

# DEEP LEARNING

## Lecture 11: Generative Models

Dr. Yang Lu

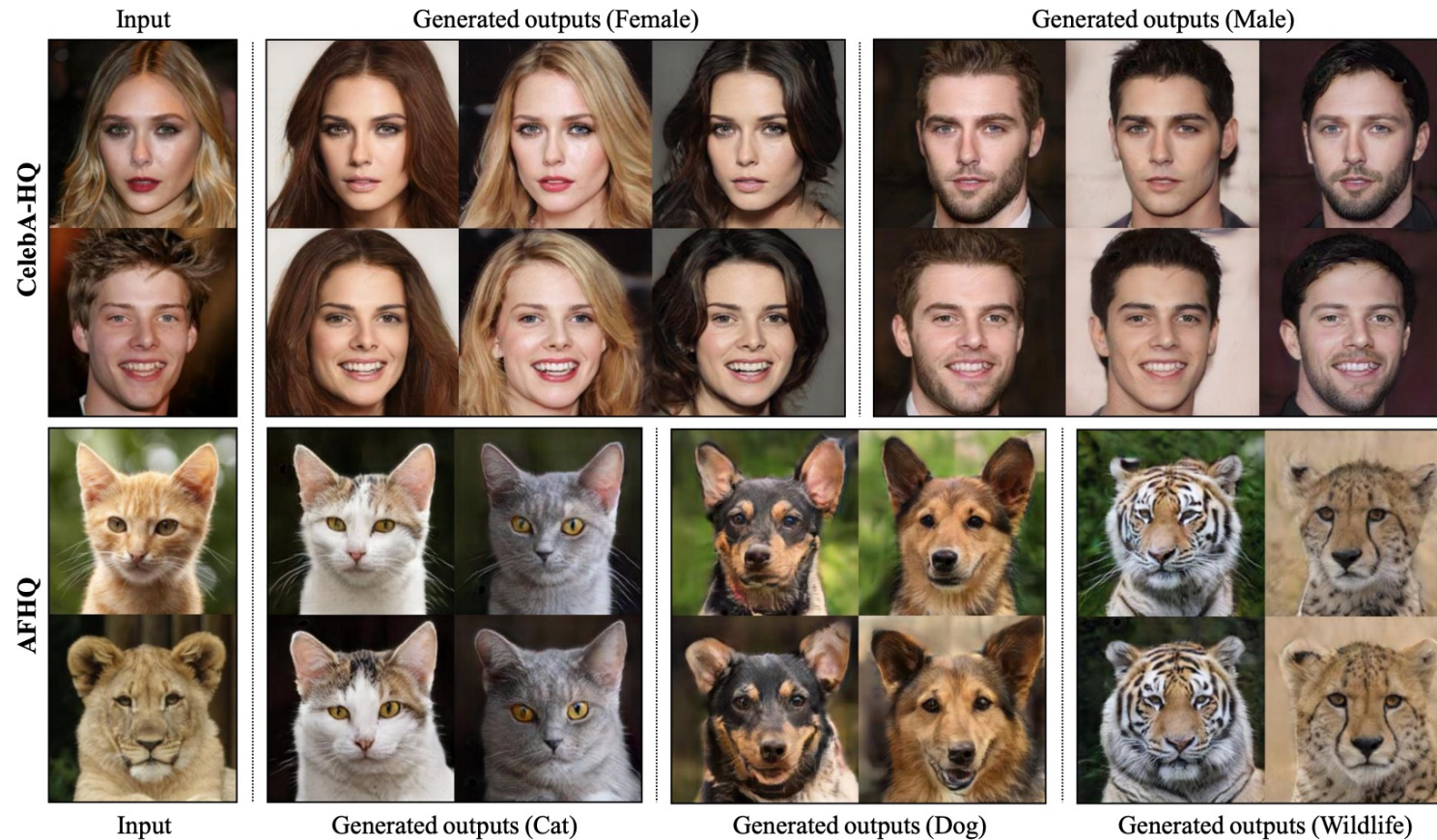
Department of Computer Science and Technology

[luyang@xmu.edu.cn](mailto:luyang@xmu.edu.cn)



# Applications of Generative Models

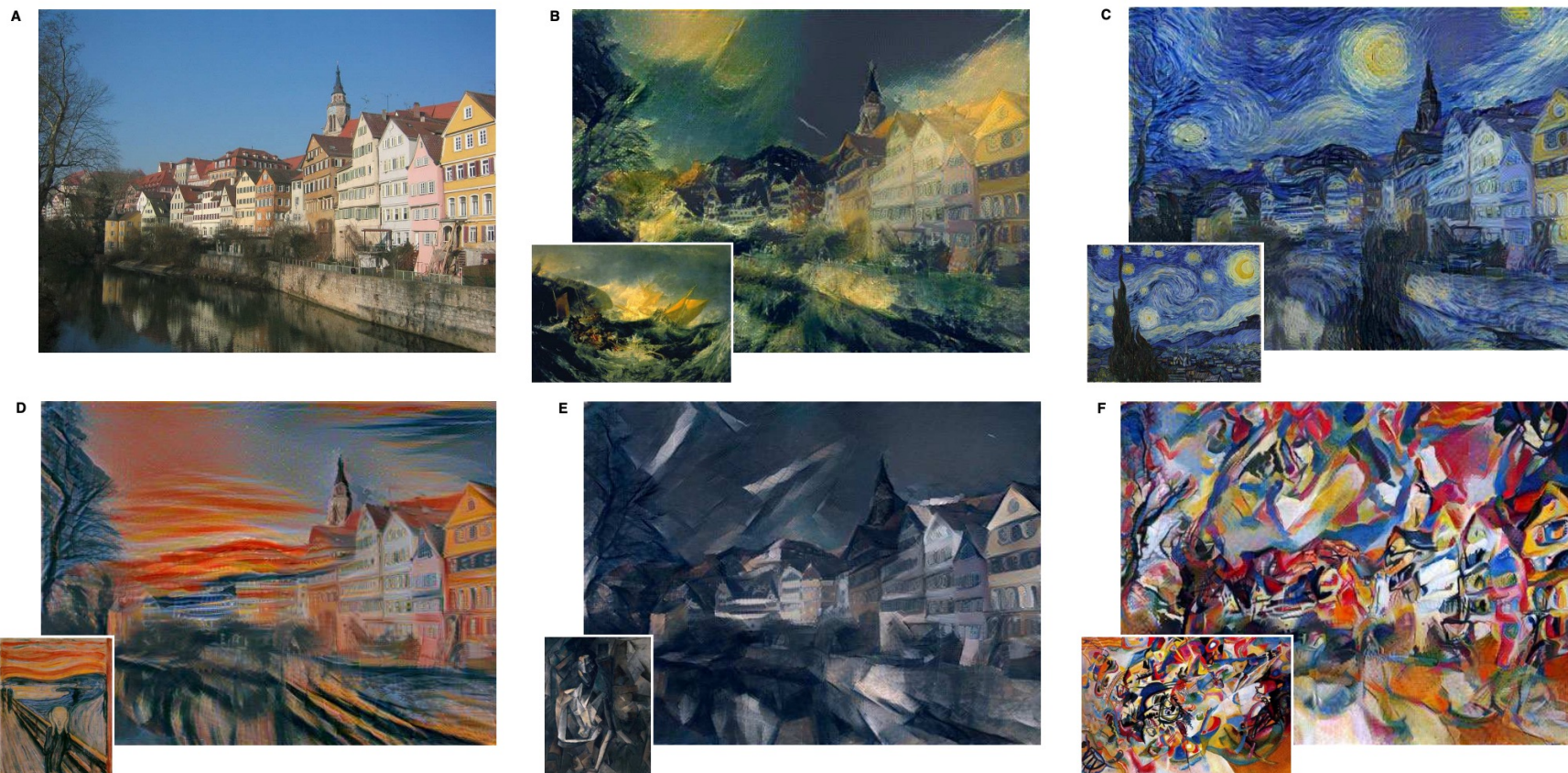
## ■ Image Translation





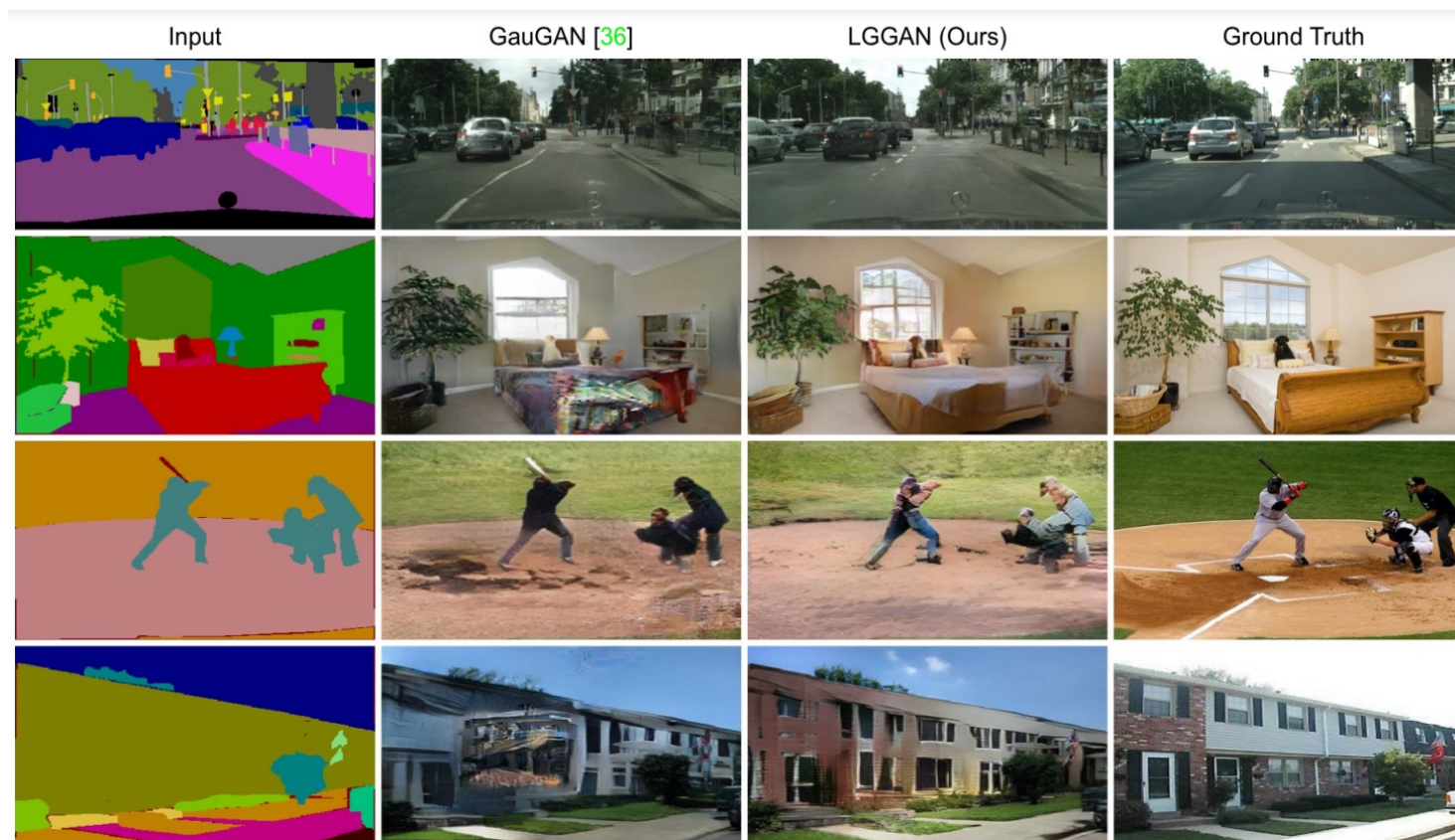
# Applications of Generative Models

## ■ Image Translation



# Applications of Generative Models

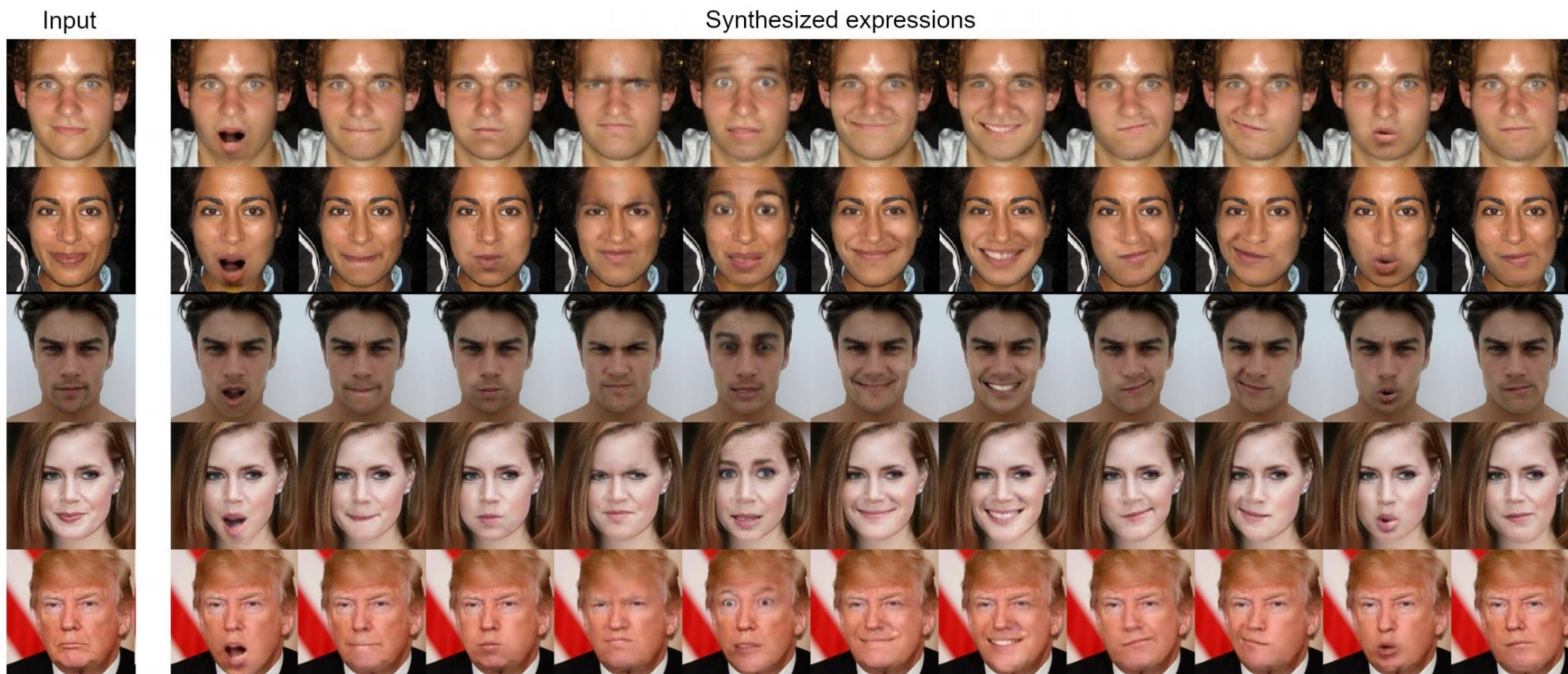
## ■ Scene Generation





# Applications of Generative Models

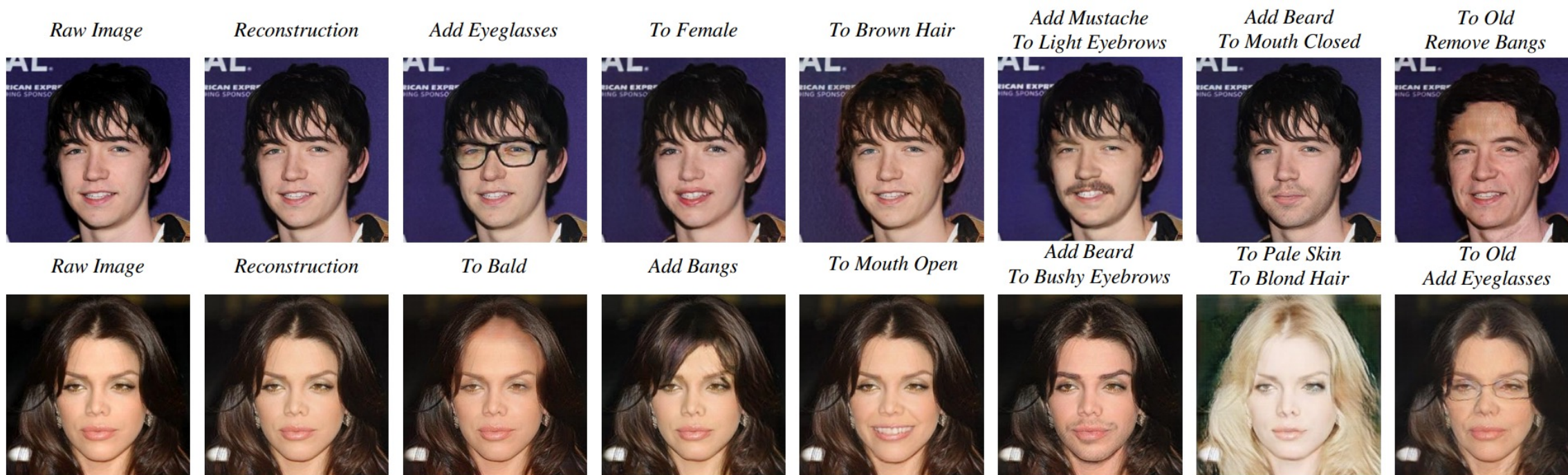
## ■ Facial Attribute Manipulation





# Applications of Generative Models

## ■ Facial Attribute Manipulation





# Applications of Generative Models

## ■ Gaze Correction



# Applications of Generative Models

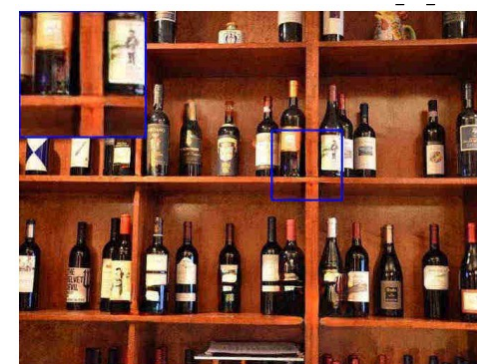
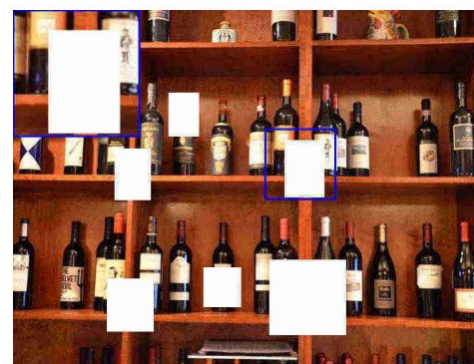
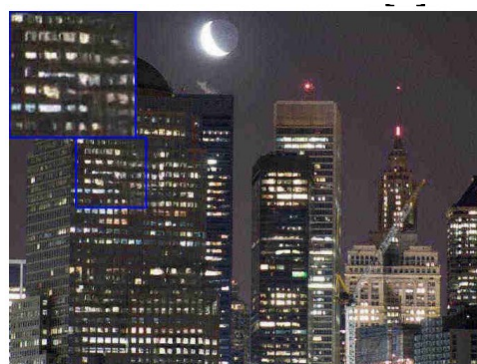
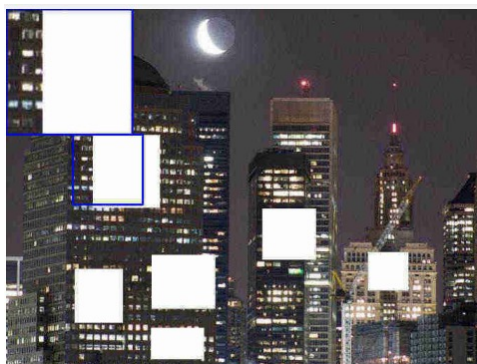
## ■ Image Animation





# Applications of Generative Models

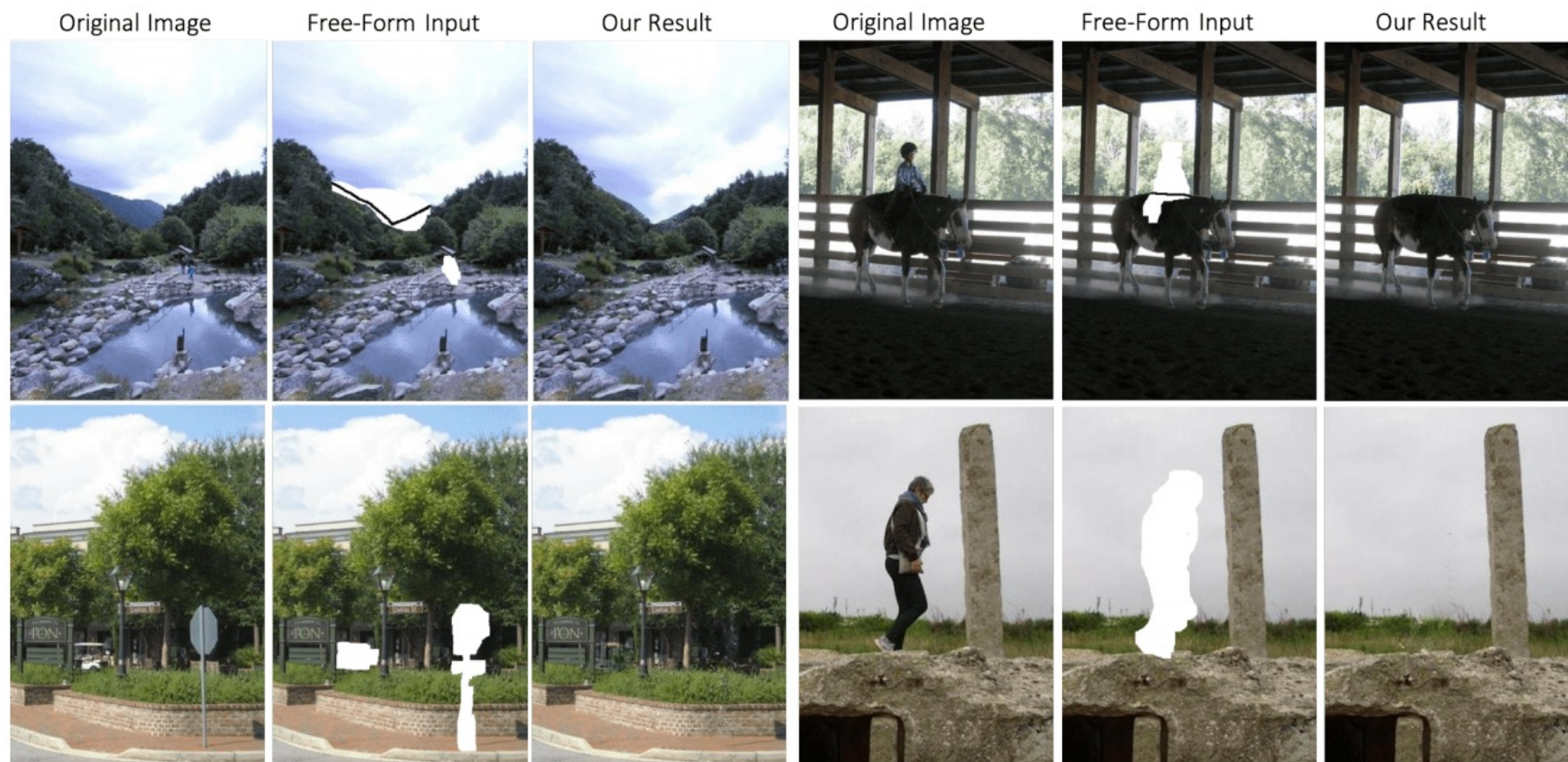
## ■ Image Inpainting





# Applications of Generative Models

## ■ Image Inpainting





# Applications of Generative Models

## ■ Image Blending



# Applications of Generative Models

## ■ Image Super-Resolution



bicubic  
(21.59dB/0.6423)



SRResNet  
(23.53dB/0.7832)



SRGAN  
(21.15dB/0.6868)



original





# Applications of Generative Models

## ■ Makeup



# Applications of Generative Models

## ■ Text-to-image generation



“a hedgehog using a calculator”



“a corgi wearing a red bowtie and a purple party hat”



“robots meditating in a vipassana retreat”



“a fall landscape with a small cottage next to a lake”





# Applications of Generative Models

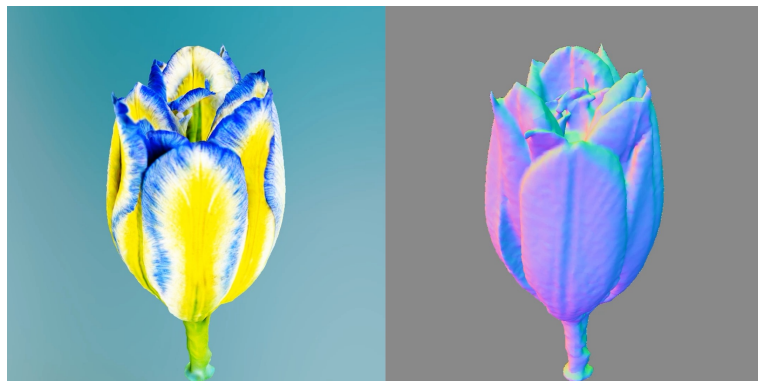
## ■ Video generation

**Text2Video-Zero:  
Text-to-Image Diffusion Models are  
Zero-Shot Video Generators**

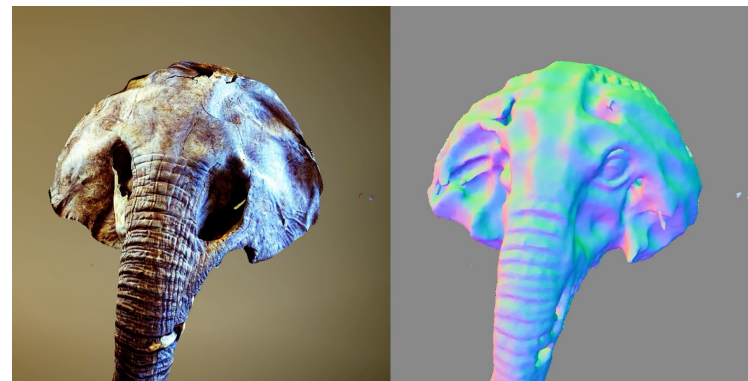


# Applications of Generative Models

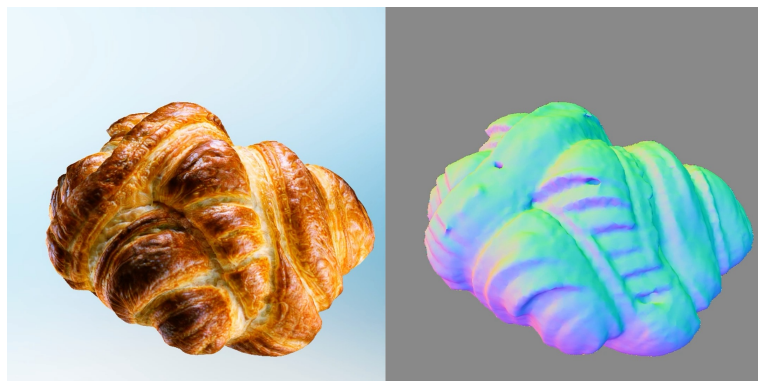
## ■ Text-to-3D generation



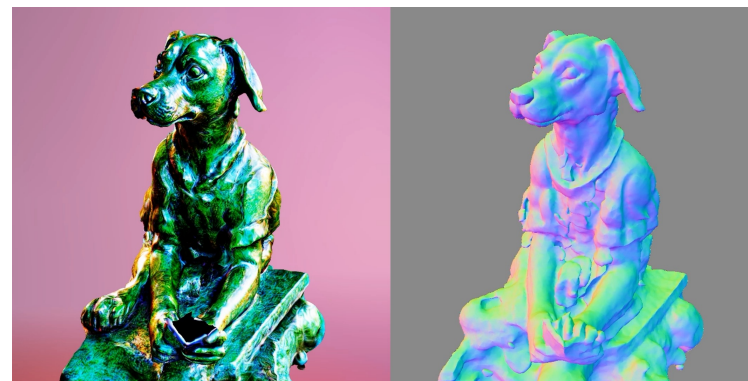
A blue tulip



An elephant skull



A delicious croissant



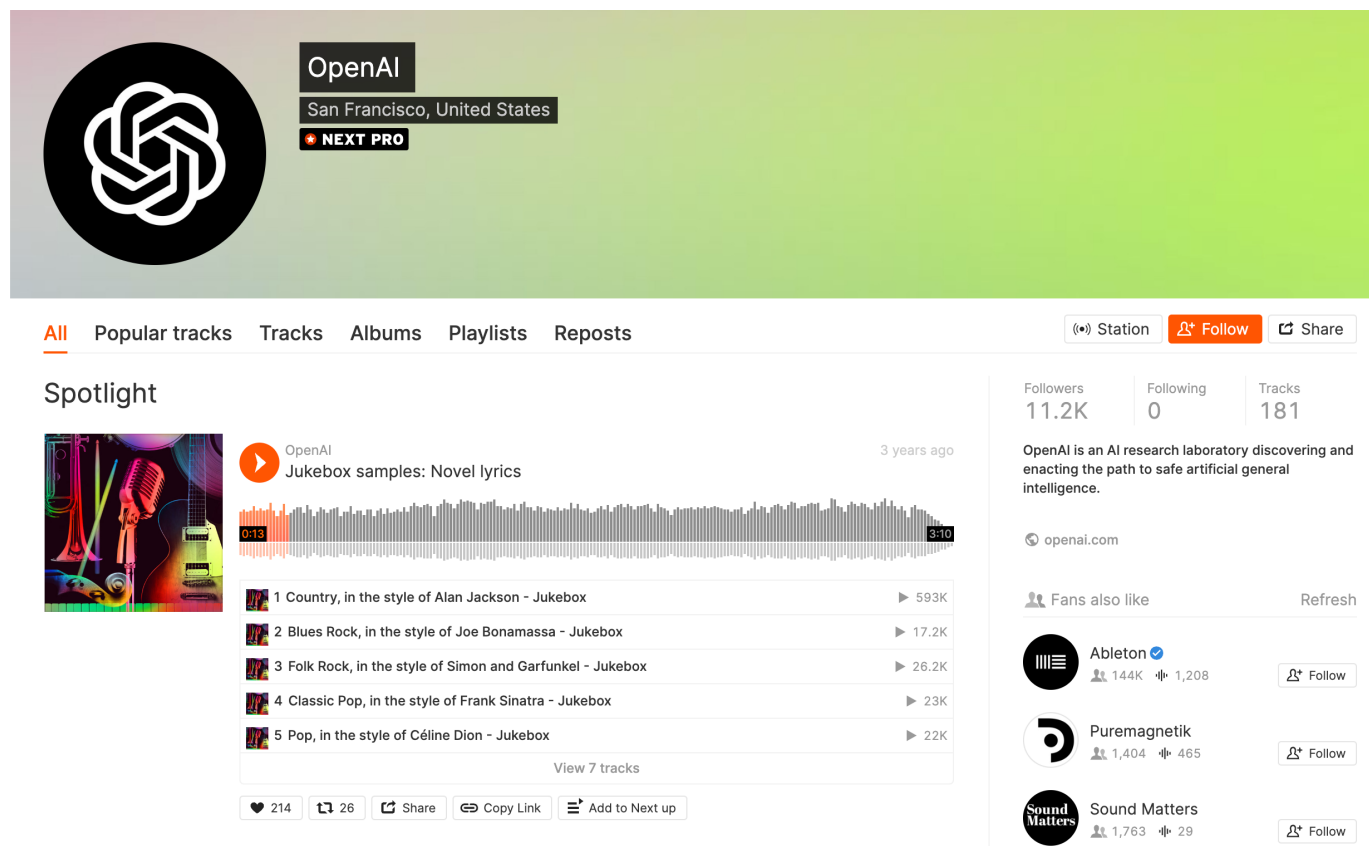
Michelangelo style statue of dog reading news on a cellphone





# Applications of Generative Models

## ■ Music Generation



The screenshot displays the OpenAI profile on SoundCloud. The profile header includes the OpenAI logo, the name "OpenAI", the location "San Francisco, United States", and a "NEXT PRO" badge. Below the header, there are tabs for "All", "Popular tracks", "Tracks", "Albums", "Playlists", and "Reposts". The "All" tab is selected. The main content area features a "Spotlight" section with a track titled "Jukebox samples: Novel lyrics" by OpenAI, uploaded 3 years ago. The track is a 3:10 audio sample with a waveform visualization. Below the track, there is a list of 5 generated tracks in a "Jukebox" format:

- 1 Country, in the style of Alan Jackson - Jukebox (593K)
- 2 Blues Rock, in the style of Joe Bonamassa - Jukebox (17.2K)
- 3 Folk Rock, in the style of Simon and Garfunkel - Jukebox (26.2K)
- 4 Classic Pop, in the style of Frank Sinatra - Jukebox (23K)
- 5 Pop, in the style of Céline Dion - Jukebox (22K)

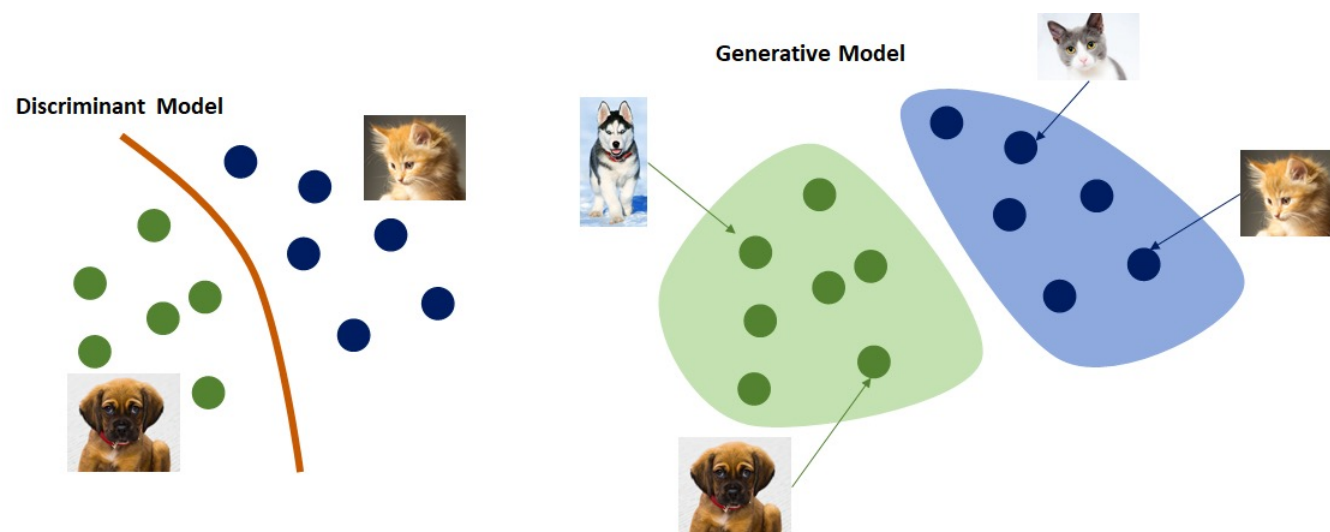
Below the list, there are options to "View 7 tracks", "Like" (214), "Repost" (26), "Share", "Copy Link", and "Add to Next up". On the right side of the profile, there are statistics: "Followers 11.2K", "Following 0", and "Tracks 181". A bio states: "OpenAI is an AI research laboratory discovering and enacting the path to safe artificial general intelligence." Below the bio, there is a link to "openai.com" and a section titled "Fans also like" with a "Refresh" button. The "Fans also like" section lists three accounts: Ableton (144K followers, 1,208 tracks), Puremagnetik (1,404 followers, 465 tracks), and Sound Matters (1,763 followers, 29 tracks).

Generated music samples: [https://soundcloud.com/openai\\_audio](https://soundcloud.com/openai_audio)



# Generative vs. Discriminative

- In machine learning, two main approaches are called the **generative approaches** and the **discriminative approaches**.
- Given an observable variable  $X$  and a target variable  $Y$ :
  - A generative model is a statistical model of the data distribution  $P(X)$  or the joint probability distribution on  $X \times Y$ :  $P(X, Y)$ .
  - A discriminative model is a model of the conditional distribution of  $Y$  given  $X$ :  $P(Y|X = x)$ .





# Discriminative Approaches

- Most supervised learning methods fall into discriminative approaches.
- Given data:  $(x, y)$ ,  $x$  is data,  $y$  is label.
- Goal: Learn a function to map  $x \rightarrow y$ , namely posterior probability  $P(Y|X = x)$ .
- Examples: Classification, regression, object detection, face recognition, sentiment classification, etc.



→ cat

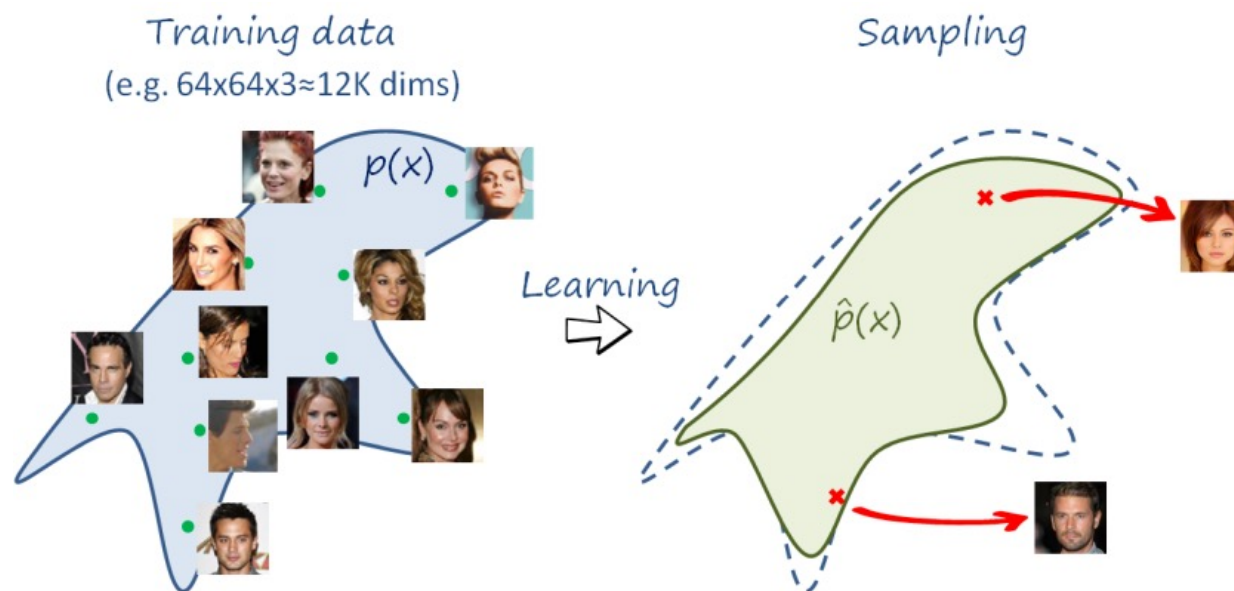


→ car



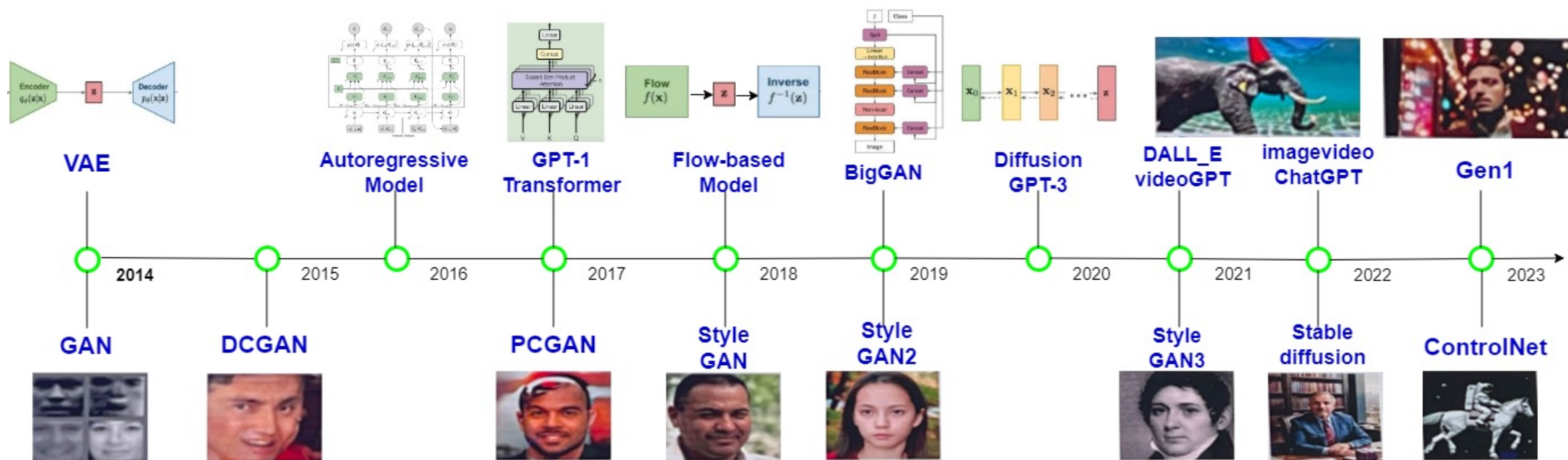
# Generative Approaches

- Given training data, generate new samples from same distribution.
- Objectives:
  1. Learn  $p_{model}(x)$  that approximates  $p_{data}(x)$ .
  2. Sample a new  $x$  from  $p_{model}(x)$ .





# Generative Models



# Outlines

- Variational Autoencoder
- Generative Adversarial Nets
- Diffusion Model



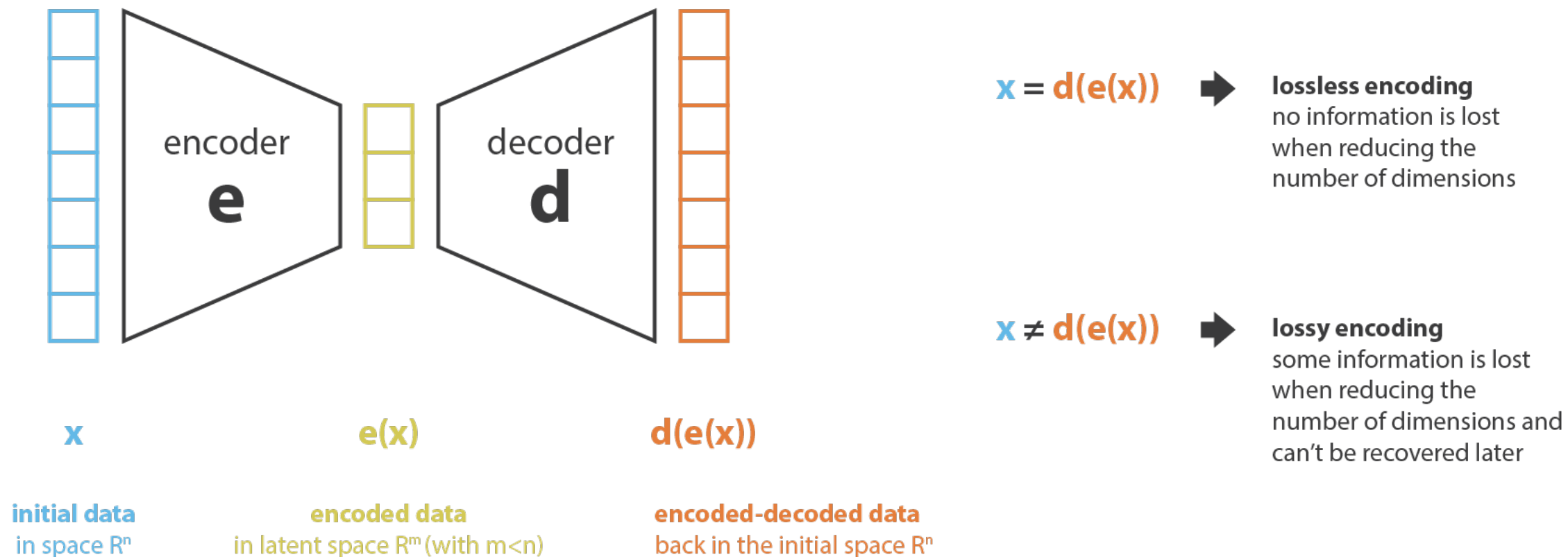




# VARIATIONAL AUTOENCODER



# Autoencoder



$$\text{loss} = \|x - \hat{x}\|^2 = \|x - d(z)\|^2 = \|x - d(e(x))\|^2$$





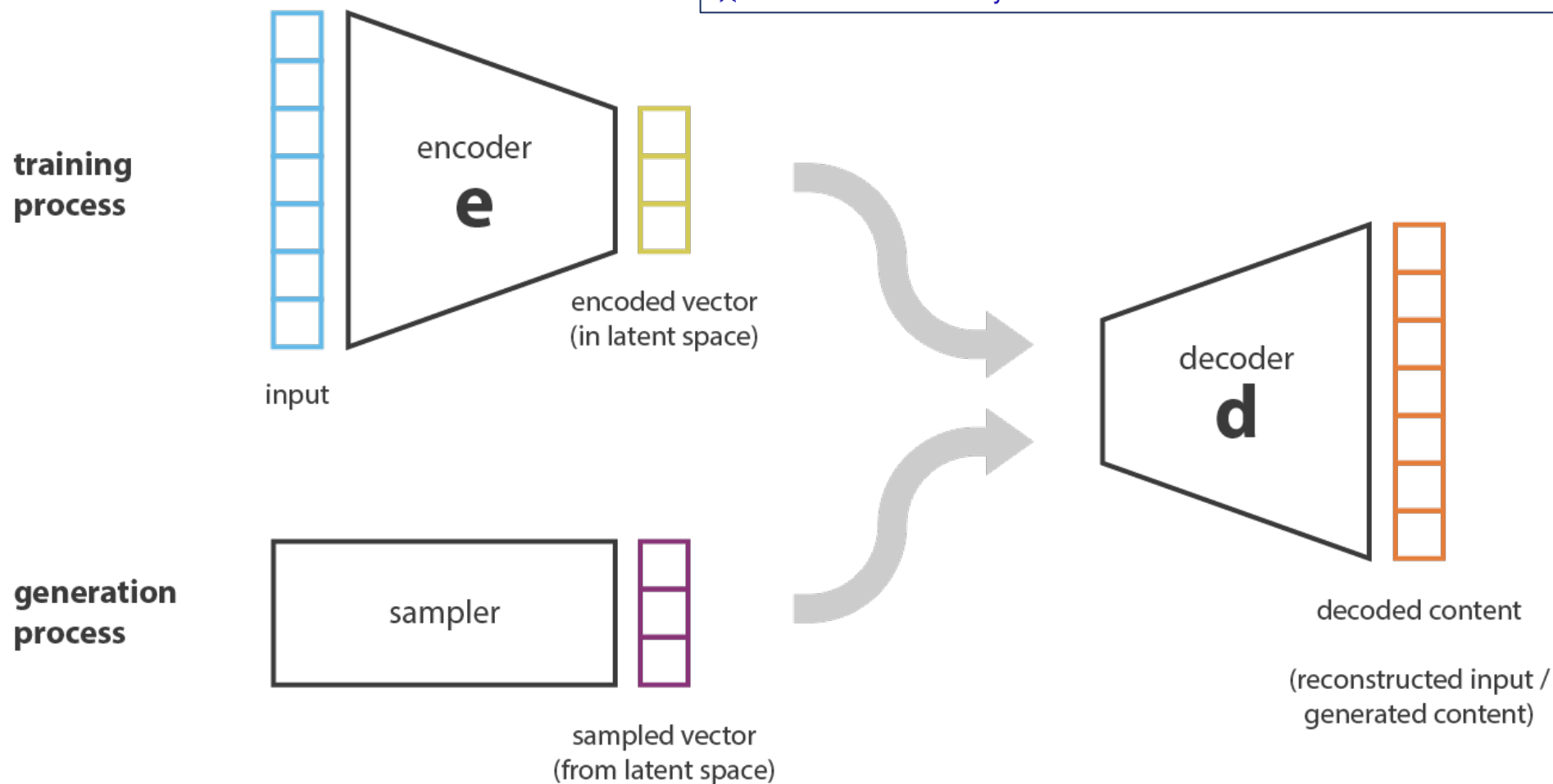
# Variational Autoencoder

## Auto-encoding variational bayes

[DP Kingma](#), [M Welling](#) - arXiv preprint arXiv:1312.6114, 2013 - [arxiv.org](#)

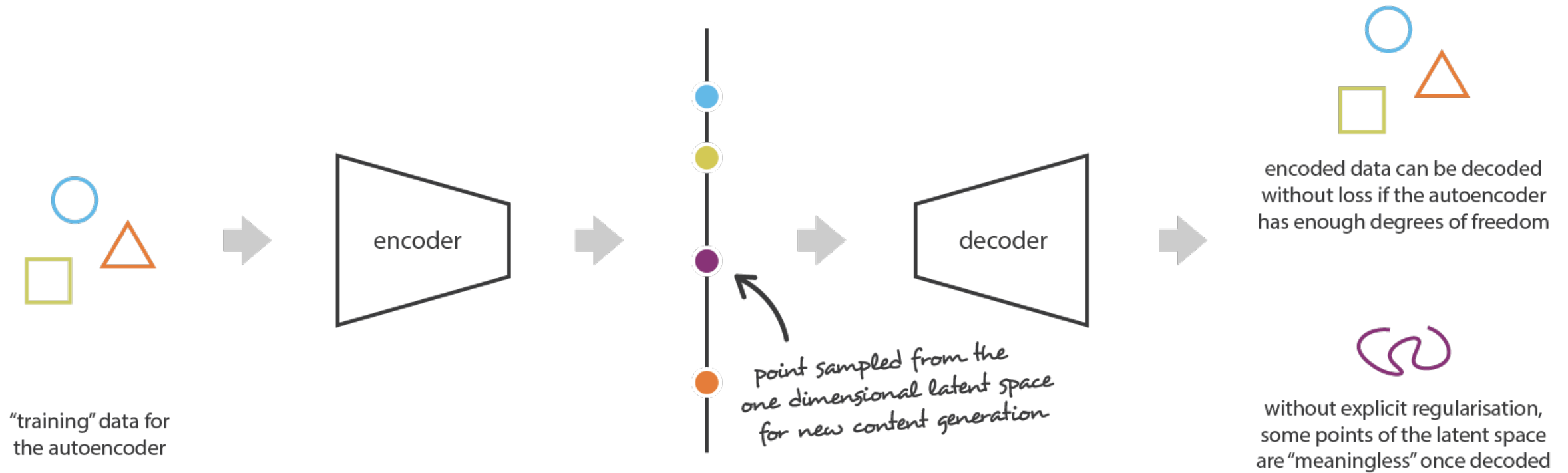
... **variational** lower bound yields a simple differentiable unbiased estimator of the lower bound; this SGVB (Stochastic Gradient **Variational Bayes**... , we propose the **AutoEncoding** VB (AEVB...

☆ Save ↀ Cite Cited by 30917 Related articles All 44 versions ↀ



# Variational Autoencoder

- The autoencoder is solely **trained to encode and decode** with as few loss as possible, no matter how the latent space is organized.





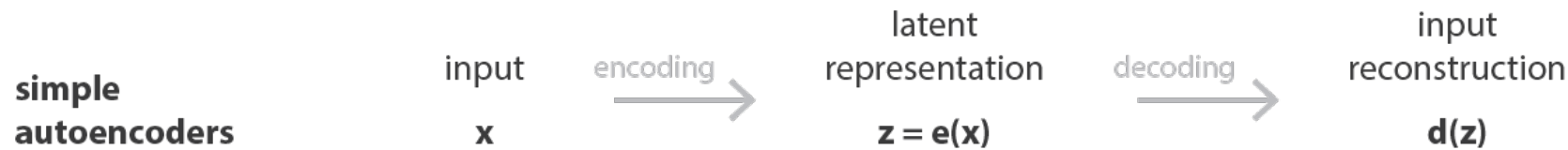
# Variational Autoencoder

- A variational autoencoder can be defined as being an autoencoder whose training is:
  - regularized to avoid overfitting;
  - ensure that the latent space has good properties that enable generative process.



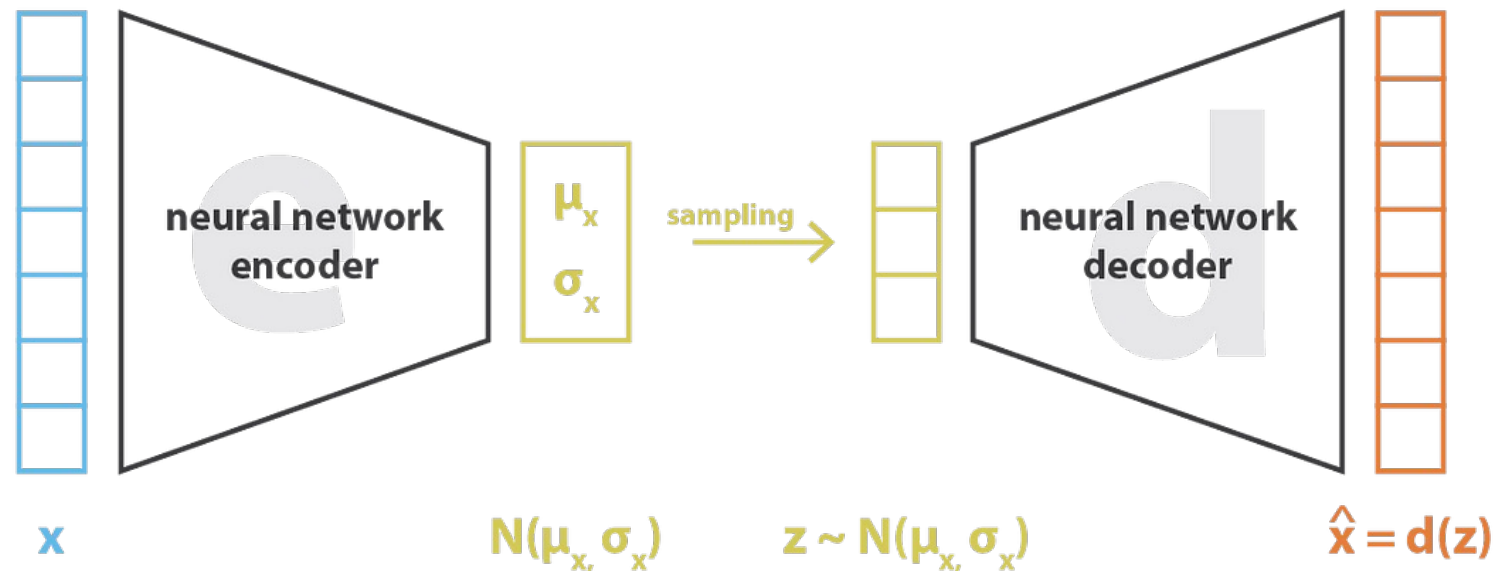
# Variational Autoencoder

- Instead of encoding an input as a single point, we encode it as a **distribution over the latent space**.



# Variational Autoencoder

- Regularize the organization of the latent space by making the distributions returned by the encoder **close to a standard normal distribution**.

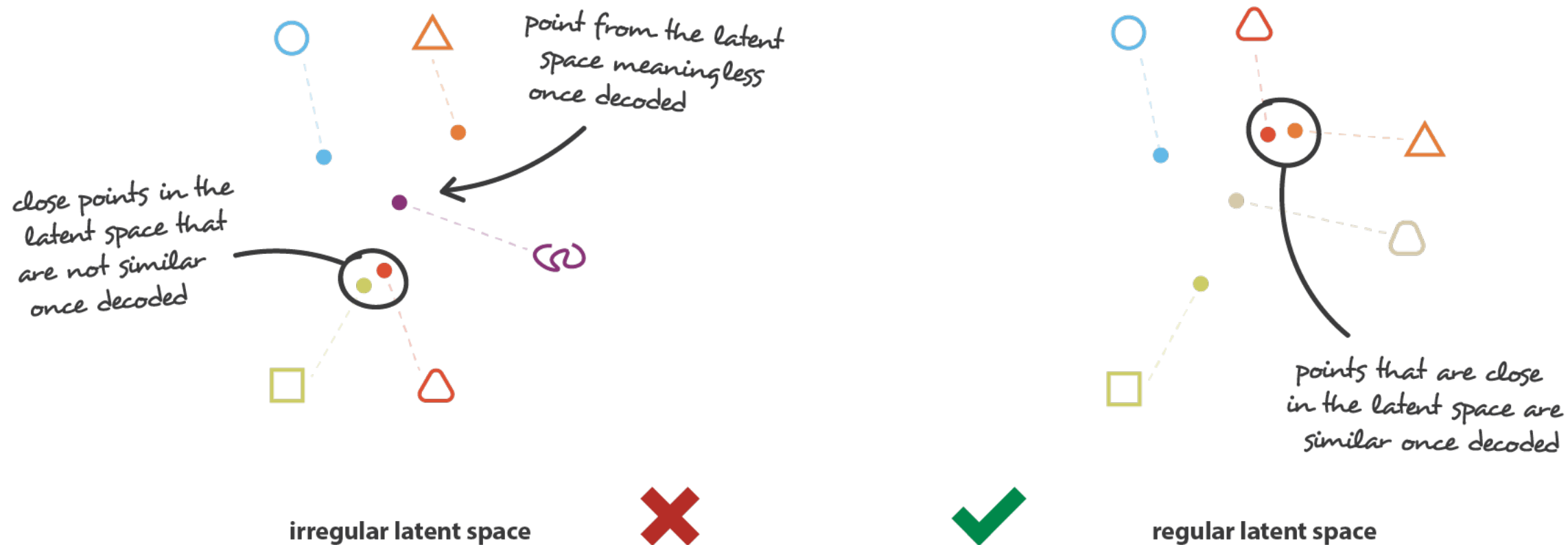


$$\text{loss} = ||x - \hat{x}||^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)] = ||x - d(z)||^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)]$$



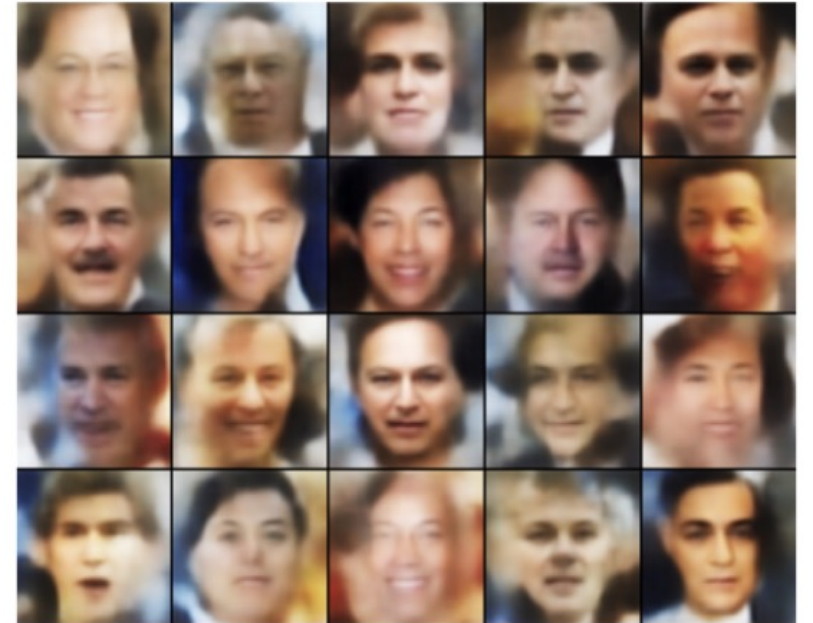


# Variational Autoencoder



# Variational Autoencoder

- The loss of the autoencoder is to minimize both:
  - the **reconstruction loss** (how similar the autoencoder's output to its input);
  - the **latent loss** (how close its hidden nodes were to a normal distribution).
- It doesn't guarantee **the quality of the generated image**. A major drawback of VAEs is the blurry outputs that they generate. VAE models tend to produce unrealistic, blurry samples.



Human faces  
generated by VAEs





GAN



# GAN

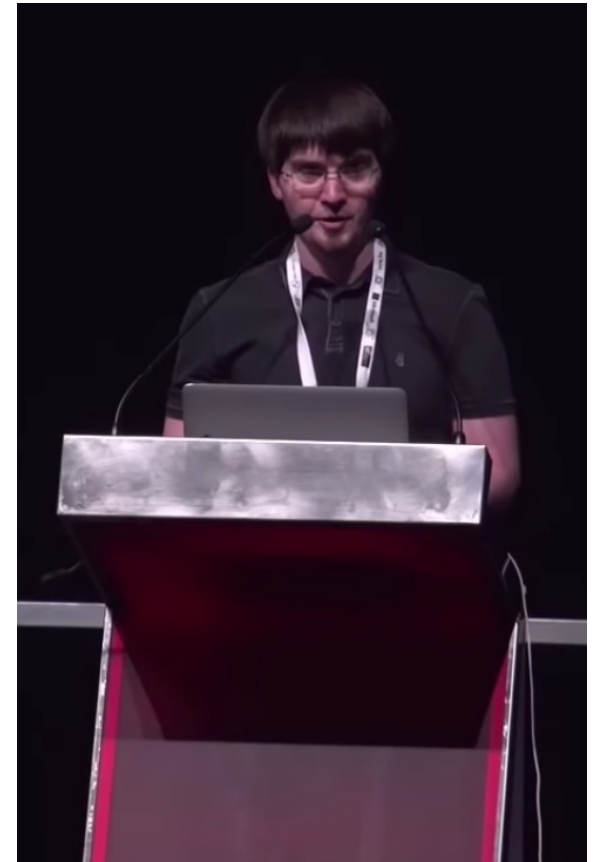
## Generative adversarial nets

[I Goodfellow, J Pouget-Abadie...](#) - Advances in neural ..., 2014 - proceedings.neurips.cc

... We propose a new framework for estimating **generative** models via **adversarial nets**, in which we simultaneously train two models: a **generative** model G that captures the data ...

☆ Save 📄 Cite Cited by 61677 Related articles All 61 versions 🔗

- GAN was proposed by Ian Goodfellow in 2014.
- Yann LeCun described GANs as “*the most interesting idea in the last 10 years in machine learning*”.
- Ian presented and explained his paper in NIPS 2016 with a 2-hour presentation.

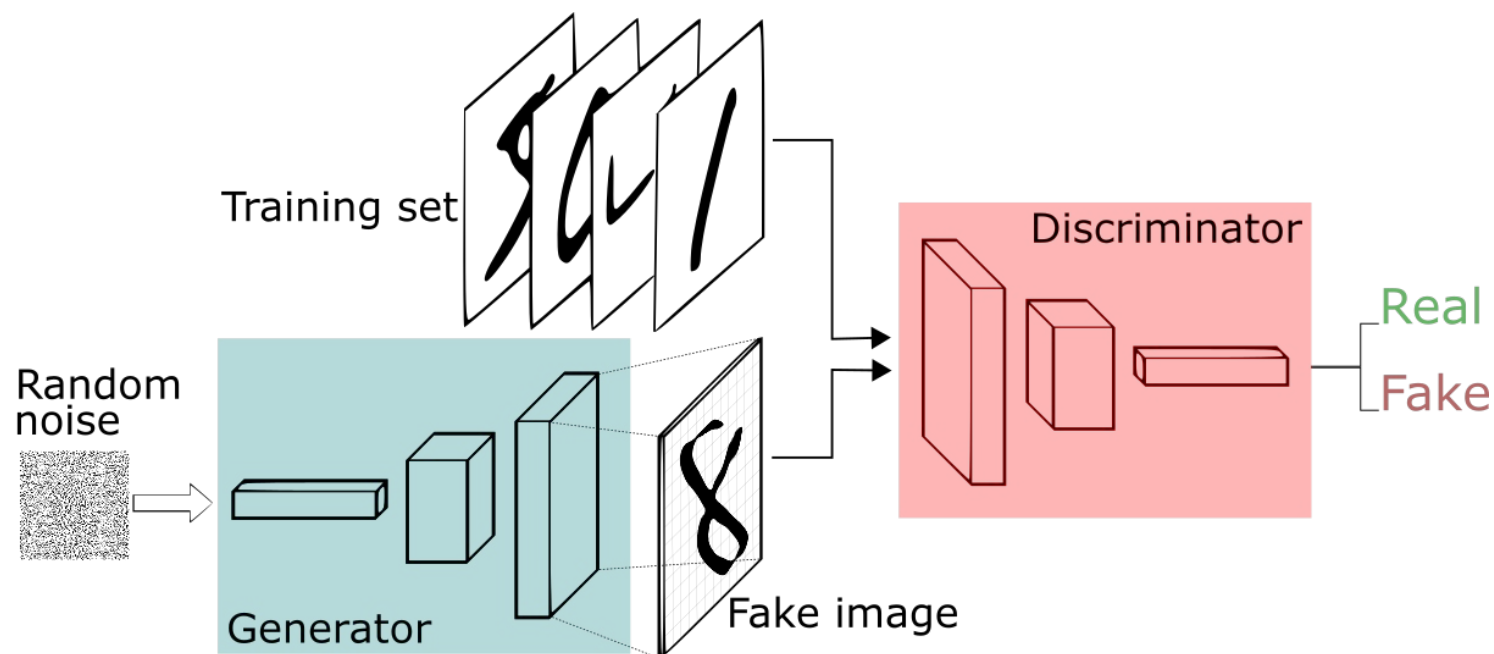


Ian in NIPS 2016



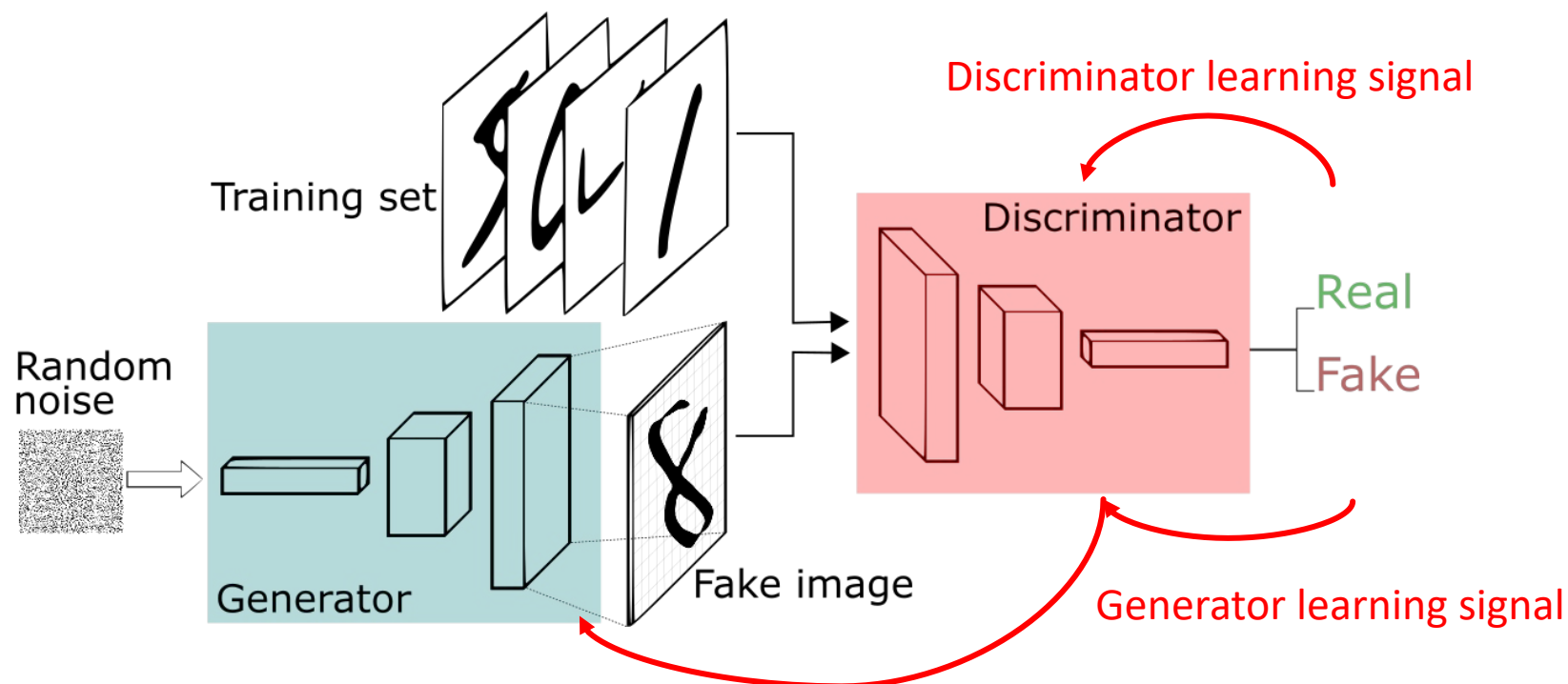
# GAN: How to Do

- GAN is composed by a **generator** and a **discriminator**. They are both neural networks.
- **Generator network**: try to fool the discriminator by generating real-looking images.
- **Discriminator network**: try to distinguish between real and fake images.



# GAN: How to Do

- Generator and discriminator tells each other where it was wrong.
  - Generator tells discriminator how I fool you.
  - Discriminator tells generator how I detect you.





# GAN: How to Learn

- Given a prior on input noise variables  $\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})$ , generator  $G_{\theta_g}(\mathbf{z})$  take  $\mathbf{z}$  as input and map it into data space.
- Discriminator  $D_{\theta_d}(\mathbf{x})$  takes the data  $\mathbf{x} \sim p_{data}$  as input and output probability that  $\mathbf{x}$  came from the real data rather than generated data  $G_{\theta_g}(\mathbf{z})$ .
- $D_{\theta_d}$  and  $G_{\theta_g}$  have different goals (1 for real, 0 for fake):
  - Generator wants:  $D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \rightarrow 1$ .
  - Discriminator wants:  $D_{\theta_d}(\mathbf{x}) \rightarrow 1, D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \rightarrow 0$ .



# GAN: How to Learn

- By maximizing the log-likelihood, the overall objective is to simultaneously train over all  $x$  with random generated  $z$ :

- train  $G_{\theta_g}$  to minimize  $\log \left( 1 - D_{\theta_d} \left( G_{\theta_g}(z) \right) \right)$ ;

- train  $D_{\theta_d}$  to maximize  $\log D_{\theta_d}(x)$  and  $\log \left( 1 - D_{\theta_d} \left( G_{\theta_g}(z) \right) \right)$ .

- In other words,  $D_{\theta_d}$  and  $G_{\theta_g}$  play the following two-player minimax game:

$$\min_{\theta_g} \max_{\theta_d} \left[ \underbrace{\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x)}_{\text{Discriminator output for real data } x} + \underbrace{\mathbb{E}_{z \sim p_z(z)} \log \left( 1 - D_{\theta_d} \left( G_{\theta_g}(z) \right) \right)}_{\text{Discriminator output for generated fake data } G_{\theta_g}(z)} \right]$$



# GAN: How to Learn

Objective:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{\mathbf{x} \sim p_{data}} \log D_{\theta_d}(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} \log \left( 1 - D_{\theta_d} \left( G_{\theta_g}(\mathbf{z}) \right) \right) \right]$$

Alternate between:

- Gradient ascent on discriminator:

$$\max_{\theta_d} \left[ \mathbb{E}_{\mathbf{x} \sim p_{data}} \log D_{\theta_d}(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} \log \left( 1 - D_{\theta_d} \left( G_{\theta_g}(\mathbf{z}) \right) \right) \right]$$

- Gradient descent on generator:

$$\min_{\theta_g} \left[ \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} \log \left( 1 - D_{\theta_d} \left( G_{\theta_g}(\mathbf{z}) \right) \right) \right]$$

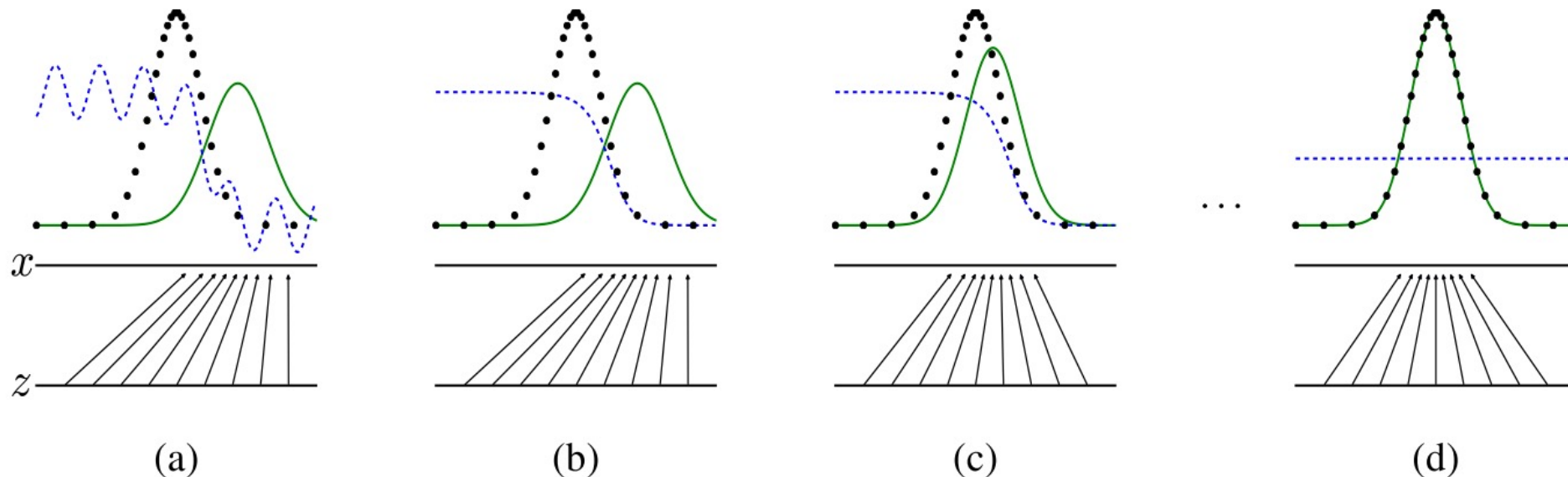




# GAN: How to Learn

- After several steps of training, if  $G$  and  $D$  have enough capacity, they will reach a point at which both cannot improve because  $p_g = p_{data}$ .
- Generator can generate real image.
- Discriminator is unable to differentiate between the two distributions, i.e.  $D_{\theta_d}(x) = 1/2$ .

----- discriminative distribution  
..... data distribution  
———— generative distribution



# GAN: Algorithm

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator,  $k$ , is a hyperparameter. We used  $k = 1$ , the least expensive option, in our experiments.

**for** number of training iterations **do**

**for**  $k$  steps **do**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Sample minibatch of  $m$  examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution  $p_{\text{data}}(x)$ .
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right].$$

**end for**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Update the generator by descending its stochastic gradient:

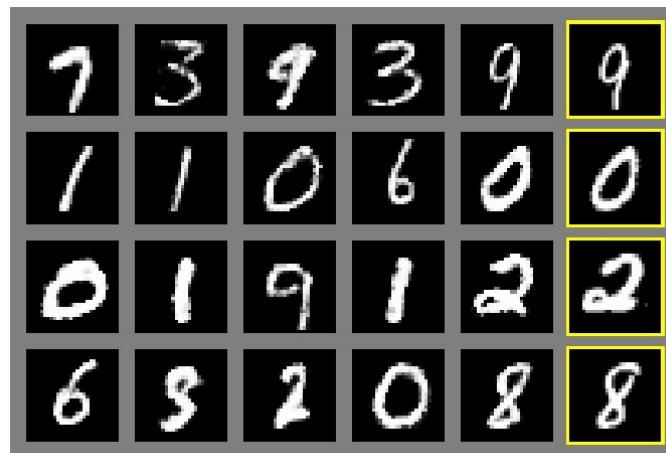
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))).$$

**end for**

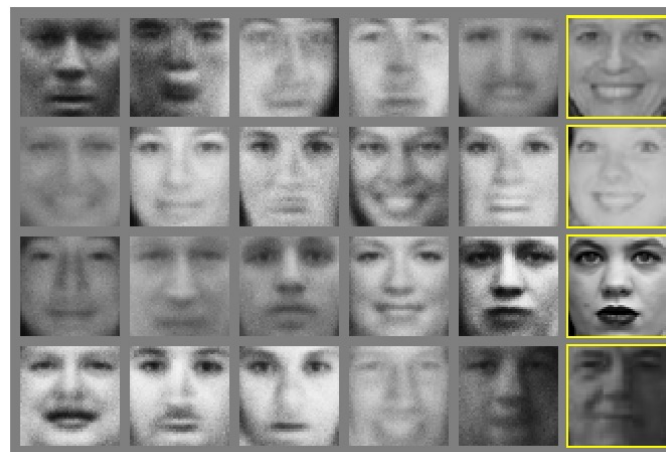
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.



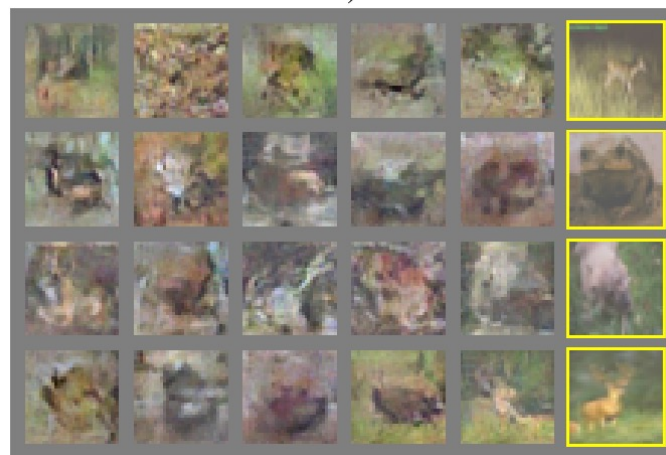
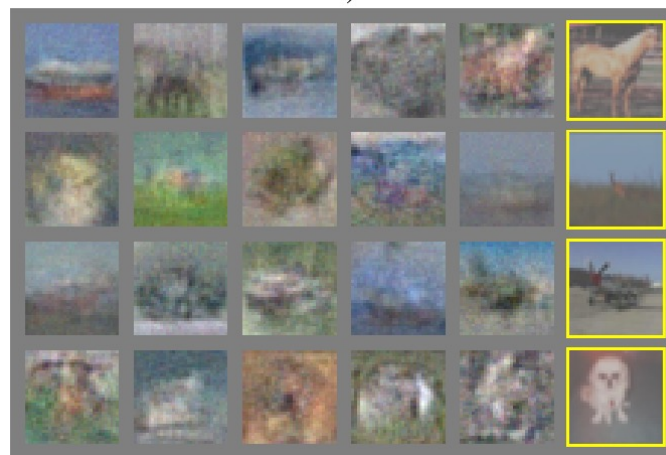
# GAN: Result



a)



b)



Rightmost column shows the nearest training example of the neighboring sample, in order to demonstrate that the model has not memorized the training set.



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)

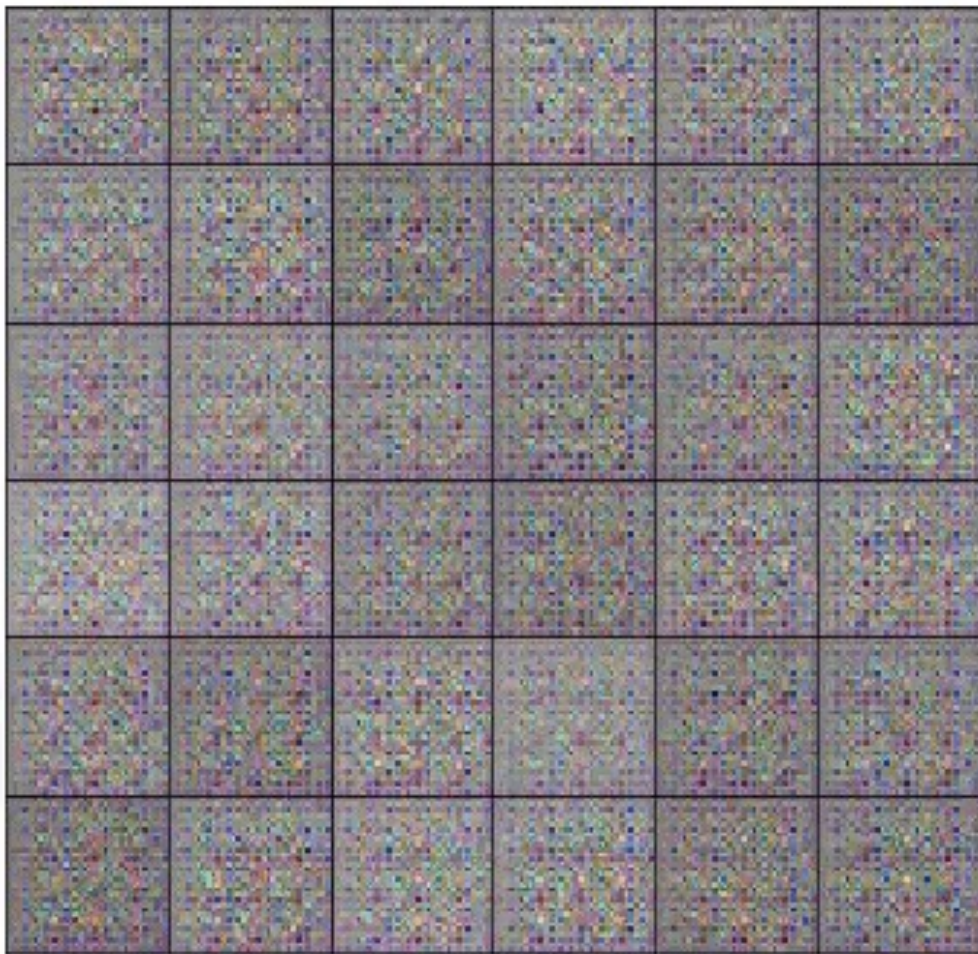


廈門大學 计算机科学与技术系

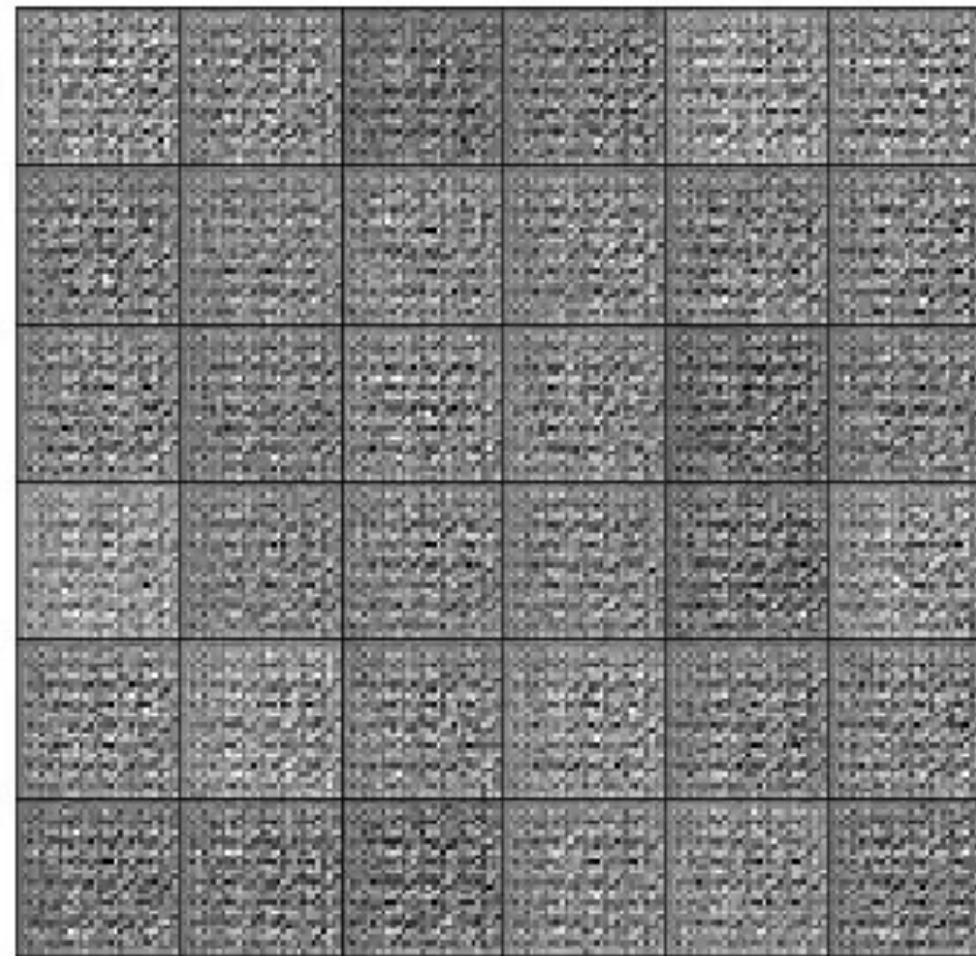
Department of Computer Science and Technology, Xiamen University



# GAN: Result



SVHNs



MNIST





# GAN Starts an Era



2014



2015



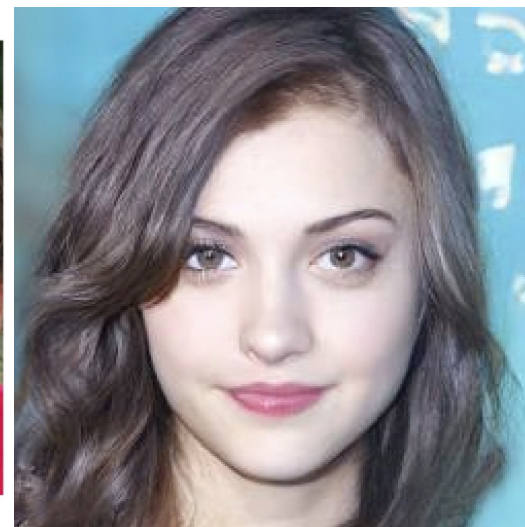
2016



2017



2018



2020



2023



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)



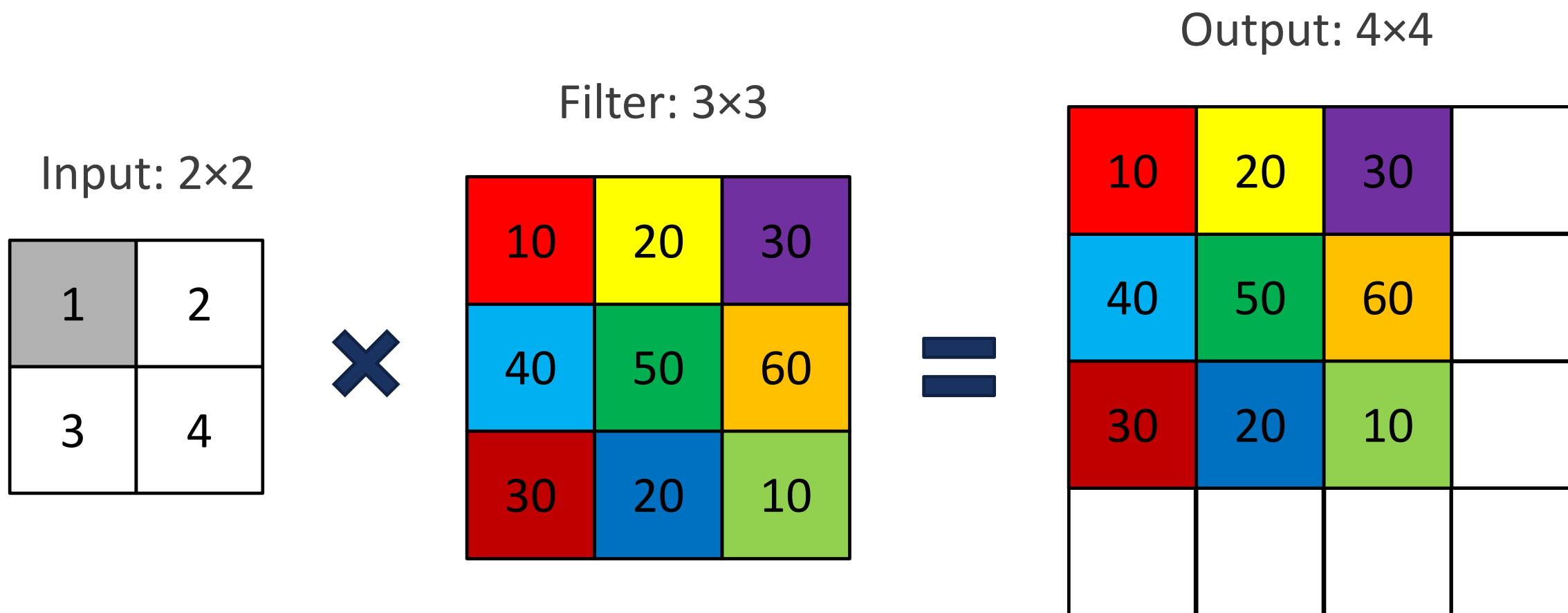
廈門大學 计算机科学与技术系

Department of Computer Science and Technology, Xiamen University

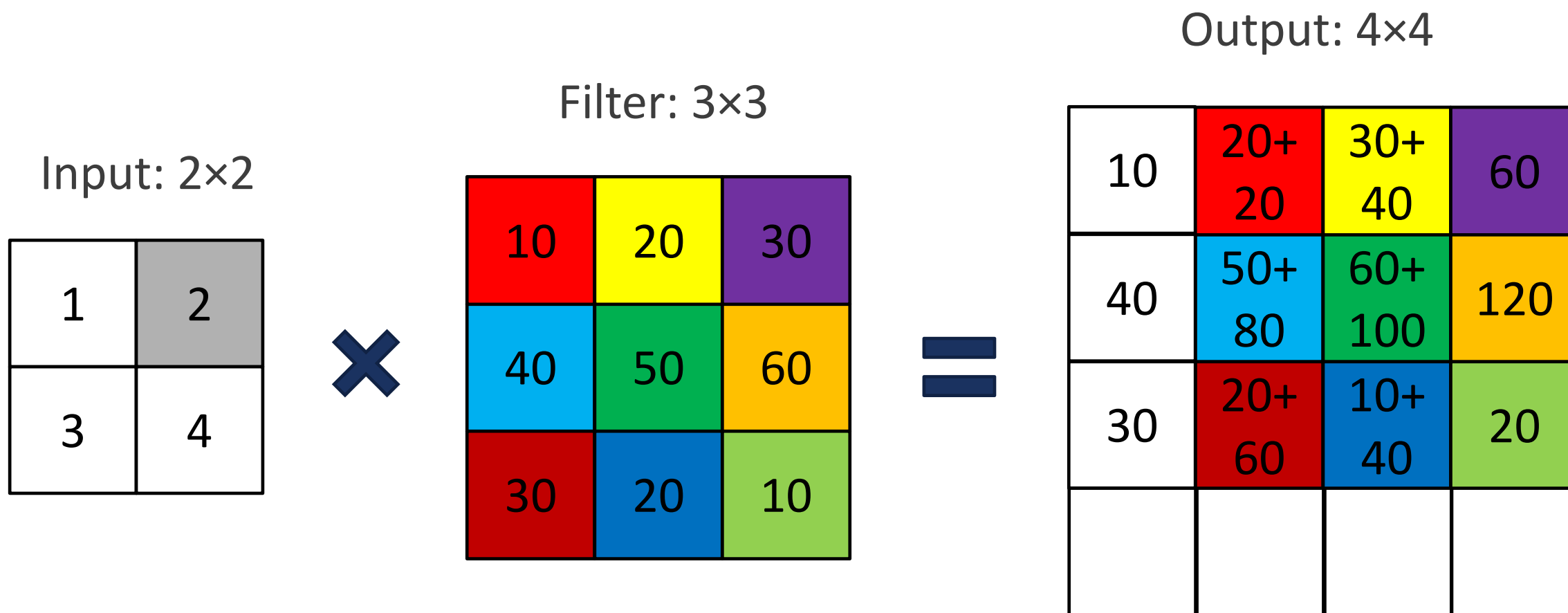
Image source: Ian Goodfellow. Samples from Goodfellow et al., 2014, Radford et al., 2015, Liu et al., 2016, Karras et al., 2017, Karras et al., 2018, Ho et al. 2020, Zhang et al. 2023

- Vanilla GAN simply uses MLP, rather than CNN in both generator and discriminator.
- CNN can be easily applied to discriminator.
- Now the problem is: how can CNN be used as a generator?
  - Pooling leads to downsampling, how to upsampling?

# Fractionally-Strided Convolutions

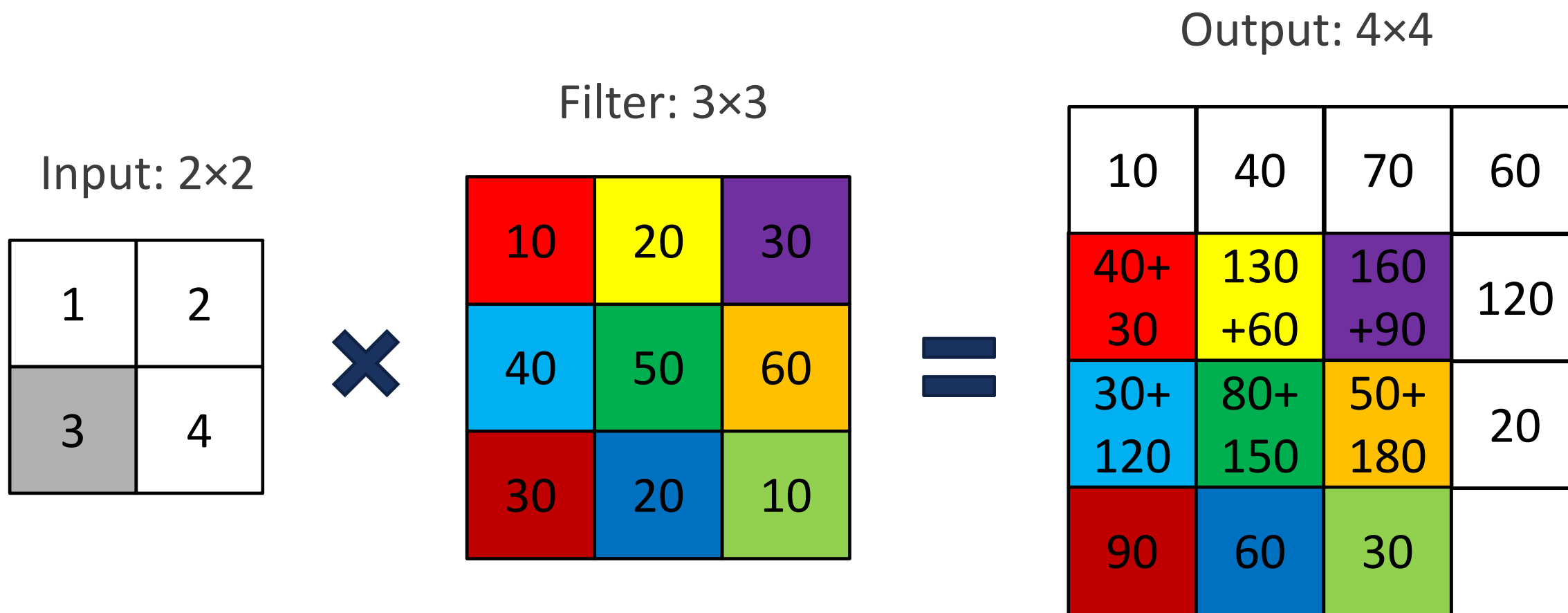


# Fractionally-Strided Convolutions

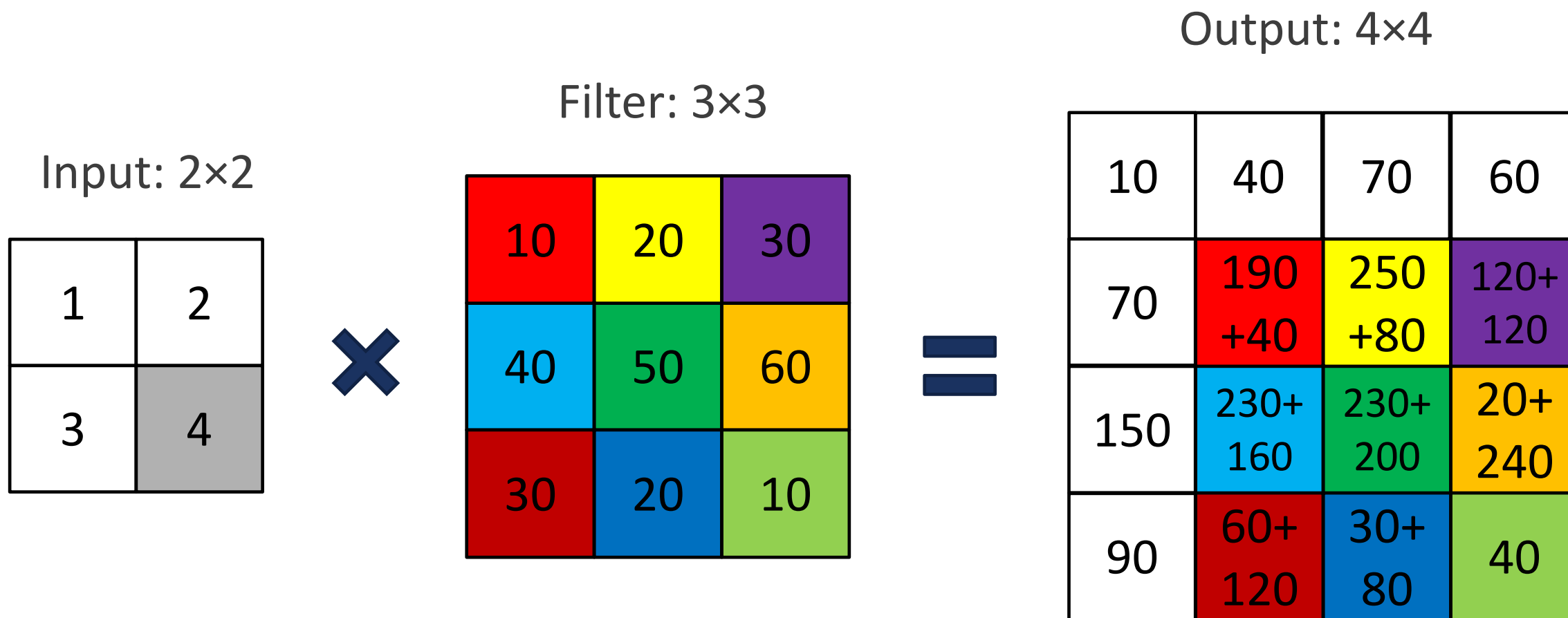




# Fractionally-Strided Convolutions

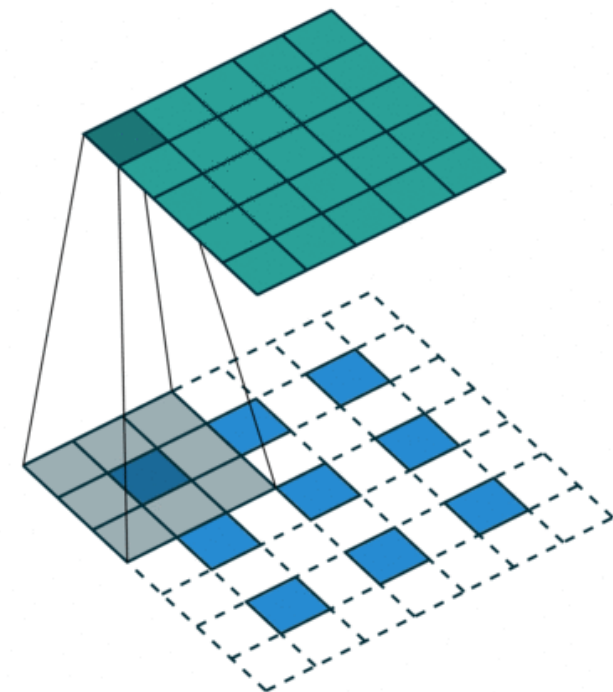


# Fractionally-Strided Convolutions



# Fractionally-Strided Convolutions

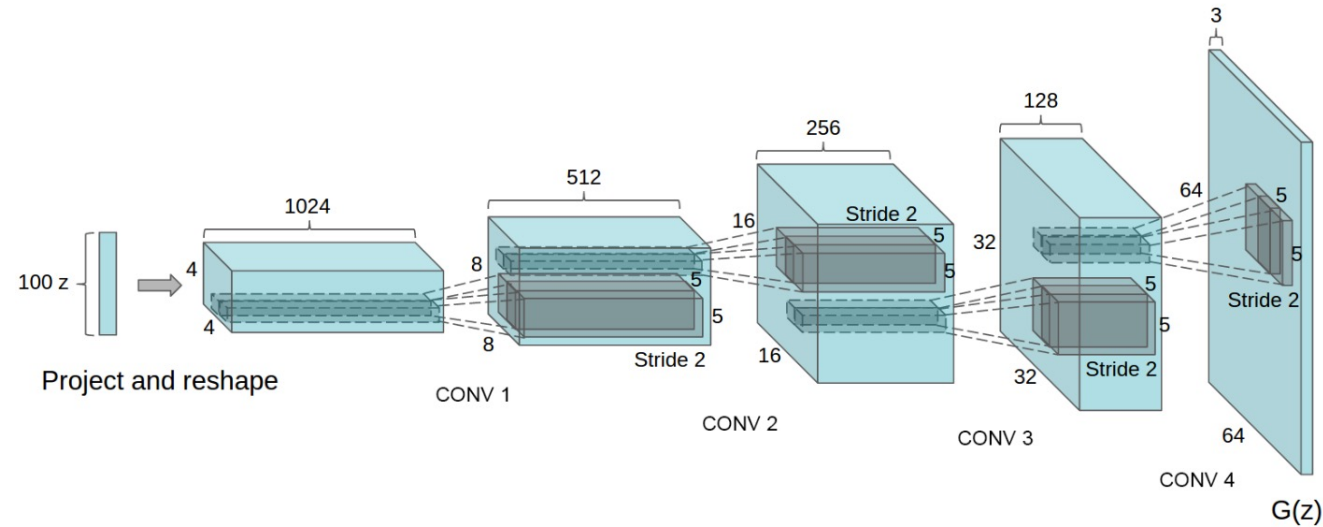
- Fractionally-strided convolutions are also called transposed convolutions.
  - PyTorch: `torch.nn.ConvTranspose2d`.
  - TensorFlow: `tf.keras.layers.Conv2DTranspose`.
- Some researchers are used to call deconvolutions. However, true deconvolutions are the inverse operation of convolution, which is not the same as fractionally-strided convolutions.



Transposed convolution with stride is equivalent to convolving with zero-padding and inserting zeros.



# DCGAN



## Architecture guidelines for stable Deep Convolutional GANs

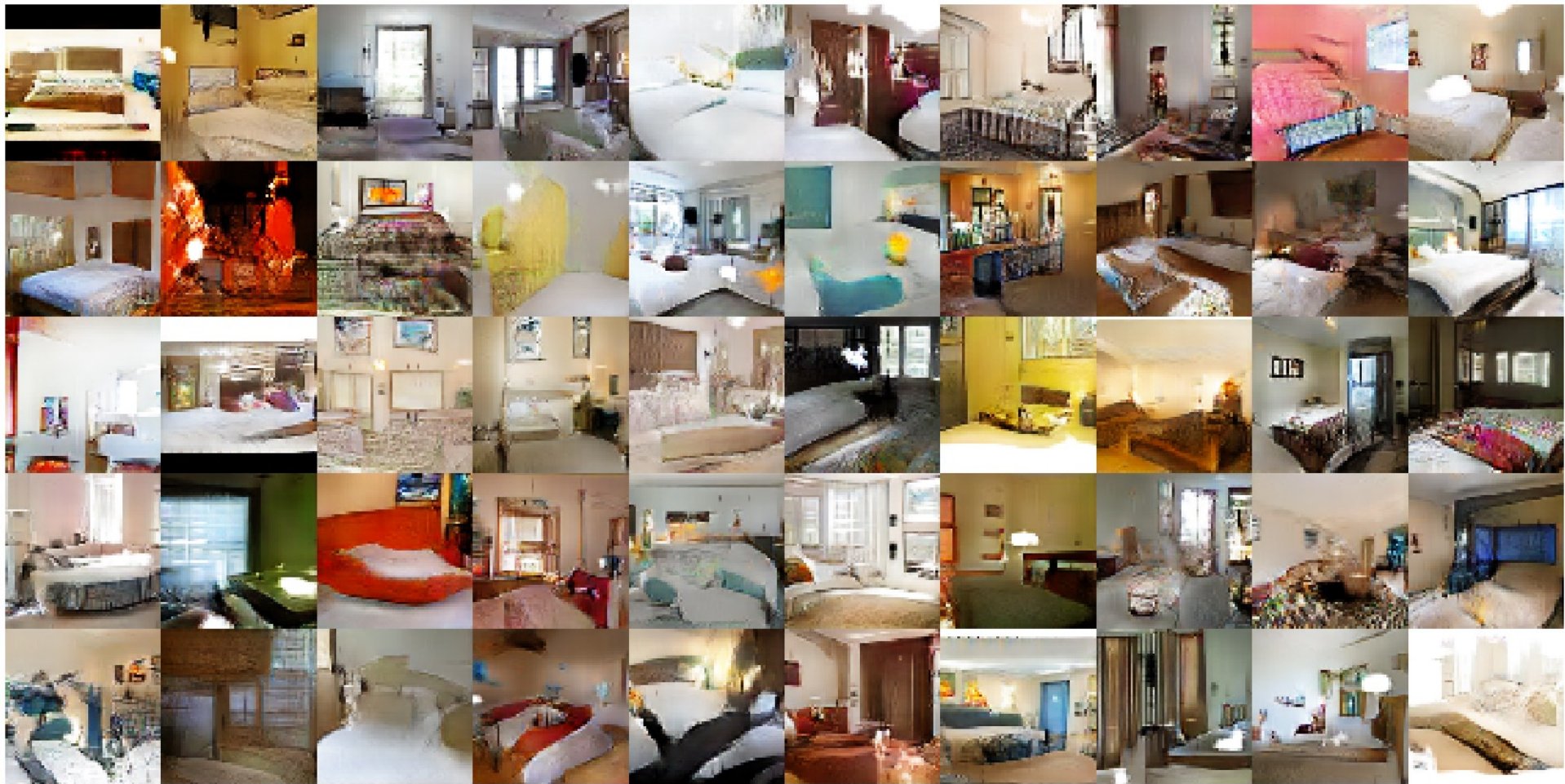
- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

stride=2 everywhere





# DCGAN: Visual Results



The generated bedrooms look very nice (at that time)!



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)



廈門大學 计算机科学与技术系

Department of Computer Science and Technology, Xiamen University

Image source: Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." arXiv preprint arXiv:1511.06434 (2015).



# DCGAN: Walking in the Latent Space

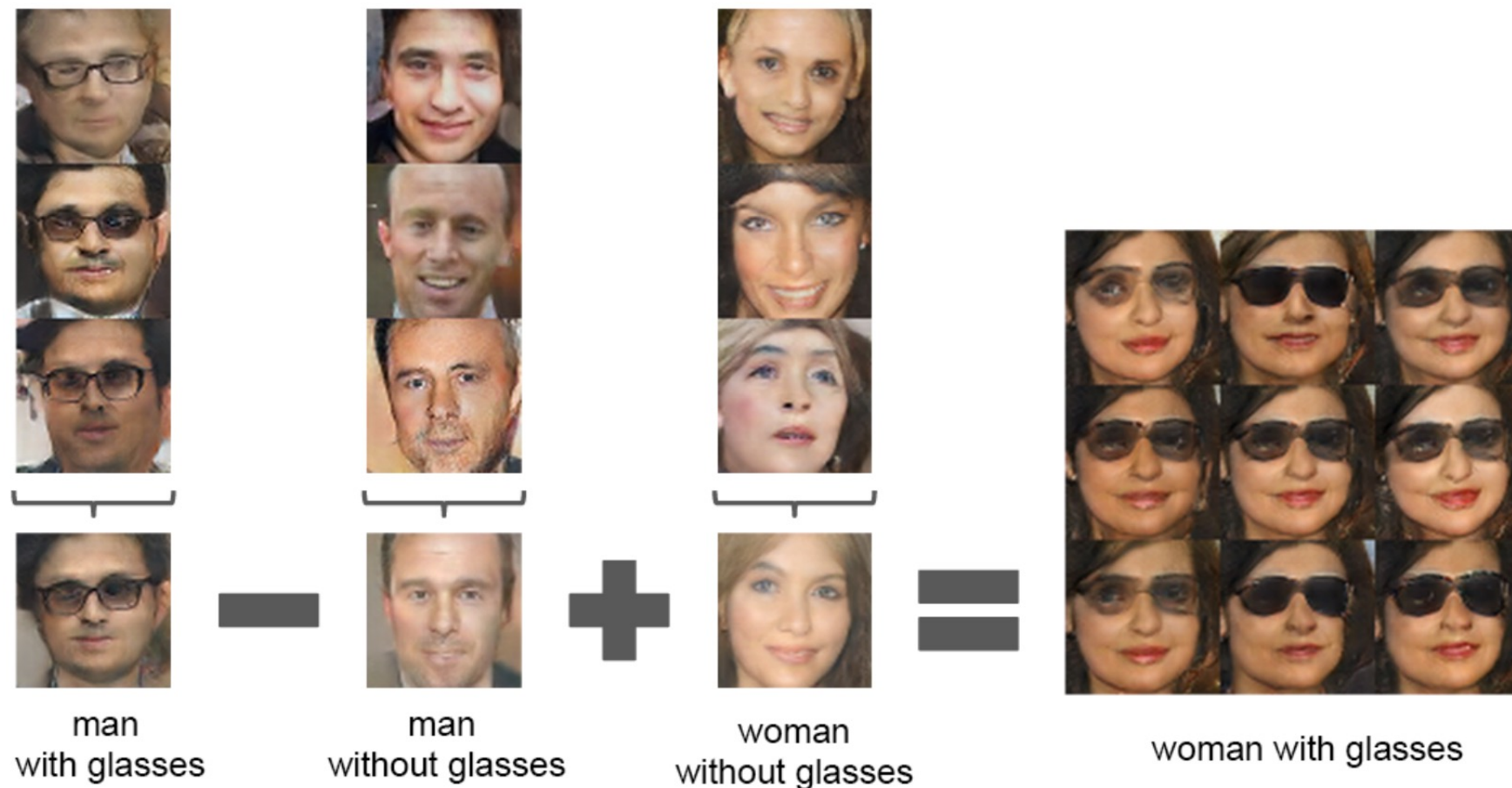
- If walking in this latent space results in semantic changes to the image generations (such as objects being added and removed), we can reason that the model has learned relevant and interesting representations.



Interpolation between a series of 9 random points in  $Z$  show that the space learned has smooth transitions.



# DCGAN: Vector Arithmetic



For each column, the  $Z$  vectors of samples are averaged. Arithmetic was then performed on the mean vectors creating a new vector  $Y$ .



## DCGAN: Use as Feature Extractor

- Train on Imagenet-1k and then use the discriminator's convolutional features from all layers.
- Maxpooling each layers representation to produce a  $4 \times 4$  spatial grid.
- These features are then flattened and concatenated to form a 28672 dimensional vector.
- A regularized linear L2-SVM classifier is trained on top of them.

Model	Accuracy	Accuracy (400 per class)	max # of features units
1 Layer K-means	80.6%	63.7% ( $\pm 0.7\%$ )	4800
3 Layer K-means Learned RF	82.0%	70.7% ( $\pm 0.7\%$ )	3200
View Invariant K-means	81.9%	72.6% ( $\pm 0.7\%$ )	6400
Exemplar CNN	84.3%	77.4% ( $\pm 0.2\%$ )	1024
DCGAN (ours) + L2-SVM	82.8%	73.8% ( $\pm 0.4\%$ )	512





- **We can't control what we generate** from the vanilla GAN.
  - Noise is the only input and it is totally random.
- How can we tell GAN what we want it to generate?
- Straightforward solution: replace data distribution by conditional distribution.

$$p(\mathbf{x}) \rightarrow p(\mathbf{x}|\mathbf{y}).$$

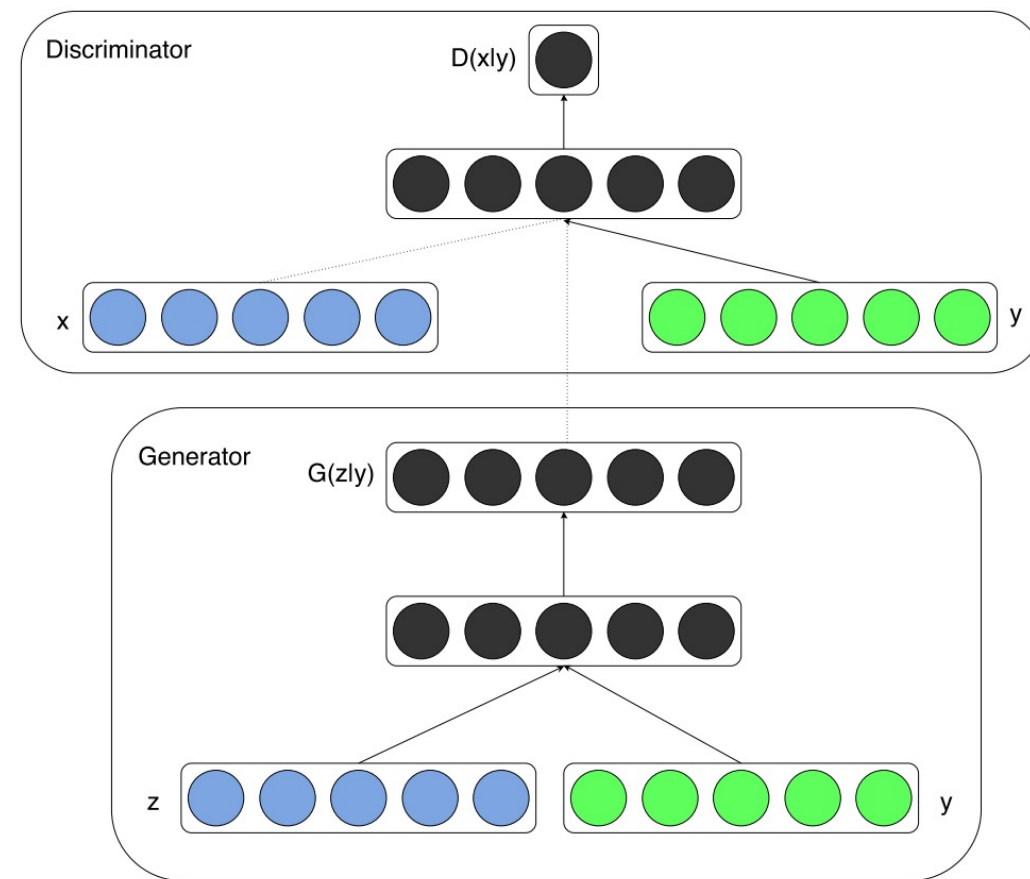
- Now, the problem becomes:
  - Generator: generate a sample for class  $y$ .
  - Discriminator: distinguish the real sample in class  $y$  and the generated sample in class  $y$ .

# CGAN

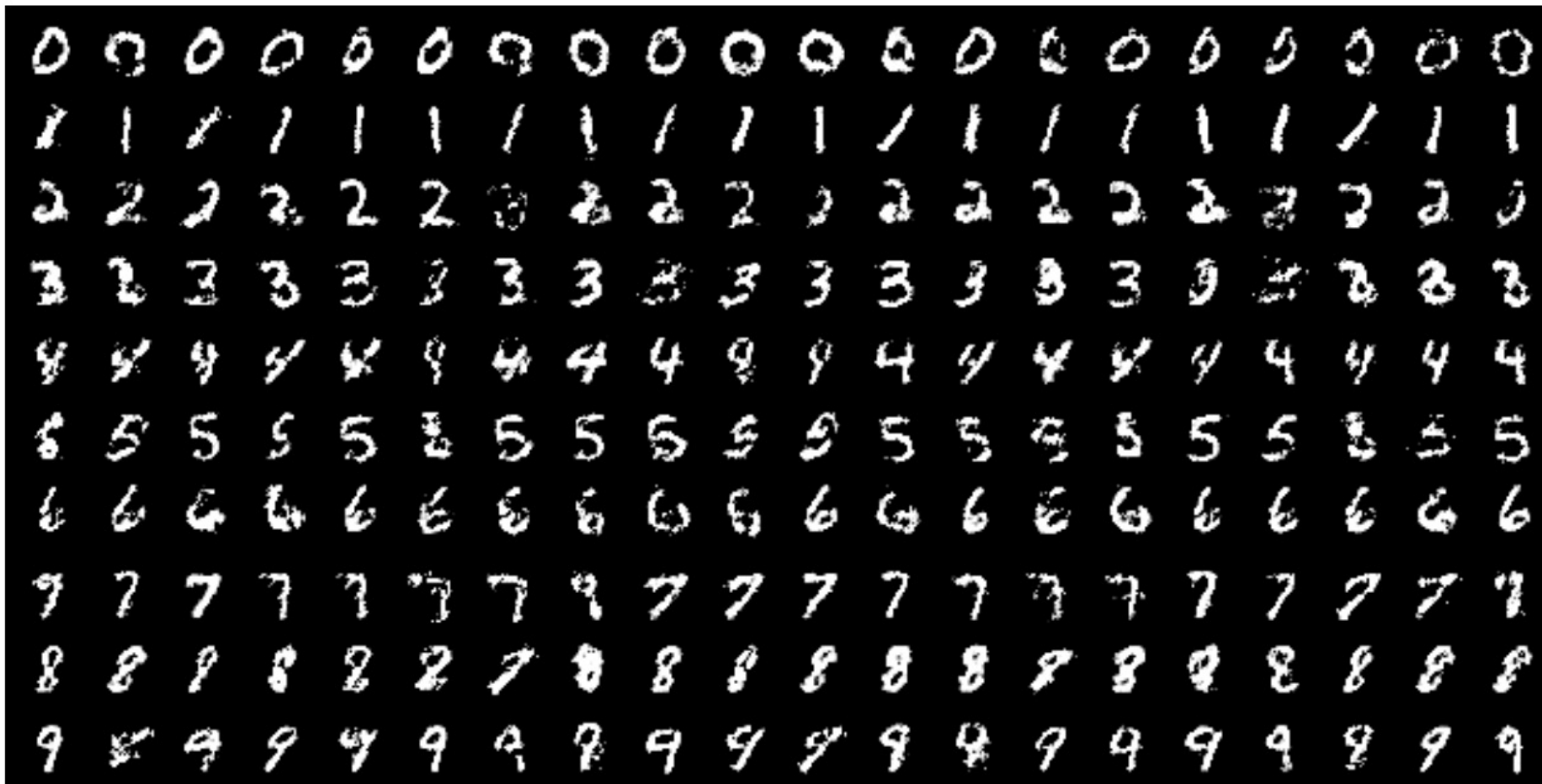
- Both generator and discriminator are conditioned on some extra information  $\mathbf{y}$ :

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x|\mathbf{y}) + \mathbb{E}_{z \sim p(z)} \log \left( 1 - D_{\theta_d} \left( G_{\theta_g}(z|\mathbf{y}) \right) \right) \right].$$

- $\mathbf{y}$  could be any kind of auxiliary information, such as class labels or data from other modalities.
- E.g. the speech of saying that class.



# CGAN

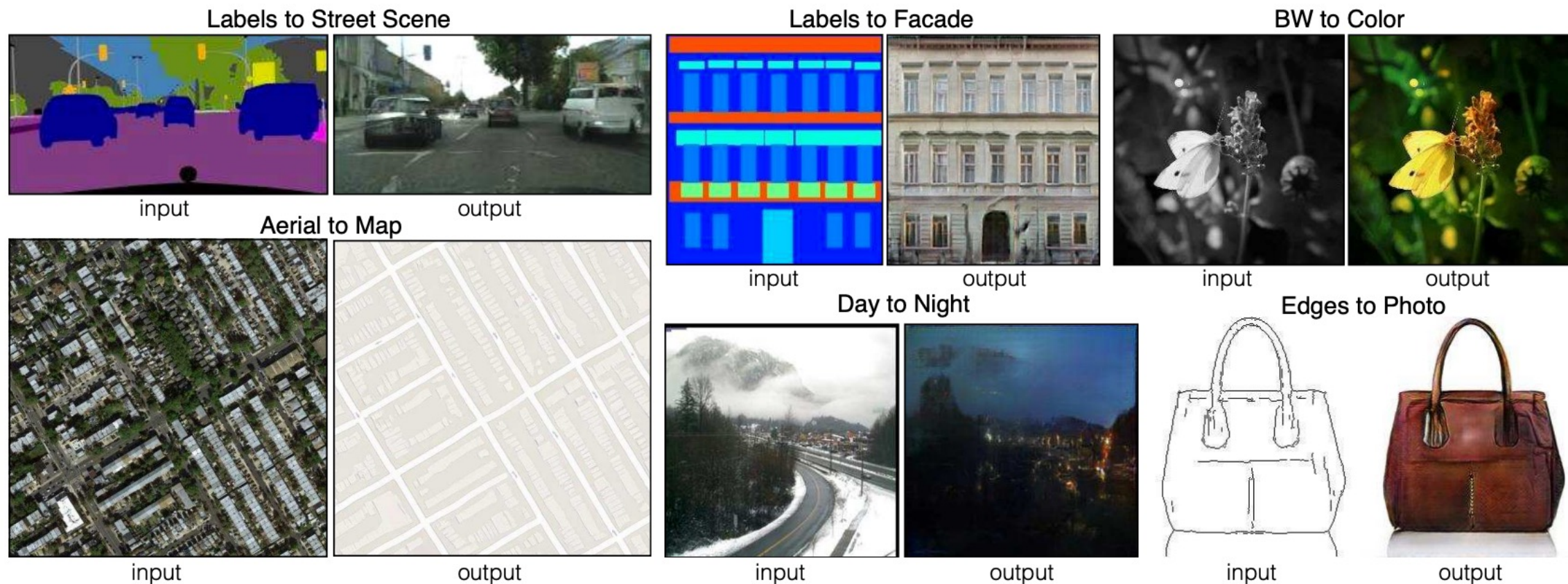


Generate MNIST digits by directly feeding one-hot class label.





- Given a pair of images, transfer the style of one image to another.

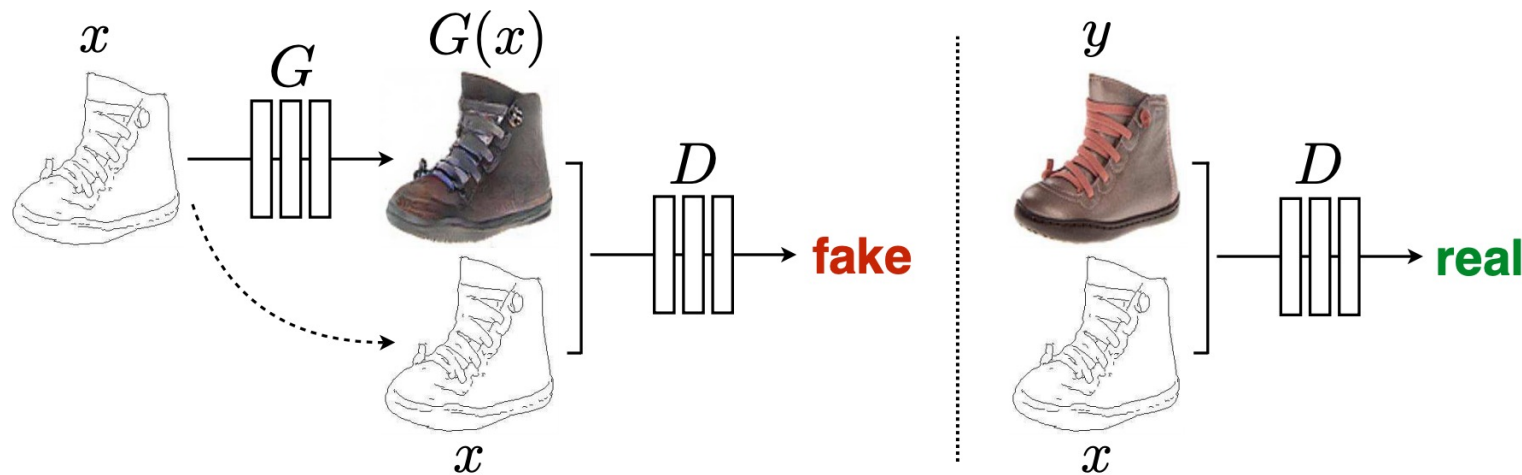




# pix2pix

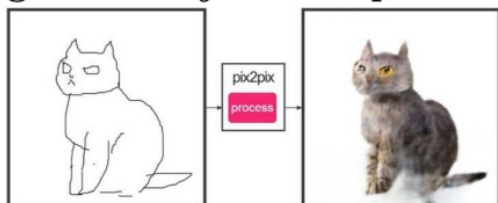
- The model is based on CGAN.
- As an improvement, the generator is tasked to not only fool the discriminator but also to be near the ground truth output.
- $L_1$  penalization is added to the loss of CGAN to encourage less blurring:

$$L_{L_1}(G) = \mathbb{E}_{x,y,z} [\|y - G(x, z)\|_1].$$

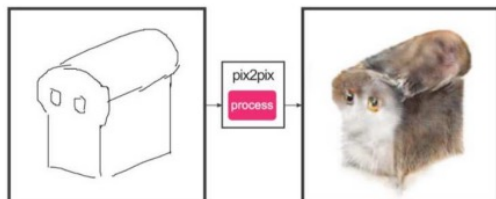


# pix2pix: More Applications

#edges2cats by Christopher Hesse



by @gods\_tail



by @ivymyt



by @vvid

Sketch→Portrait



by Mario Klingemann

“Do as I do”



by Brannon Dorsey

Depth→Streetview



by Jasper van Loenen

Palette generation



by Jack Qiao

Background removal



by Kaihu Chen

Sketch → Pokemon



by Bertrand Gondouin



# CycleGAN

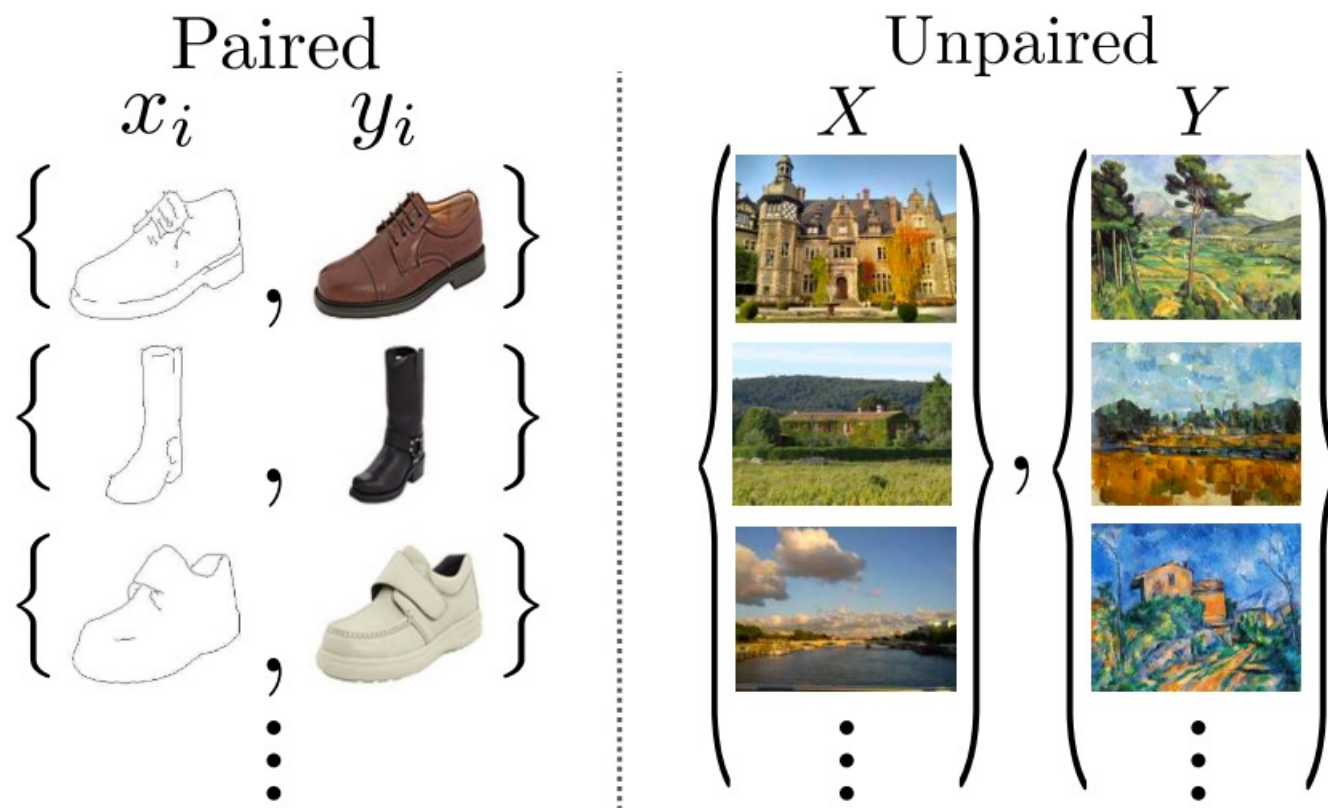
Unpaired image-to-image translation using cycle-consistent adversarial networks

JY Zhu, T Park, P Isola, AA Efros - Proceedings of the IEEE ..., 2017 - openaccess.thecvf.com

... losses on domains X and Y yields our full objective for **unpaired image-to-image** translation. ... We first compare our approach against recent methods for **unpaired image-to-image** ...

☆ Save 剪 Cite Cited by 20082 Related articles All 26 versions 》》

- Paired examples can be expensive to obtain.
- Can we translate from  $X \leftrightarrow Y$  in an unsupervised manner?



Paired vs. unpaired examples

# CycleGAN

- Two generators:

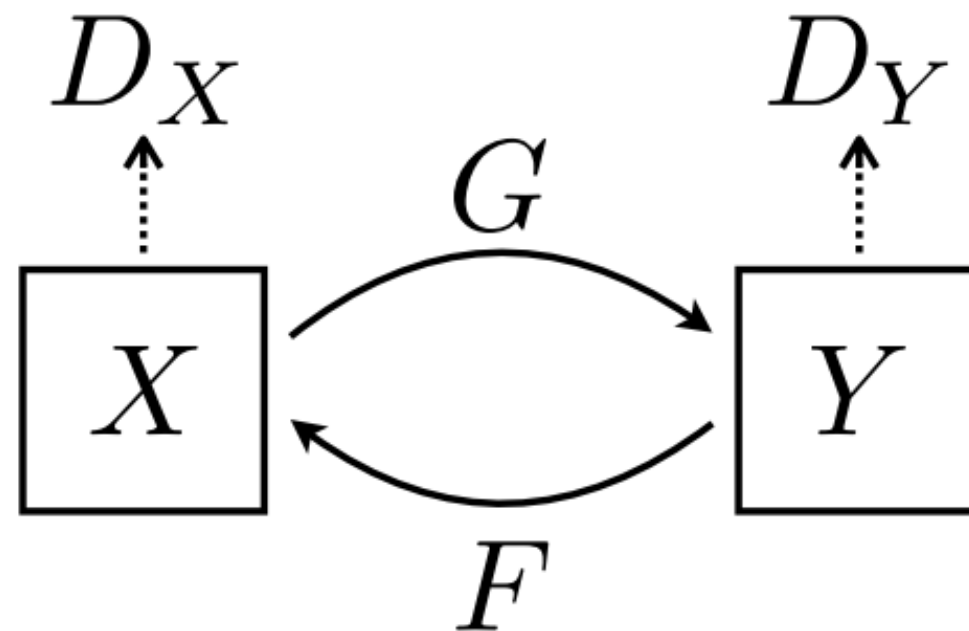
- $G: X \rightarrow Y;$

- $F: Y \rightarrow X.$

- Two discriminators:

- $D_X$  aims to distinguish between images  $\{x\}$  and translated images  $\{F(y)\};$

- $D_Y$  aims to discriminate between  $\{y\}$  and  $\{G(x)\}.$

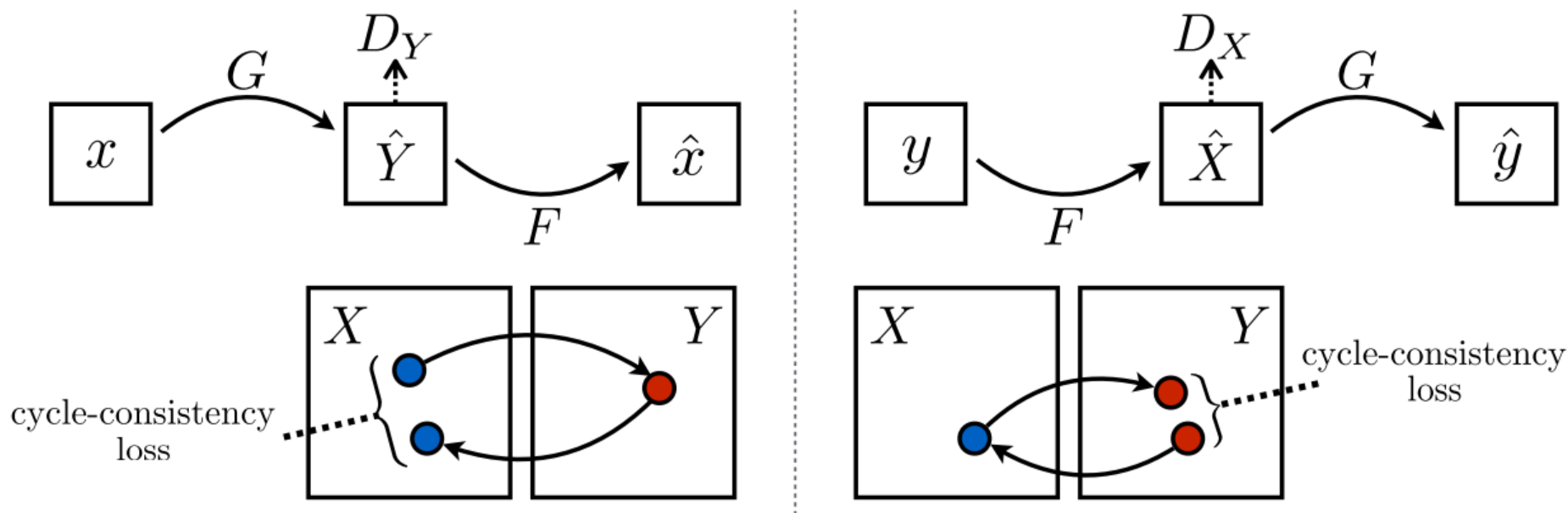




# CycleGAN

- If we can go from  $X$  to  $\hat{Y}$  via  $G$ , then it should also be possible to go from  $\hat{Y}$  back to  $X$  via  $F$ .
- Cycle consistency loss is added to the original adversarial loss:

$$L_{\text{cyc}}(G, F) = \mathbb{E}_x \left[ \|F(G(x)) - x\|_1 \right] + \mathbb{E}_y \left[ \|G(F(y)) - y\|_1 \right].$$



# CycleGAN

Failed case:



horse → zebra

Monet ↔ Photos



Monet → photo

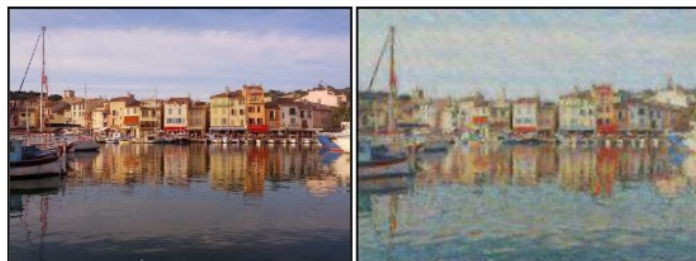


photo → Monet

Zebras ↔ Horses

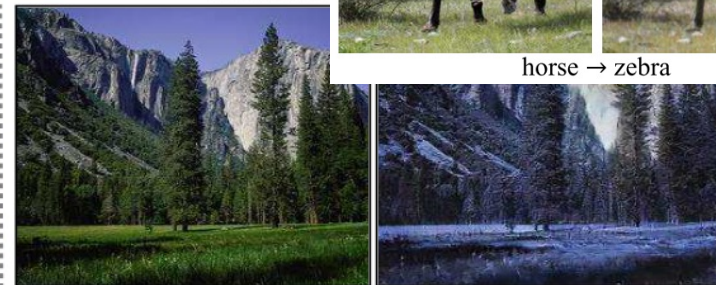


zebra → horse

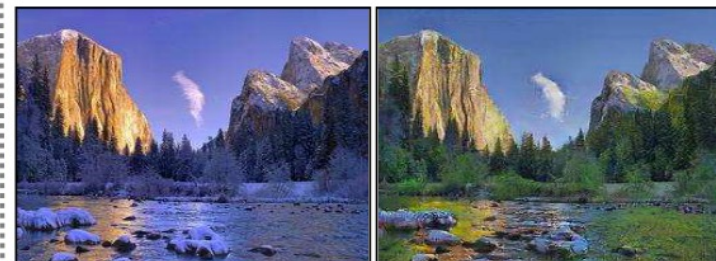


horse → zebra

Summer



summer → winter



winter → summer



Photograph



Monet



Van Gogh



Cezanne



Ukiyo-e





# DIFFUSION MODEL



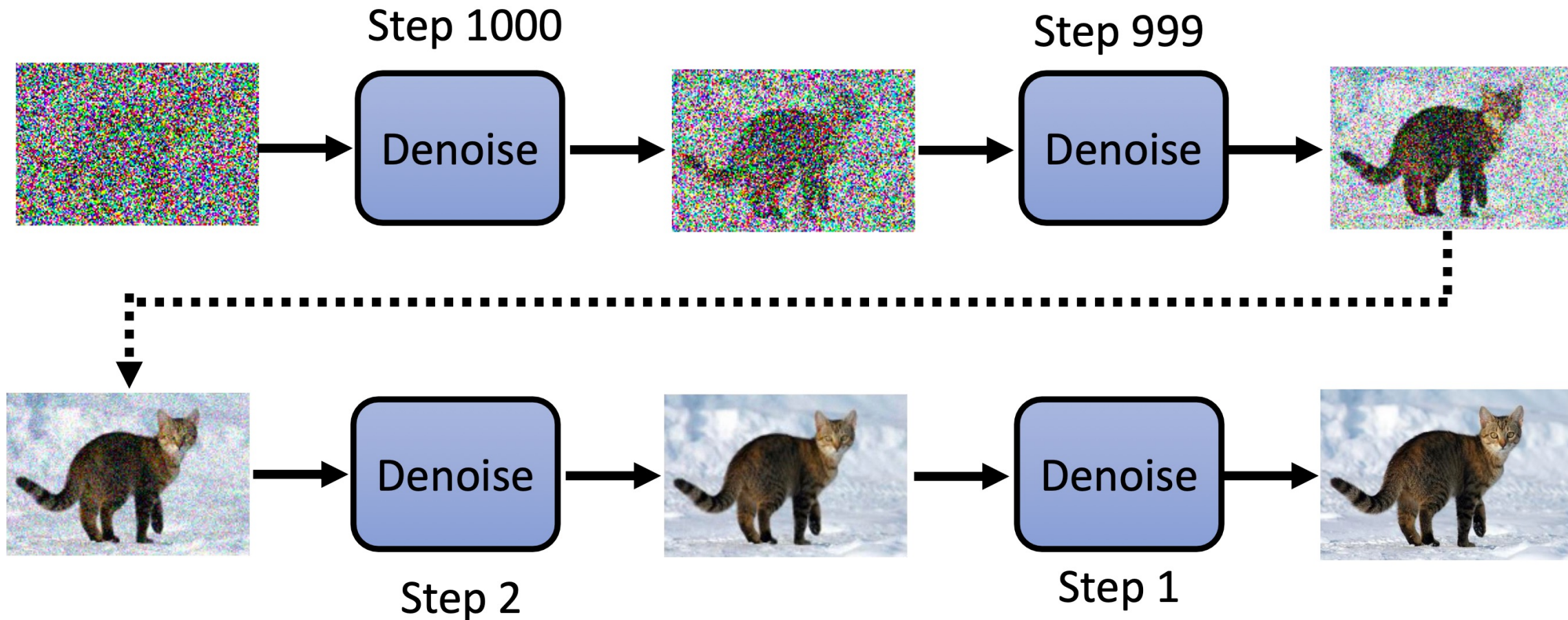
# Diffusion Model

## Denoising diffusion probabilistic models

[J Ho, A Jain, P Abbeel](#) - Advances in neural information ..., 2020 - proceedings.neurips.cc

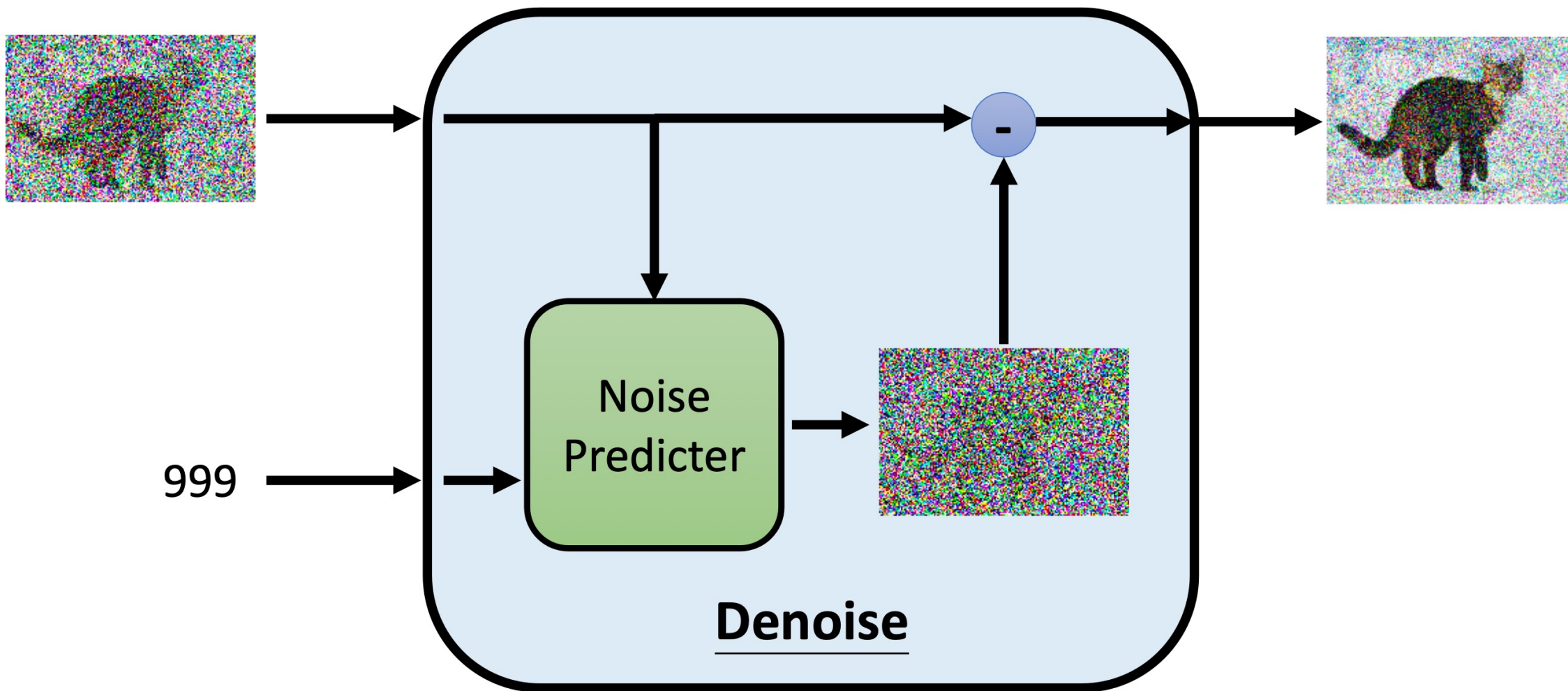
... This paper presents progress in **diffusion probabilistic models** [53]. A **diffusion probabilistic model** (which we will call a “**diffusion model**” for brevity) is a parameterized Markov chain ...

☆ Save 77 Cite Cited by 4825 Related articles All 6 versions ⌕

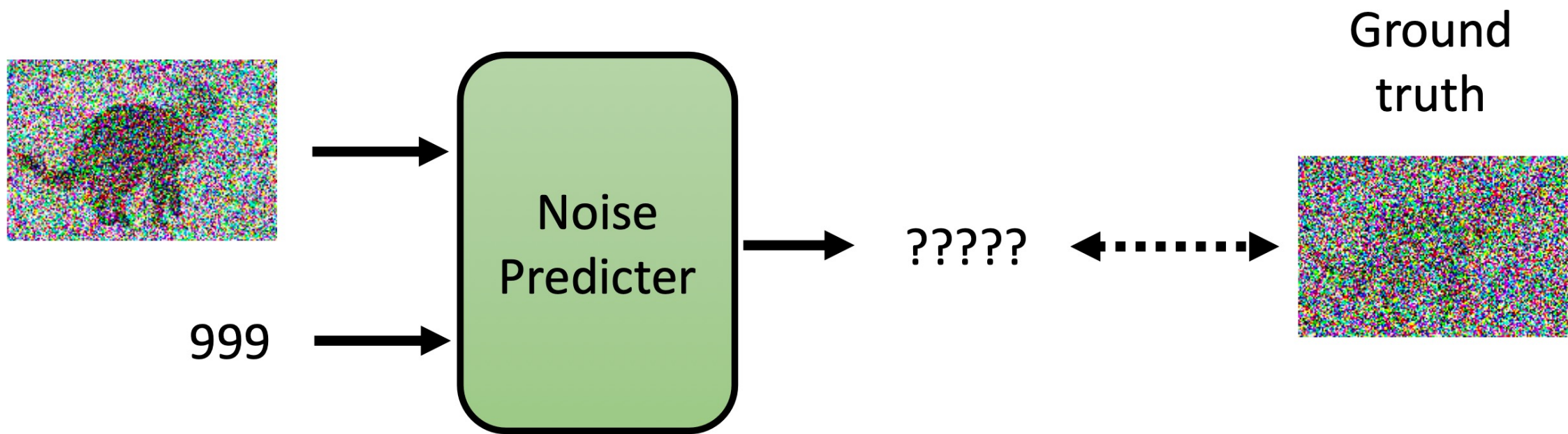




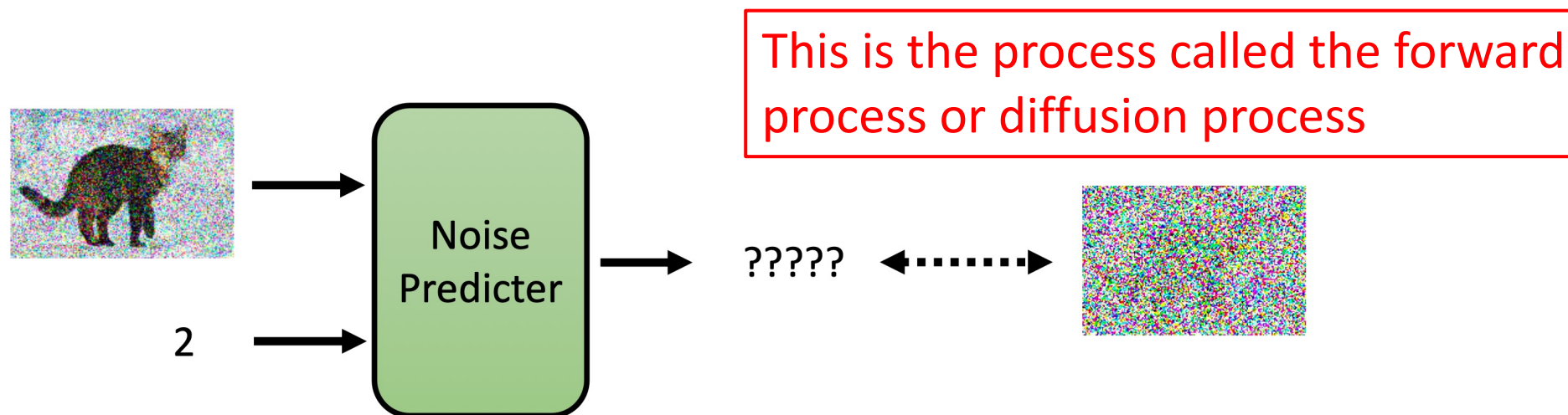
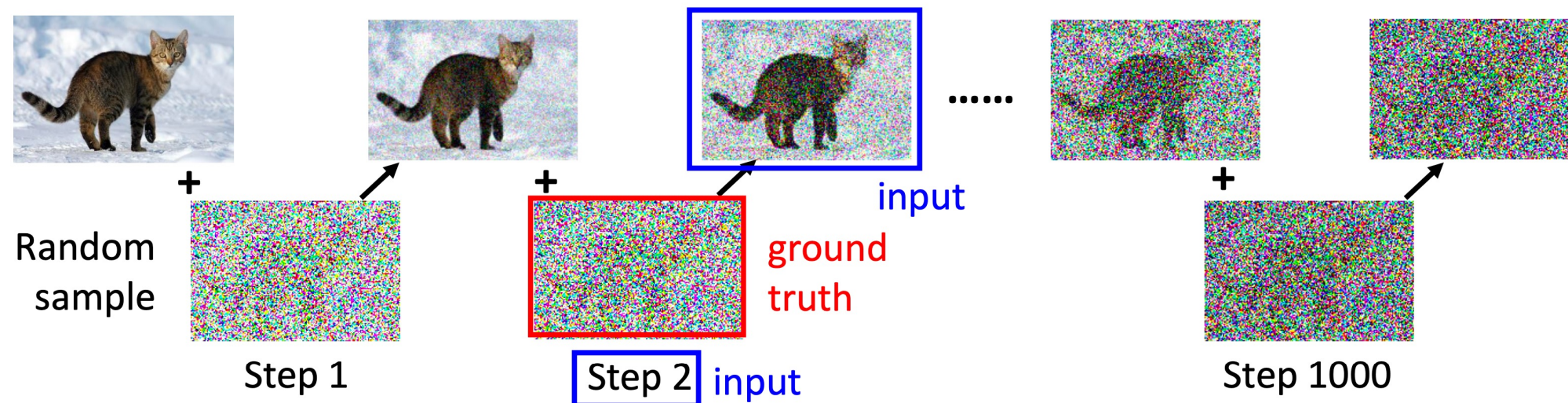
# Diffusion Model



# Diffusion Model

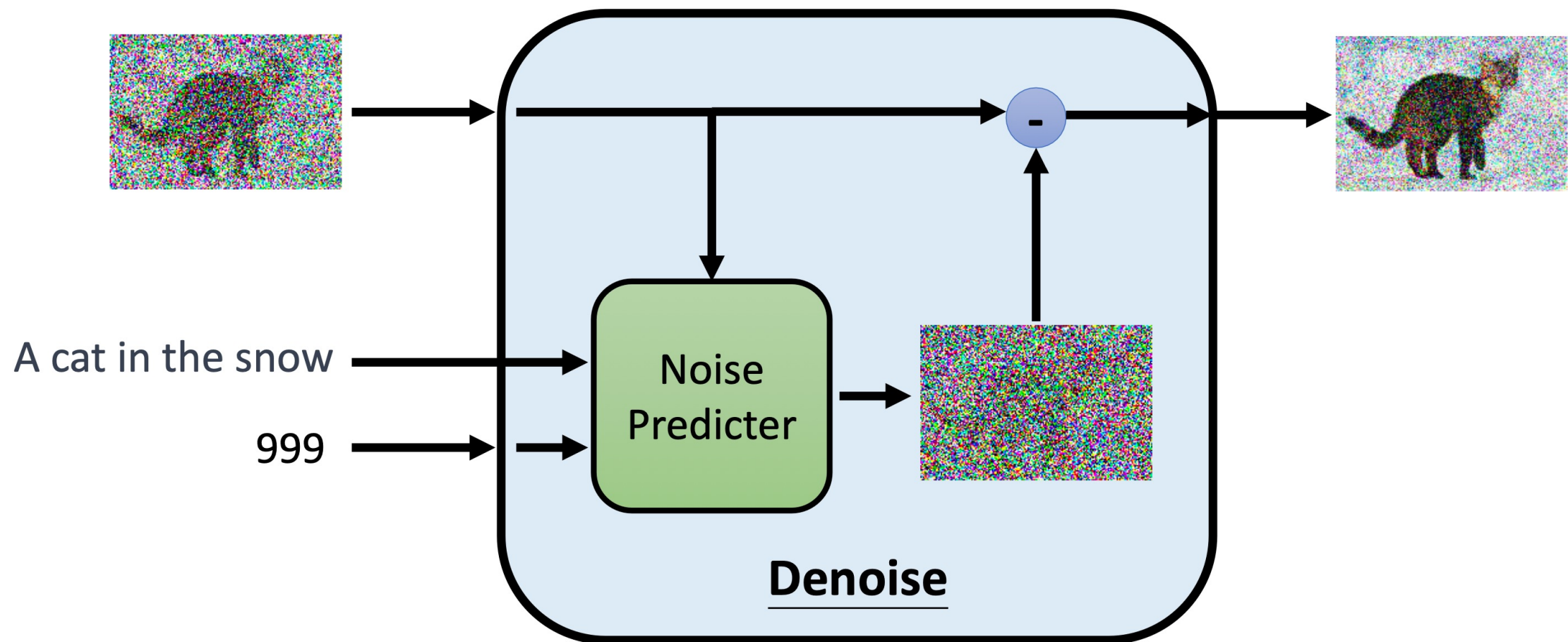


# Diffusion Model





# Diffusion Model





# Diffusion Model

## Algorithm 1 Training

- 1: **repeat**
- 2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4:  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on  
$$\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$$
- 6: **until** converged

Add ground  
truth noise

Ground truth noise

Denoise model

## Algorithm 2 Sampling

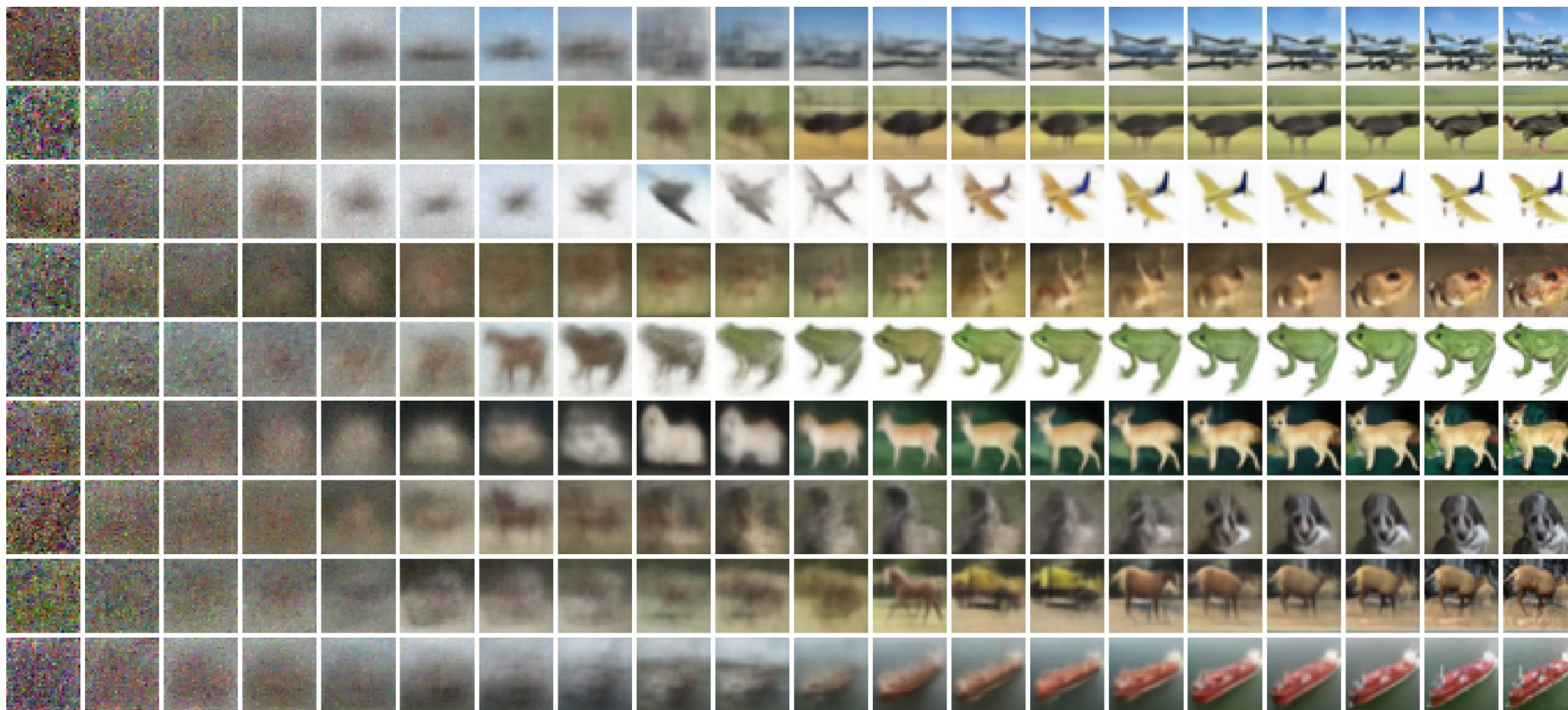
- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for**  $t = T, \dots, 1$  **do**
- 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
- 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return**  $\mathbf{x}_0$

Subtract  
predicted noise

Denoise model



# Diffusion Model



Unconditional CIFAR10 progressive generation



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)



廈門大學 计算机科学与技术系

Department of Computer Science and Technology, Xiamen University

Image source: Ho, Jonathan, Ajay Jain, and Pieter Abbeel. "Denoising diffusion probabilistic models." Advances in neural information processing systems 33 (2020): 6840-6851.



# Diffusion Model



Latent  $\mathbf{x}_{750}$

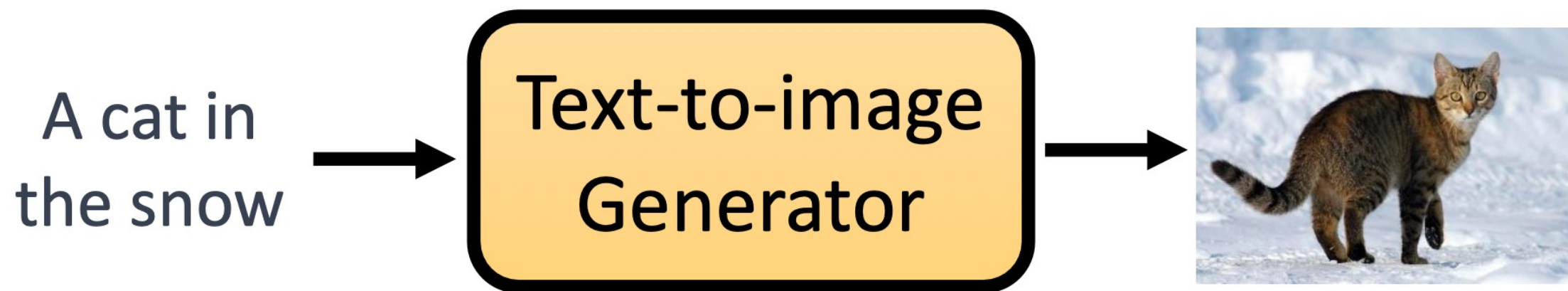


Decodings  $\mathbf{x}_0 \sim p_{\theta}(\mathbf{x}_0 | \mathbf{x}_{750})$



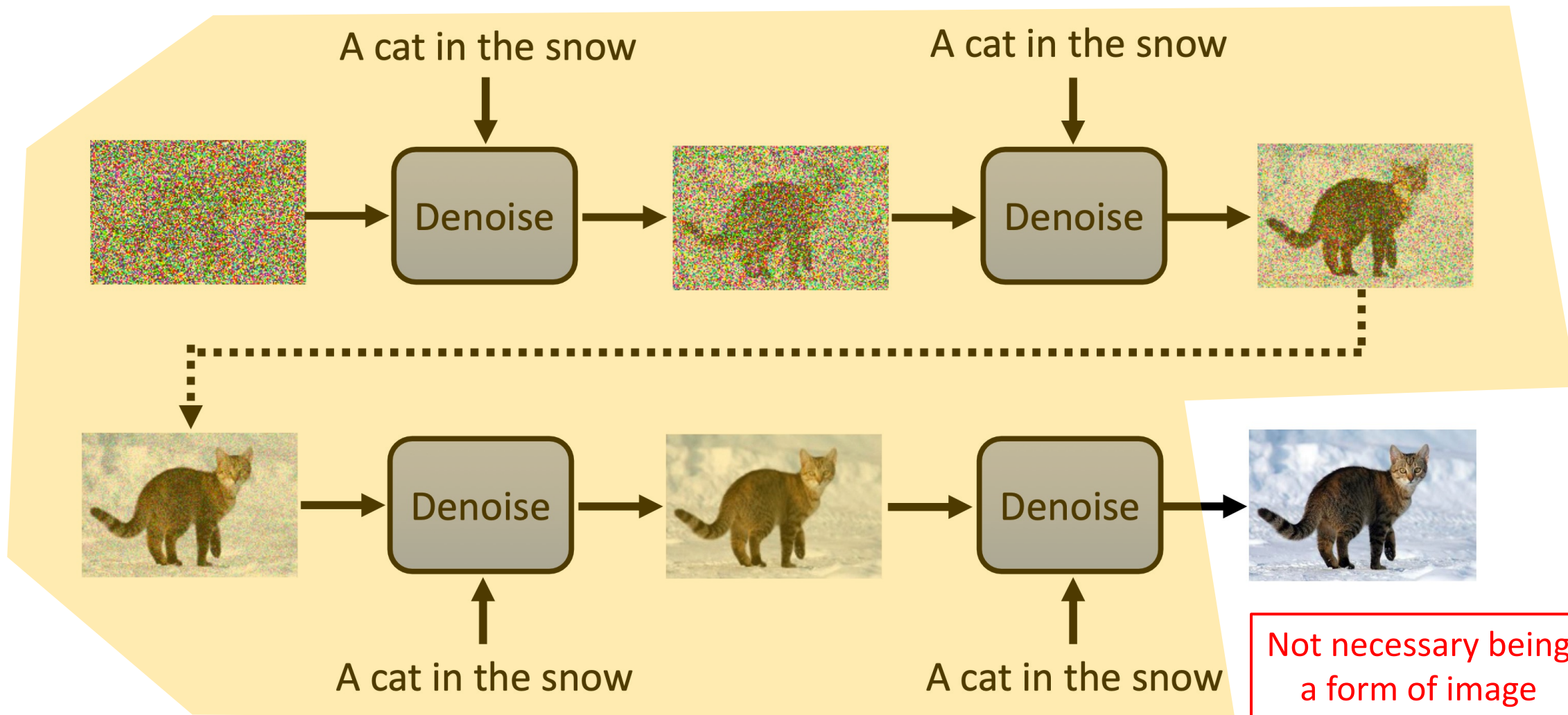
# Diffusion Model

- How can we generate desired image?





# Diffusion Model



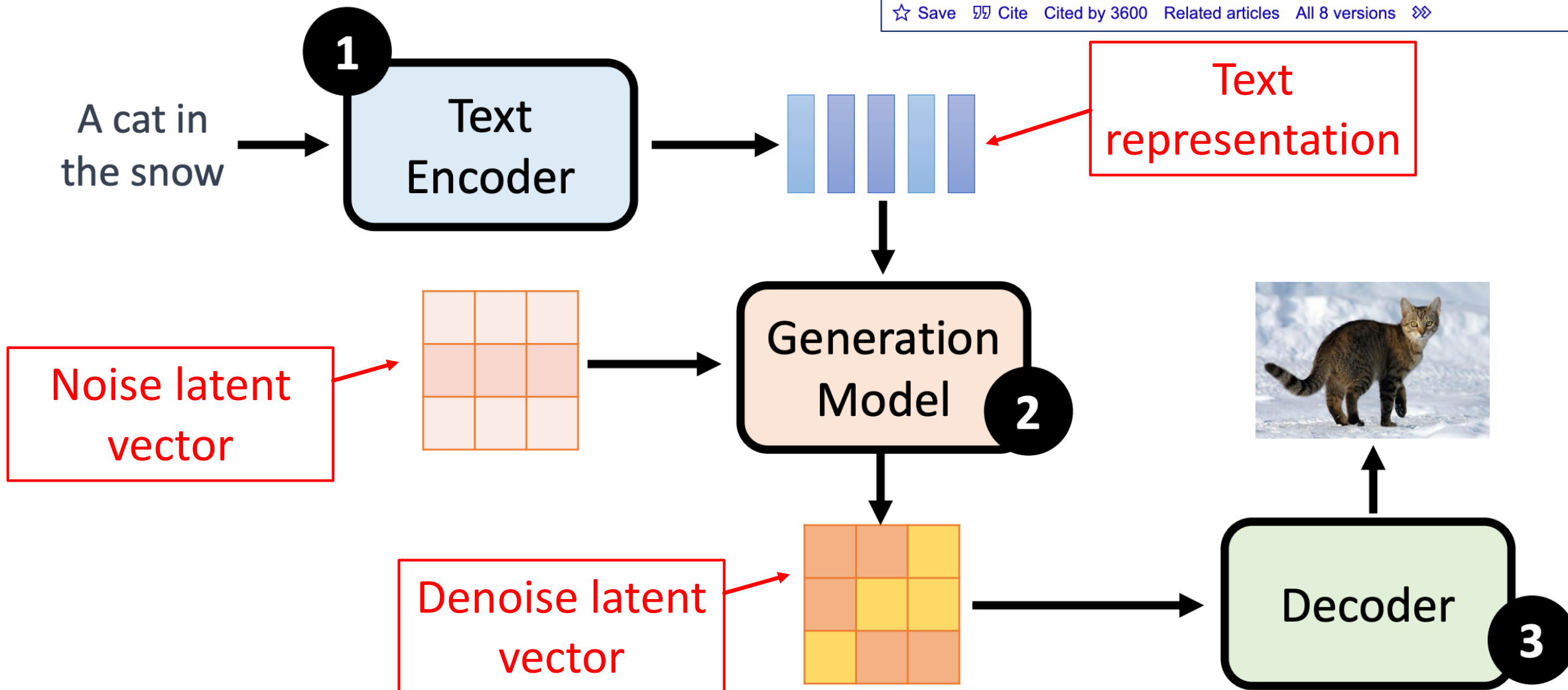
# Latent Diffusion

## High-resolution image synthesis with latent diffusion models

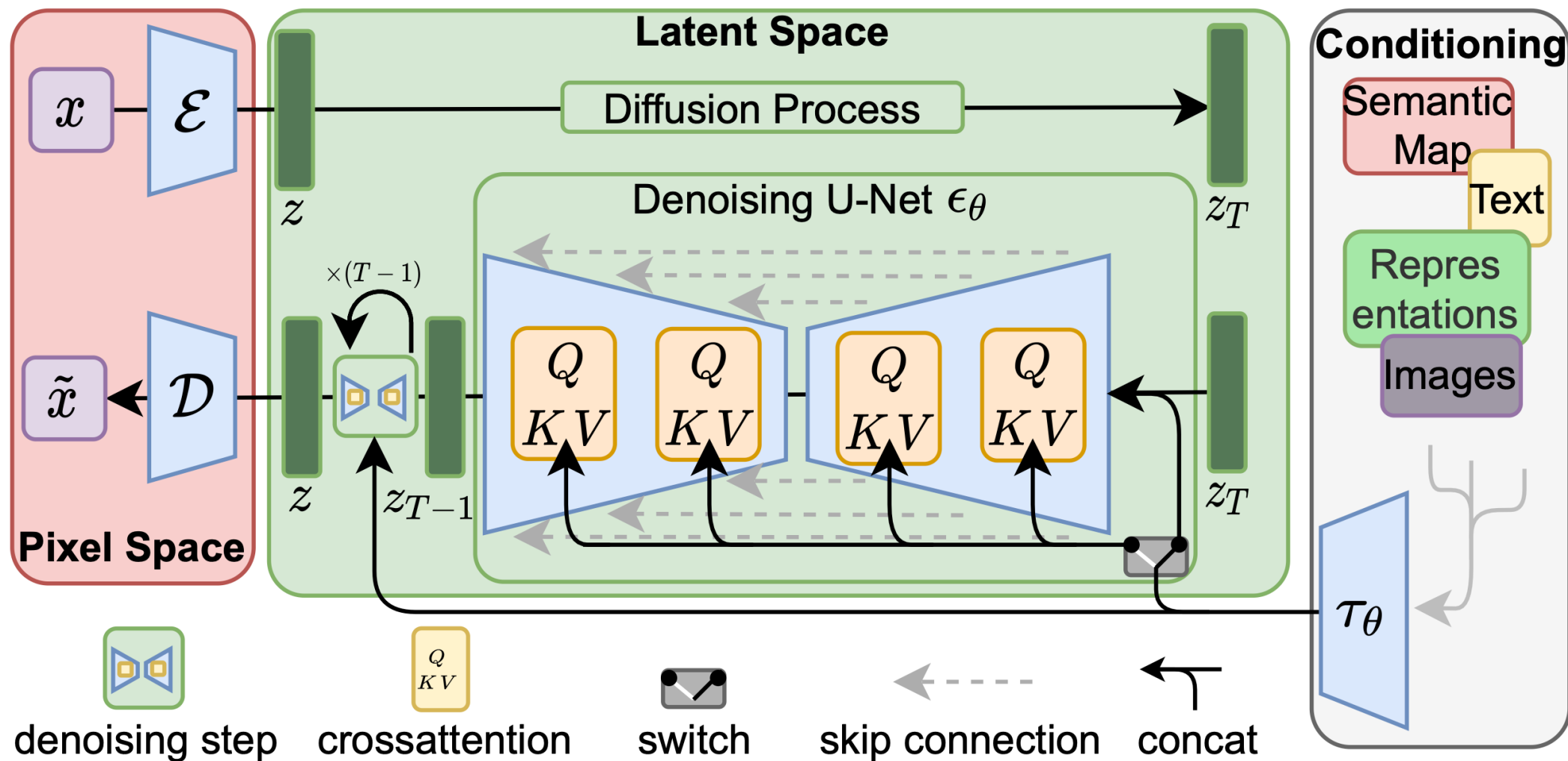
R Rombach, A Blattmann, D Lorenz... - Proceedings of the ..., 2022 - openaccess.thecvf.com

... To lower the computational demands of training diffusion models towards **high-resolution image synthesis**, we observe that although diffusion models allow to ignore perceptually ...

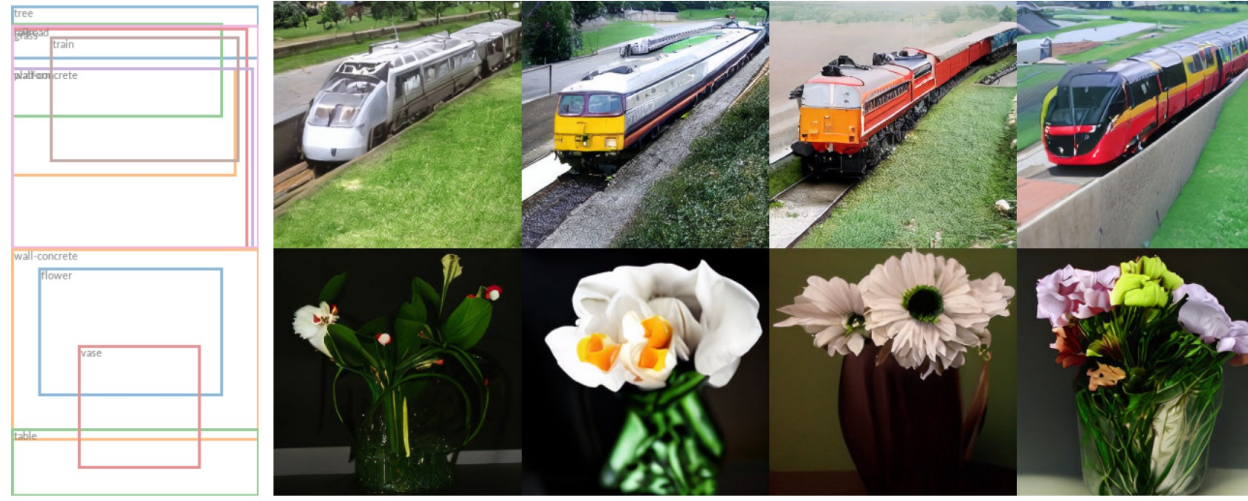
☆ Save 77 Cite Cited by 3600 Related articles All 8 versions »



# Latent Diffusion



# Latent Diffusion



Layout-to-image  
synthesis on COCO

*"A street sign that reads  
'Latent Diffusion'"*



*"An oil painting  
of a space shuttle"*



Text-to-image *LDM*  
model for user-  
defined text prompts



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)



廈門大學 计算机科学与技术系

Department of Computer Science and Technology, Xiamen University



# Latent Diffusion



## Semantic synthesis of landscape images



廈門大學信息學院(特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)

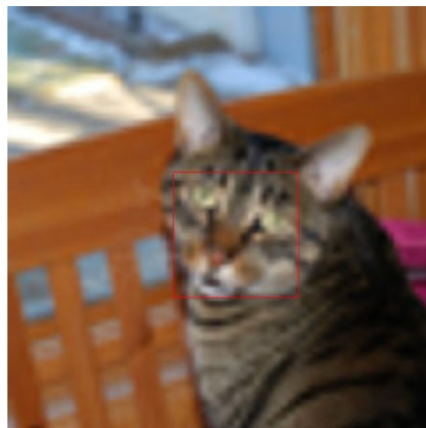


廈門大學 计算机科学与技术系

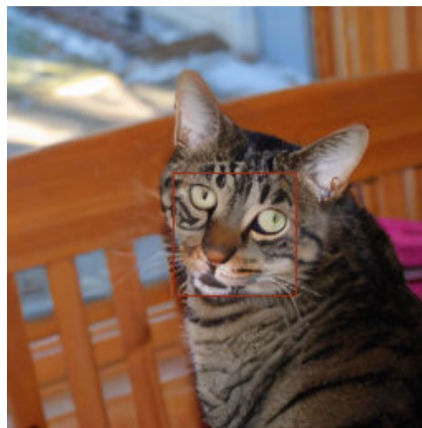
Department of Computer Science and Technology, Xiamen University

# Latent Diffusion

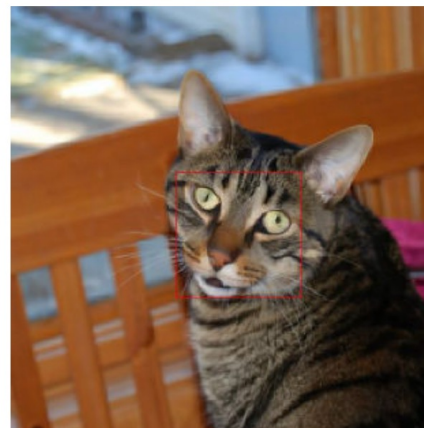
bicubic



*LDM-SR*



SR3



ImageNet 64→256 super-resolution





# Latent Diffusion

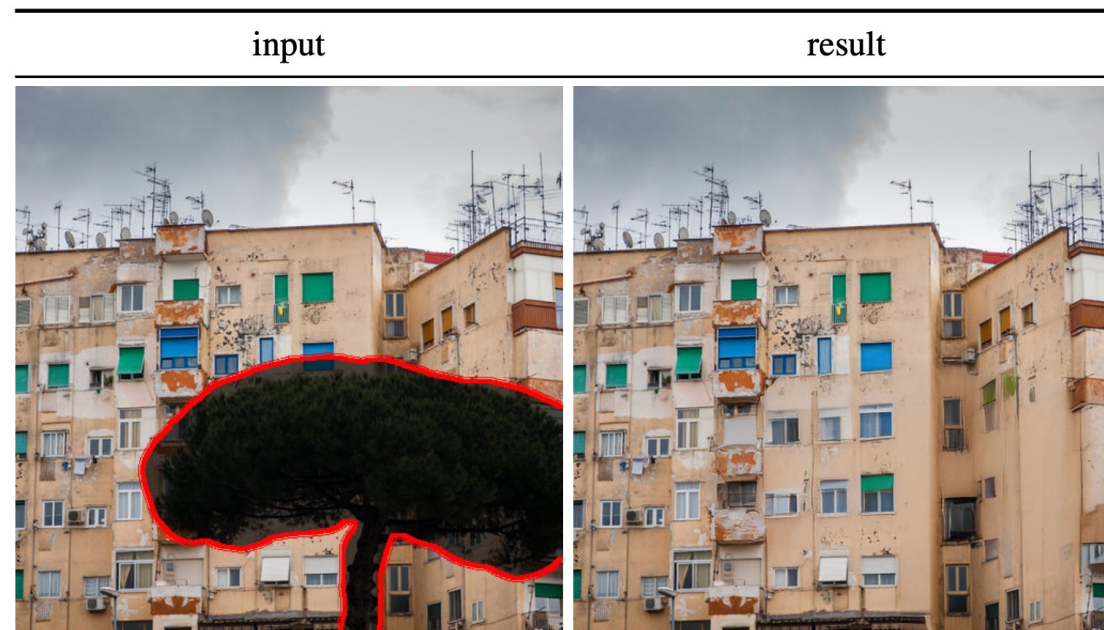
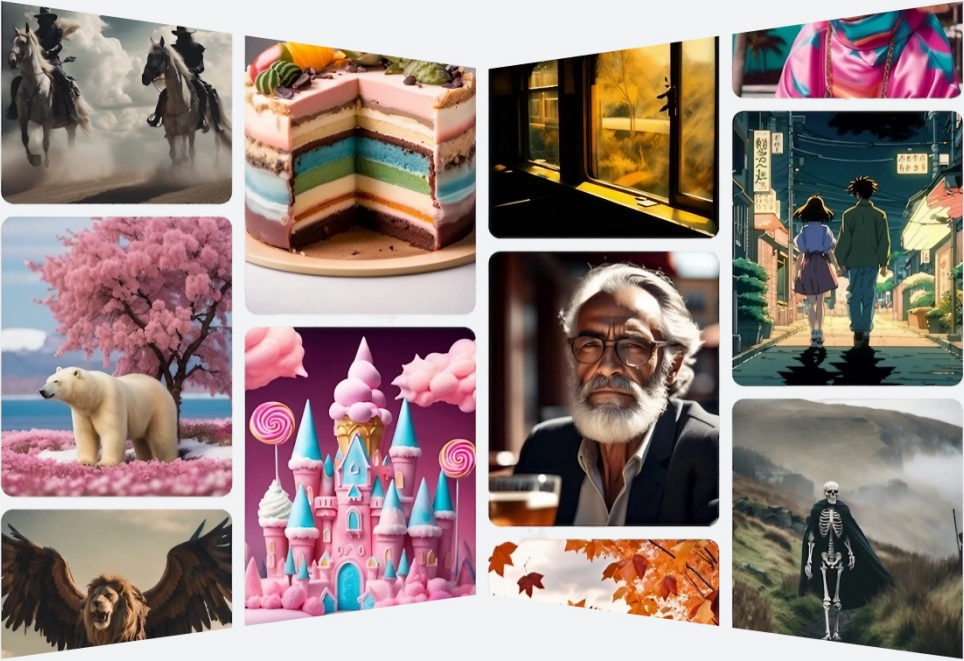


Image inpainting with latent diffusion



# Stable Diffusion

- Stable diffusion is developed by researchers from the CompVis Group at Ludwig Maximilian University of Munich and Runway with a compute donation by Stability AI and training data from non-profit organizations.



## Stable Diffusion

Get involved with the fastest growing open software project. Download and join other developers in creating incredible applications with Stable Diffusion as a foundation model.

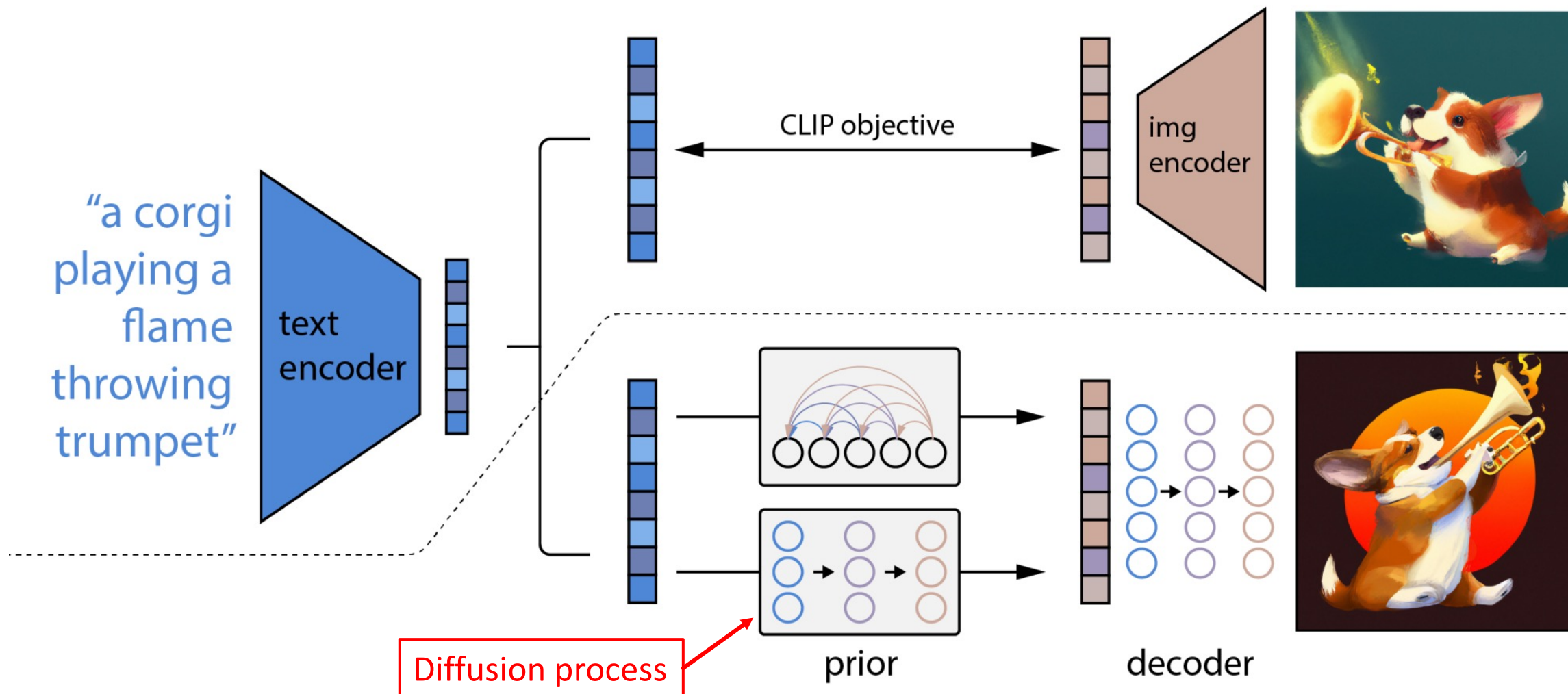
[Try Stable Diffusion](#)

[Download Code](#)

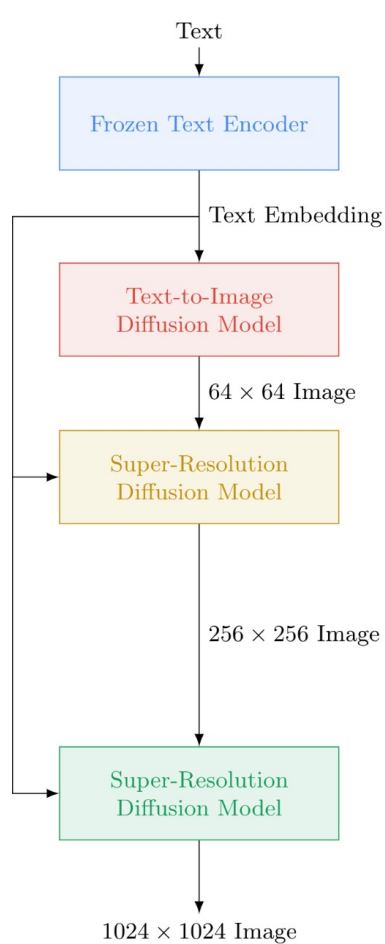




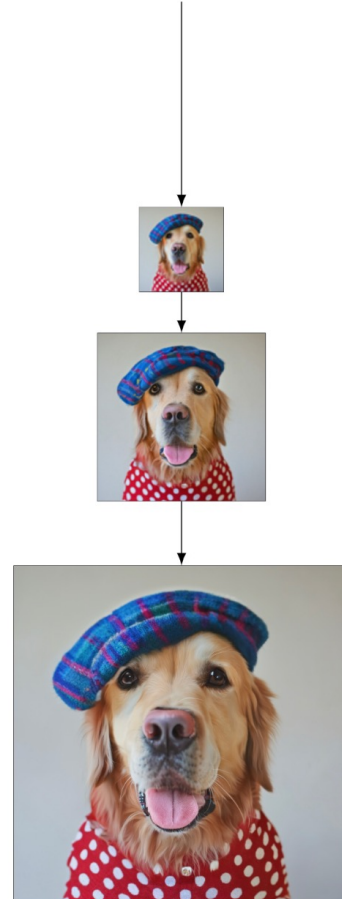
# DALL·E



# Imagen



“A Golden Retriever dog wearing a blue checkered beret and red dotted turtleneck.”



Sprouts in the shape of text 'Imagen' coming out of a fairytale book.



A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.



A high contrast portrait of a very happy fuzzy panda dressed as a chef in a high end kitchen making dough. There is a painting of flowers on the wall behind him.



Teddy bears swimming at the Olympics 400m Butterfly event.



A cute corgi lives in a house made out of sushi.



A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.



A brain riding a rocketship heading towards the moon.



A dragon fruit wearing karate belt in the snow.



A strawberry mug filled with white sesame seeds. The mug is floating in a dark chocolate sea.





# Midjourney



These fake images of the Pope and Donald Trump have gone viral in recent weeks, both were created on Midjourney v5.



# ControlNet

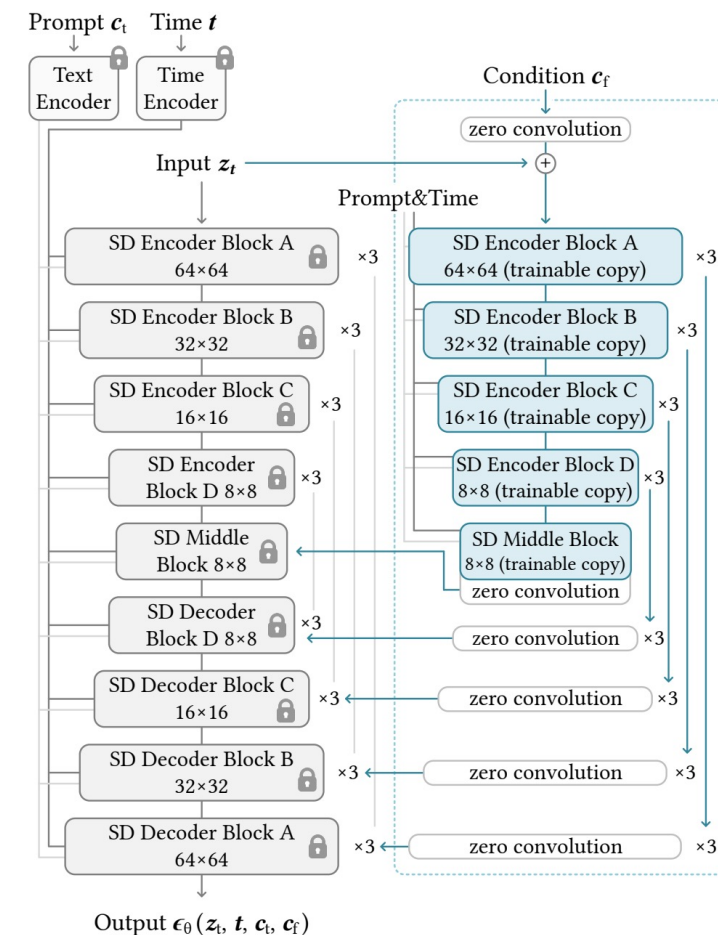
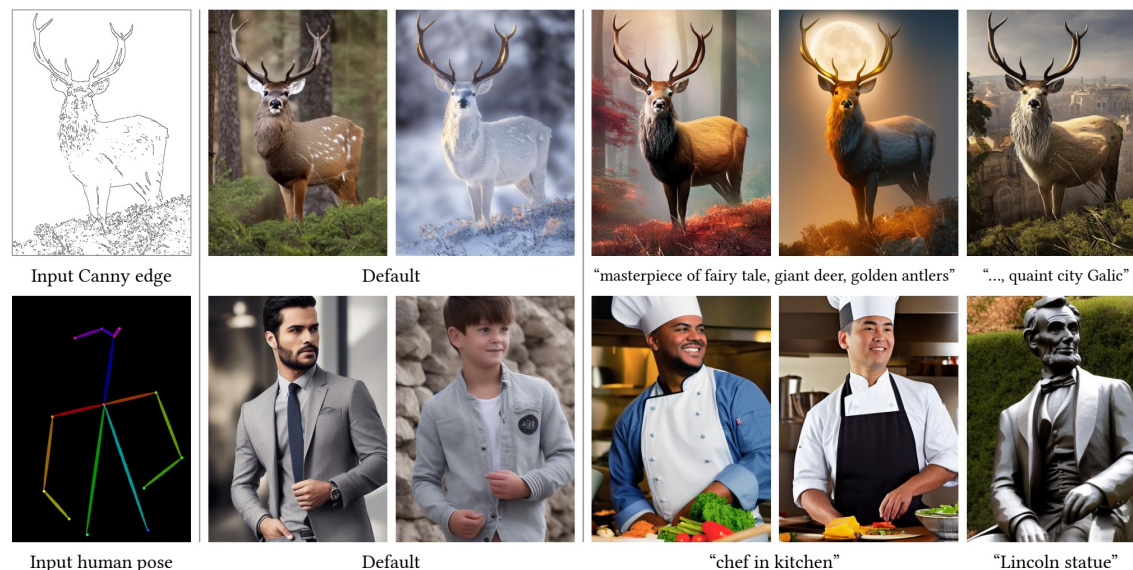
## Adding conditional control to text-to-image diffusion models

L Zhang, A Rao, M Agrawala - Proceedings of the IEEE/CVF ..., 2023 - openaccess.thecvf.com

... Learning **conditional controls** for large text-to-image diffusion ... network architecture that learns **conditional controls** for large ... for learning diverse **conditional controls**. The trainable copy ...

☆ Save 77 Cite Cited by 440 Related articles All 3 versions 》

- A new concept called **zero convolution layers**, with weights initialized to zeros so that they progressively grow during the training.
- This architecture ensures that **harmful noise is not added to the deep features** of the large diffusion model.



(a) Stable Diffusion

(b) ControlNet



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)



廈門大學 计算机科学与技术系

Department of Computer Science and Technology, Xiamen University



# Conclusion

After this lecture, you should know:

- What is a generative model?
- How generator and discriminator improve each other?
- How does transposed convolution work?
- How to design application specific loss to train with adversarial loss?
- What is the main principle of diffusion model?



# Suggested Reading

- Adversarial Nets Papers
- Tips and tricks to make GANs work
- 令人拍案叫绝的Wasserstein GAN
- 李宏毅讲Diffusion Model



# Assignment 4

- Assignment 4 is released. The deadline is 18:00, 2nd December.



# Thank you!

- Any question?
- Don't hesitate to send email to me for asking questions and discussion. 😊

