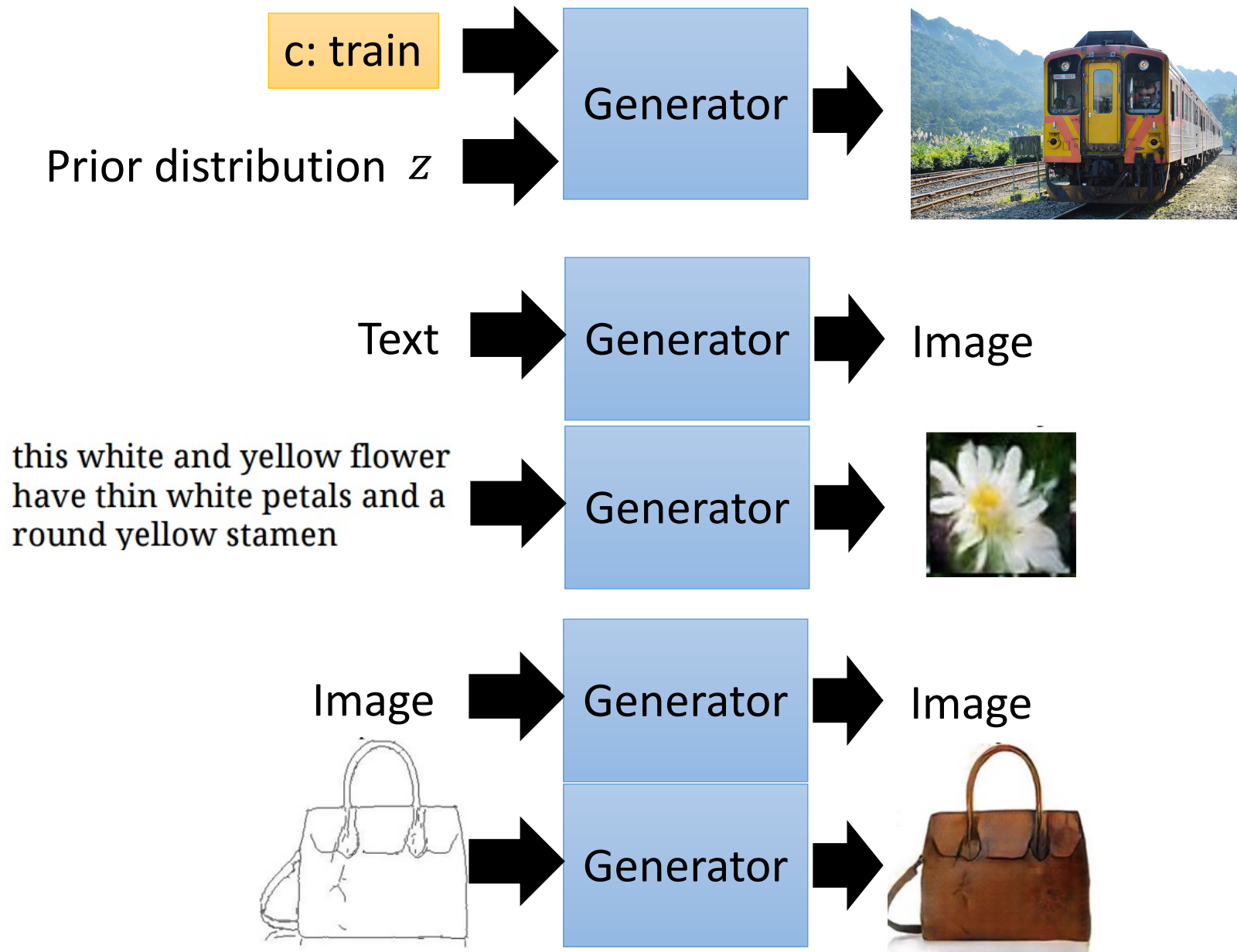# Generative models
## Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. **Conditional GANs**

# Generative models
# Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
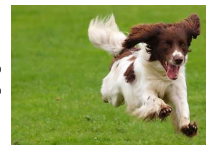5. **Conditional GANs**
   1. **Principle**

# Motivation



c: train

Prior distribution $z$

Generator

Text → Generator → Image

this white and yellow flower have thin white petals and a round yellow stamen → Generator

Image → Generator → Image

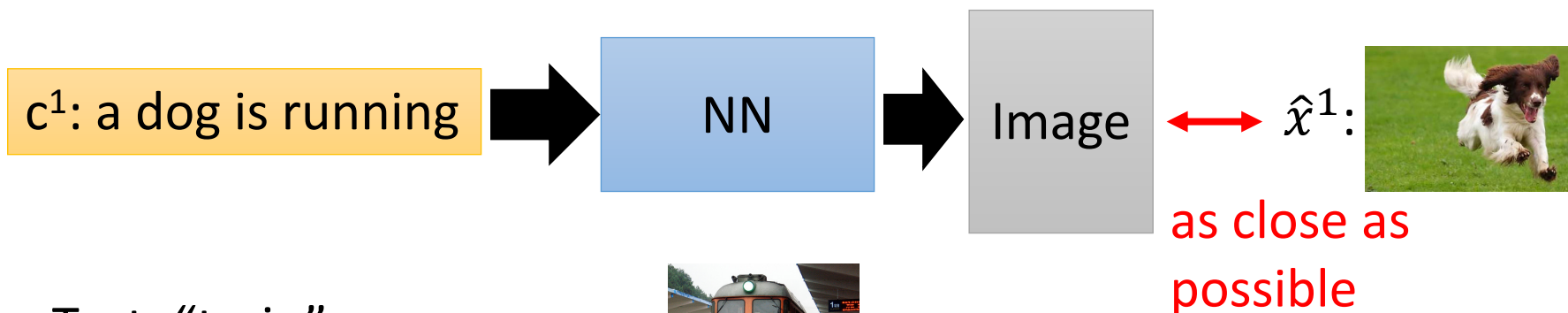# Conditional GAN

$c^1$: a dog is running    $\hat{x}^1$: 

$c^2$: a bird is flying    $\hat{x}^2$: 

- ***Text to image*** by traditional supervised learning

$c^1$: a dog is running → NN → Image ↔ $\hat{x}^1$: 

as close as possible

Text: "train"

Target of NN output

A blurry image!

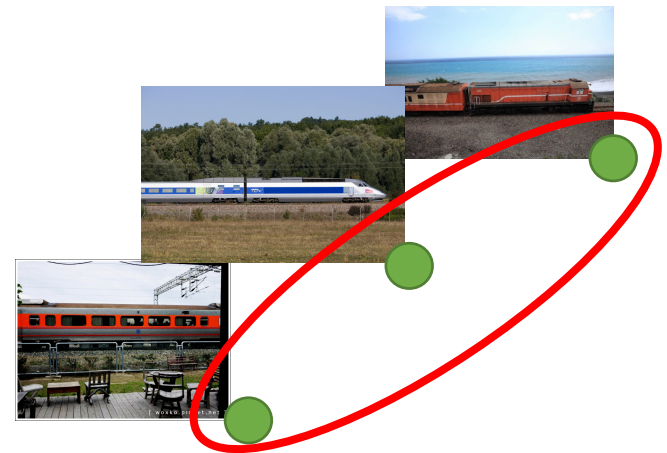# Conditional GAN



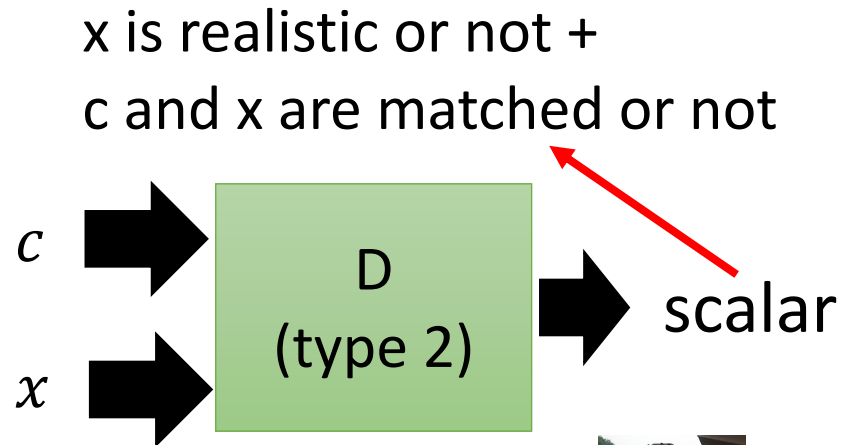c: train → G → Image $\quad x = G(c,z)$

Prior distribution $z$ →

It is a distribution

Approximate the distribution of real data

Text: "train"

# Conditional GAN

c: train → G

Prior distribution $z$ →

G → Image    $x = G(c,z)$

x is realistic or not

$x$ → D (type 1) → scalar

Positive example: 

Negative example: 

x is realistic or not +
c and x are matched or not

$c$ → D (type 2)

$x$ → D (type 2) → scalar

Positive example:    (train ,  )

Negative example:    (train ,  )

Extra neg

(cat ,  )

# Conditional GAN (cGAN model)

GAN

$$V(G, D) = \mathbb{E}_{x \sim P_{data}}[log D(x)] + \mathbb{E}_{x \sim P_G}\big[log\big(1 - D(x)\big)\big]$$
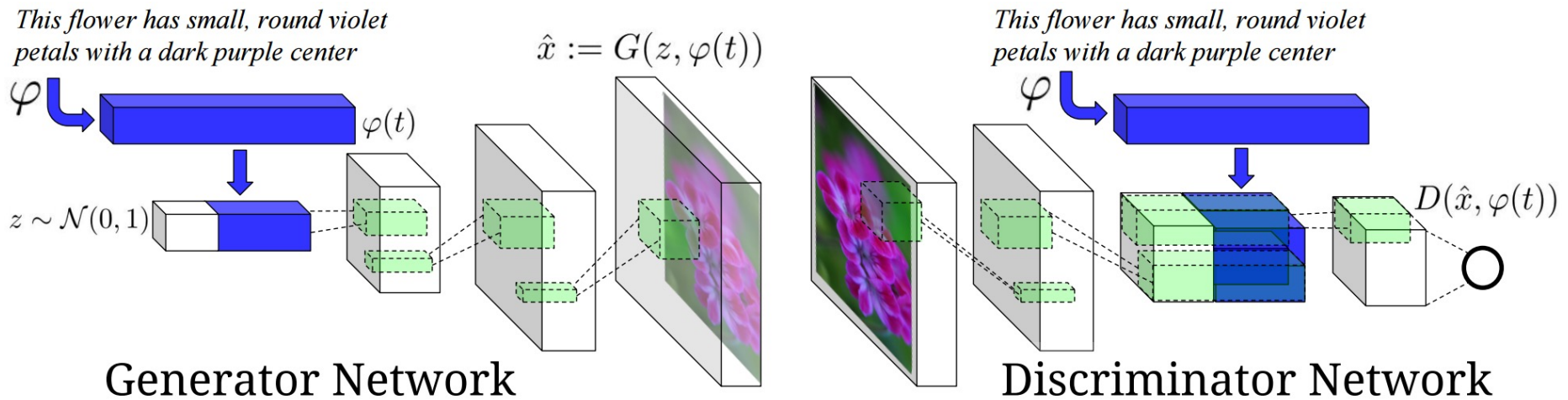
$$G^* = arg \min_G \max_D V(G, D)$$

cGAN

$$\min_G \max_D \left( \mathbb{E}_{\mathbf{x,y} \sim p_{data}(\mathbf{x,y})}[\log D(\mathbf{x,y})] + \mathbb{E}_{\mathbf{y} \sim p_{\mathbf{y}}, \mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})}[\log(1 - D(G(\mathbf{z,y}), \mathbf{y}))] \right)$$

# Generative models
# Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
   1. Principle
   2. **Text2Image**

# Text2Image: architecture example



- Positive samples:
  - real image + right texts
- Negative samples:
  - fake image + right texts
  - Real image + wrong texts

# Text2Image results



this small bird has a pink breast and crown, and black primaries and secondaries.

this magnificent fellow is almost all black with a red crest, and white cheek patch.

the flower has petals that are bright pinkish purple with white stigma

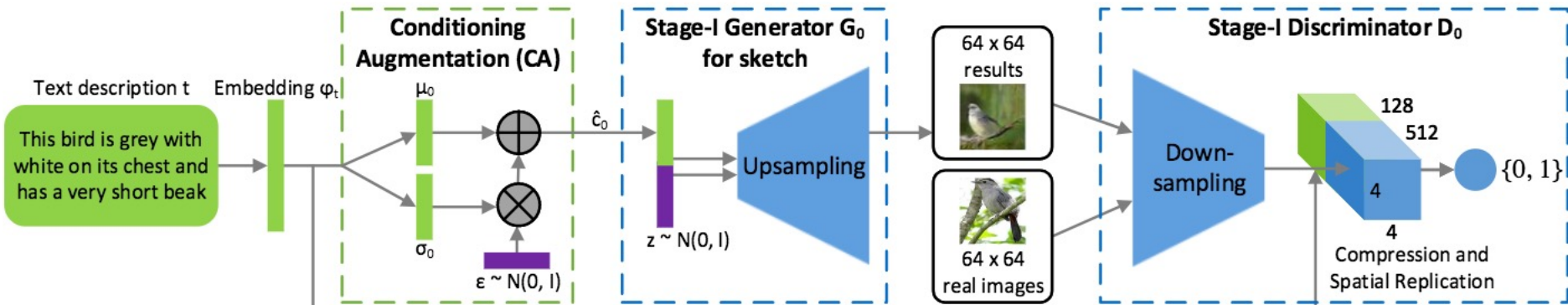this white and yellow flower have thin white petals and a round yellow stamen

# Text2Image results

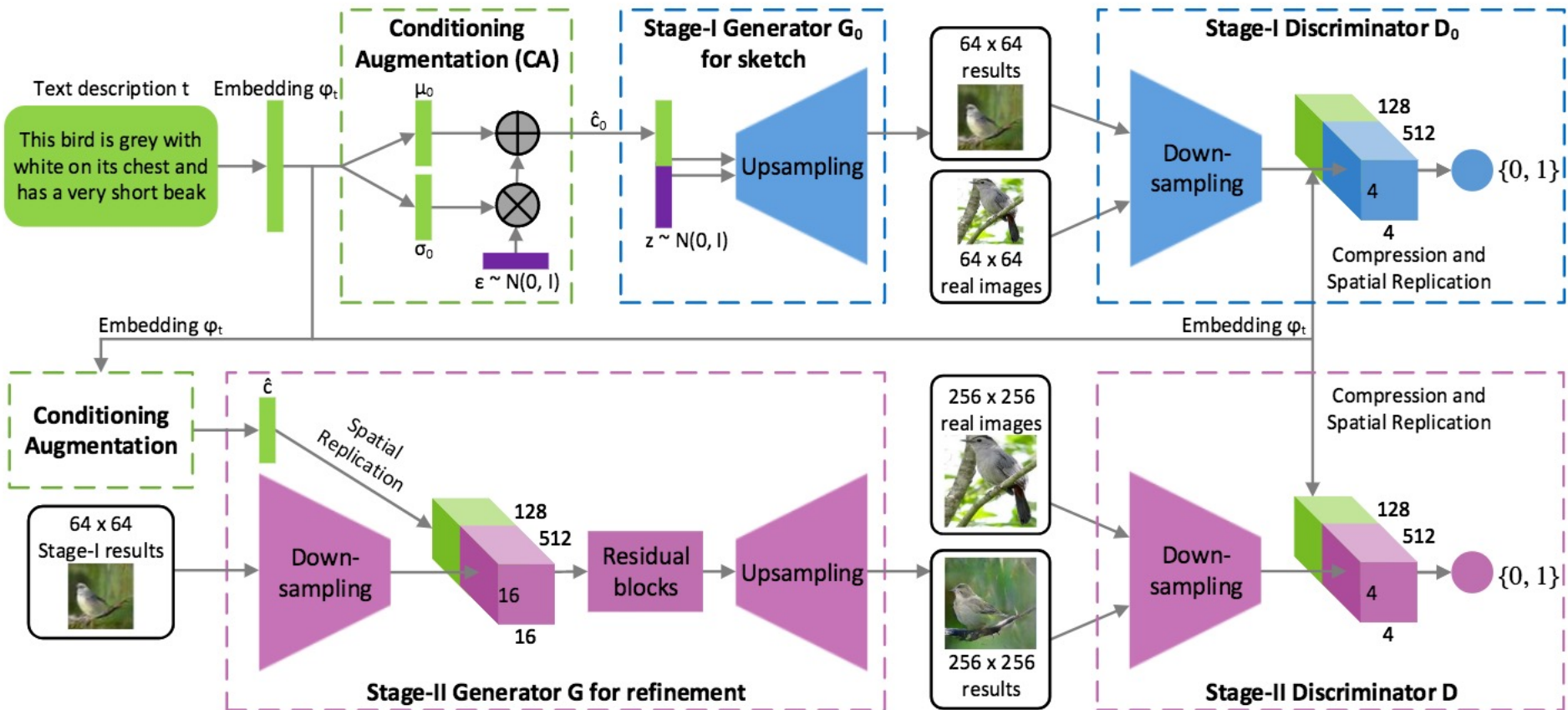| Caption | Image |
|---|---|
| this flower has white petals and a yellow stamen |  |
| the center is yellow surrounded by wavy dark purple petals |  |
| this flower has lots of small round pink petals |  |

# Text2Image: architecture example (2)

## Generating higher resolution images (from 64 to 256)

# Text2Image: architecture example (2)

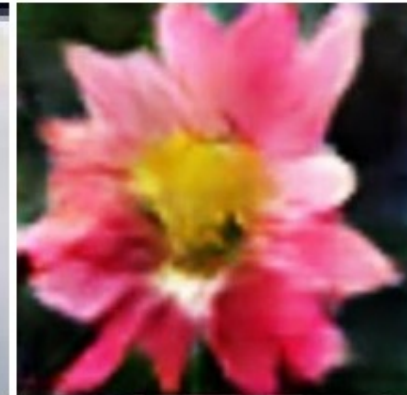## Generating higher resolution images (from 64 to 256)

# StackGAN results



This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face

This bird is white with some black on its head and wings, and has a long orange beak

This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments

(a) Stage-I images
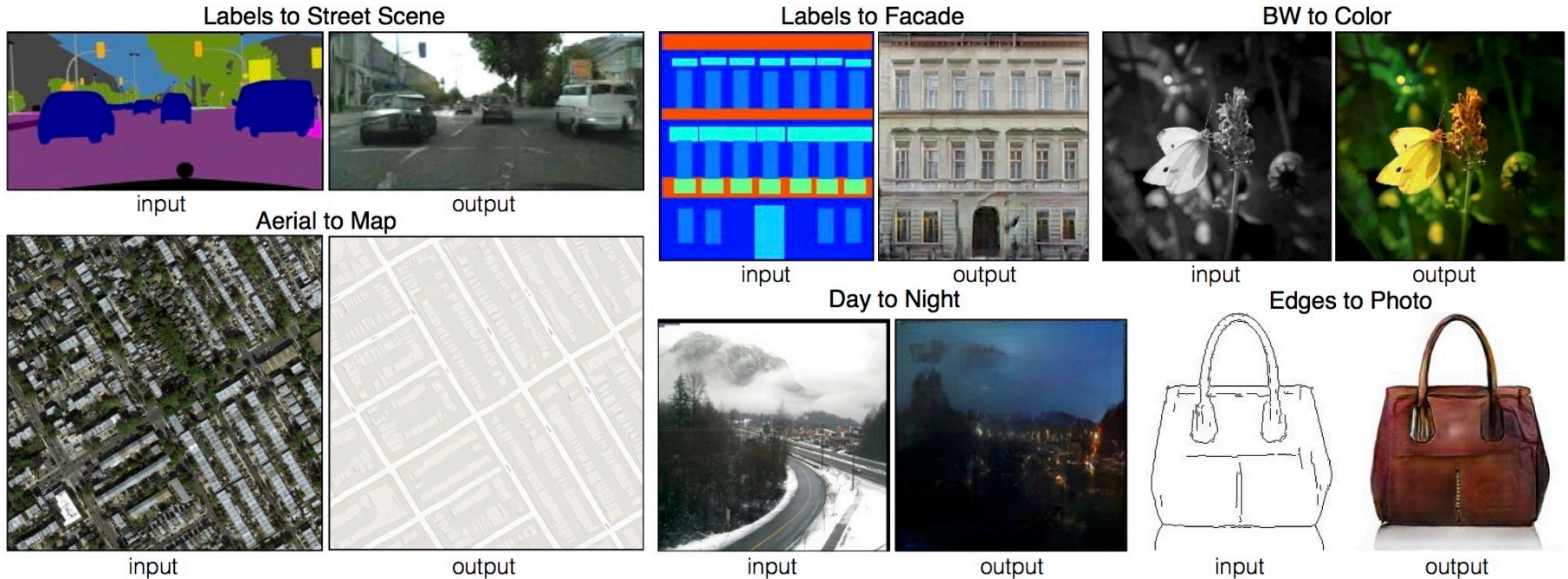
(b) Stage-II images

[Zhang et al. 2016]

# Generative models
## Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
    1. Principle
    2. Text2Image
    3. **Image2Image**

# Image-based Conditional GAN



- Conditioned on an image of different modality
- Image-to-Image Translation => pix2pix

# Image-to-image pix2pix
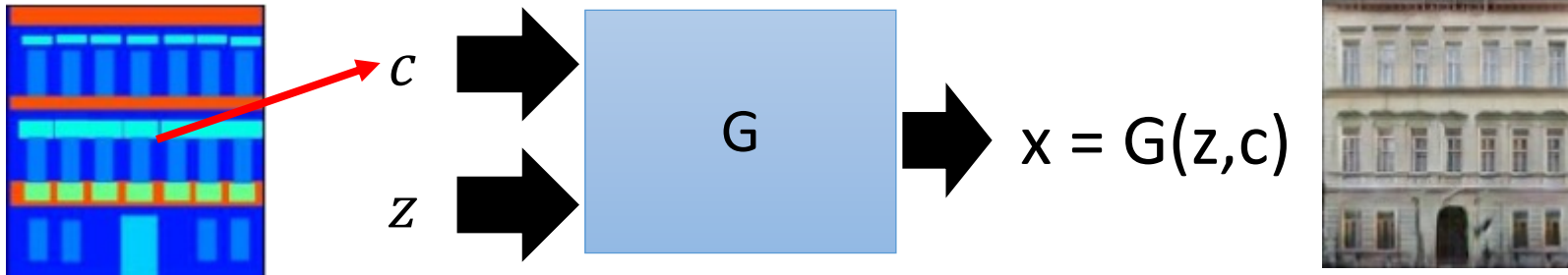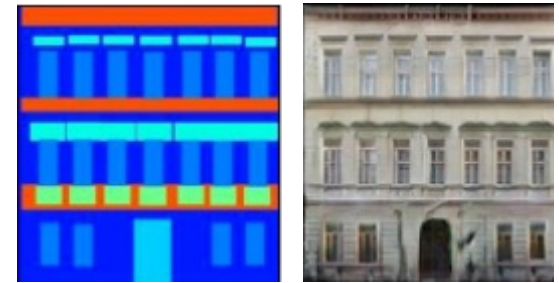


$c$

$z$

G

$x = G(z,c)$

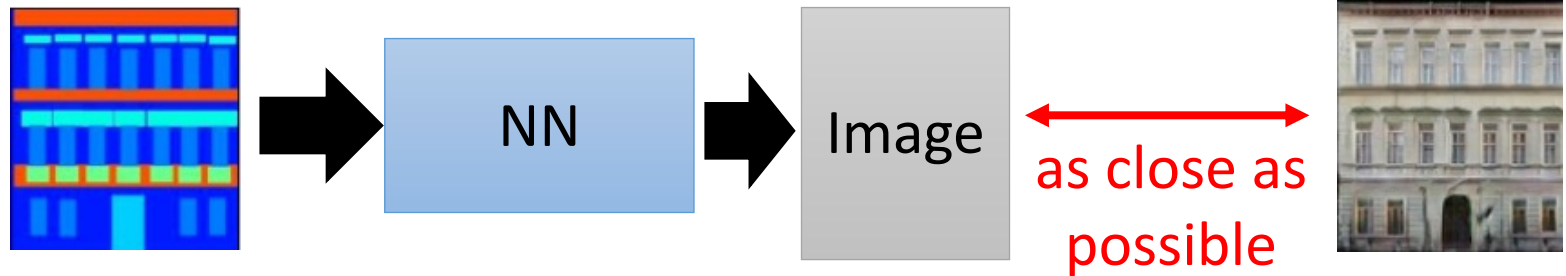# Image-to-image pix2pix



- Traditional supervised approach

 → NN → Image ↔ as close as possible 

Testing:

  It is blurry because it is the average of several images. 

input        close

# Image-to-image

- Conditional GAN

Supplementary loss (close)

Image $x = G(z,c)$

$z$

G

Image

D

scalar

Testing:

input    close    GAN    GAN + close    GT

# Positive examples

## Real or fake pair?



# Negative examples

## Real or fake pair?



**G** tries to synthesize fake
images that fool **D**

**D** tries to identify the fakes

# Label2Image



| Input | Ground truth | L1 | cGAN | L1 + cGAN |

# Edges2Image

# Pix2pixHD



*16 x 16 x c3*

*4 x 4 x c1* *8 x 8 x c2*

latent vector

$g_{gen}$ r $g^1$ r $g^2$ r

*16 x 16 x c3'* *8 x 8 x c2'* *4 x 4 x c1'*

Critic-loss function

$d^{k-2}$ $d^{k-1}$ $d^k$

16 x 16

8 x 8

## *Coarse-to-fine Generator*



Residual blocks

Residual blocks

↓ 2

## *Multi-scale Discriminators*



$D_3$

$D_2$

$D_1$

real              synthesized

## *Robust Objective*



$D_i$              $D_i$

real

*match*

*match*

real              synthesized

Semantic Map

pix2pix

CRN

Ours

# Improving Segmentation2Image strategy?

Limitation of the approach:


pix2pixHD

Directly feed the semantic layout as input to the deep network, which is processed through stacks of convolution, normalization, and nonlinearity layers.

However, this is suboptimal as the normalization layers tend to "wash away" semantic information in input semantic segmentation masks.

# Improving Segmentation2Image strategy

Proven effective for recent generative adversarial networks such as StyleGAN

Can we do the same for conditional GAN? **Conditional Normalization Layers?**



(a) Traditional

(b) Style-based generator

# Improving Segmentation2Image strategy

Recall: Adaptive instance normalization



$$\text{AdaIN}(x, y) = \sigma(y) \left( \frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$
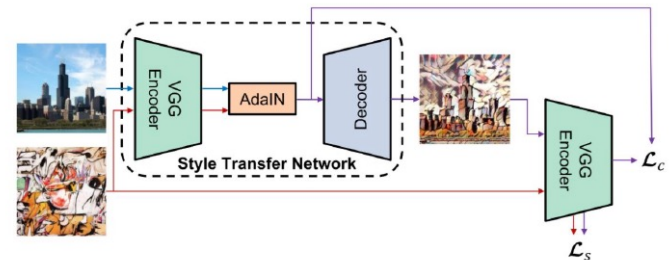
# SPADE block= spatially-adaptive denormalization:
Same idea but per class c over each channel i (N=batch size)

$$\gamma_{c,y,x}^i(\mathbf{m}) \frac{h_{n,c,y,x}^i - \mu_c^i}{\sigma_c^i} + \beta_{c,y,x}^i(\mathbf{m})$$

$$\mu_c^i = \frac{1}{N H^i W^i} \sum_{n,y,x} h_{n,c,y,x}^i$$

$$\sigma_c^i = \sqrt{\frac{1}{N H^i W^i} \sum_{n,y,x} (h_{n,c,y,x}^i)^2 - (\mu_c^i)^2}$$



**SPADE paper = [Semantic Image Synthesis with Spatially-Adaptive Normalization CVPR 2019]**

# SPADE Generator



SPADE ResBlk

pix2pixHD

# SPADE Generator



Better preserve semantic information against common normalization layers

# SPADE results



| Label | Ground Truth | Multi-modal results |

# Spade and follow-up approaches



| Label map | Ground truth | SPADE | CC-FPSE | OASIS |

# Editing with *conditional-or-structured-latent* GANs



**BlobGan** [Epstein *et al.* ECCV'22]

Blob appearance features
$\psi$

$\phi = c_x, c_y, s, a, \theta$
Blob spatial parameters

$z \sim \mathcal{N}(0, I)$

Layout network

$G$

$\frac{s}{\sqrt{a}}$   $s\sqrt{a}$   $x$   $\theta$

# Editing with *conditional-or-structured-latent* GANs



| Original blobs | Original image | Edited blobs | Edited image |

Blob 30

Blob 29

Blob 30

Blob 29

# Editing with *conditional-or-structured-latent* GANs

Example of Counterfactual optimization for editing

# Generative models
# Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
    1. Principle
    2. Text2Image
    3. Image2Image
    4. **Inpainting and general missing data encoder**

# Inpainting task

- Complete the missing part



Inpainting

# Inpainting as unsupervised learning with GAN loss



Reconstruct missing pixels by decoding using context
Loss defined on the predicted patch and the real one (known at training time)

# First proposition -- Architecture

- Architecture: Encoder/Fully connected/Decoder



- DC-GAN for inpainting task
- **Input:** $227 \times 227 \times 3$ image
- **Output:** encoder context features $(6 \times 6 \times 256)$

# Channel-wise fully-connected layer

- **Input / output:** 6 × 6 × 256 channels
- **First layer:** Channel-wise fully-connected
  (each 6 × 6 input connected to the corresponding 6 × 6 output)
- **Second layer:** Stride 1 convolution to mix channels



# Decoder

- **Architecture:** Same as DC-GAN: 5 up-convolutional layers
  ("deconv" + ReLU)
- **Input:** decoder context features 6 × 6 × 256
- **Output:** 227 × 227 × 3 image

# Training: Masking the images

- **How to define the mask ?**

  ▶ Center region of the image
  ▶ Random regions (chosen solution)
  ▶ Random segmentation mask from VOC (said to be equivalent to random regions)

- **Formal definition:** Defined by a mask $\hat{M} \in \{0,1\}^{227 \times 227}$ with 1 if the pixel should be masked

# Training: Loss - Overview

- Trained completely from scratch to fill-up the masked areas
- **Problem:** multiple plausible solutions
- **Solution:** combining 2 losses:
  - ▶ $\mathcal{L}_{rec}$ **L2 reconstruction loss:** learn the structure of the missing region (average multiple modes in prediction)
  - ▶ $\mathcal{L}_{adv}$ **Adversarial loss:** make it look real (pick a mode from the distribution)

$$\min_F \mathcal{L} = \lambda_{rec}\mathcal{L}_{rec} + \lambda_{adv}\mathcal{L}_{adv}$$

$$\mathcal{L}_{rec}(x) = \left\| \hat{M} \odot \left( x - F((1 - \hat{M}) \odot x) \right) \right\|_2$$

$$\mathcal{L}_{adv} = \max_D \mathbb{E}_{x \in \mathcal{X}} \left[ \log(D(x)) + \log \left( 1 - D(F((1 - \hat{M}) \odot x)) \right) \right]$$

- Rq: The encoder-decoder is the generator, D is a CNN

# Results

**Dataset:** StreetView Paris and ImageNet



| Image | Ours($L2$) | Ours(Adv) | Ours(L2+Adv) |

# Semantic inpainting - Qualitative results

# Generalizing inpainting: missing data encoder



(1) inpainting

(2) reverse intpainting

(3) colorization

(4) random extrapolation

(5) random extrapolation + colorization



original image

$x^R$   $x^R + w^R$

generator

completed image

# Results

# Generative models

## Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
    1. Principle
    2. Text2Image
    3. Image2Image
    4. Inpainting and general missing data encoder
    5. **Learning unpaired Transformation**

# Unpaired Transformation

***paired data***



Transform an object from one domain to another ***without paired data***



photo → van Gogh

Domain X

Domain Y

# Cycle GAN

https://arxiv.org/abs/1703.10593
https://junyanz.github.io/CycleGAN/

Domain X

Domain Y

Domain X

Become similar
to domain Y

$G_{X \rightarrow Y}$

Not what we want

ignore input

$D_Y$ → scalar

Input image
belongs to
domain Y or not

Domain Y

Cycle GAN

Domain X

Domain Y

as close as possible

Lack of information for reconstruction

$G_{X \to Y}$

$G_{Y \to X}$

$D_Y$ → scalar

Input image belongs to domain Y or not

Domain Y

Cycle GAN

Domain X
Domain Y

as close as possible

$G_{X \to Y}$

$G_{Y \to X}$

scalar: belongs to domain X or not

$D_X$

$D_Y$

scalar: belongs to domain Y or not

$G_{Y \to X}$

$G_{X \to Y}$

as close as possible

# Results -- Cycle GAN



photo → van Gogh

Domain X

Domain Y

Monet ⟳ Photos

Monet → photo

photo → Monet

Zebras ⟳ Horses

zebra → horse

horse → zebra

Summer ⟳ Winter

summer → winter

winter → summer

# GANs: works in progress

A lot of things to better understand, to use, adapt, …



Patch co-occurrence discriminator $D_{\text{patch}}$
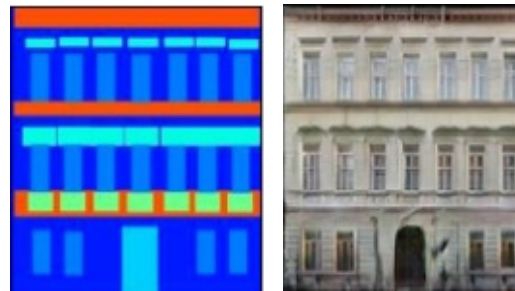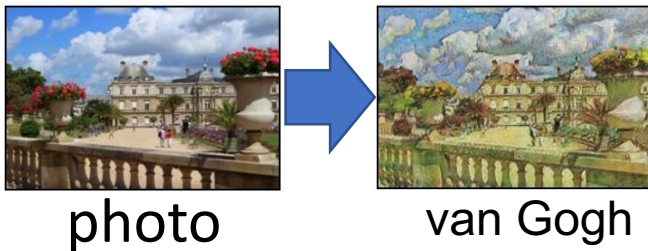
# Generative models

## Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
6. **Diffusion models**

# Generative Modelling with Diffusion models



Gaussian Distribution

Noising Process

Generative Process

Distribution of real images

# Generative Modelling with Diffusion models

*DDPM: Denoising Diffusion Probabilistic Models*
*In context with other genenrative Models:*

# Generative Modelling with Diffusion models

*DDPM: Denoising Diffusion Probabilistic Models*

$$\mathbf{x}_t = \sqrt{1 - \beta_t}\,\mathbf{x}_{t-1} + \sqrt{\beta_t}\,\epsilon \qquad \mathbf{x}_t = \sqrt{\alpha_t}\,\mathbf{x}_0 + \sqrt{1 - \alpha_t}\,\epsilon$$



Training:

*The reverse process refers to learning a model* $\epsilon_\theta$ *which approximates the noise added at a given timestep t:*

$$\mathcal{L} = \mathbb{E}_{\mathbf{x_0}, t, \epsilon}\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|_2^2$$

Sampling:

$$\hat{\mathbf{x}}_0 = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \sqrt{1 - \alpha_t}\cdot\epsilon_\theta(\mathbf{x}_t, t)),$$

$$\mathbf{x}_{t-1} = \frac{(\alpha_{t-1} - \alpha_t)\sqrt{\alpha_{t-1}}}{\alpha_{t-1}(1 - \alpha_t)}\hat{\mathbf{x}}_0 + \frac{(1 - \alpha_{t-1})\sqrt{\alpha_t}}{(1 - \alpha_t)\sqrt{\alpha_{t-1}}}\mathbf{x}_t + \sigma_t\mathbf{z}$$

# Generative Modelling with Diffusion models

## *DDPM: Denoising Diffusion Probabilistic Models*

**Algorithm 1** Training

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

**Algorithm 2** Sampling

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

# U-Net for Diffusion models

# Text-conditional Diffusion models

**● Training:**



"A cute cat"

Add noises

Predict (U-Net)

*Predict denoised images*

**● Sampling:**



Gaussian noise

"A cute dog"

Linear Interpolation

Predict (U-Net)

Output

*Repeat N times*

*Generated Image*

📖Denoising Diffusion Probabilistic Models, *Jonathan Ho, Ajay Jain, Pieter Abbeel*, NeurIPS 2020

# Text-conditional Diffusion models

# Text-conditional Diffusion models



Image

Encoder    Latent Vector    Forward    Noisy Latent Vector

Decoder

Generated Image

Perceptual Compression

Denoising Network to recover the original latent vector

Reverse Diffusion Step

Pretrained De-Noising Autoencoder with skip connections for attention

Semantic Compression

Text/Image Transformer

**Latent Space**

$z_{t-1}$ ..... $q(z_t|z_{t-1})$ ..... $z_t$

Diffusion Process

Denoising U-Net $\epsilon_\theta$

$\times (T-1)$

$Q$ $KV$    $Q$ $KV$    $Q$ $KV$    $Q$ $KV$

$z$    $z_{T-1}$

$z_{t-1}$ ..... $p_\theta(z_{t-1}|z_t)$ ..... $z_t$

$z$

$z_T$

$z_T$

$\mathcal{E}$

$x$

$\tilde{x}$

$\mathcal{D}$

**Pixel Space**

$\tau_\theta$

**Conditioning**

Semantic Map

Text

Repres entations

Images

denoising step    crossattention    switch    skip connection    concat

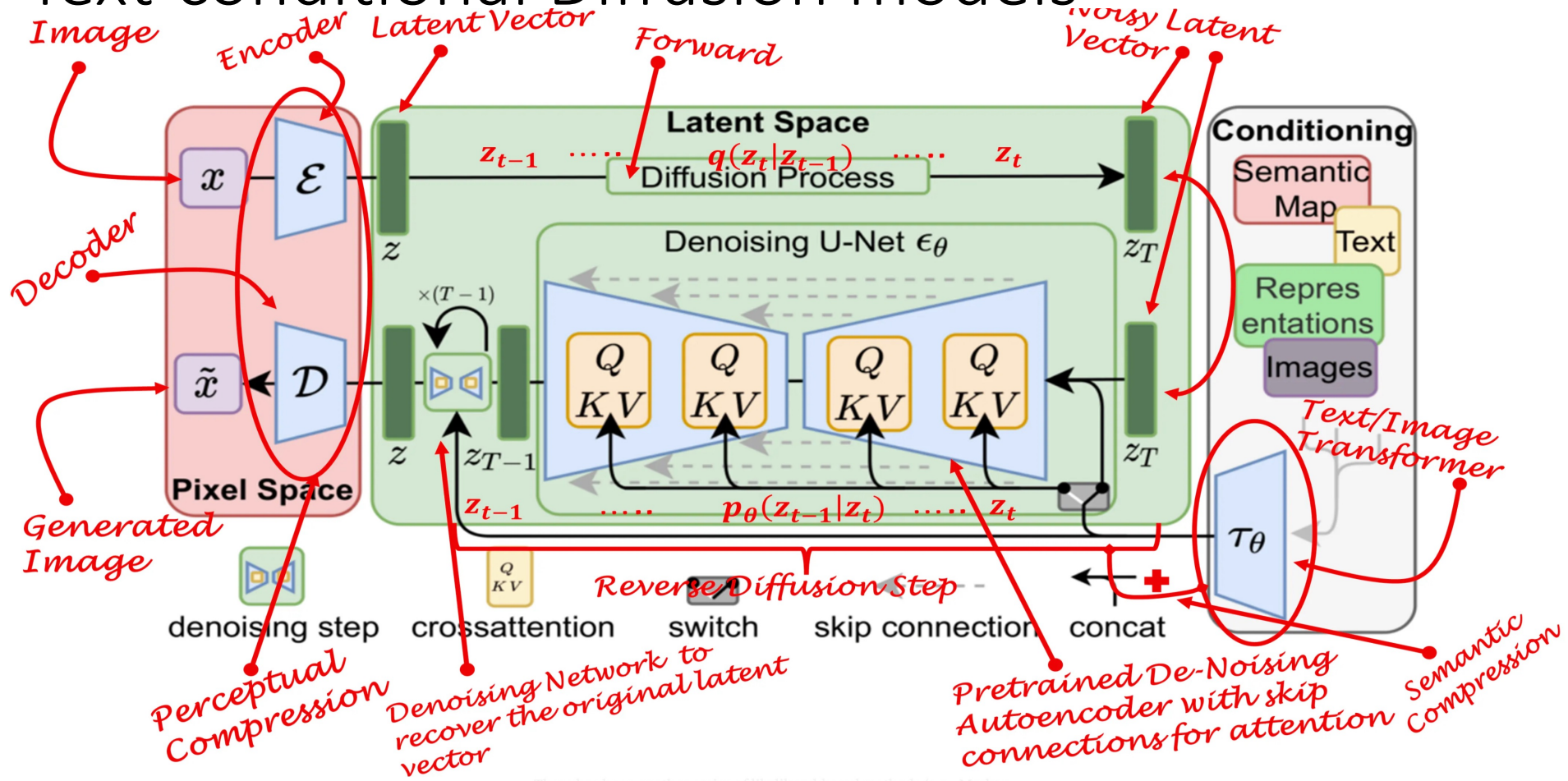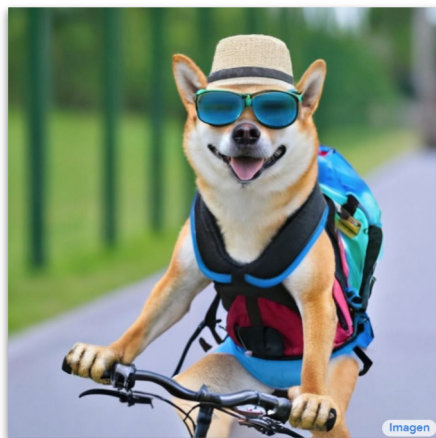# Image Generation using Diffusion Models



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.

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

📖Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding, *Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, Mohammad Norouzi,* NeurIPS 2022

# Generation using Diffusion Models

And the (next) big thing is …

Text to Video



Transition      Transition      Animation      Prediction