

PAT 464/564 (Winter 2026)

Generative AI for Music & Audio Creation

Lecture 14: Diffusion Models

Instructor: Hao-Wen Dong

Representative Types of Deep Generative Models

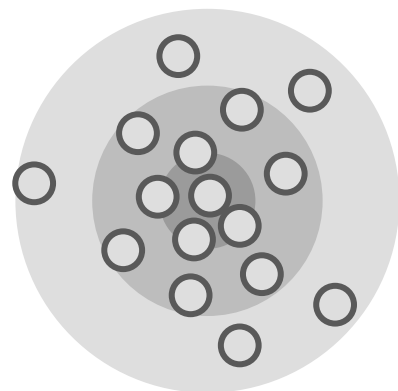
- **Deep autoregressive models**
 - Recurrent neural network (RNN)
 - Long short-term memory (LSTM)
 - Transformer model
- **Deep latent variable models**
 - Variational autoencoder (VAE)
 - Generative adversarial network (GAN)
 - Diffusion model **Today's topic!**
 - Flow-based model
- *And many others...*

Deep Latent Variable Models

Deep Latent Variable Models

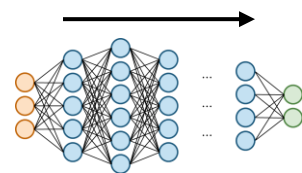
- **Intuition:** Learn to map a known distribution to the data distribution

Known distribution

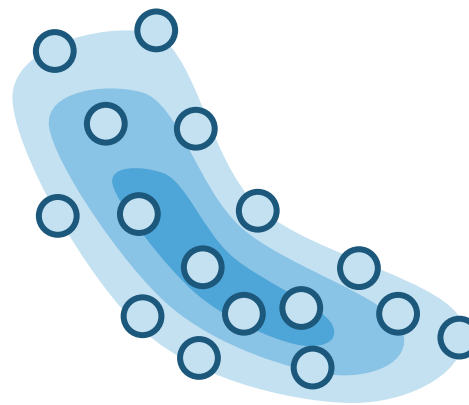


$P(z)$

$P(x | z)$



Data distribution

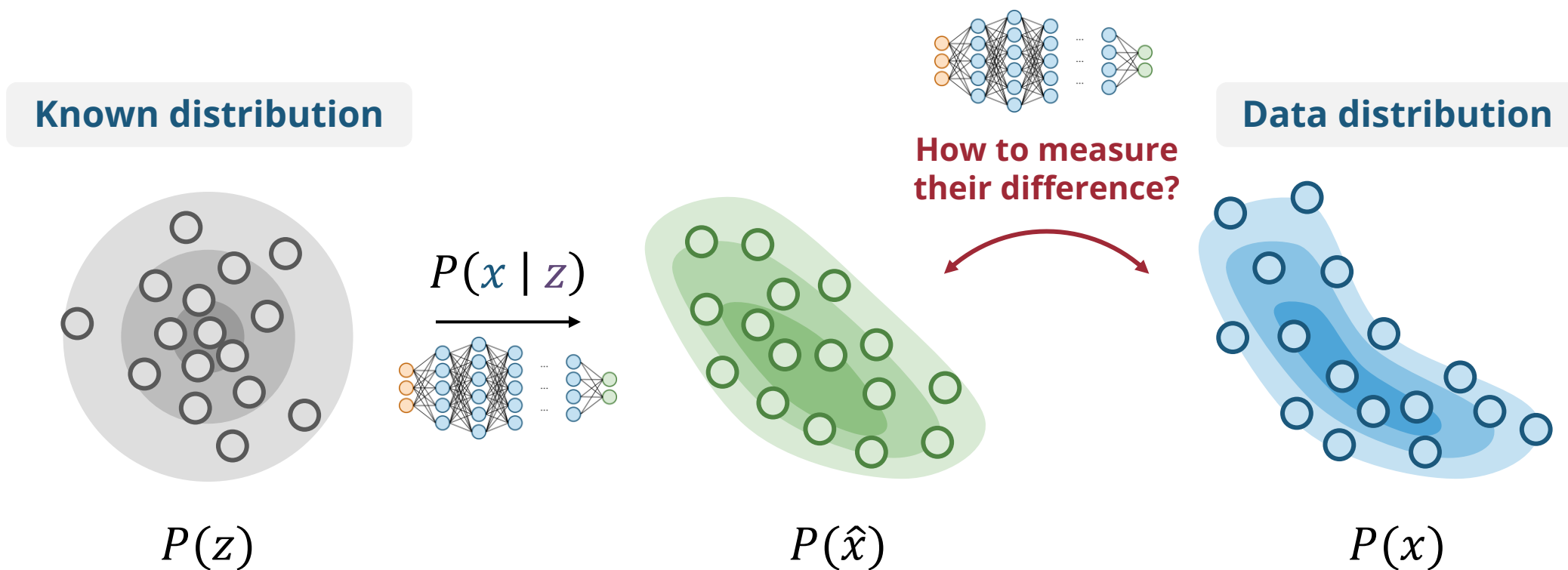


$P(x)$

$$P(x) = P(z) P(x | z)$$

Deep Latent Variable Models

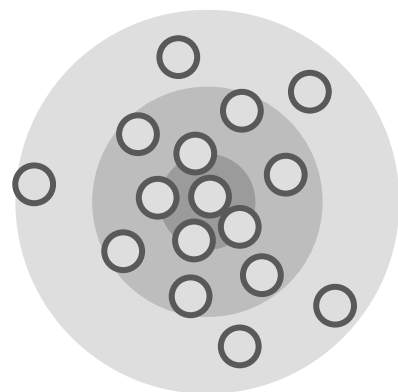
- **Intuition:** Learn to map a known distribution to the data distribution



Deep Latent Variable Models

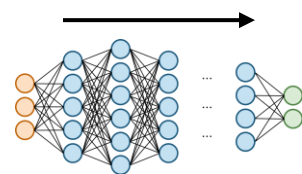
- **Intuition:** Learn to map a known distribution to the data distribution

Known distribution

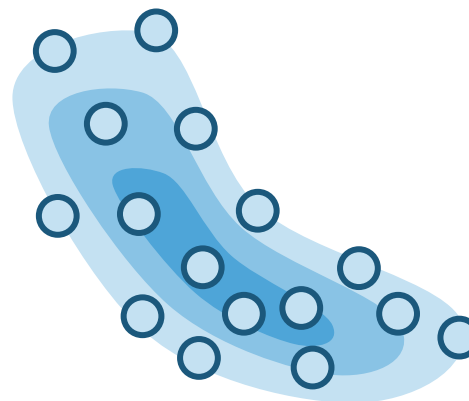


$P(z)$

$P(x | z)$



Data distribution



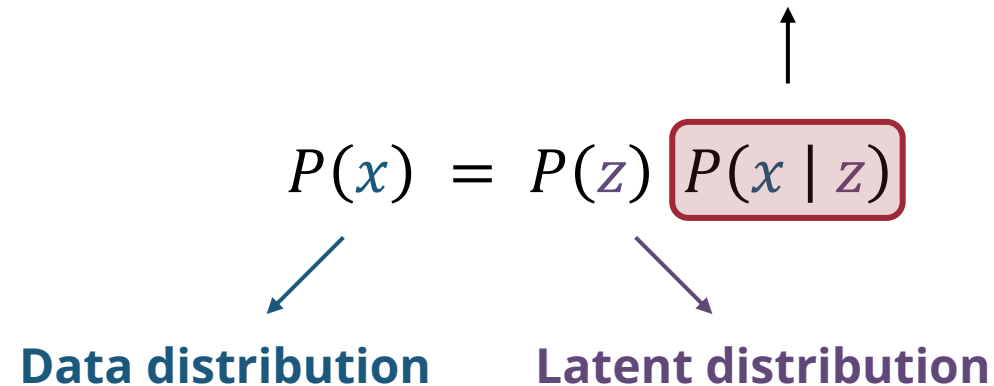
$P(x)$

$$P(x) = P(z) P(x | z)$$

Deep Latent Variable Models

- **Intuition:** Learn to map a known distribution to the data distribution

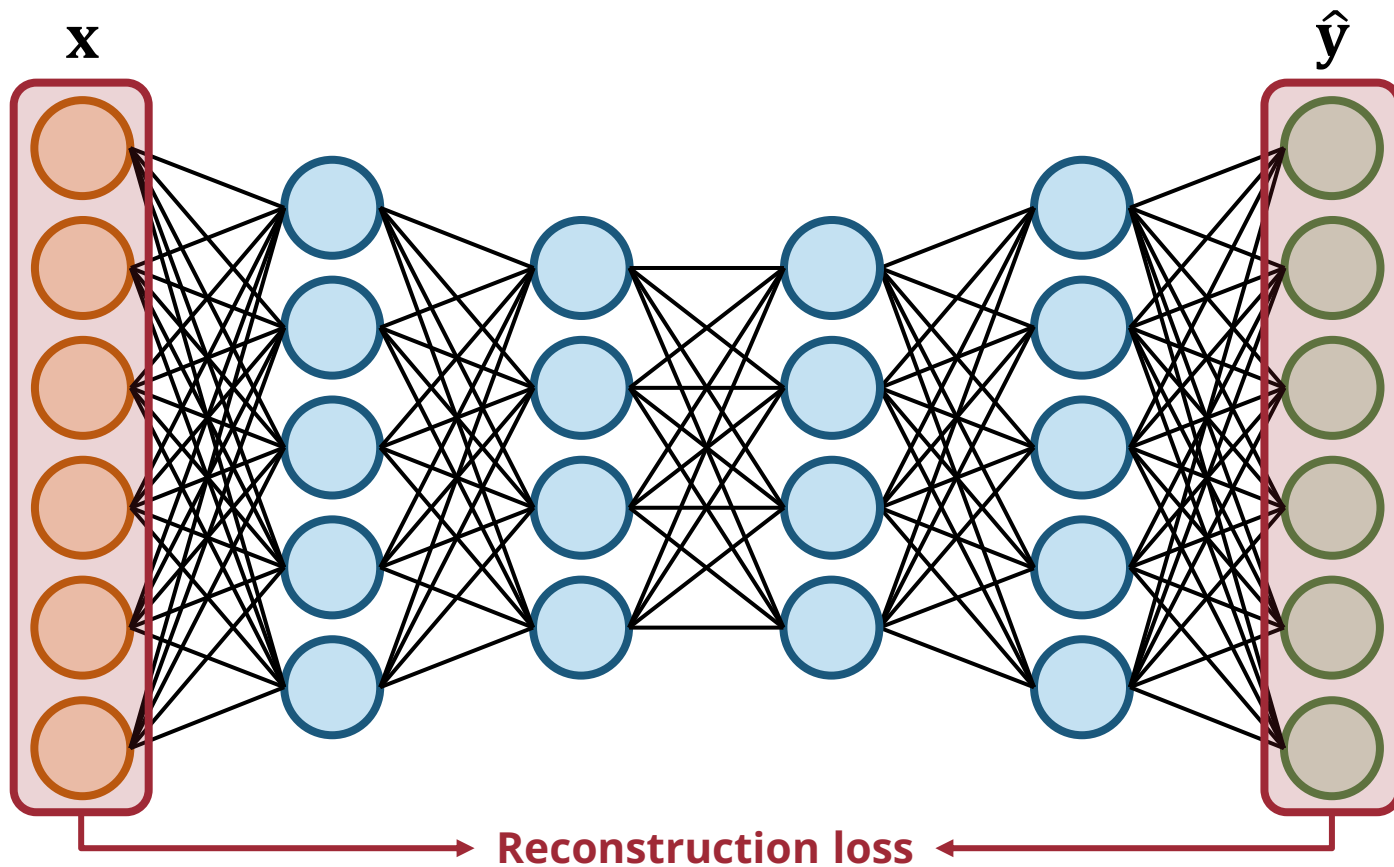
What we want the model to learn!



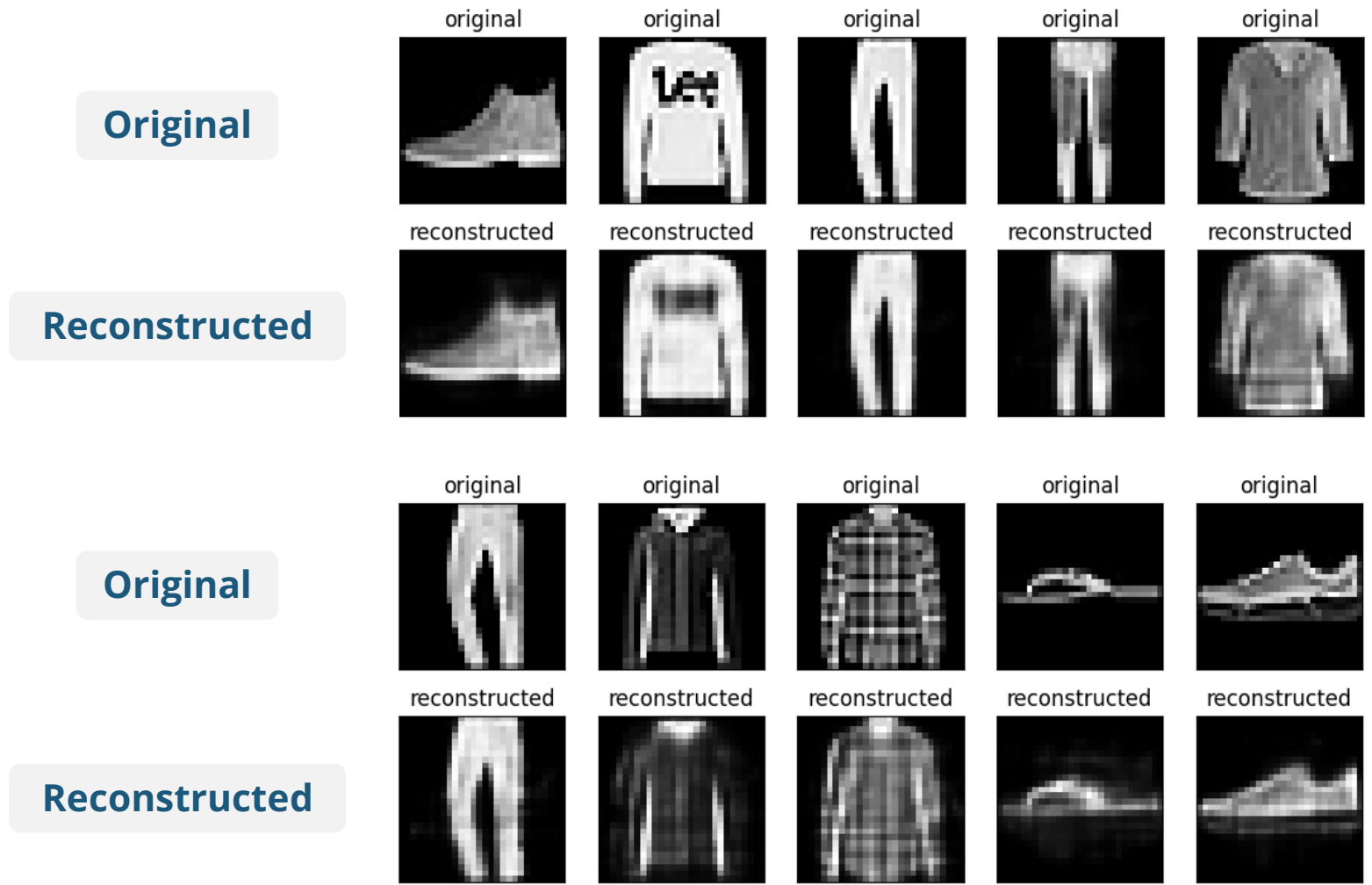
Diffusion Models

Autoencoders

- A neural network where the **input and output are the same**

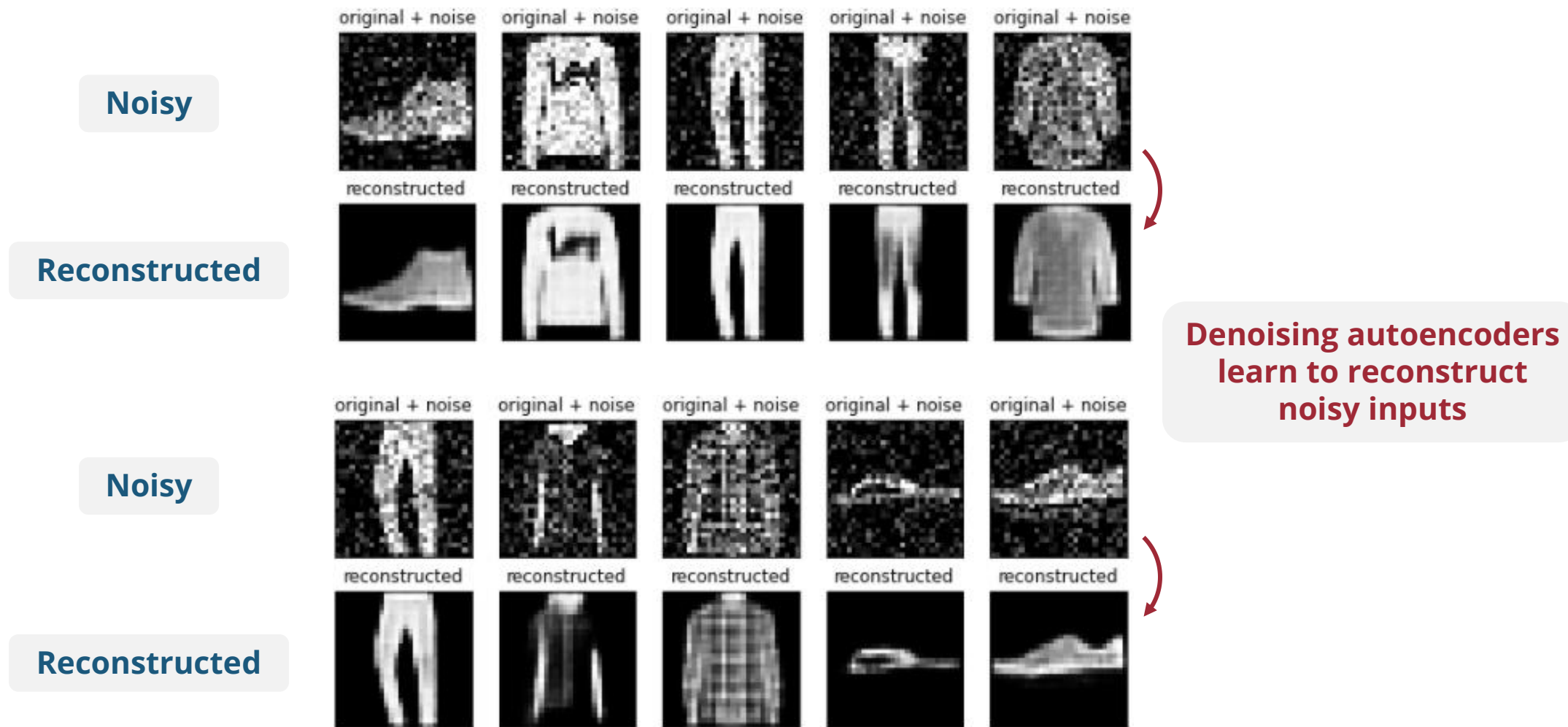


Autoencoders: Reconstruction Examples



(Source: tensorflow.org)

Denoising Autoencoders



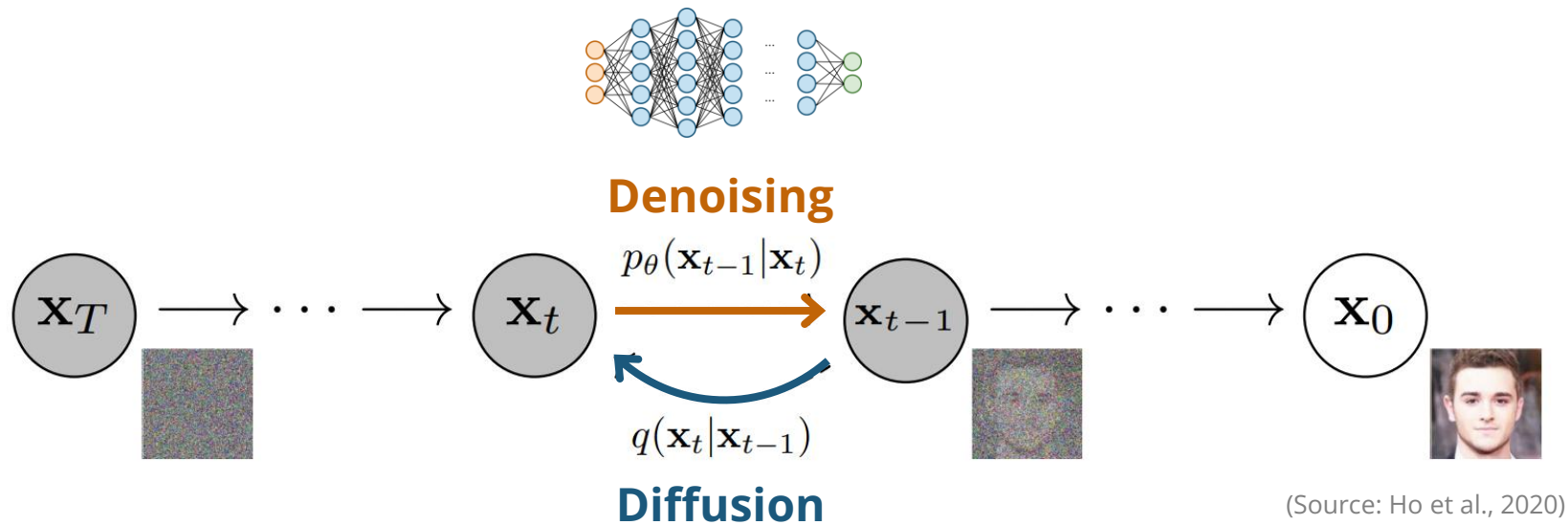
(Source: tensorflow.org)

tensorflow.org/tutorials/generative/autoencoder

Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol, "Extracting and Composing Robust Features with Denoising Autoencoders," *ICML*, 2008.

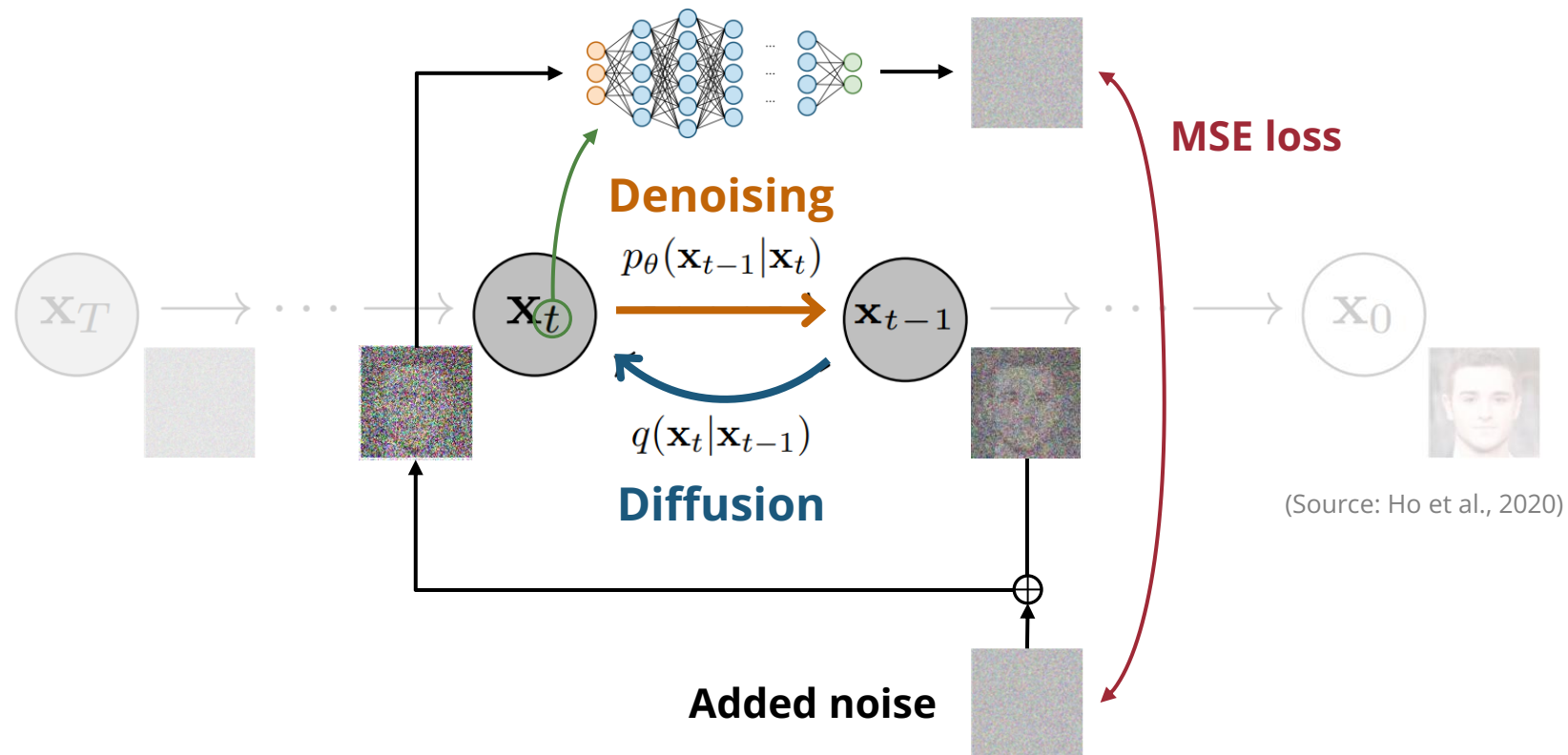
Diffusion Models (Ho et al., 2020)

- Intuition:** Many denoising autoencoders stacked together



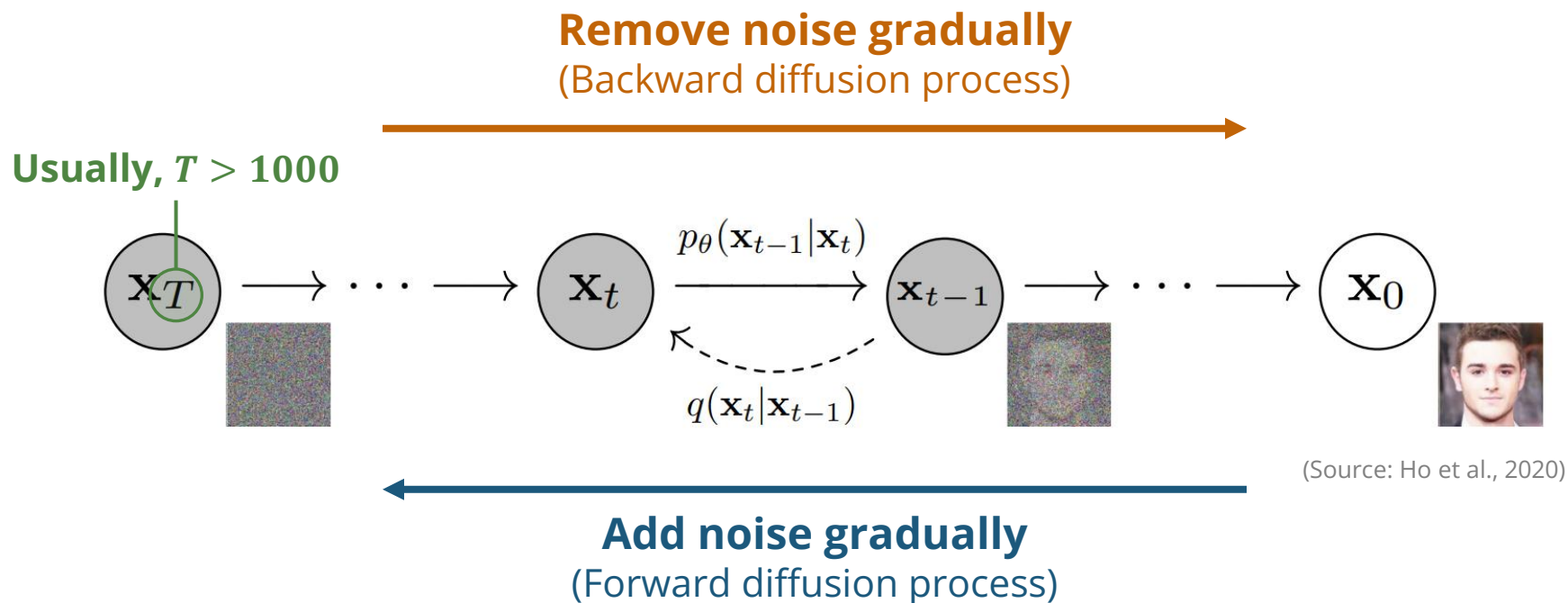
Diffusion Models: Training

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

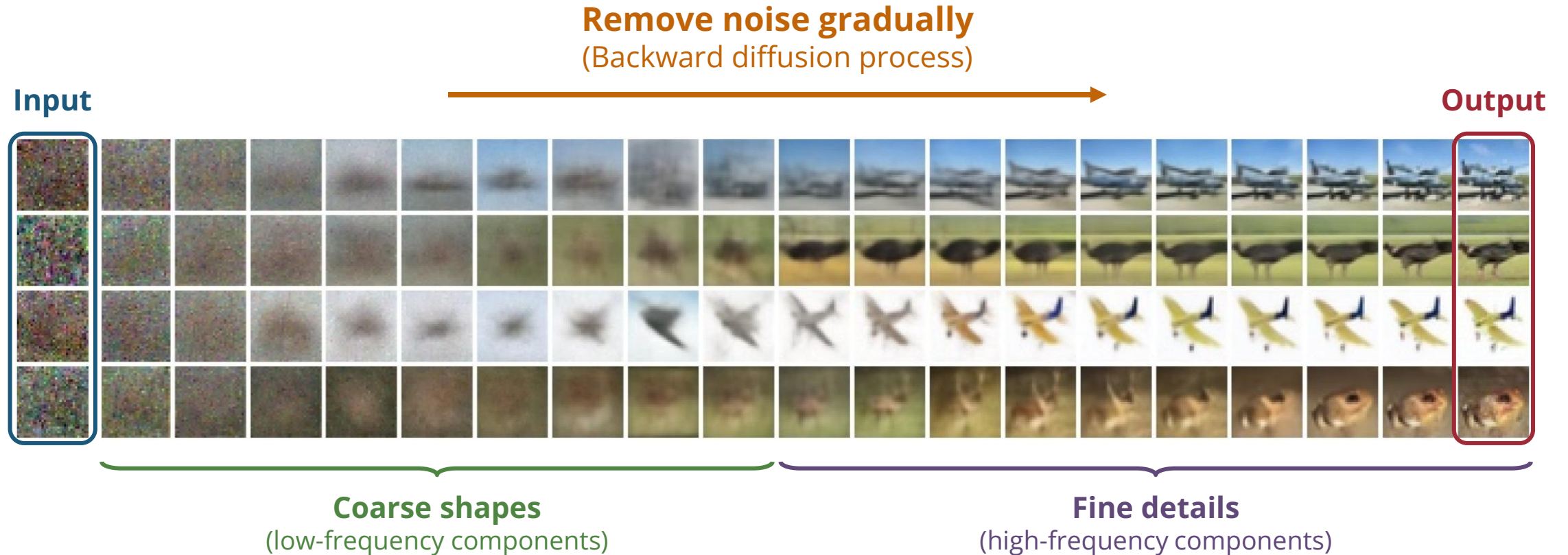


Diffusion Models (Ho et al., 2020)

- Intuition:** Many denoising autoencoders stacked together



Diffusion Models: Generation

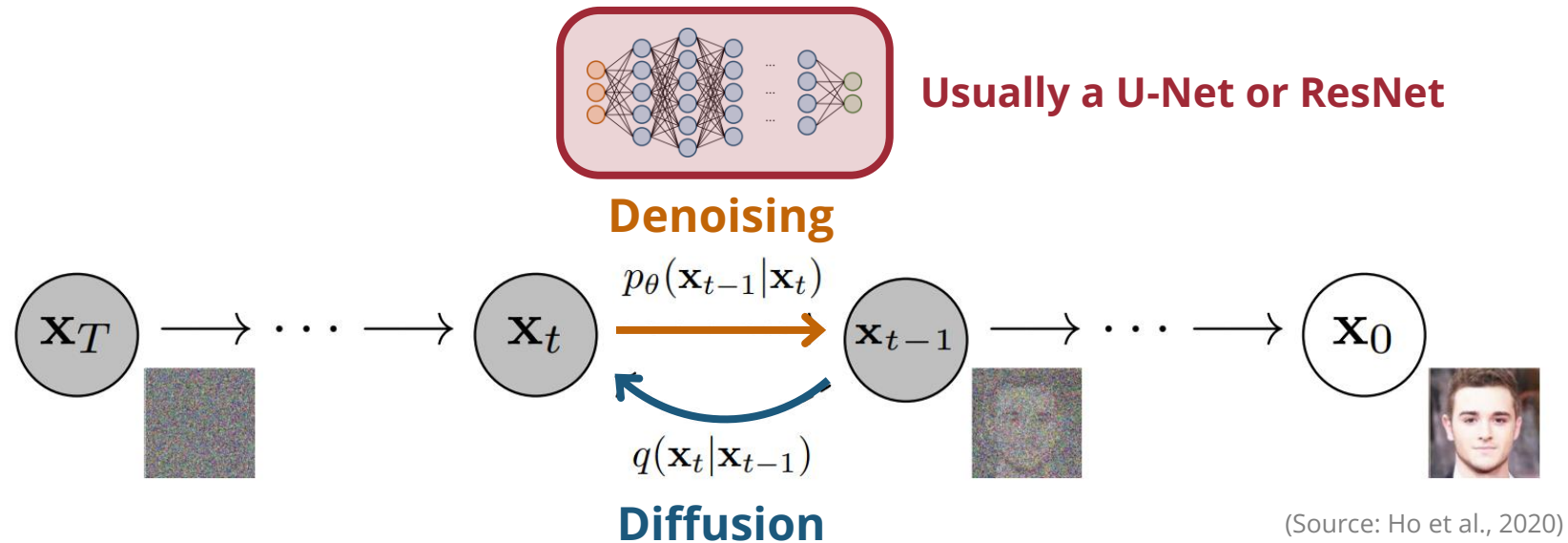


(Source: Ho et al., 2020)

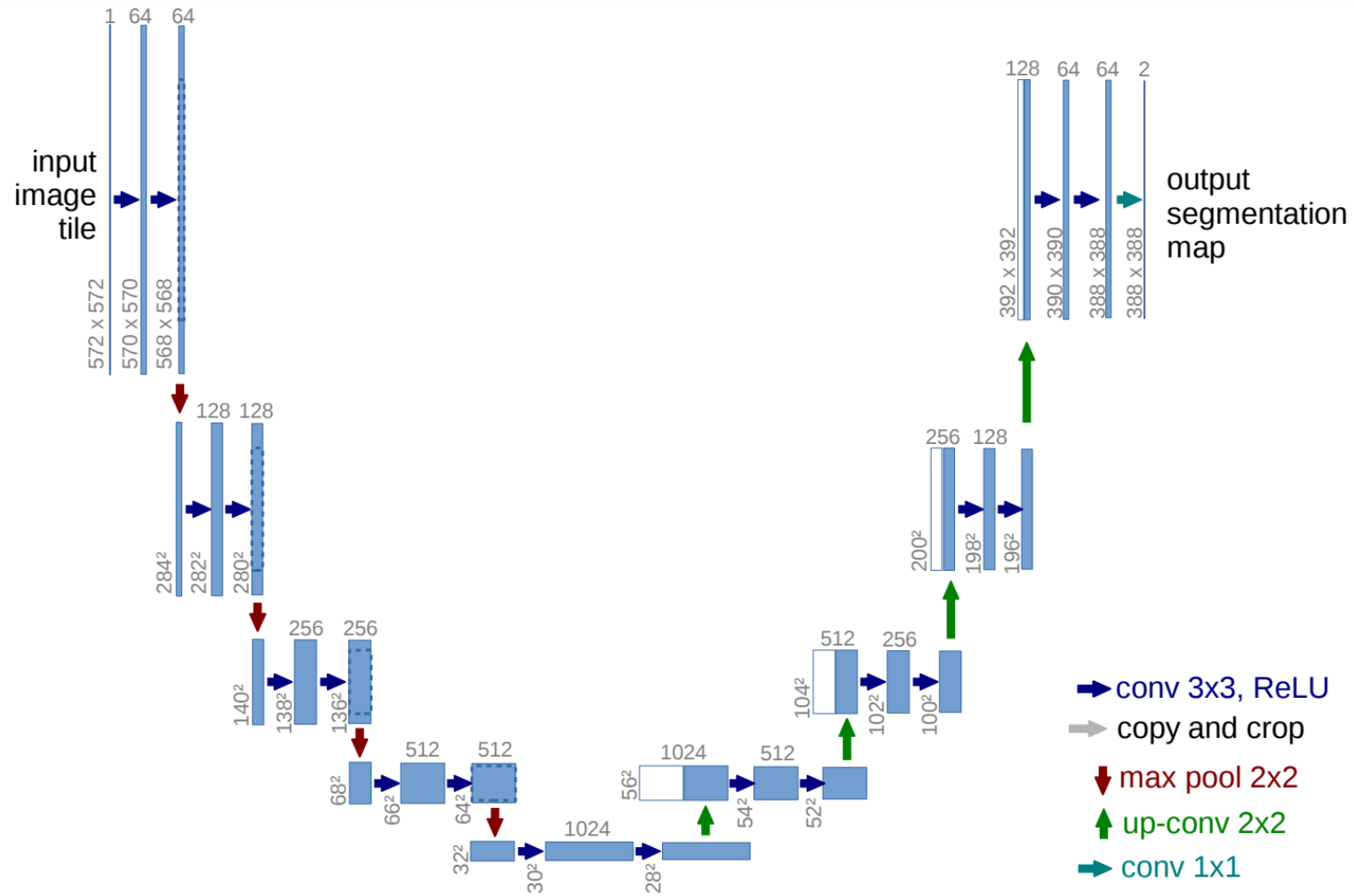
U-Nets & ResNets

Diffusion Models (Ho et al., 2020)

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

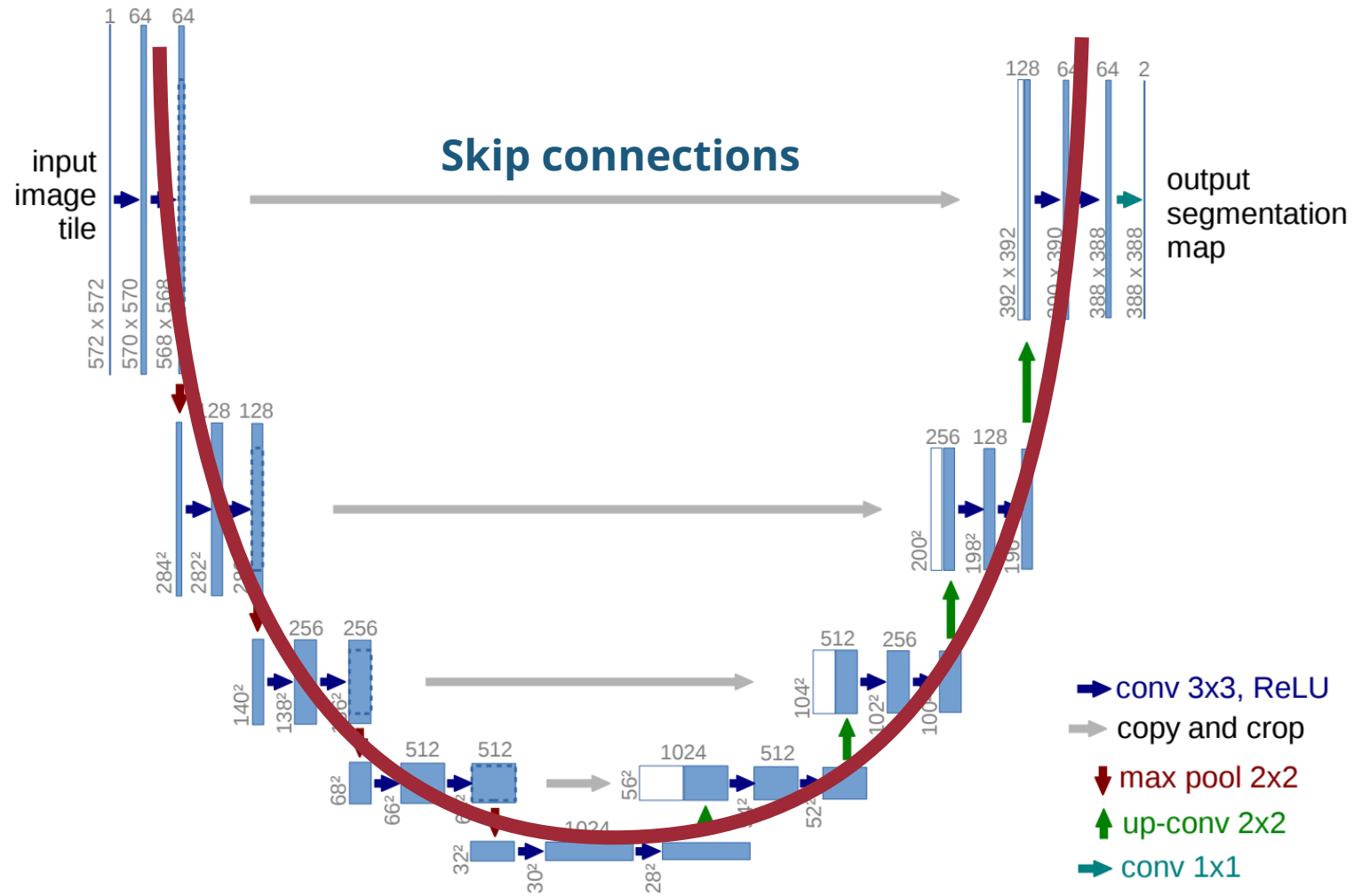


U-Net (Ronneberger et al., 2015)



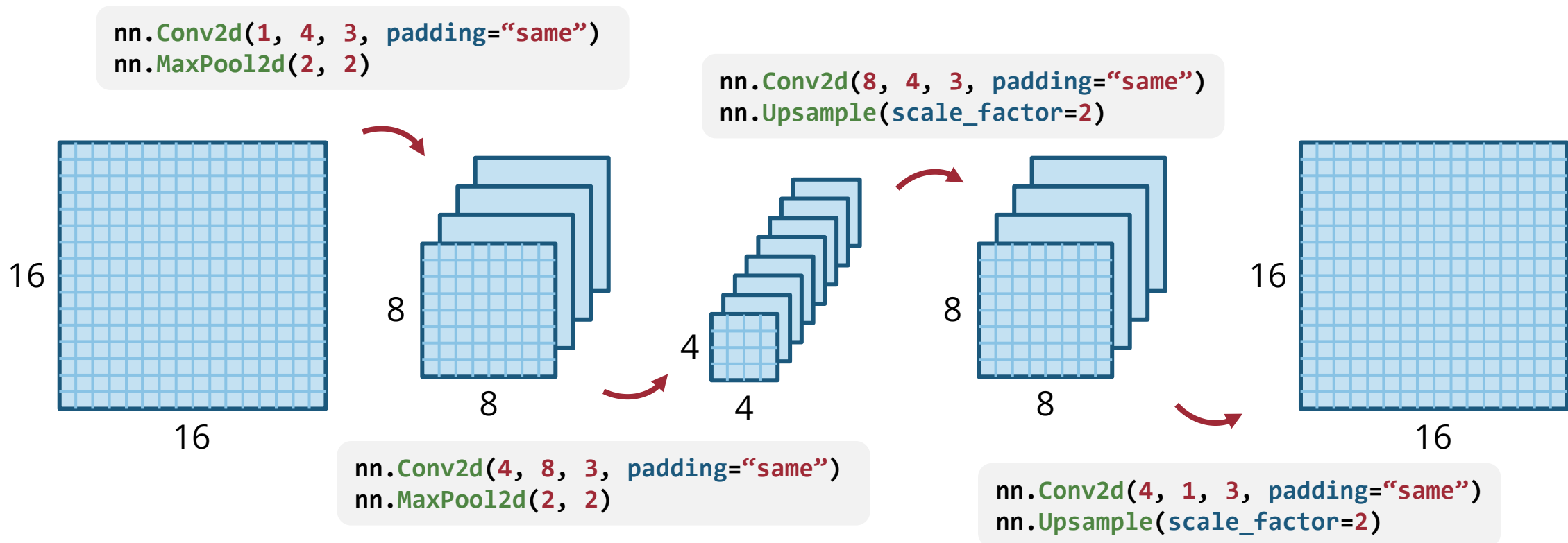
(Source: Ronneberger et al., 2015)

U-Net (Ronneberger et al., 2015)

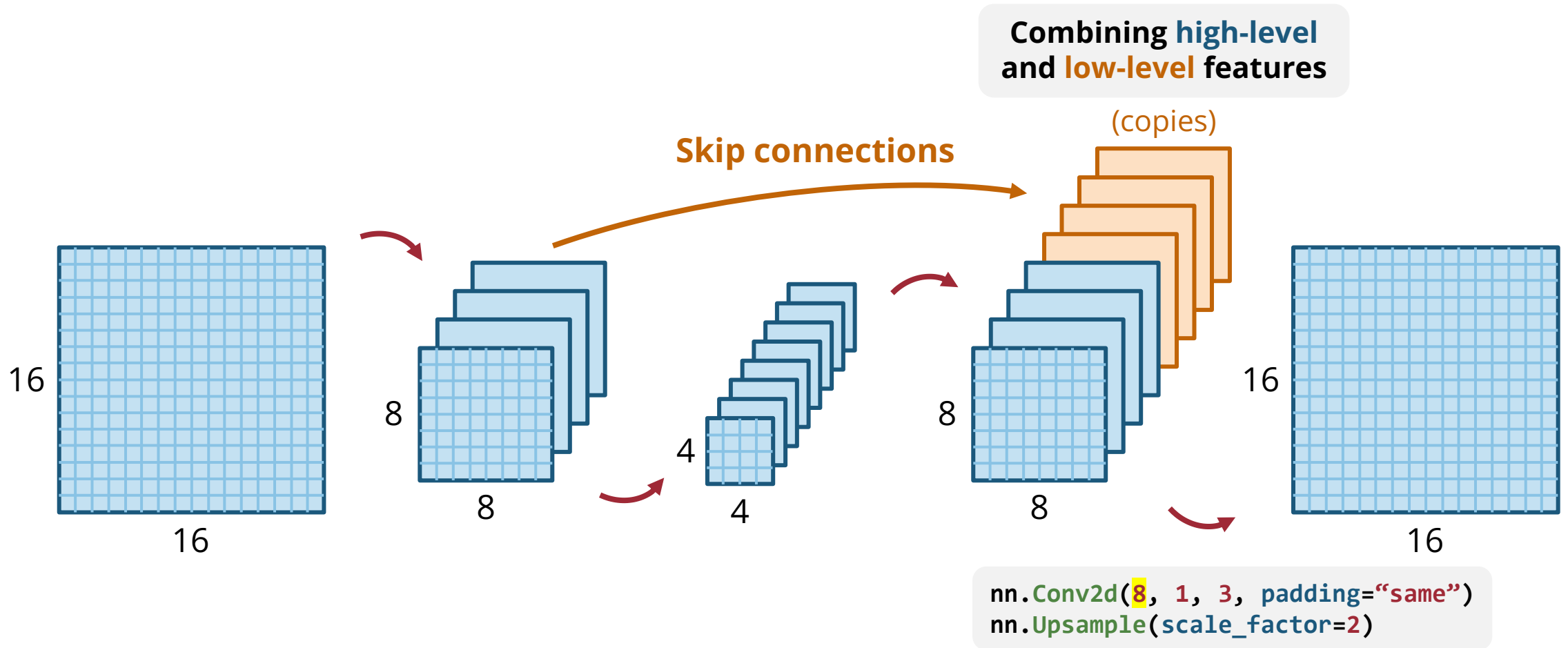


(Source: Ronneberger et al., 2015)

A Toy Example of U-Net

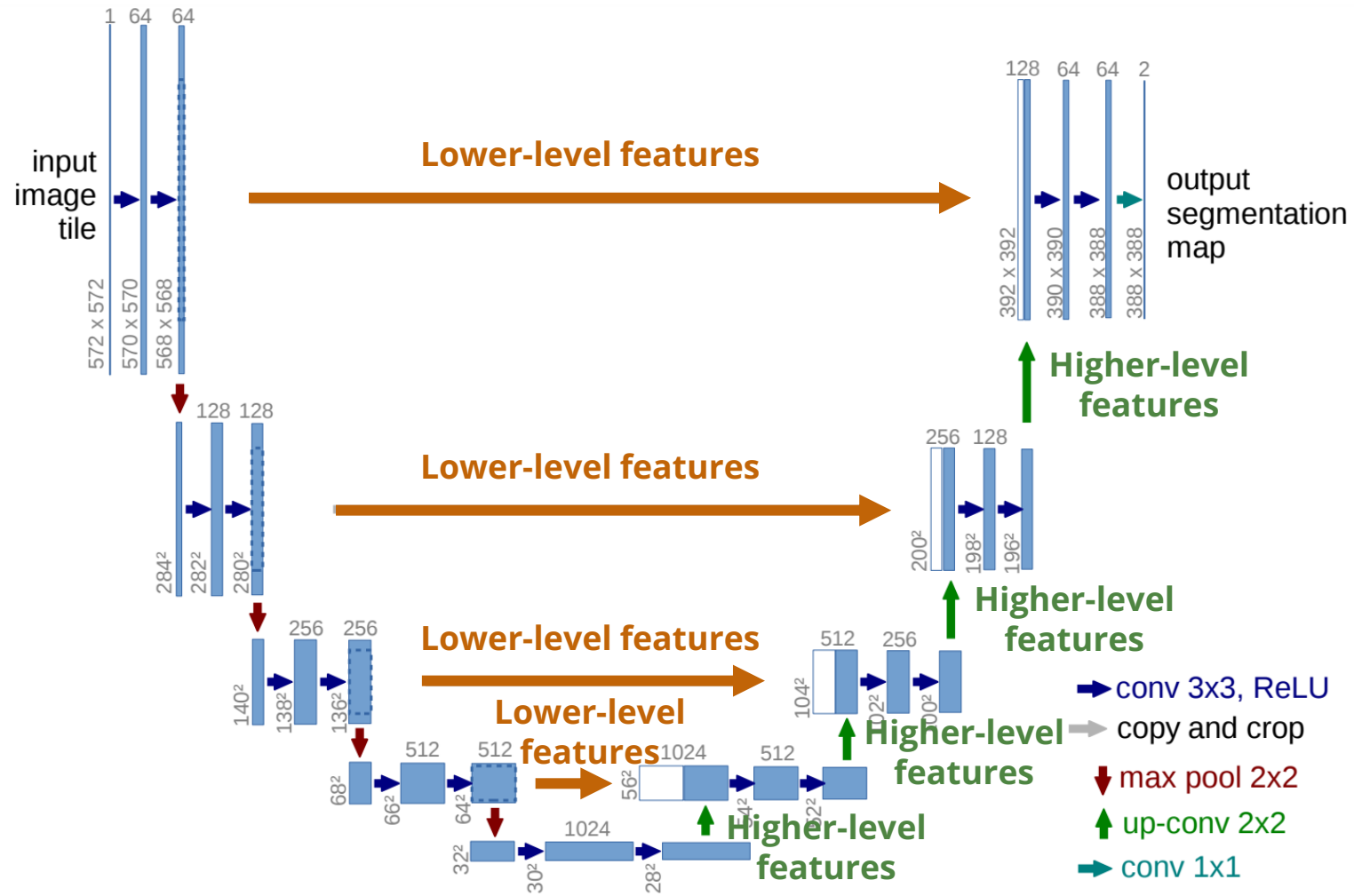


A Toy Example of U-Net



U-Nets are useful when the inputs and outputs have the same shape!

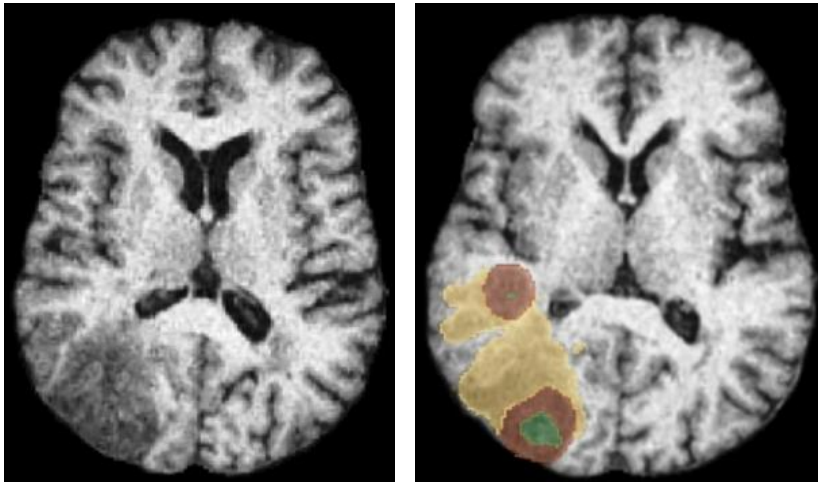
U-Net (Ronneberger et al., 2015)



(Source: Ronneberger et al., 2015)

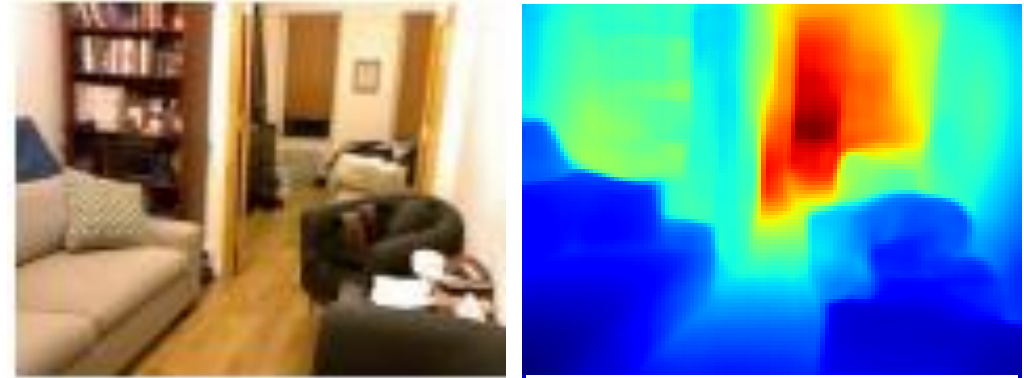
Applications of U-Nets

Tumor Segmentation



(Source: Kharaji et al., 2024)

Depth Estimation



(Source: Barakat, 2018)

Image Segmentation



(Source: Kirillov et al., 2023)

Omar Barakat, "Depth estimation with deep Neural networks part 1," *Medium*, January 11, 2018

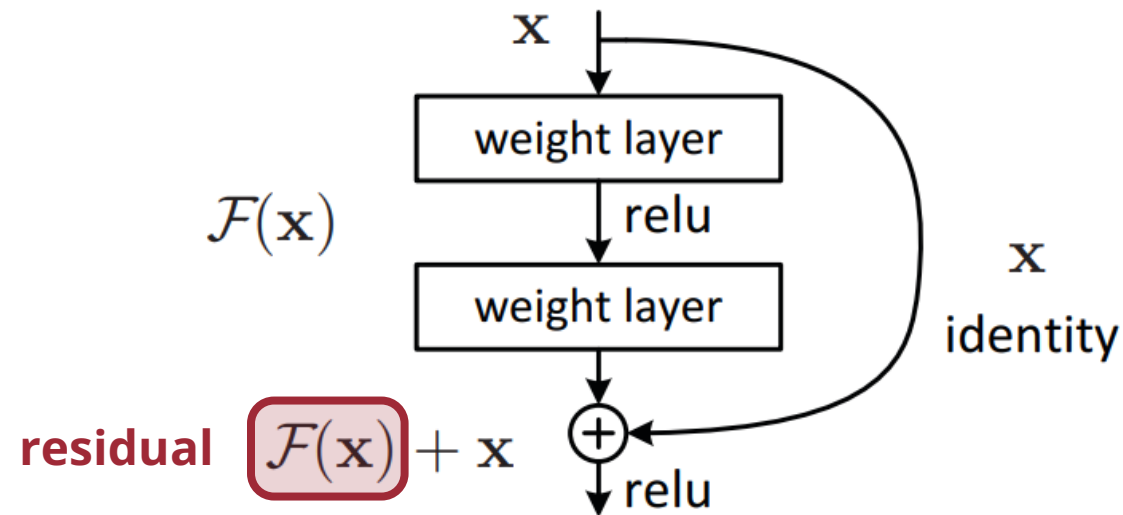
Mona Kharaji, Hossein Abbasi, Yasin Orouskhani, Mostafa Shomalzadeh, Foad Kazemi, and Maysam Orouskhani, "nnU-Net for Brain Tumor Segmentation," *Neuroscience Informatics*, 2024.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick, "Segment Anything," *ICCV*, 2023.

ResNet: Residual Neural Network (He et al., 2016)

- **Intuition:** Learn how to **update** the input x to $x + \Delta x$

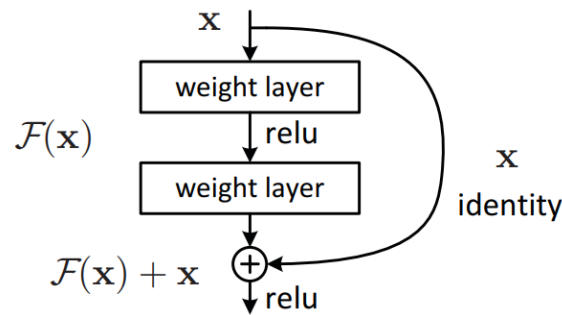
Δx
↓
residual



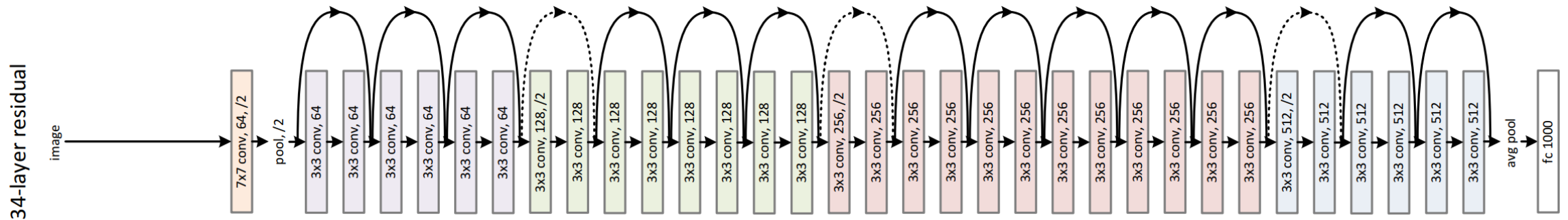
(Source: He et al., 2016)

ResNet: Residual Neural Network (He et al., 2016)

- **Intuition:** Learn how to **update** the input x to $x + \Delta x$



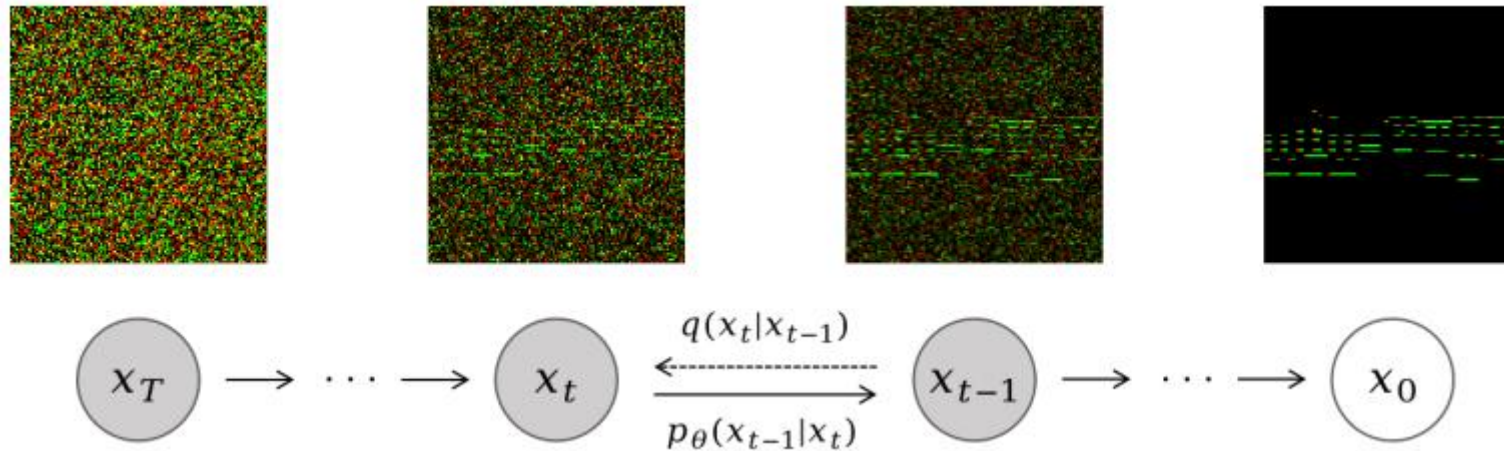
(Source: He et al., 2016)



(Source: He et al., 2016)

Generating Music using Diffusion Models

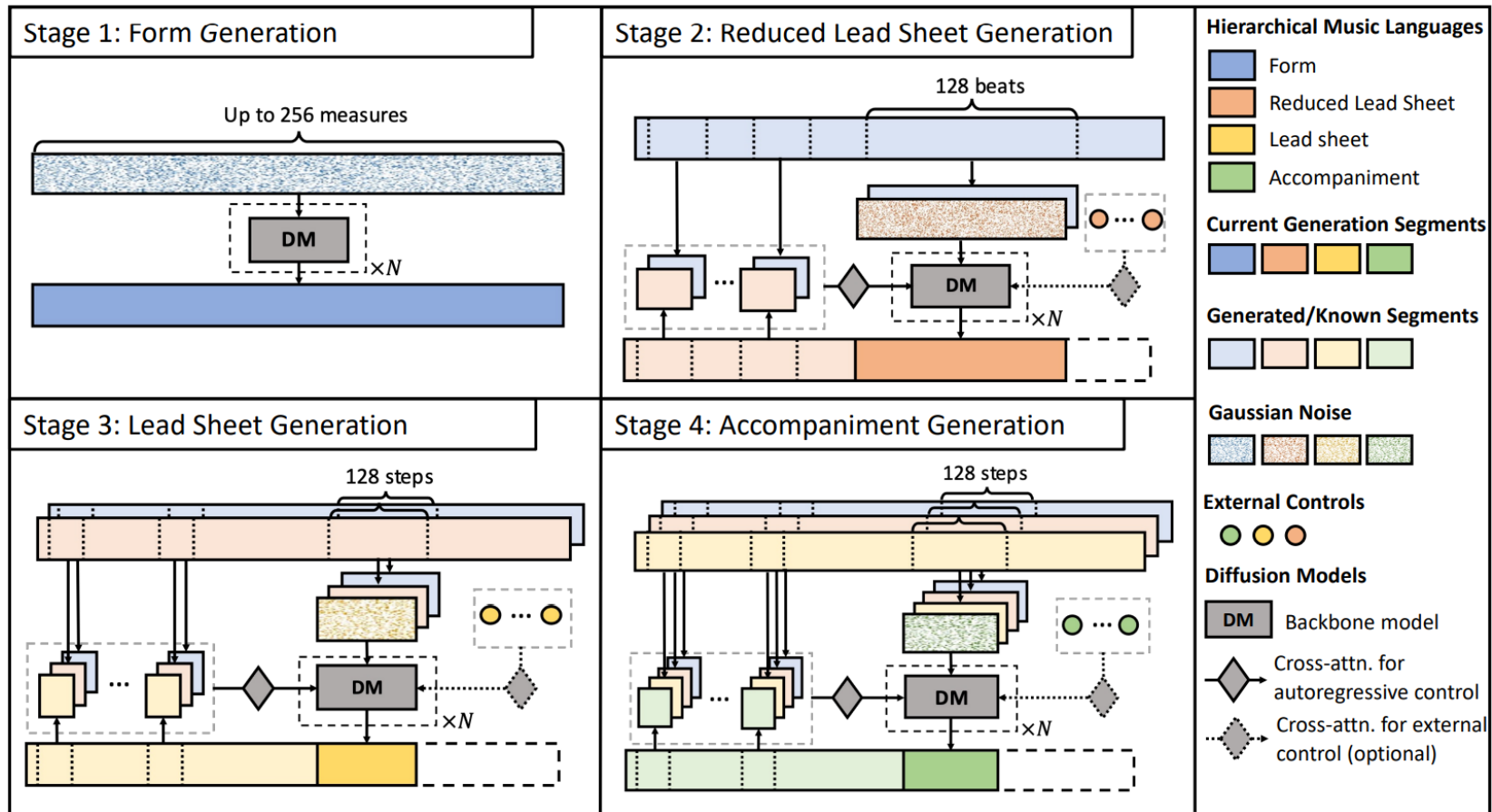
Polyffusion (Min et al., 2023)



(Source: Min et al., 2023)

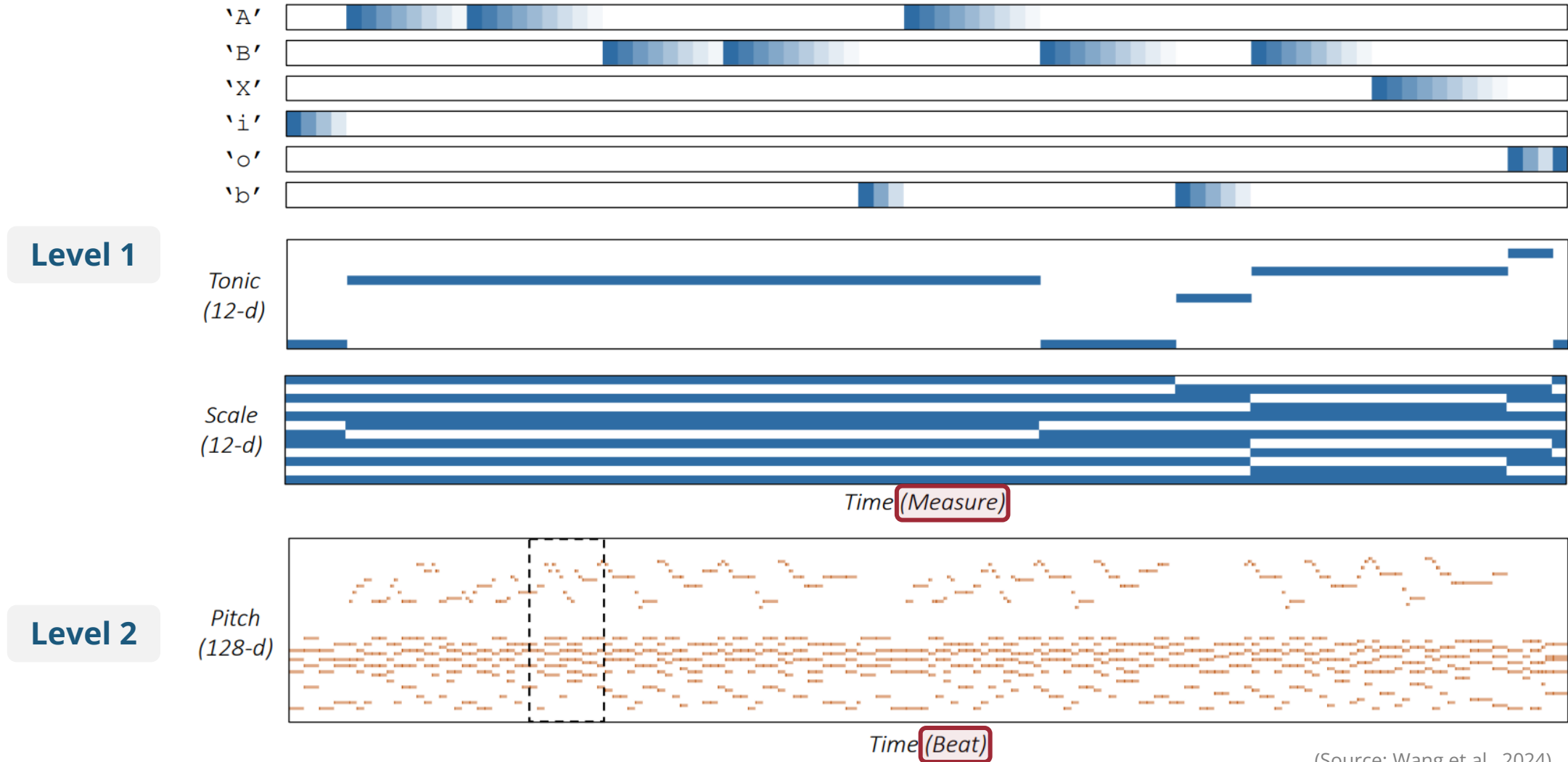
polyffusion.github.io

Cascaded Diffusion Models (Wang et al., 2024)



(Source: Wang et al., 2024)

Cascaded Diffusion Models (Wang et al., 2024)

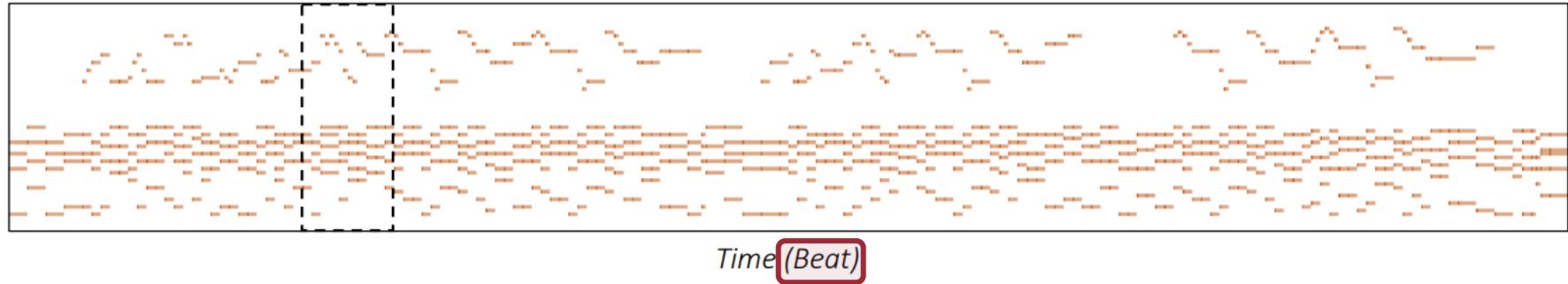


(Source: Wang et al., 2024)

Cascaded Diffusion Models (Wang et al., 2024)

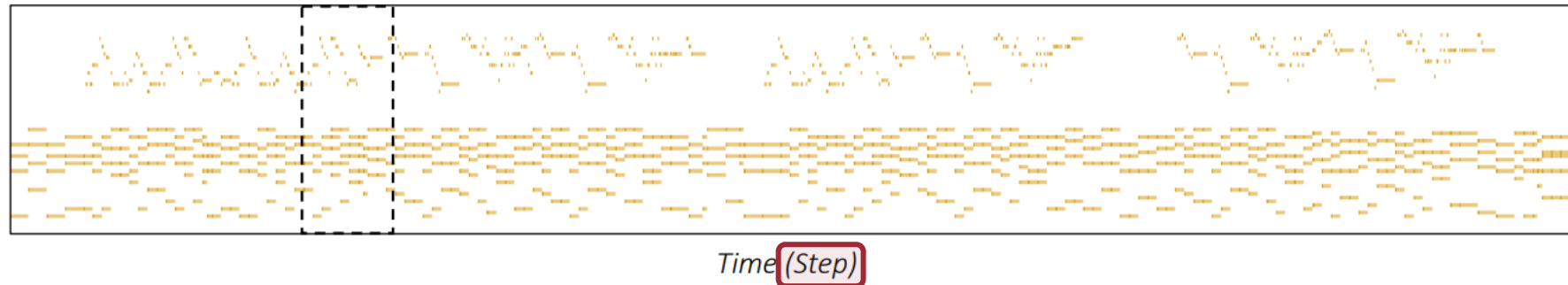
Level 2

Pitch
(128-d)



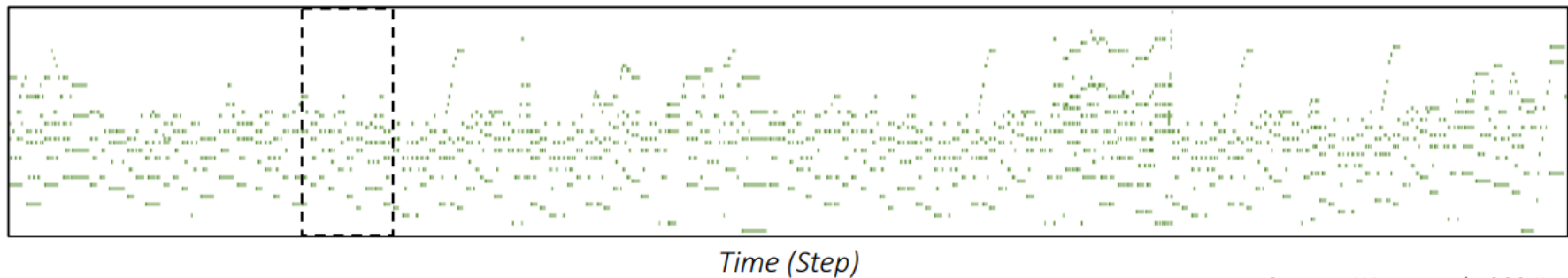
Level 3

Pitch
(128-d)



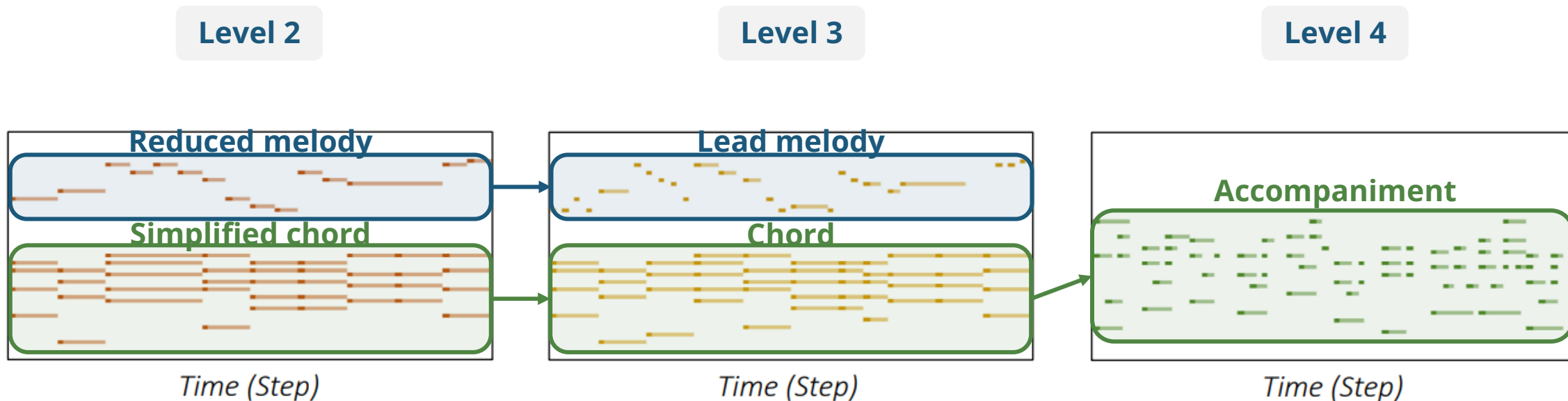
Level 4

Pitch
(128-d)



(Source: Wang et al., 2024)

Cascaded Diffusion Models (Wang et al., 2024)



(Source: Wang et al., 2024)

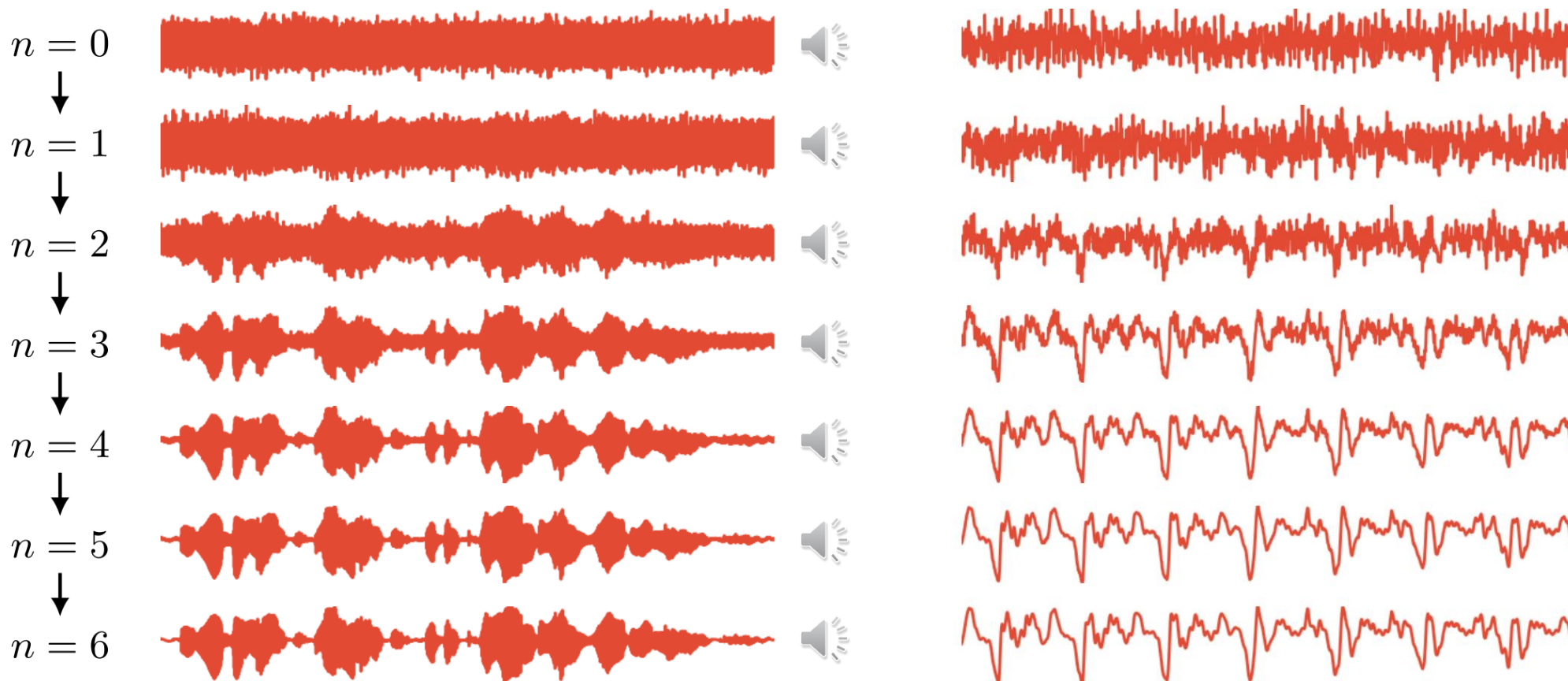
wholesonggen.github.io

Diffusion Models for Audio

WaveGrad: Diffusion for Waveforms (Chen et al., 2021)

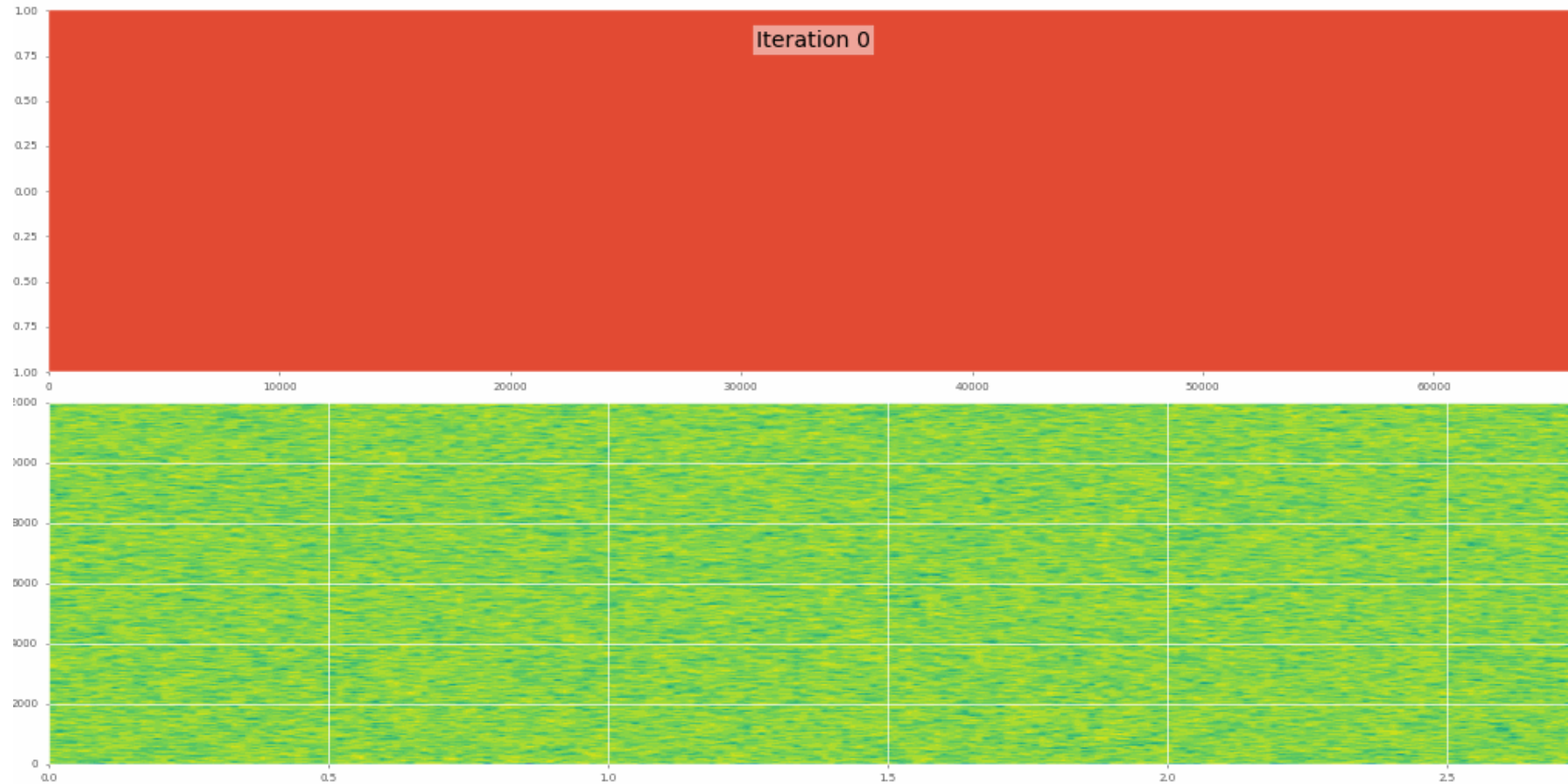
Text: Here are the match lineups for the Colombia Haiti match.

Zoom in



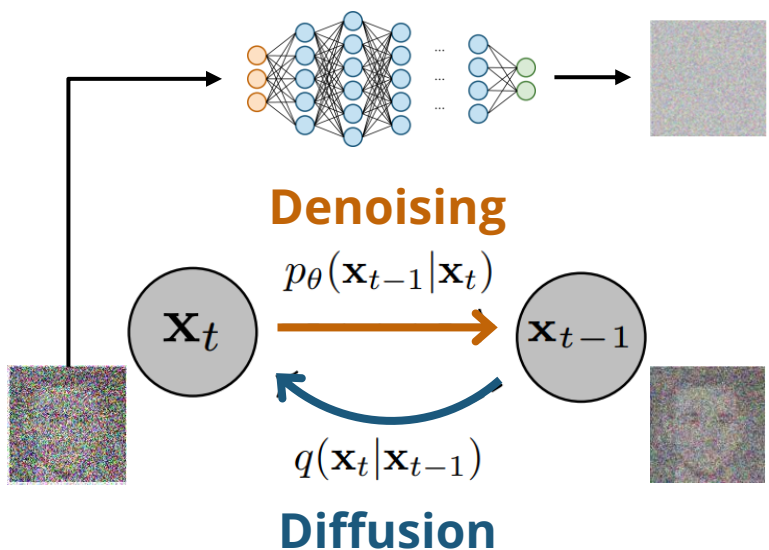
(Source: Chen et al., 2021)

WaveGrad: Diffusion for Waveforms (Chen et al., 2021)

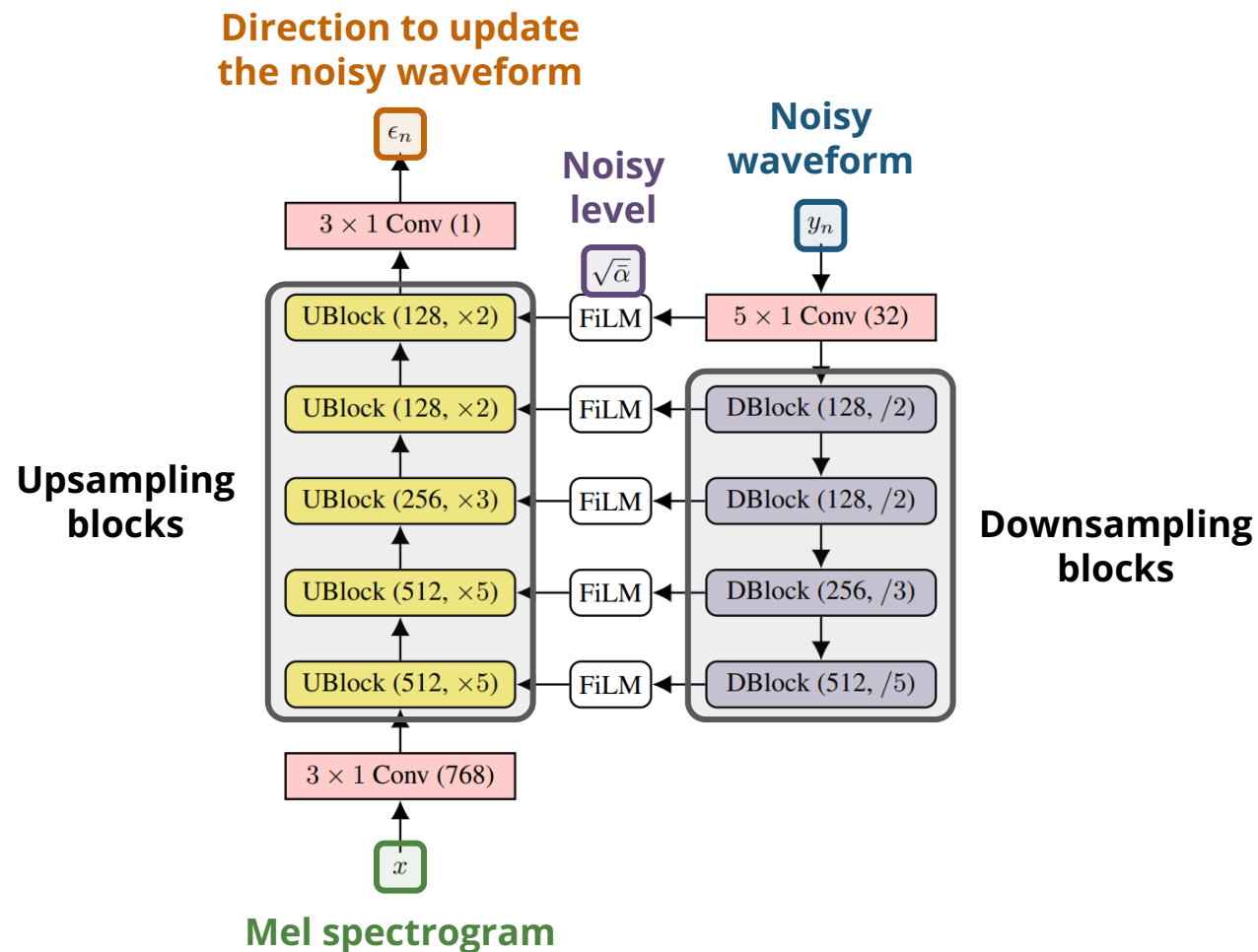


(Source: Chen et al., 2021)

WaveGrad: Diffusion for Waveforms (Chen et al., 2021)



(Source: Ho et al., 2020)

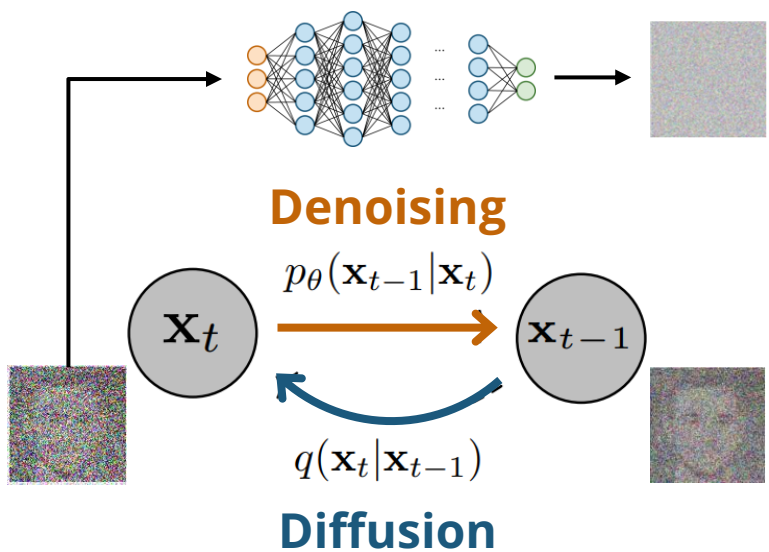


(Source: Chen et al., 2021)

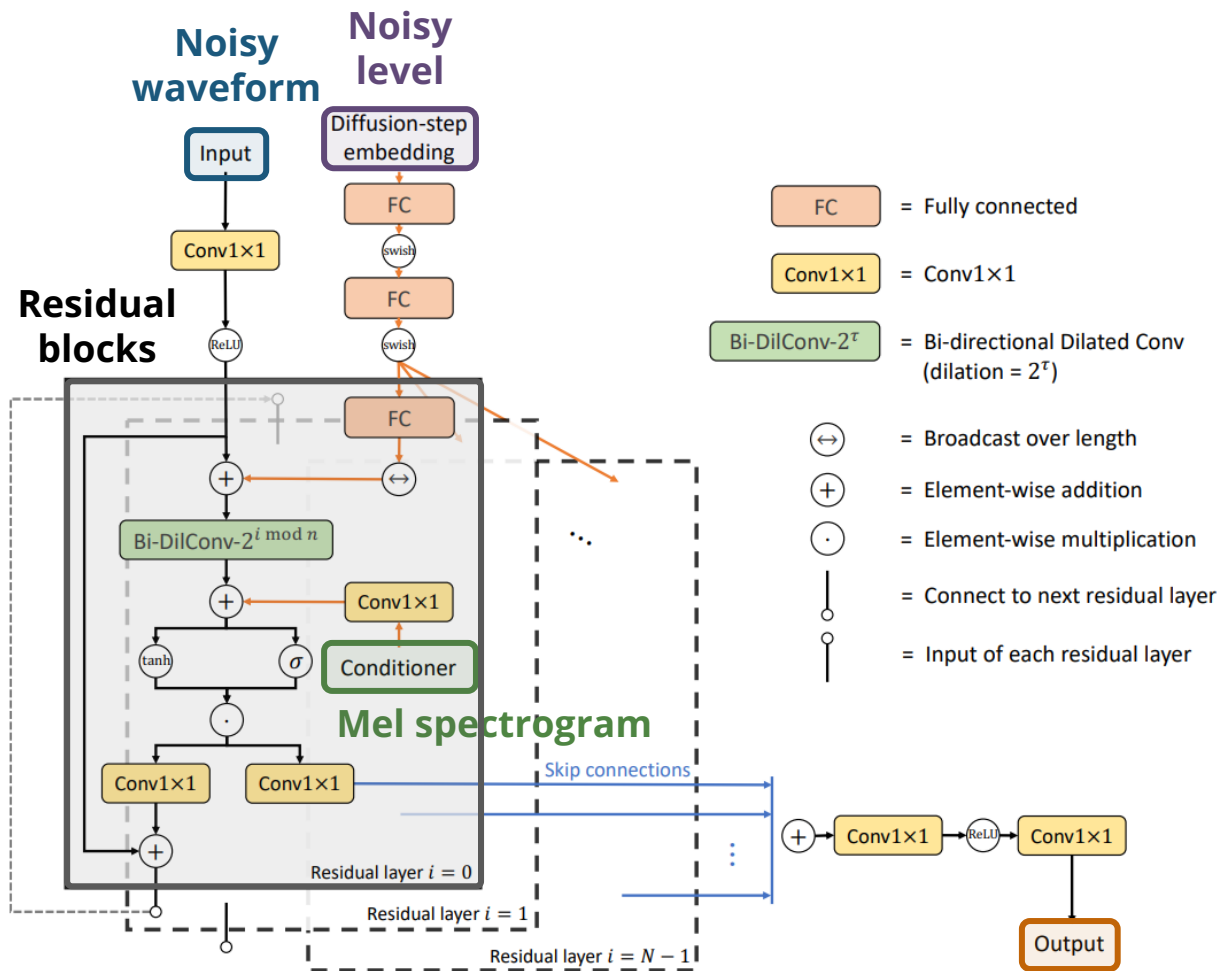
Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising Diffusion Probabilistic Models," *NeurIPS*, 2020.

Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan, "WaveGrad: Estimating Gradients for Waveform Generation," *ICLR*, 2021.

DiffWave: Diffusion Model for Waveforms (Kong et al., 2021)



(Source: Ho et al., 2020)



Direction to update the noisy waveform

(Source: Kong et al., 2021)

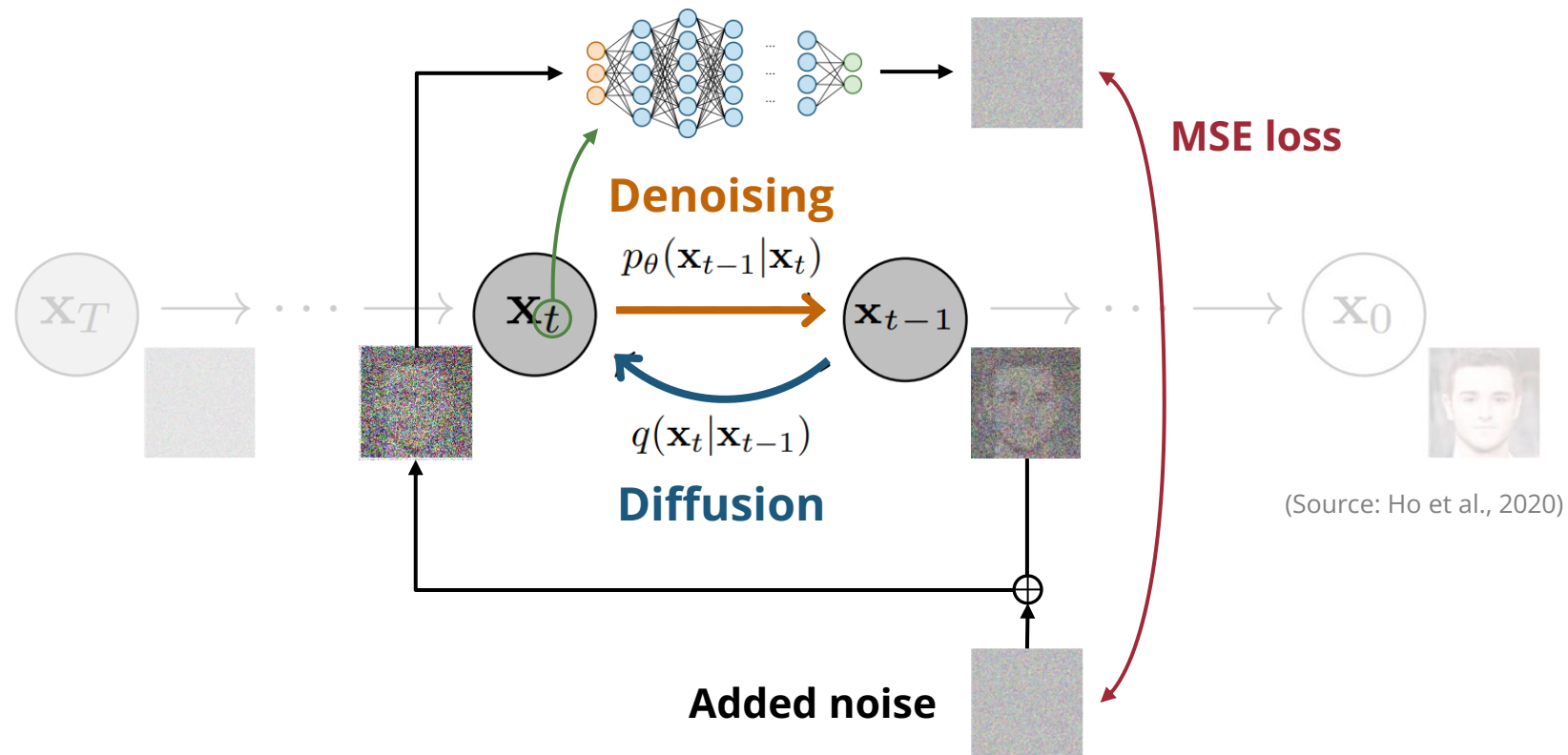
Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising Diffusion Probabilistic Models," *NeurIPS*, 2020.

Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan, "WaveGrad: Estimating Gradients for Waveform Generation," *ICLR*, 2021.

Efficient Diffusion Models

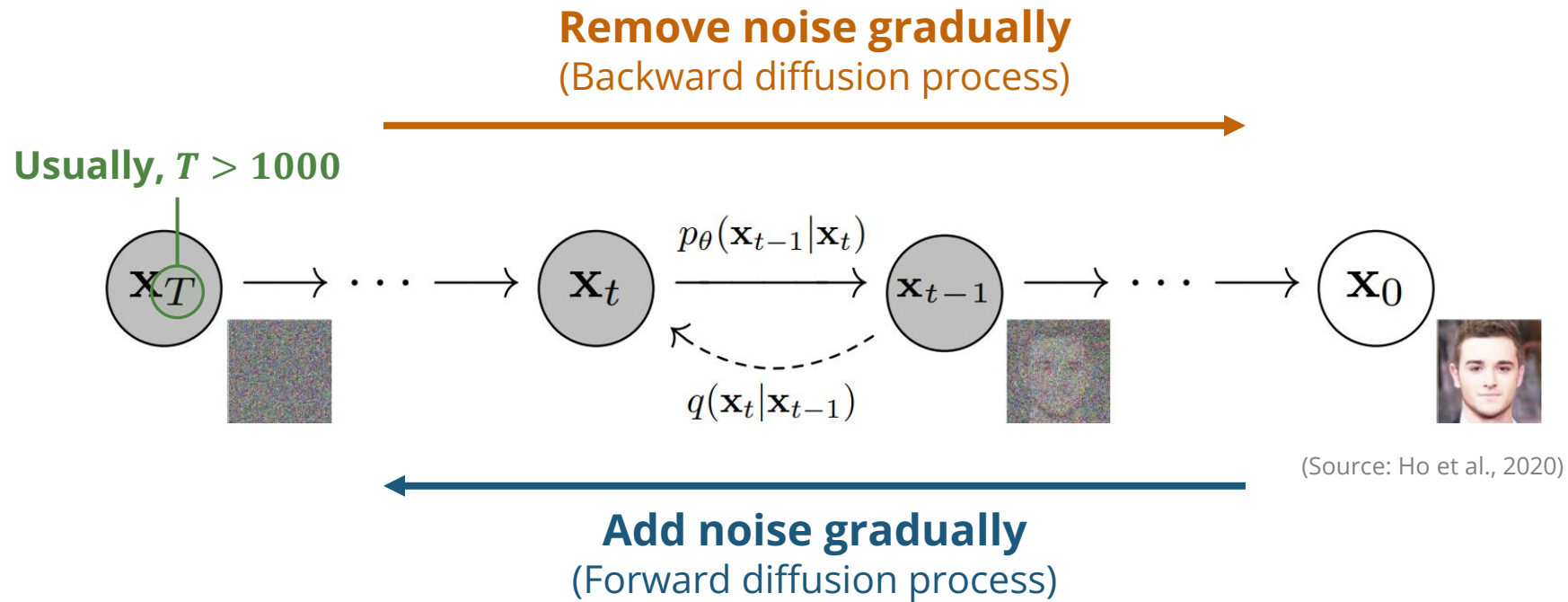
Diffusion Models: Training

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



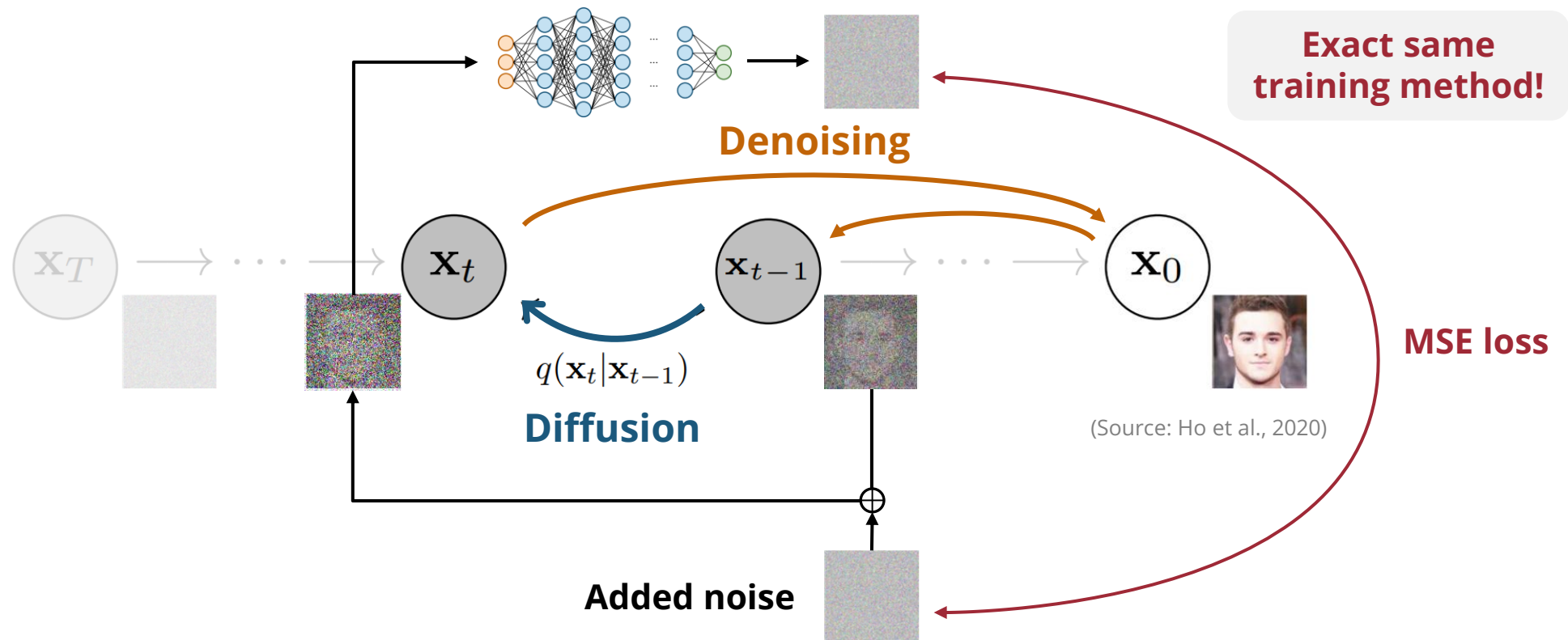
Diffusion Models (Ho et al., 2020)

- Intuition:** Many denoising autoencoders stacked together



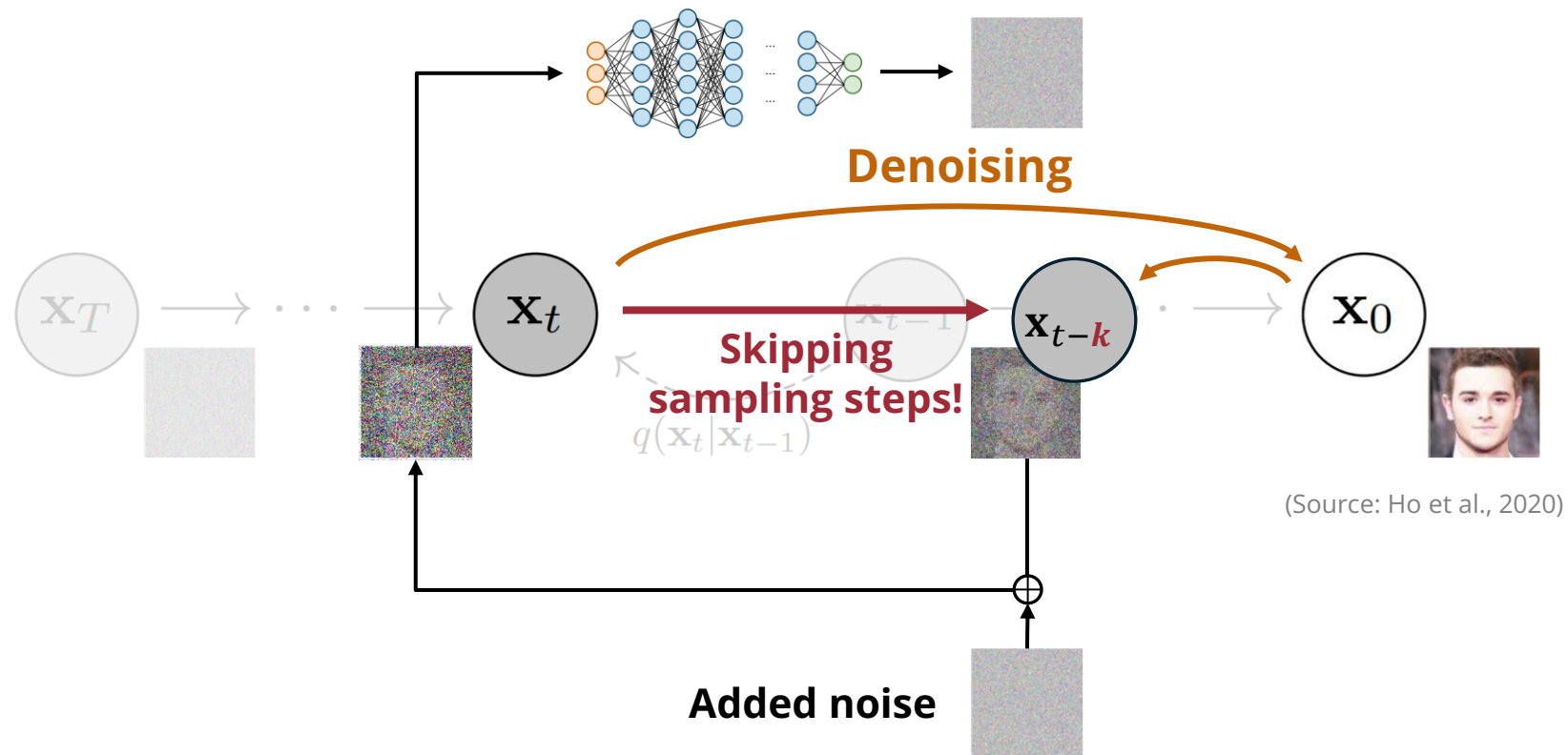
Fast Sampling for Diffusion Models

- **Intuition:** Skip some sampling steps



Fast Sampling for Diffusion Models

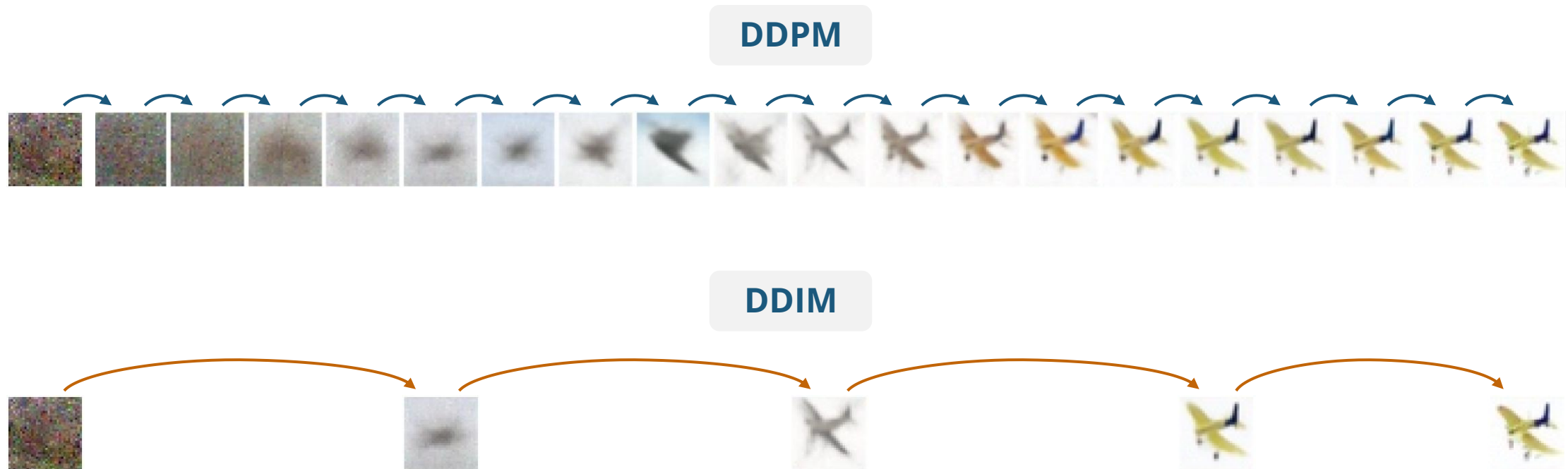
- **Intuition:** Skip some sampling steps



Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising Diffusion Probabilistic Models," *NeurIPS*, 2020.
Jiaming Song, Chenlin Meng, and Stefano Ermon, "Denoising Diffusion Implicit Models," *ICLR*, 2021.

Fast Sampling for Diffusion Models

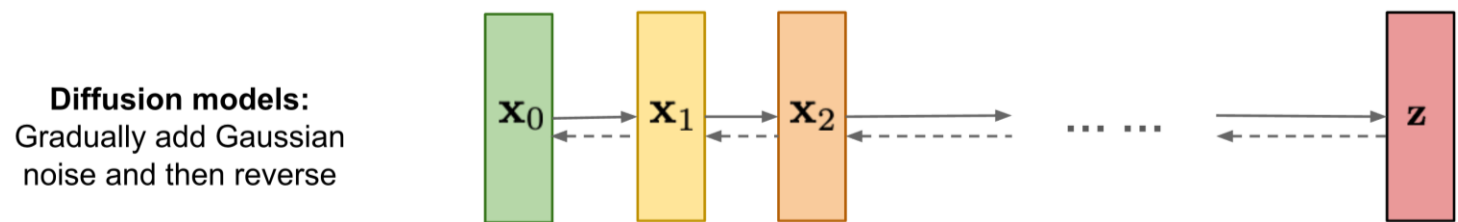
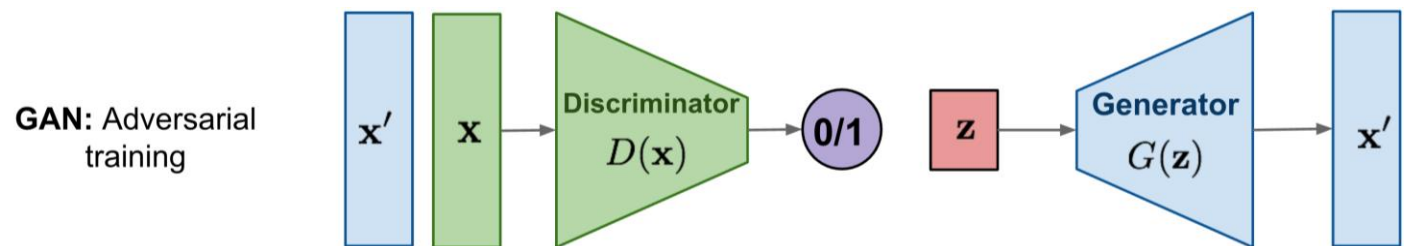
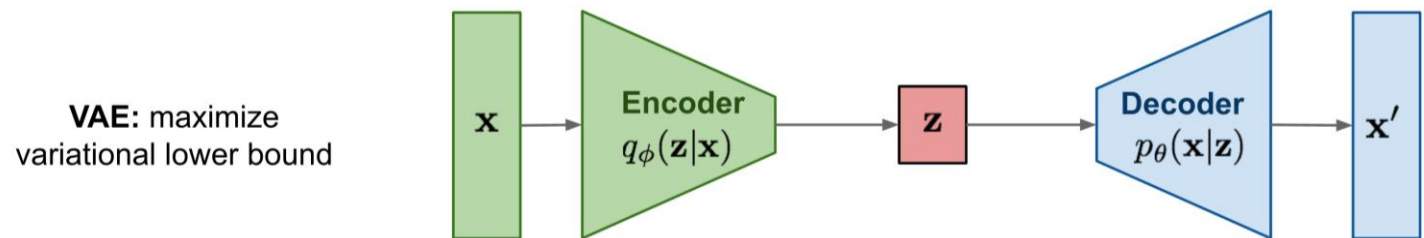
- **Intuition:** Skip some sampling steps



(Source: Ho et al., 2020)

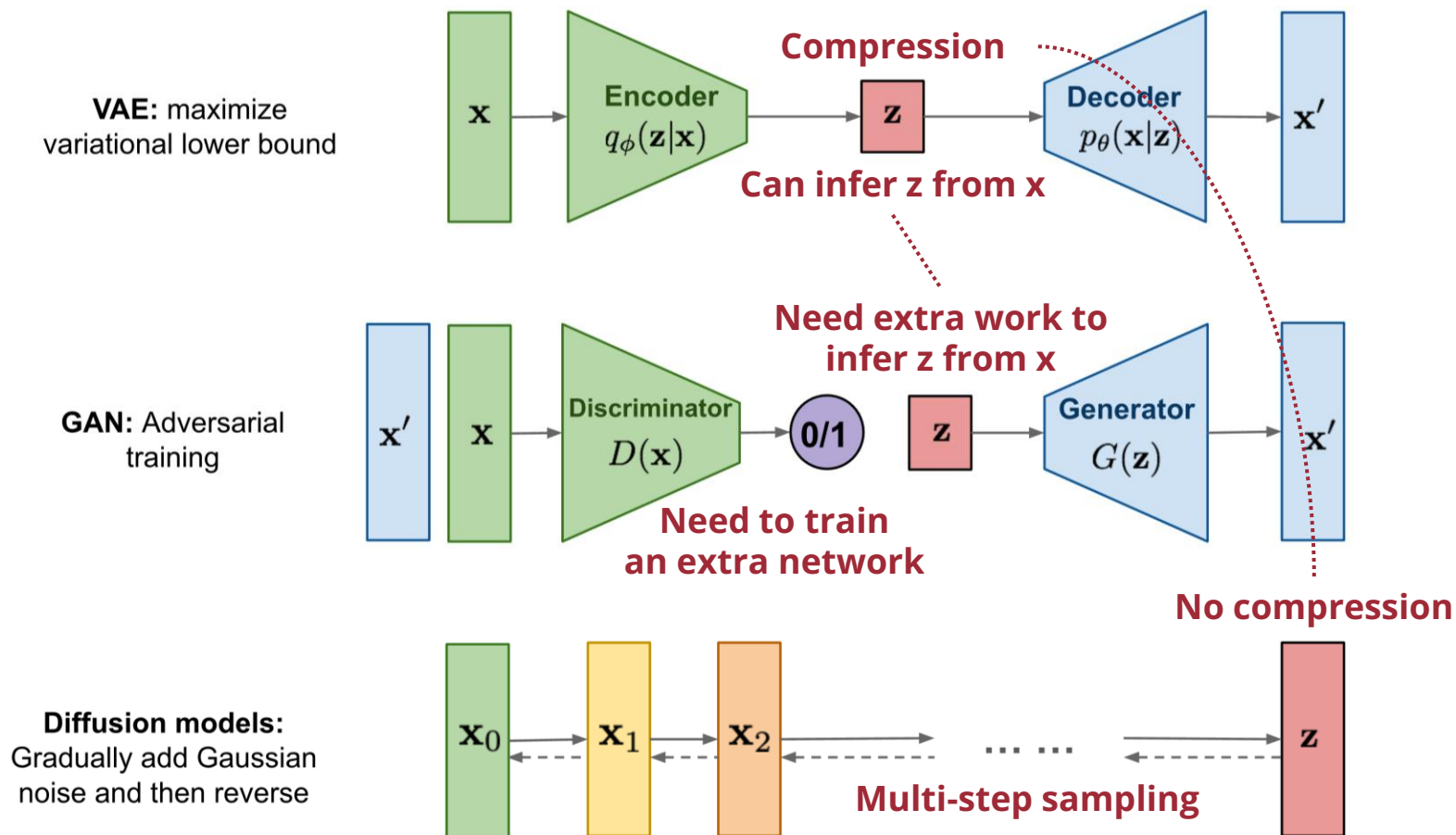
Comparison of Deep Generative Models

Comparison of Deep Latent Variable Models



(Source: Weng, 2021)

Comparison of Deep Latent Variable Models



(Source: Weng, 2021)

Network Architectures vs. Training Frameworks

Network architectures

Multilayer perceptron (MLP)
Convolutional neural networks (CNNs)
Recurrent neural networks (RNNs)
Transformers
ResNets
U-Nets
⋮

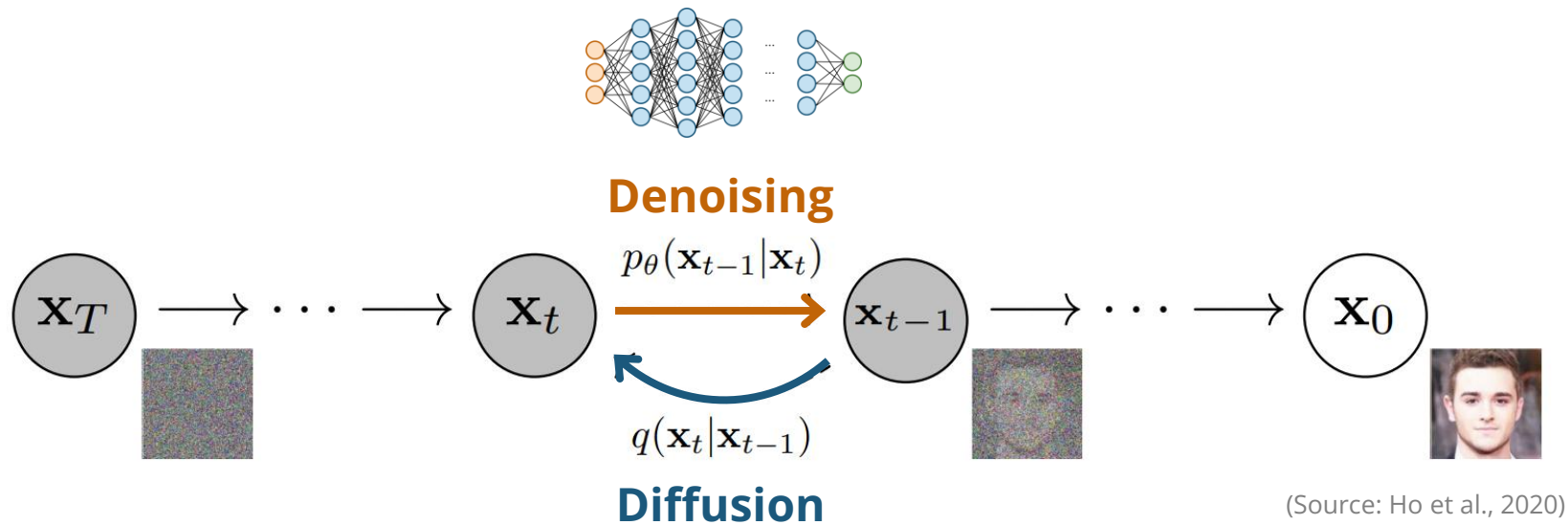
Training frameworks

Autoregressive
Autoencoders
Variational autoencoders (VAEs)
Generative adversarial networks (GANs)
Diffusion models
Consistency models
⋮

Recap

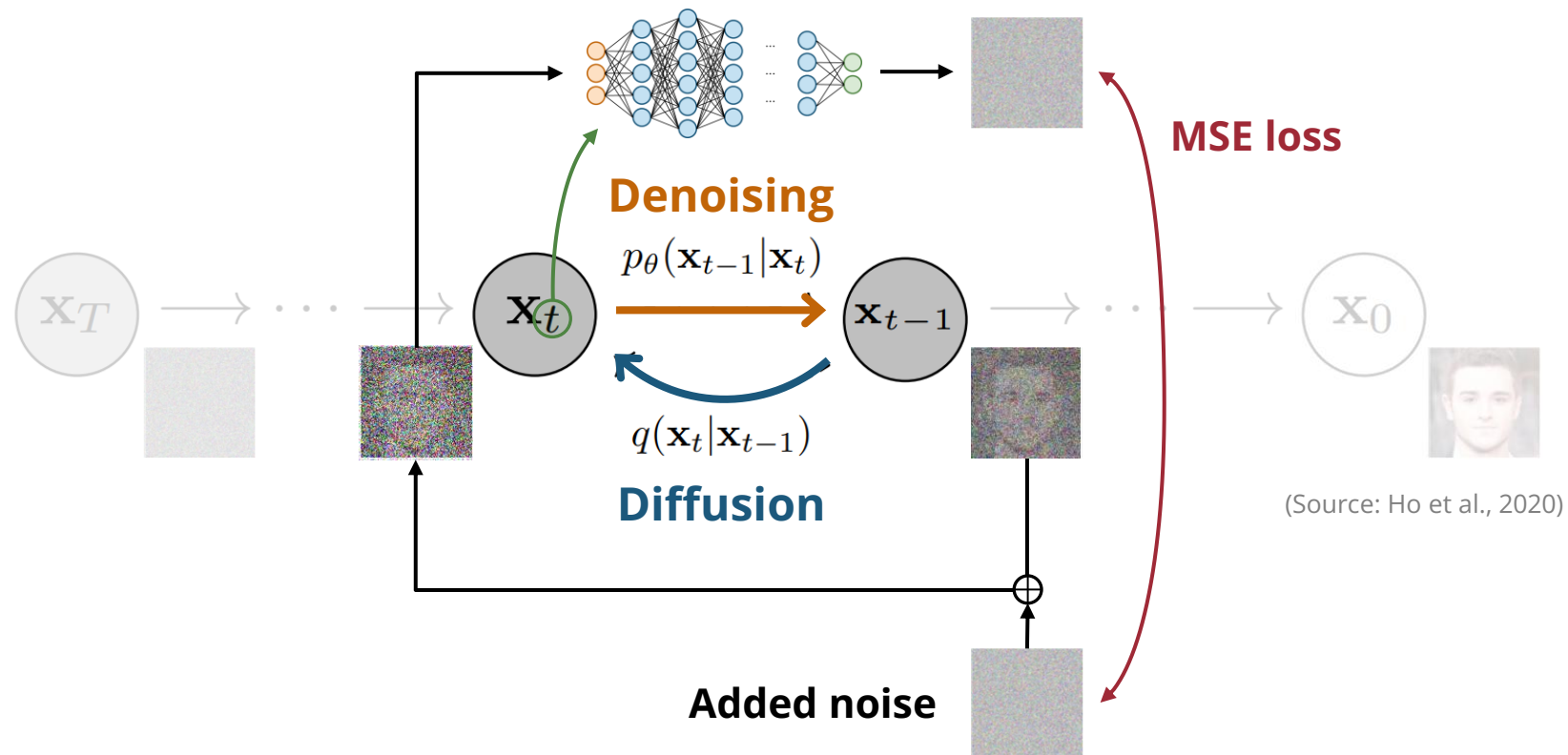
Diffusion Models (Ho et al., 2020)

- Intuition:** Many denoising autoencoders stacked together



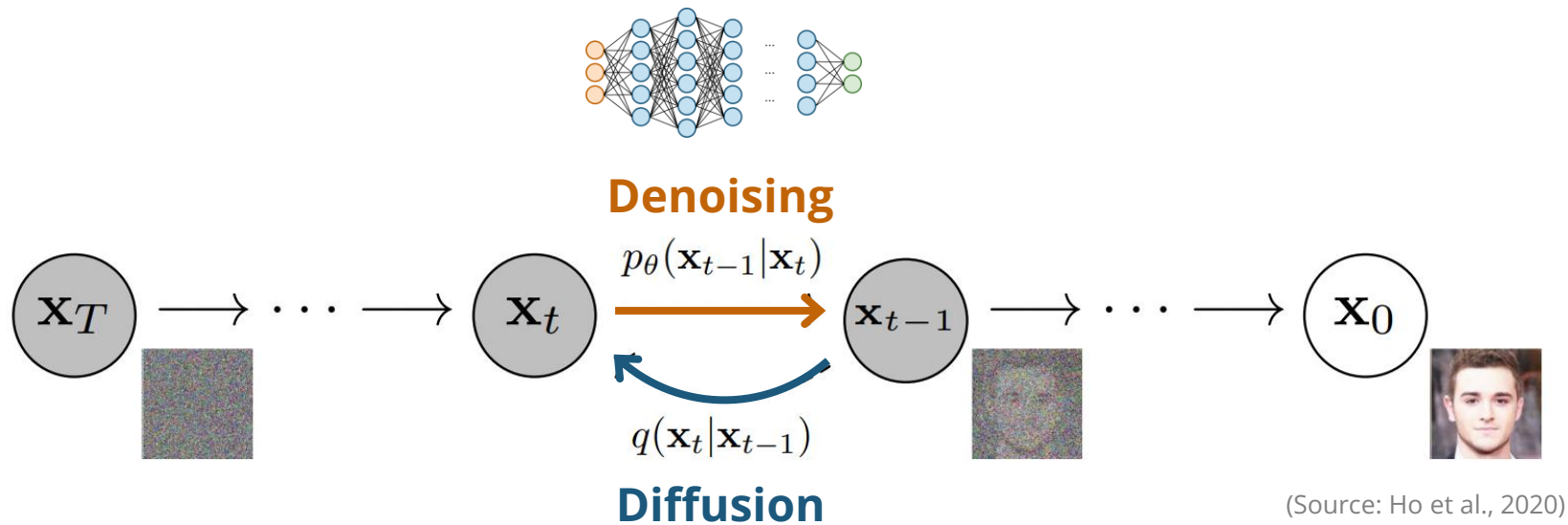
Diffusion Models: Training

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

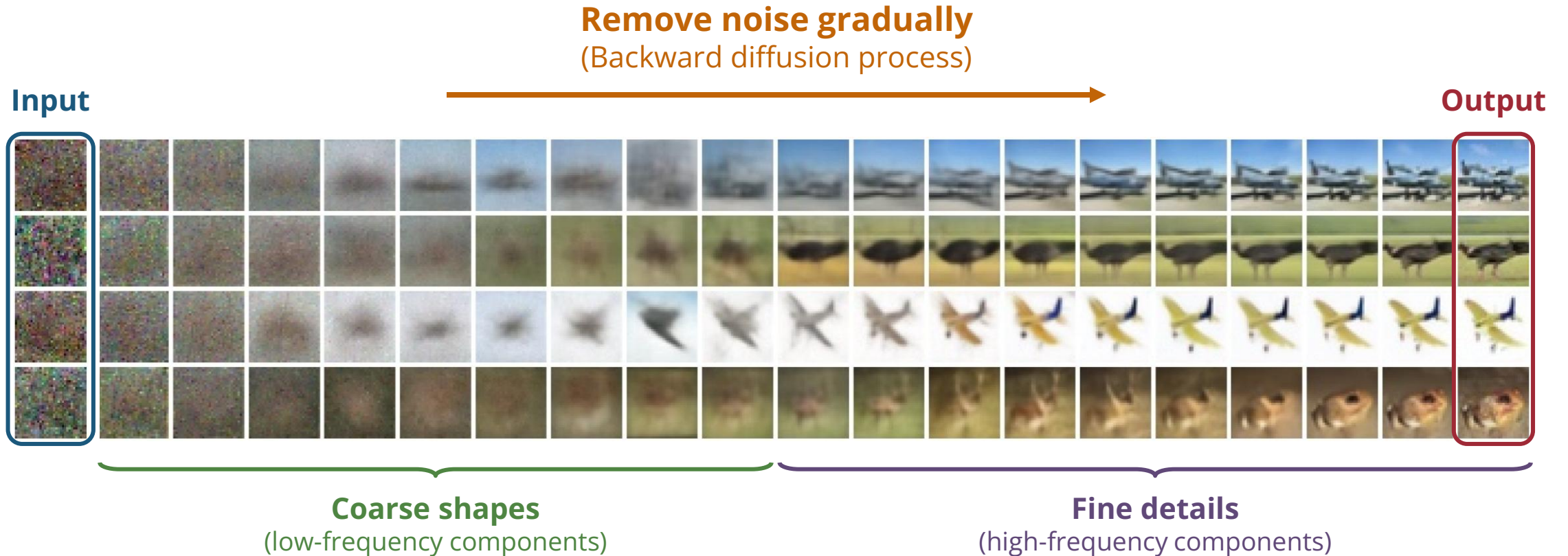


Diffusion Models (Ho et al., 2020)

- Intuition:** Many denoising autoencoders stacked together

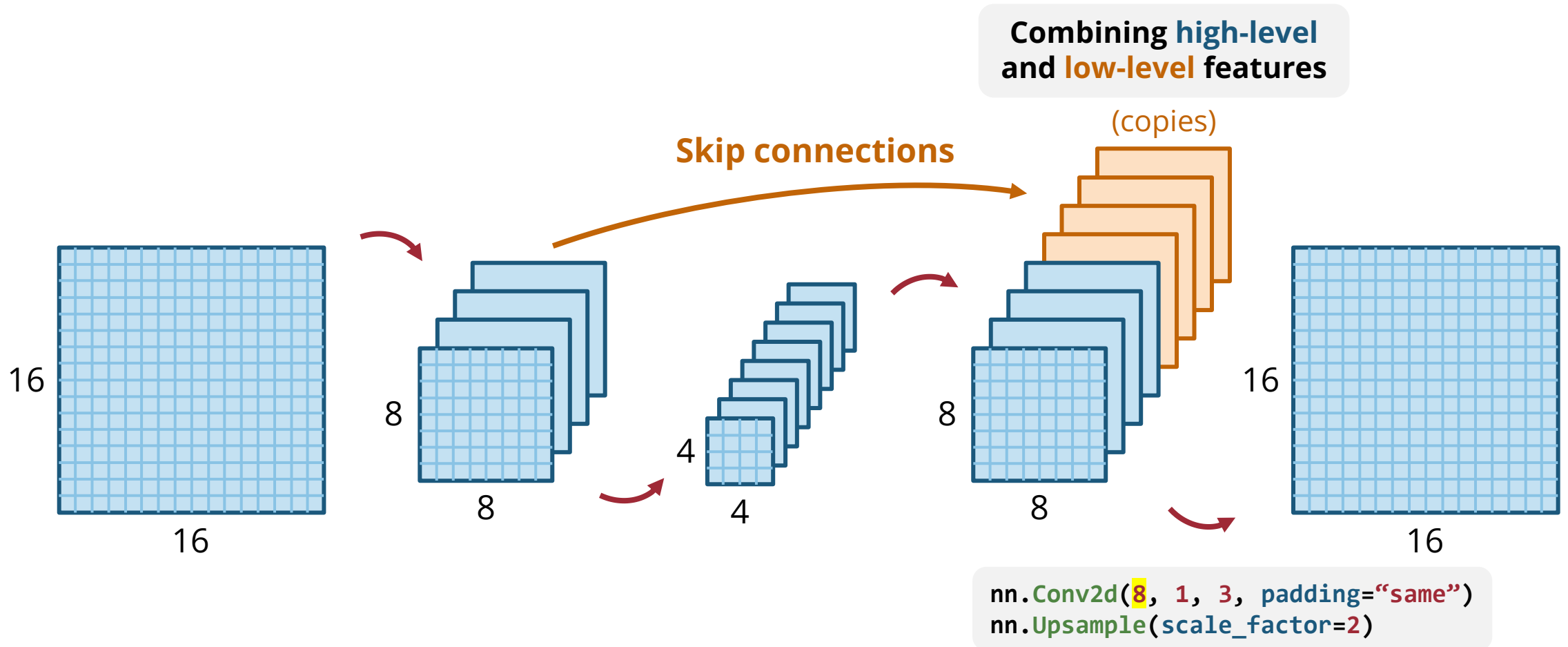


Diffusion Models: Generation



(Source: Ho et al., 2020)

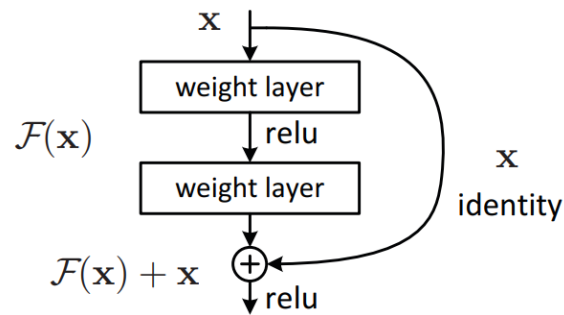
A Toy Example of U-Net



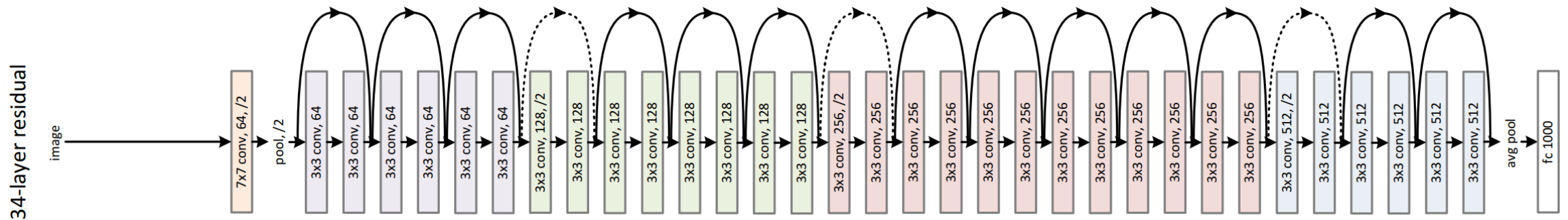
U-Nets are useful when the inputs and outputs have the same shape!

ResNet: Residual Neural Network (He et al., 2016)

- **Intuition:** Learn how to **update** the input x to $x + \Delta x$

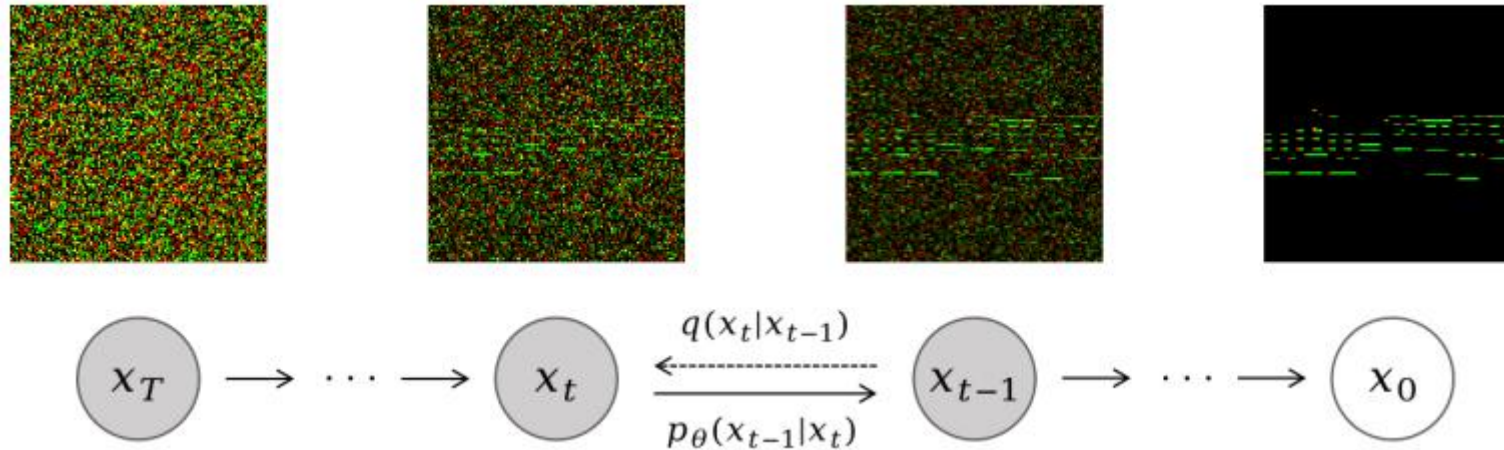


(Source: He et al., 2016)



(Source: He et al., 2016)

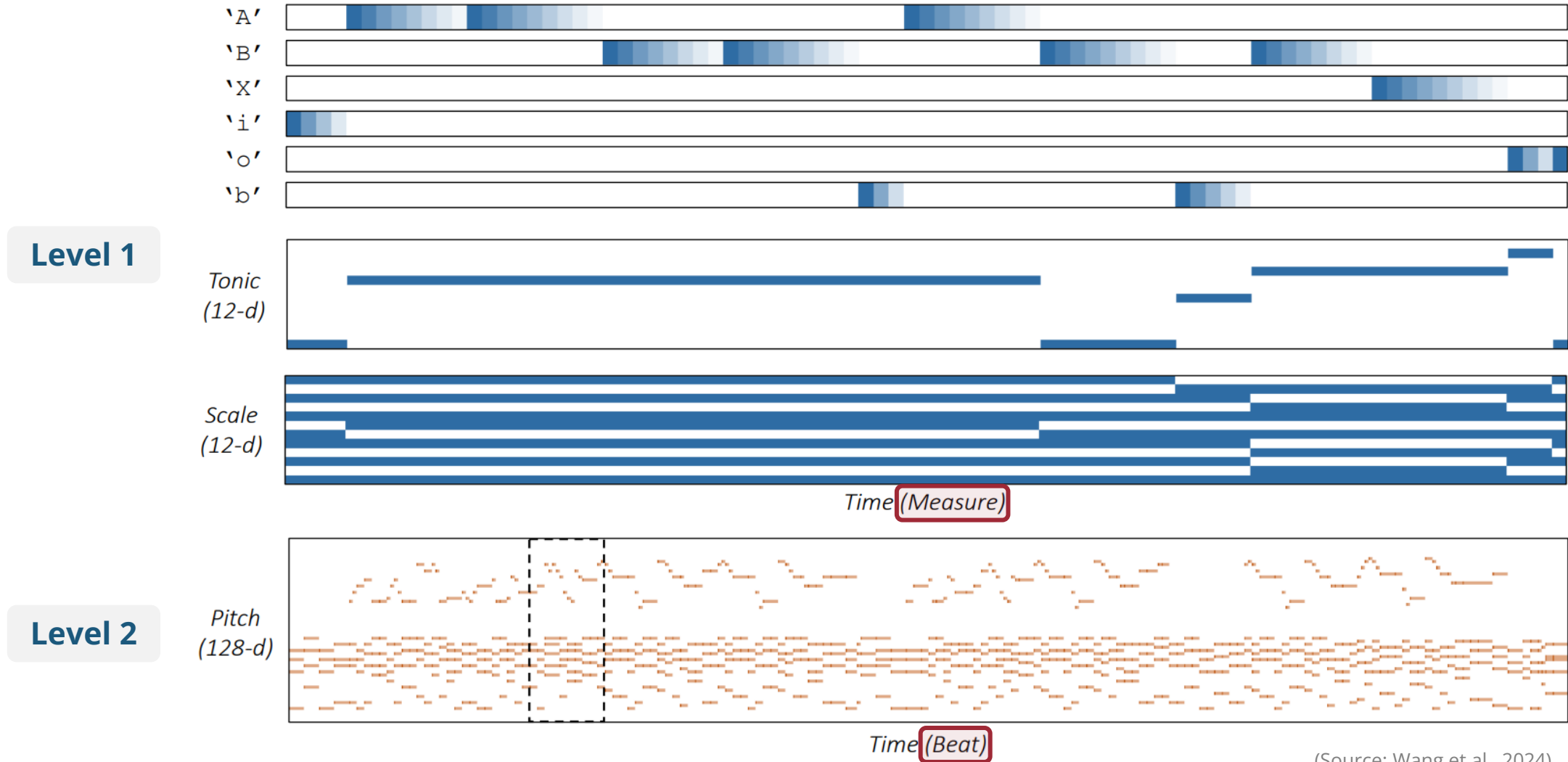
Polyffusion (Min et al., 2023)



(Source: Min et al., 2023)

polyffusion.github.io

Cascaded Diffusion Models (Wang et al., 2024)

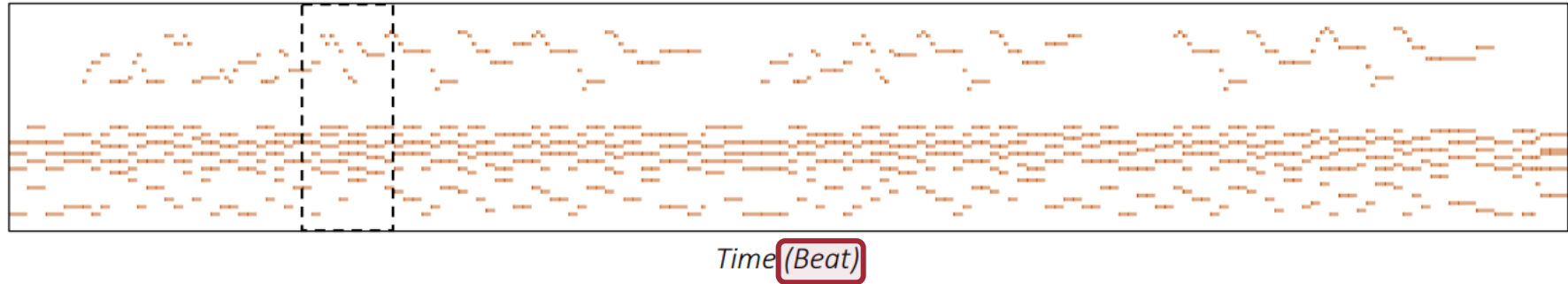


(Source: Wang et al., 2024)

Cascaded Diffusion Models (Wang et al., 2024)

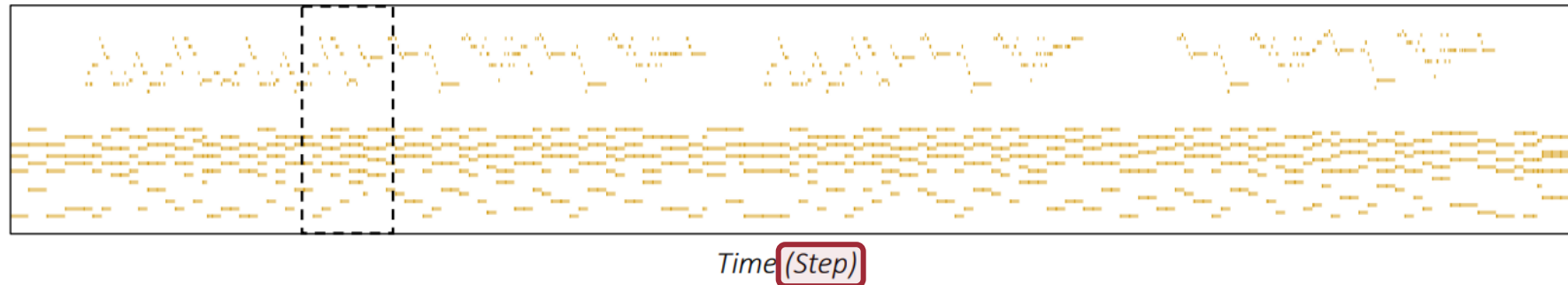
Level 2

Pitch
(128-d)



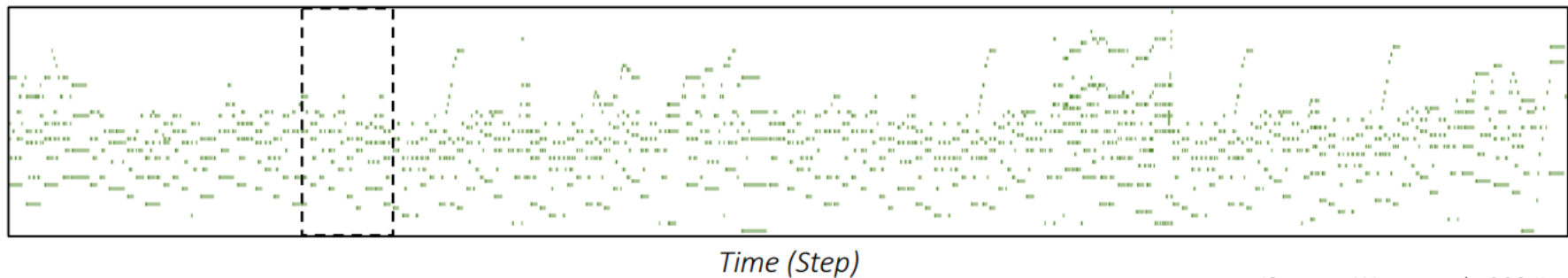
Level 3

Pitch
(128-d)



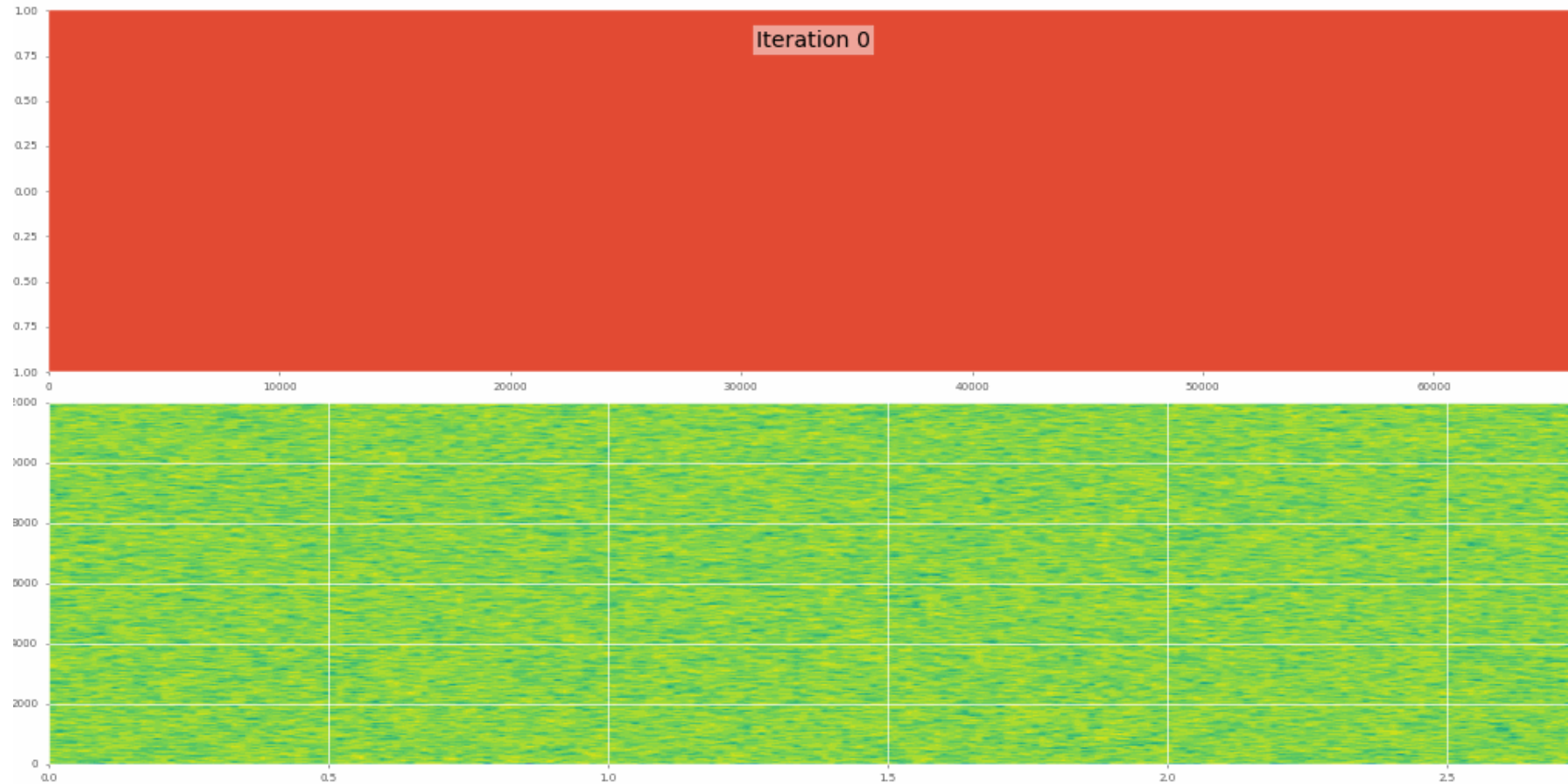
Level 4

Pitch
(128-d)



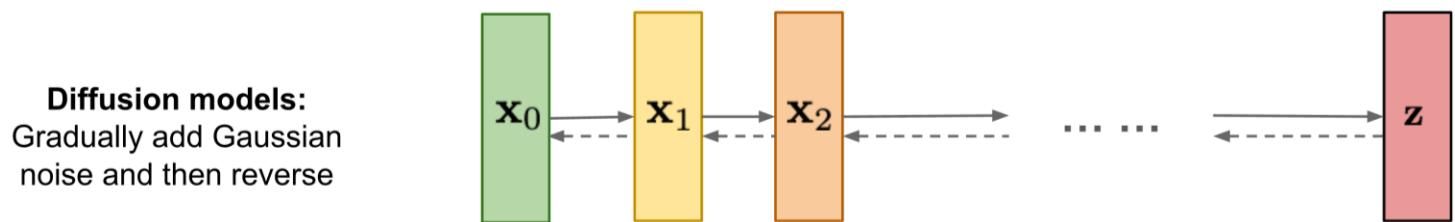
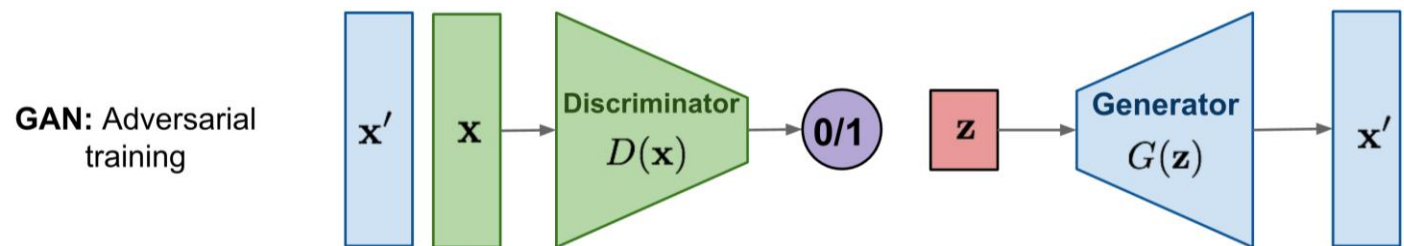
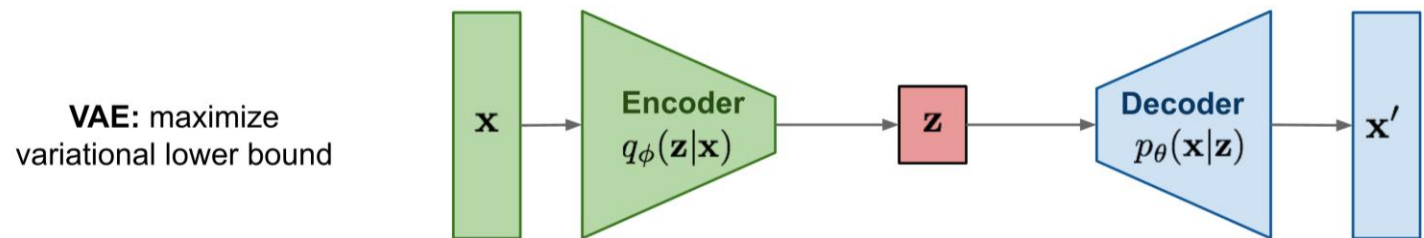
(Source: Wang et al., 2024)

WaveGrad: Diffusion for Waveforms (Chen et al., 2021)



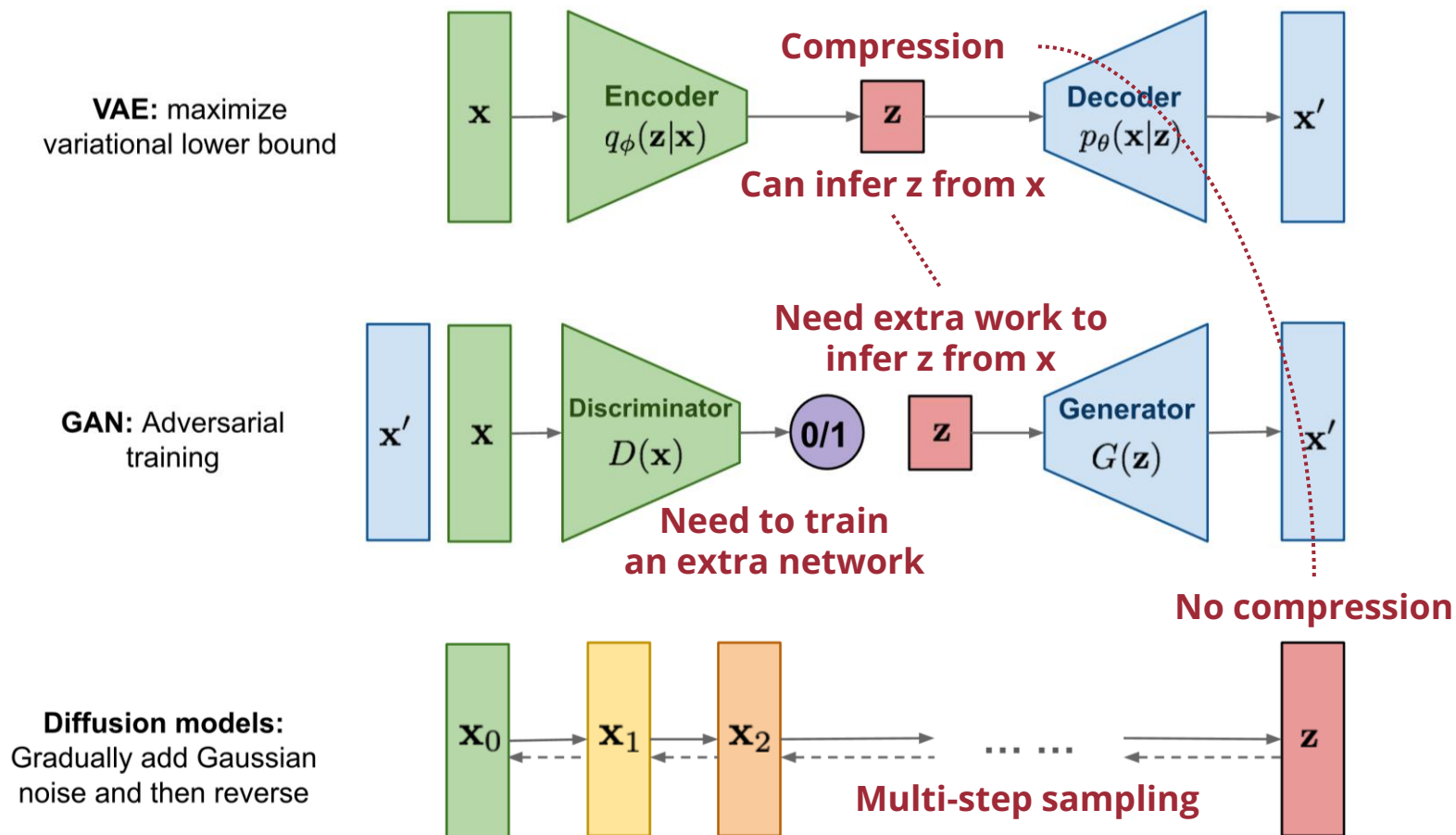
(Source: Chen et al., 2021)

Comparison of Deep Latent Variable Models



(Source: Weng, 2021)

Comparison of Deep Latent Variable Models



(Source: Weng, 2021)

Network Architectures vs. Training Frameworks

Network architectures

Multilayer perceptron (MLP)

Convolutional neural networks (CNNs)

Recurrent neural networks (RNNs)

Transformers

ResNets

U-Nets

⋮

Training frameworks

Autoregressive

Autoencoders

Variational autoencoders (VAEs)

Generative adversarial networks (GANs)

Diffusion models

Consistency models

⋮

Next Lecture

Latent Diffusion Models

