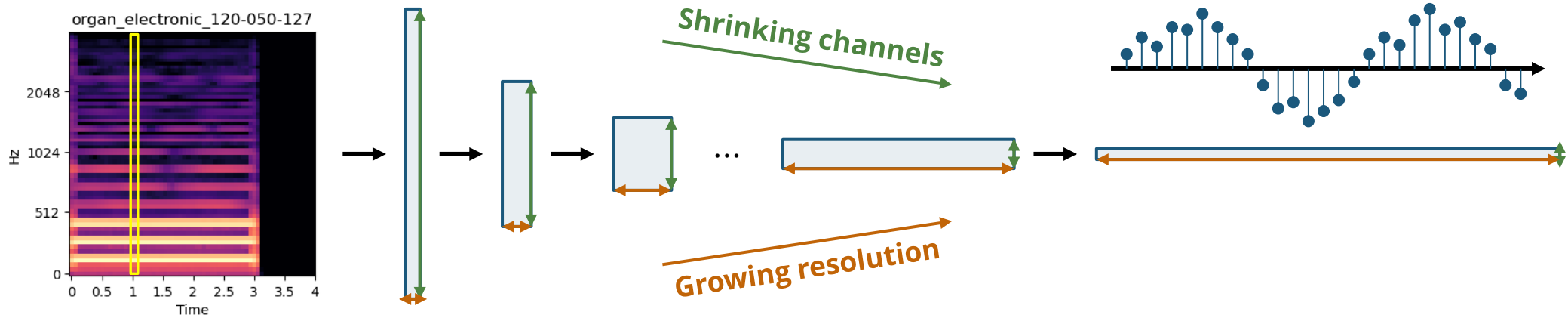
*

1	-1	-1
-1	1	-1
-1	-1	1

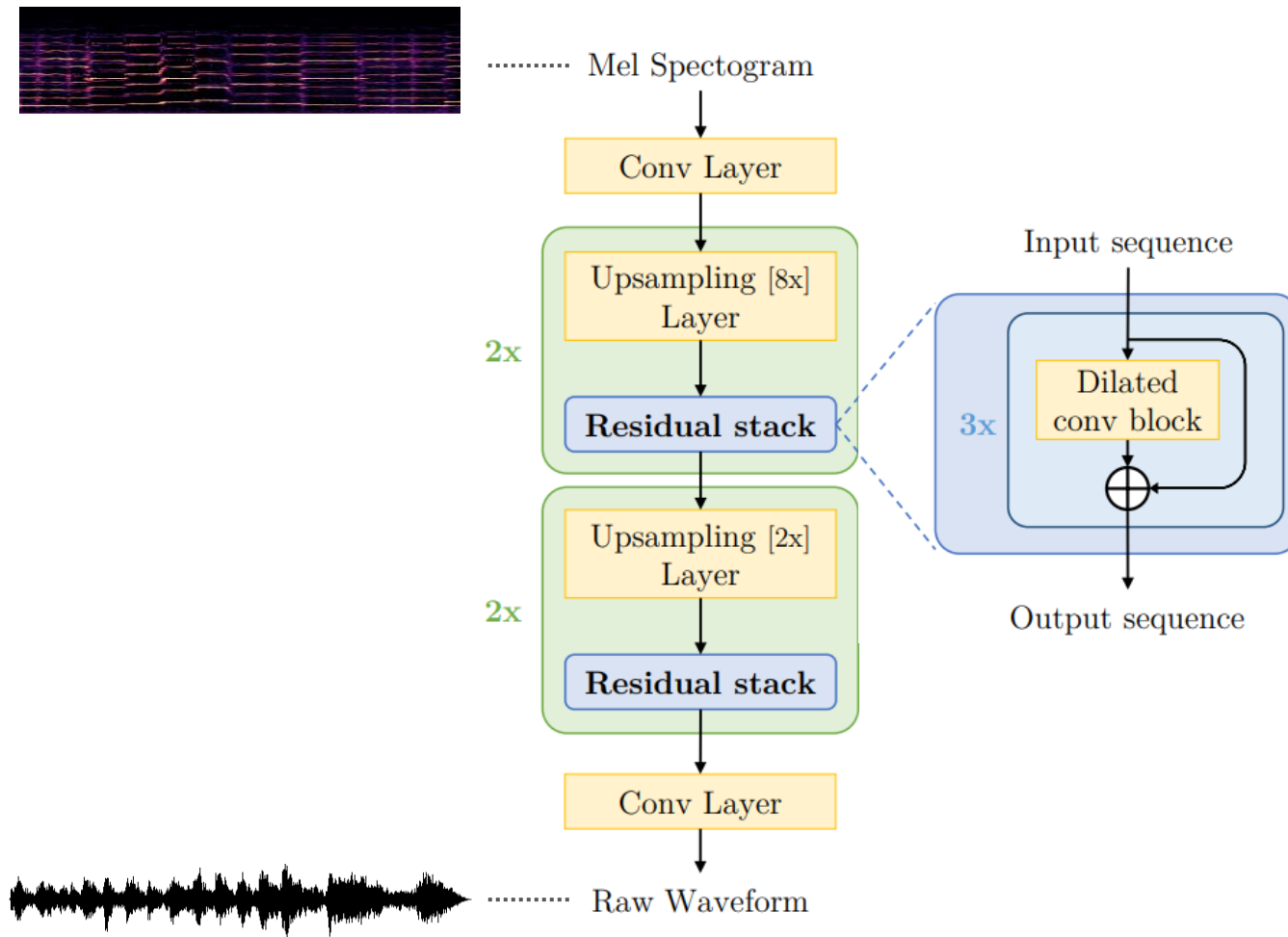
=

1	0	0	1
0	4	-2	0
0	-2	4	0
1	0	0	1

Transposed Convolution for Vocoders



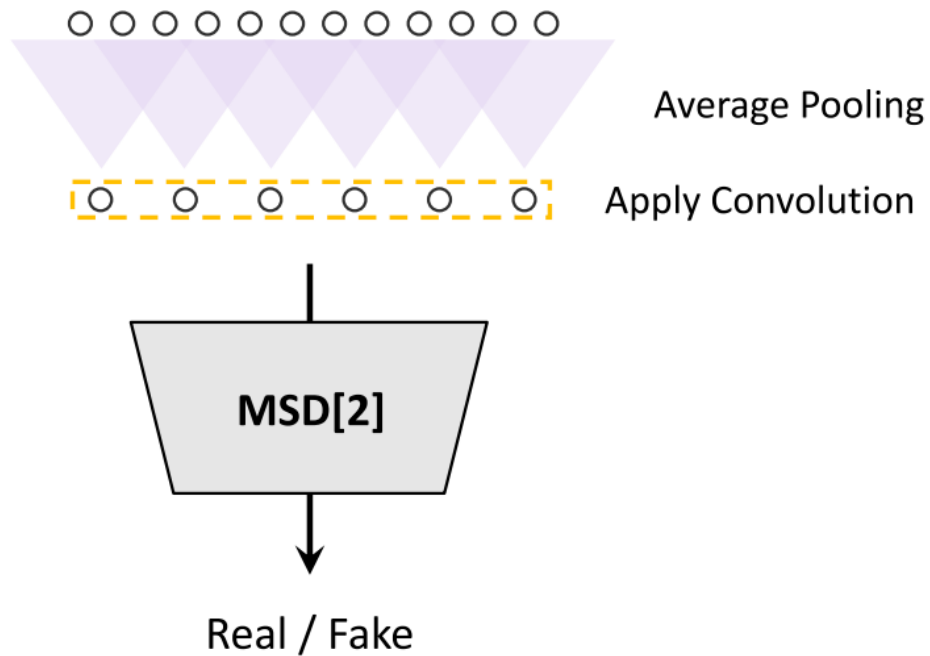
Example: MelGAN (Kumar et al., 2019)



(Source: Kumar et al., 2019)

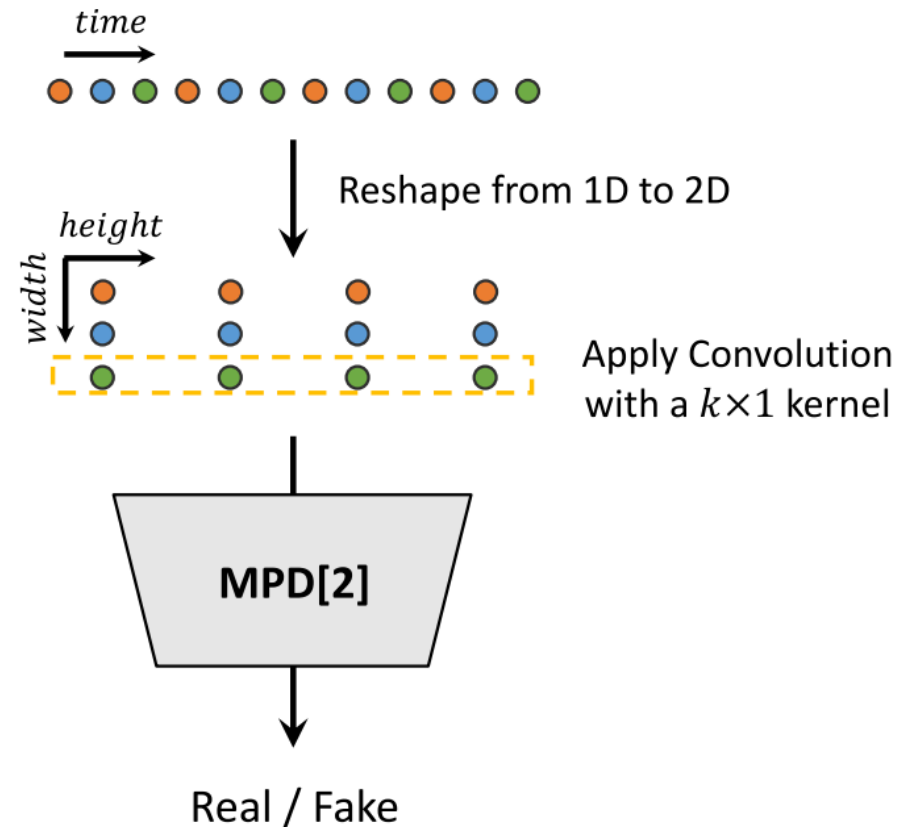
Example: MelGAN (Kumar et al., 2019)

Multi-scale discriminator



(Source: Kong et al., 2019)

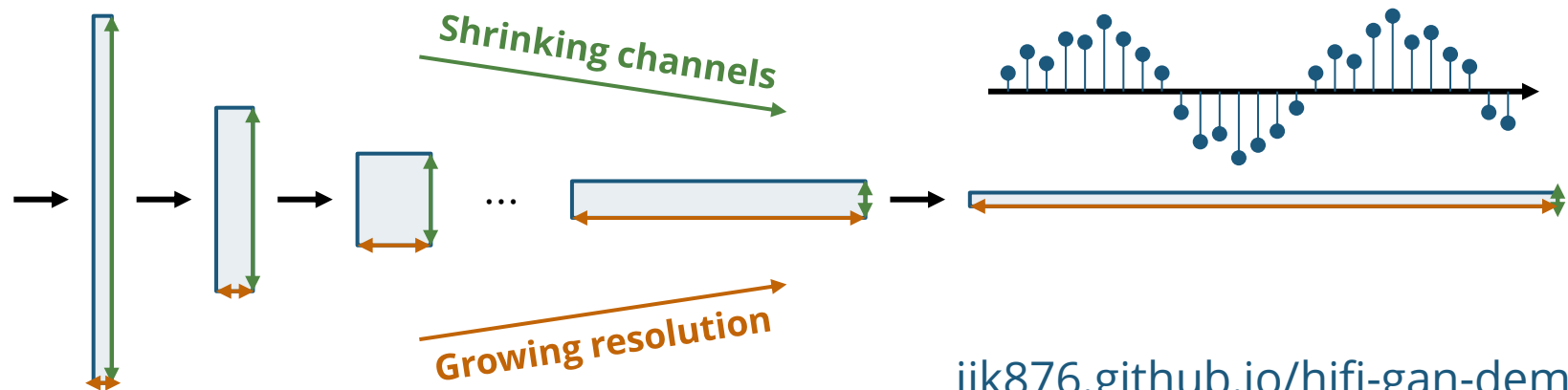
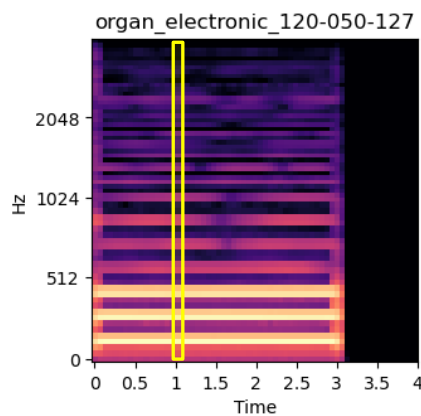
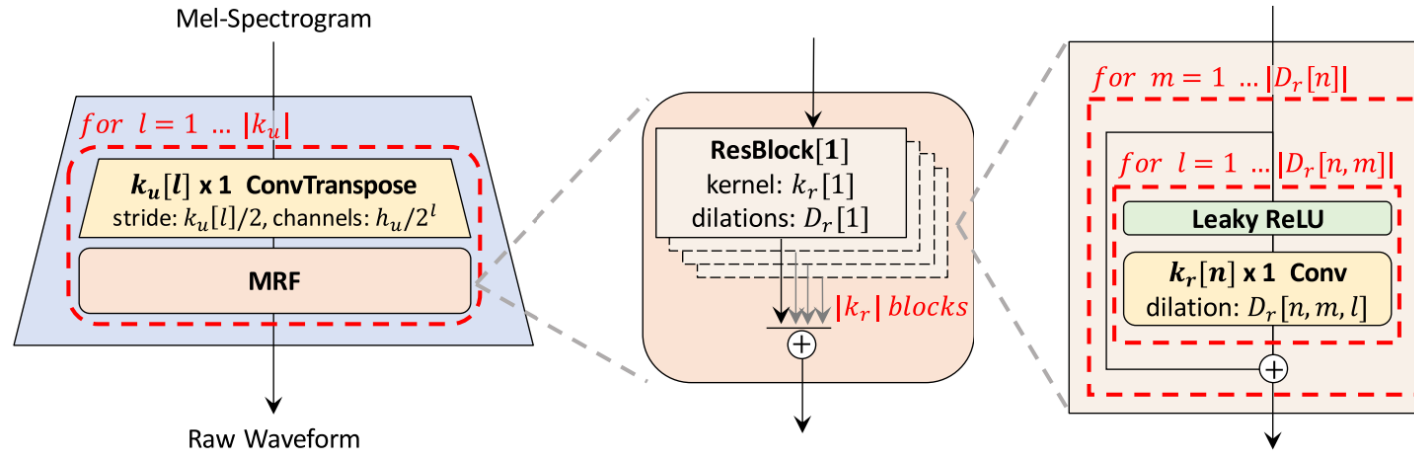
Multi-period discriminator



Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre de Brebisson, Yoshua Bengio, and Aaron Courville, "[MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis](#)," *NeurIPS*, 2019.

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

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



jik876.github.io/hifi-gan-demo

Optional Reading

- A very nice blog on “**Generating music in the waveform domain**” by Sander Dieleman: sander.ai/2020/03/24/audio-generation

(Recap) Four Paradigms



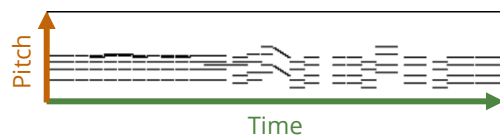
Symbolic music generation

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



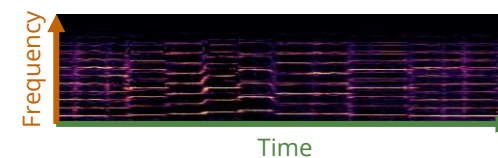
Audio-domain music generation

Time series-based

Image-based



Waveform



Spectrogram

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