

## Module 09 – Project Class

# Text-Guided Image Generation Using Conditional Flow Matching

[Code - Data](#)

Nguyen Quoc Thai  
MSc in Computer Science

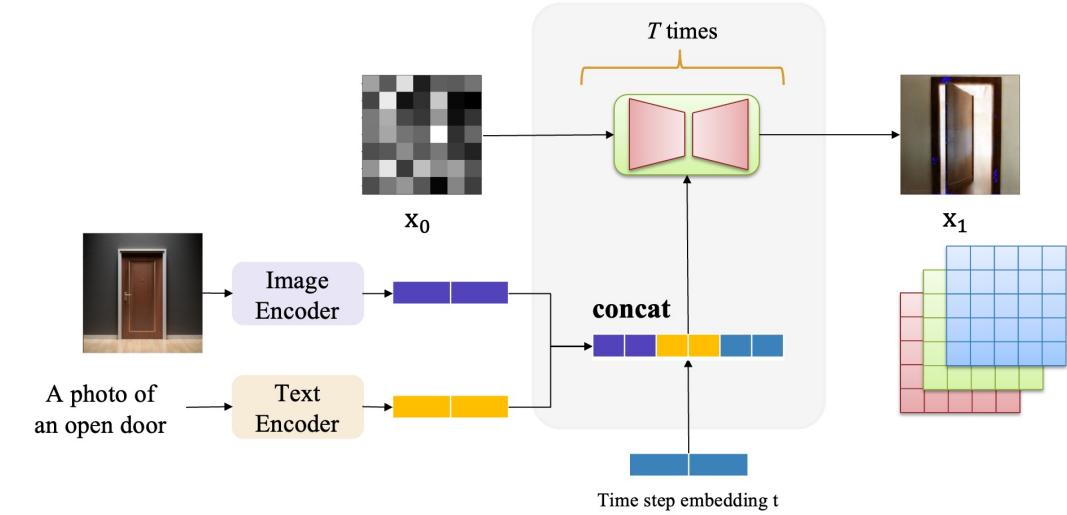
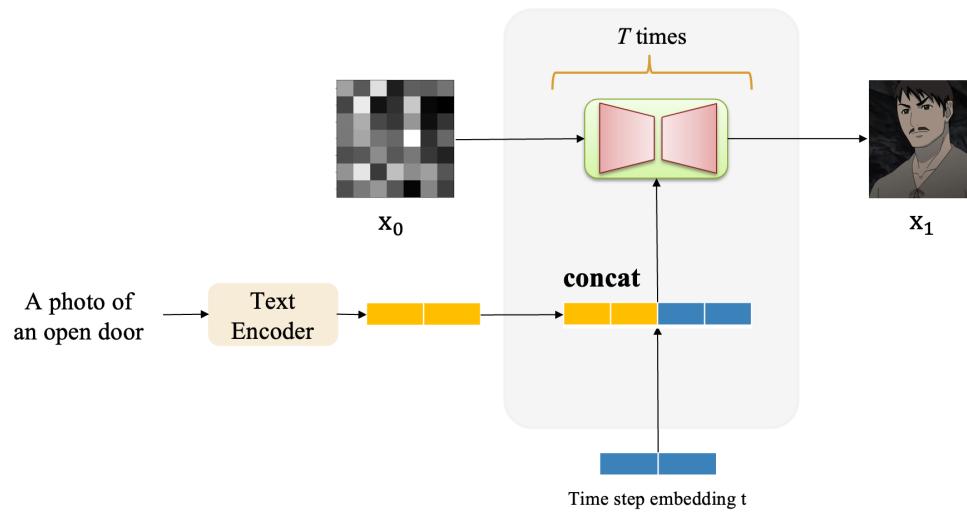
# Objectives

## Conditional Flow Matching (Review)

- ❖ Flow Matching Model
- ❖ Training & Sampling
- ❖ Conditional Flow Matching

## Text-Guided Image Generation

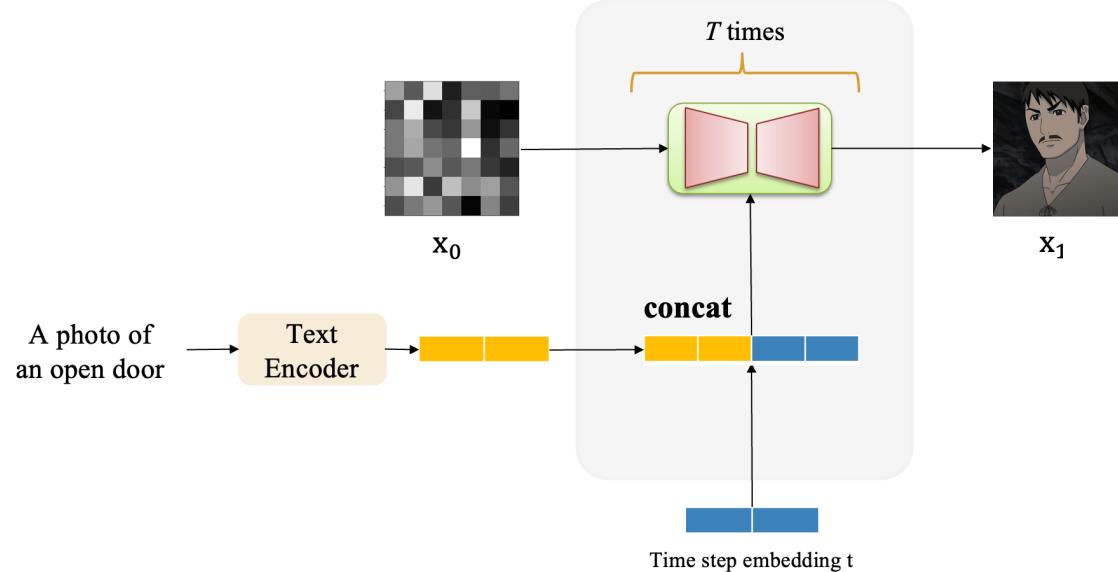
- ❖ Text-To-Image Generation
- ❖ Text-Guided Image Generation
- ❖ TedBench Dataset



# Outline

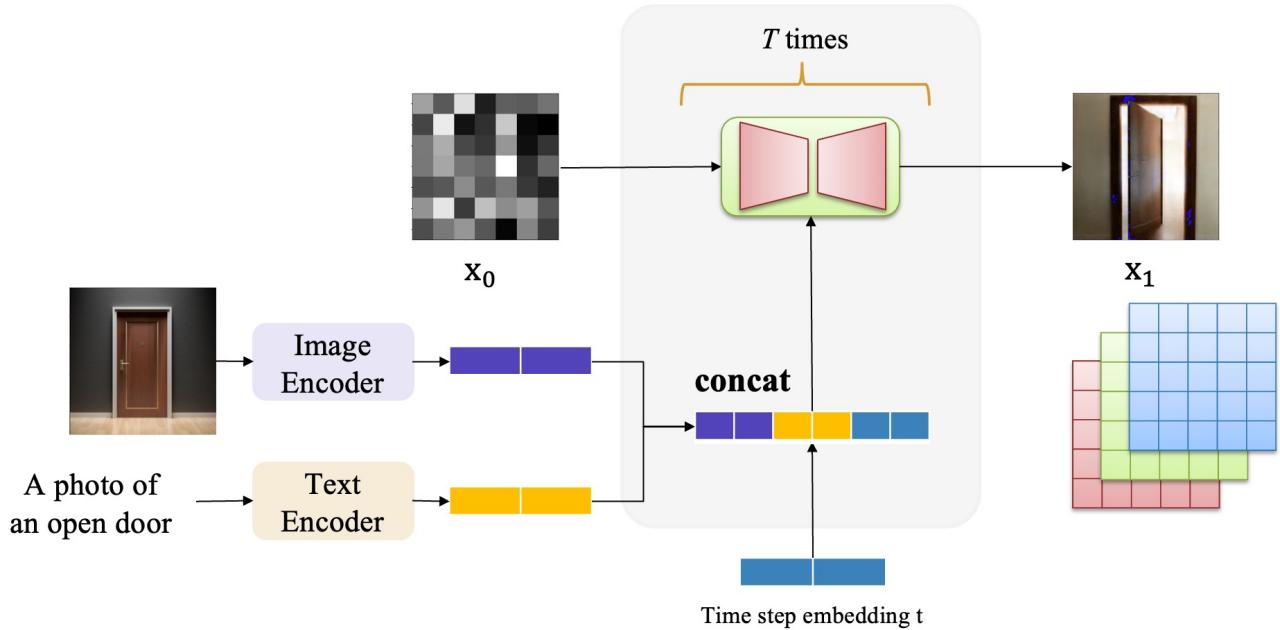
SECTION 1

## Conditional Flow Matching



SECTION 2

## Text-Guided Image Generation

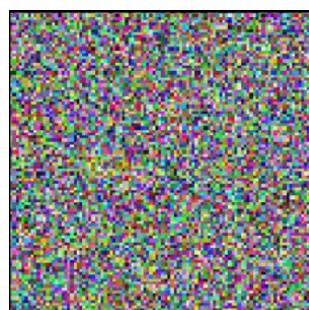


# Conditional Flow Matching

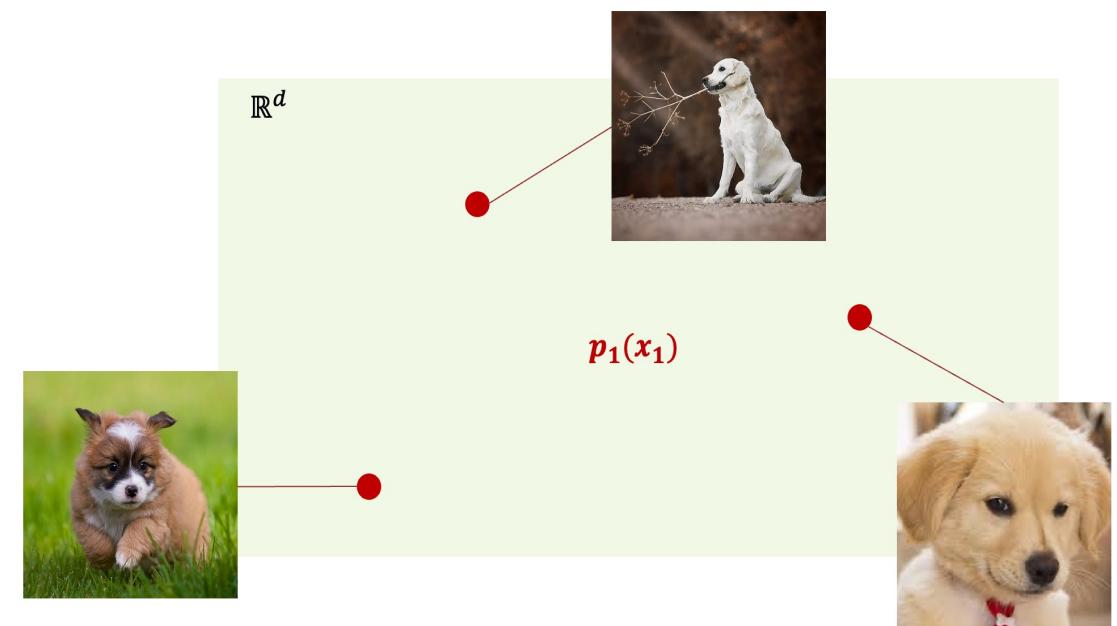
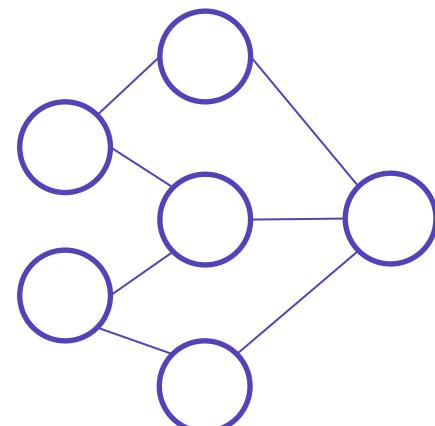


## The Generative Modeling

- Generate new samples from an underlying distribution (the training data)



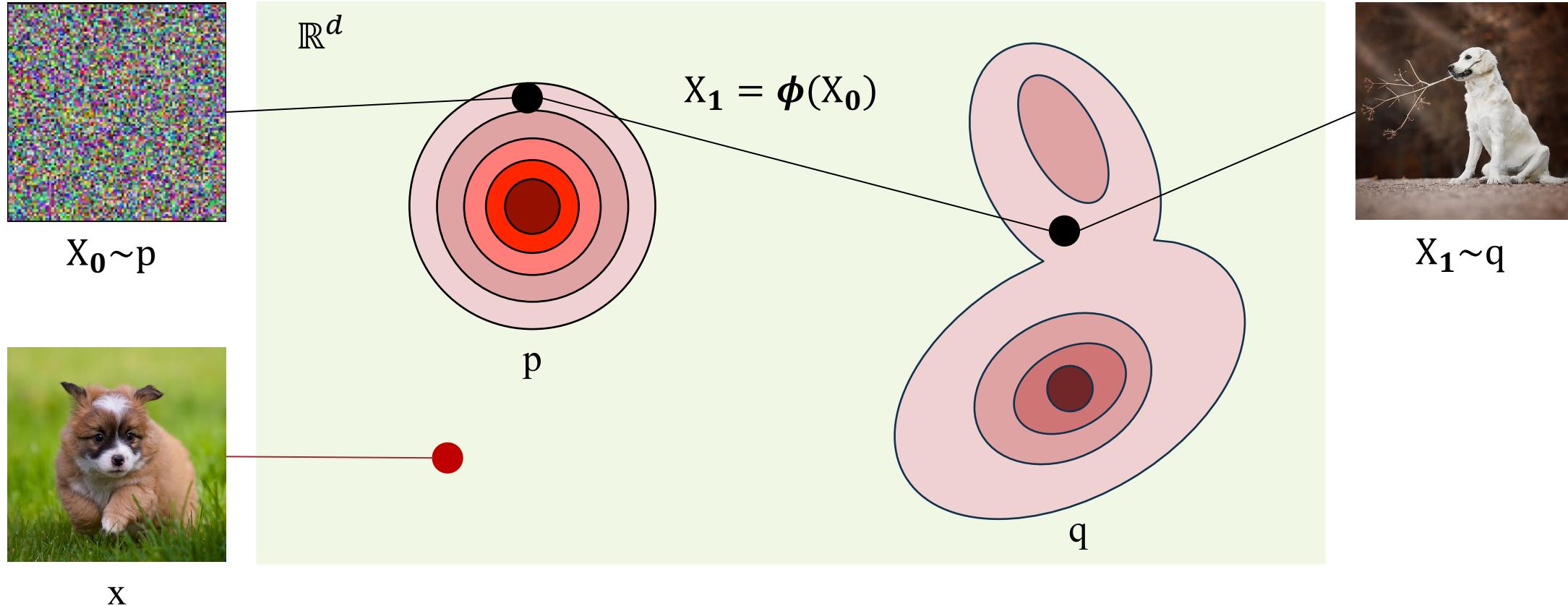
$p_0(x_0)$



# Conditional Flow Matching



## The Generative Modeling



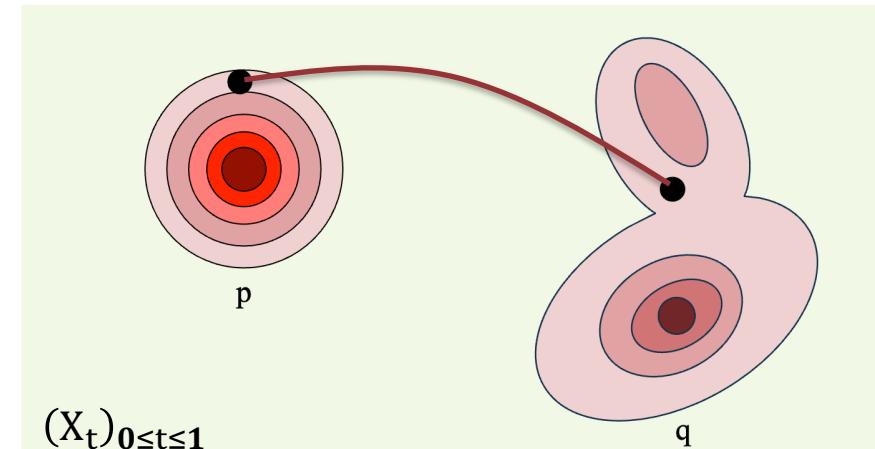
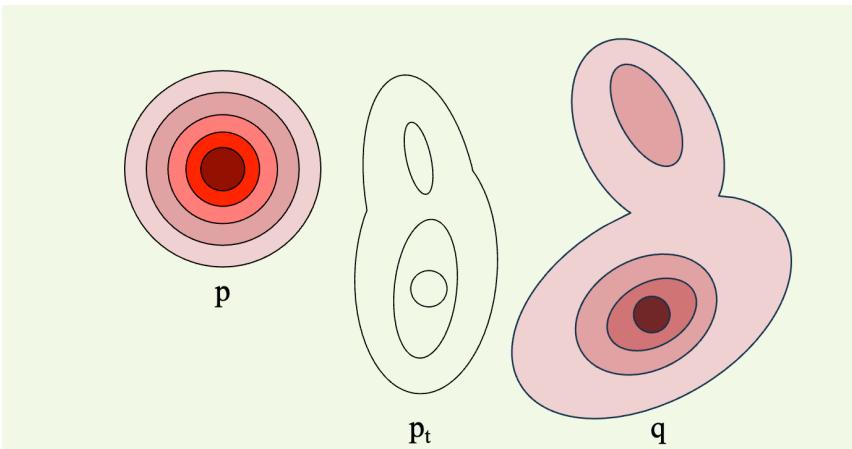
# Conditional Flow Matching



## Flow Matching

- A scalable method to train flow generative models
- Train by regressing a velocity, sample by following the velocity

$$\begin{aligned} & (X_t)_{0 \leq t \leq 1} \\ & X_{t+h} = \phi_{t+h|t}(X_t) \\ & X_t \sim p_t \end{aligned}$$



Flow Matching

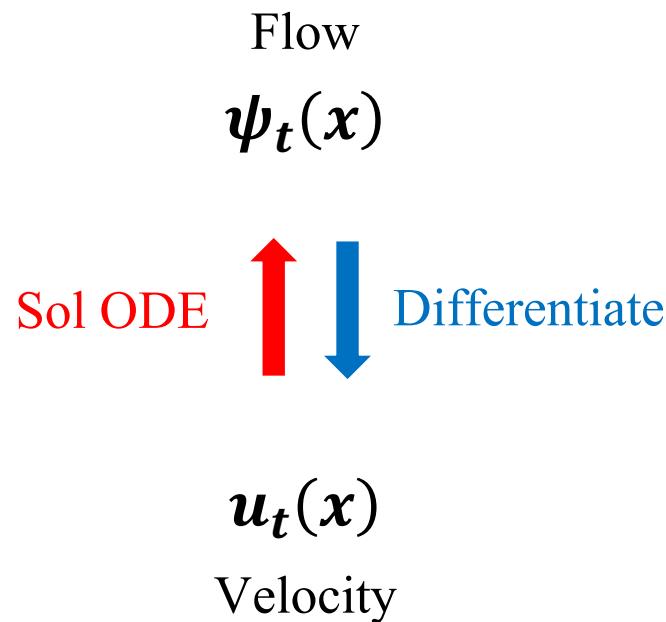
# Conditional Flow Matching



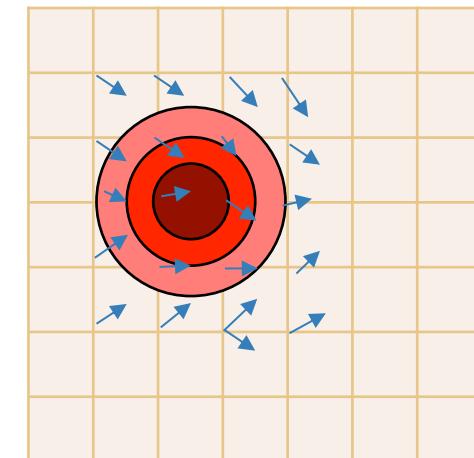
## Velocity in Flow Matching

- Velocity: the difference between the noisy and clean data, the rate of transformation between the two states:

$$\mathbf{v}_t = \mathbf{x}_1 - \mathbf{x}_0$$



Velocity  $u_t$  generates  $p_t$  if  
 $X_t = \psi_t(X_0) \sim p_t$

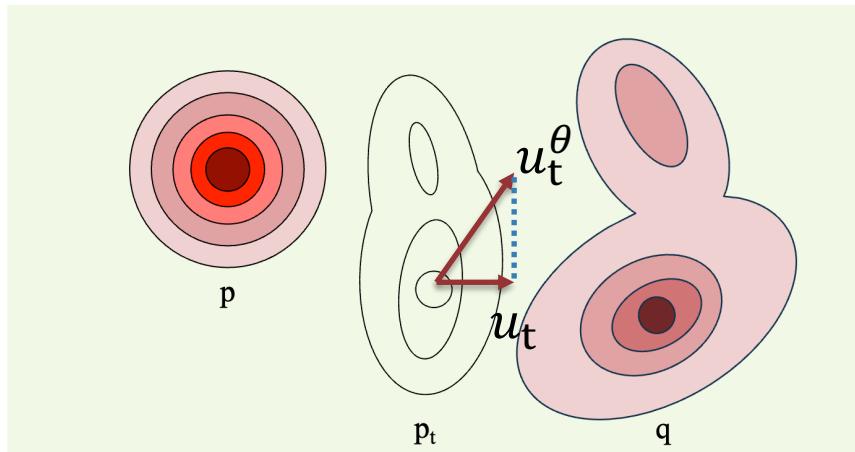


# Conditional Flow Matching

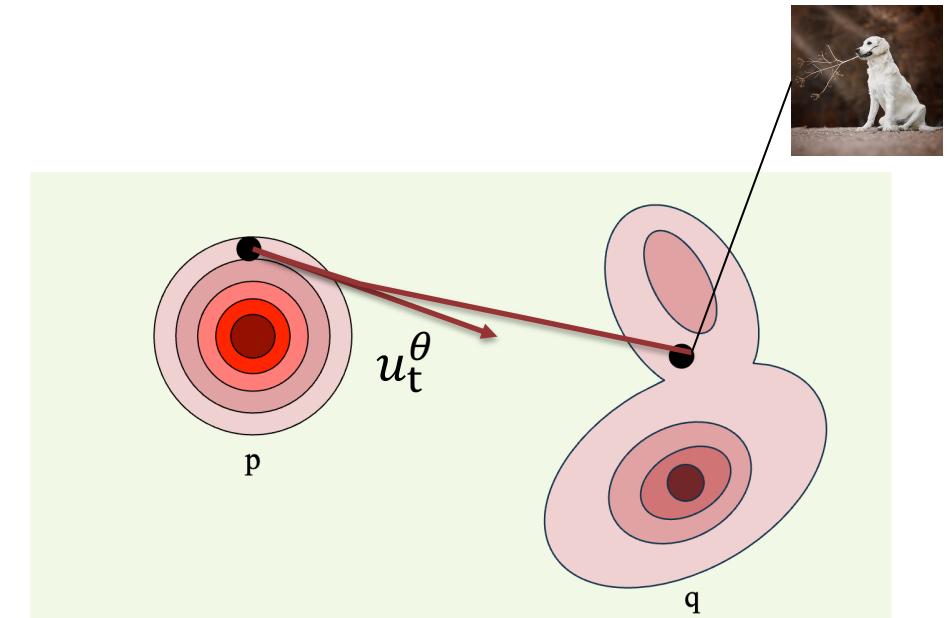


## Flow Matching

Train a velocity generating  $p_t$  with  
 $p_0 = p$  and  $p_1 = q$



Sample from  $X_0 \sim p$



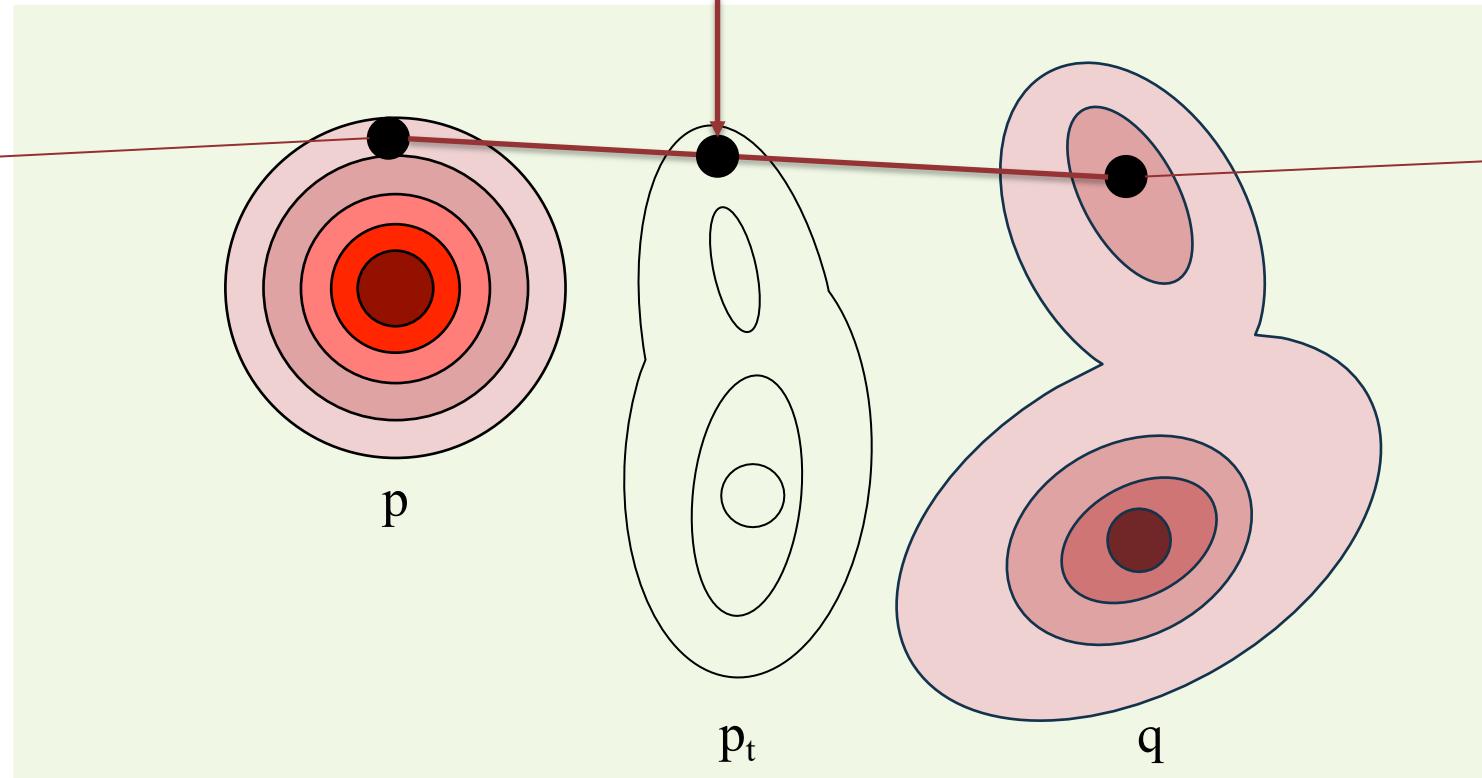
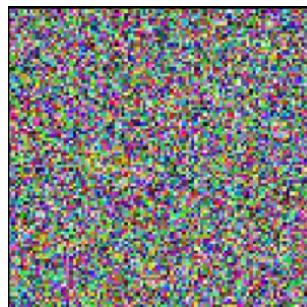
# Conditional Flow Matching



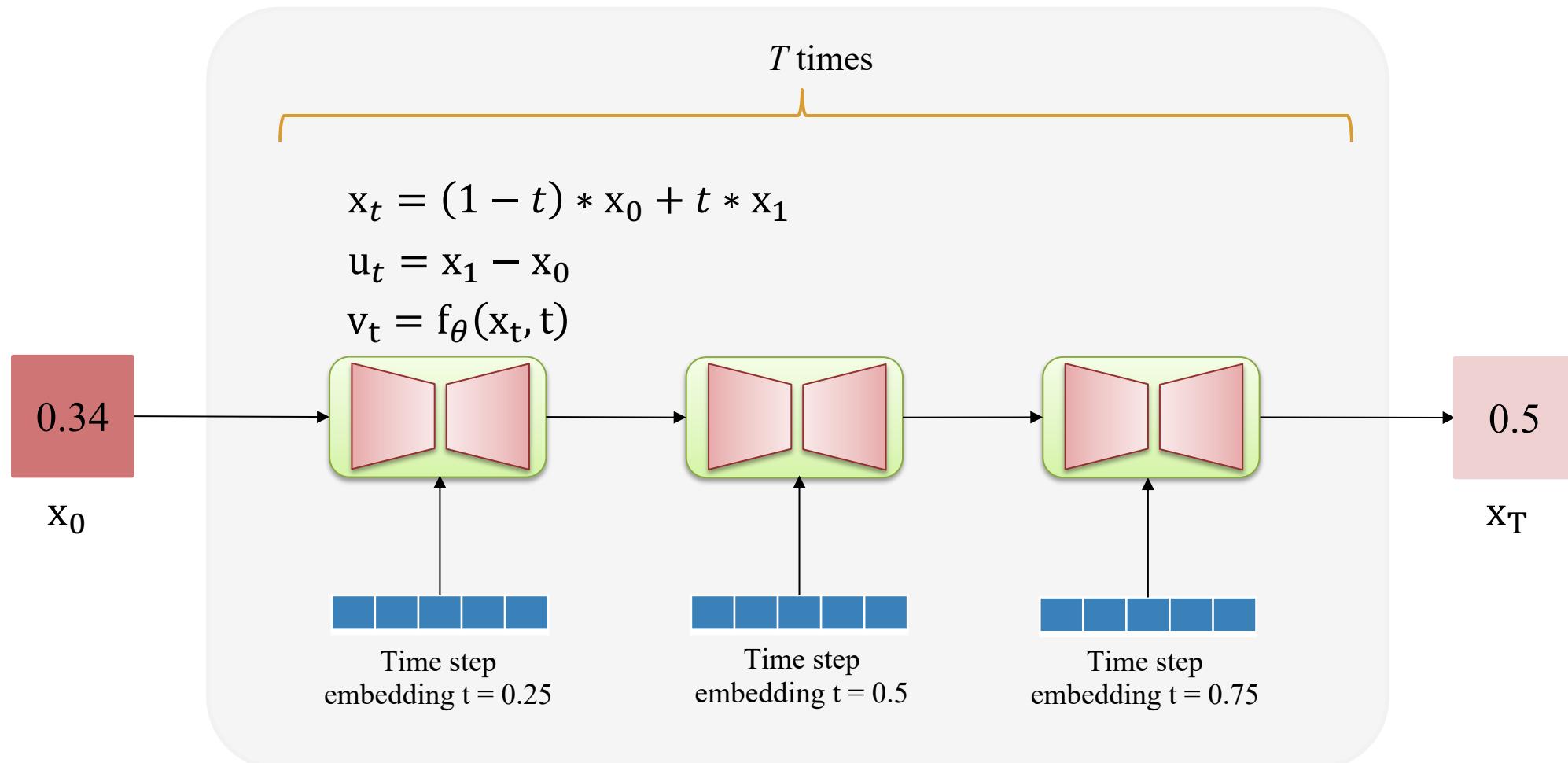
## Training in Flow Matching

$$\mathbb{E}_{t, X_0, X_1} \|u_t^\theta(X_t) - (X_1 - X_0)\|$$

$$X_t = (1 - t) * X_0 + t * X_1$$



# Conditional Flow Matching



# Conditional Flow Matching



## Sampling in Flow Matching

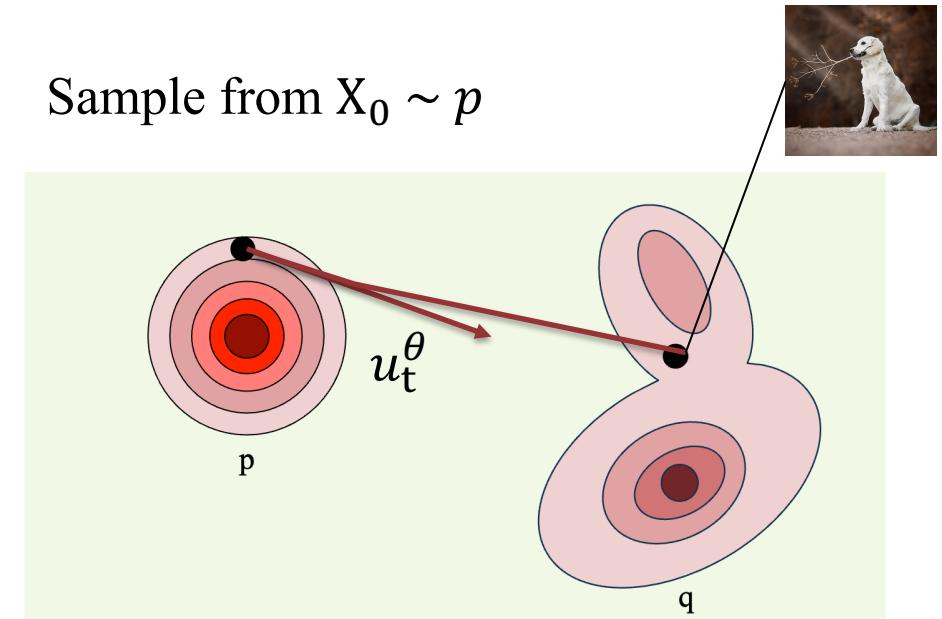
- Euler's method for solving ordinary differential equations (ODEs)

$$\frac{dy}{dt} = f_\theta(y, t)$$

$$y(t + \Delta t) \approx y(t) + \left( \frac{dy(t)}{dt} \right) \Delta t$$

$$y(t + \Delta t) \approx y(t) + f_\theta(y, t) \Delta t$$

Sample from  $X_0 \sim p$



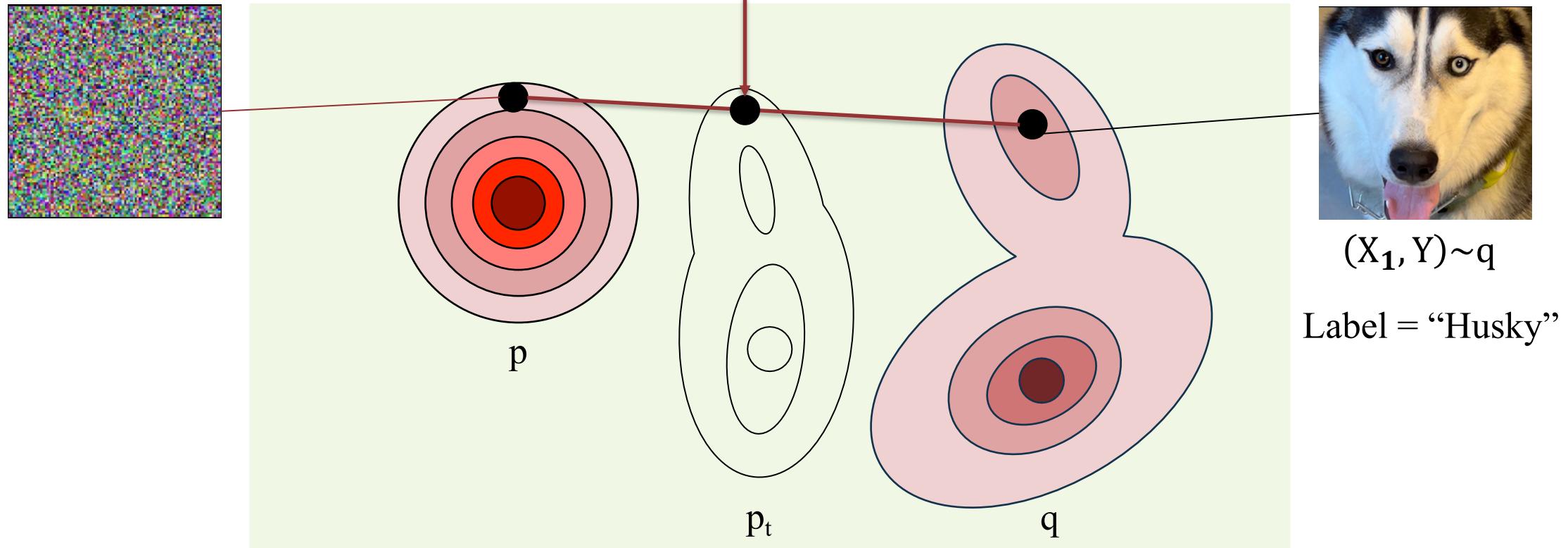
# Conditional Flow Matching



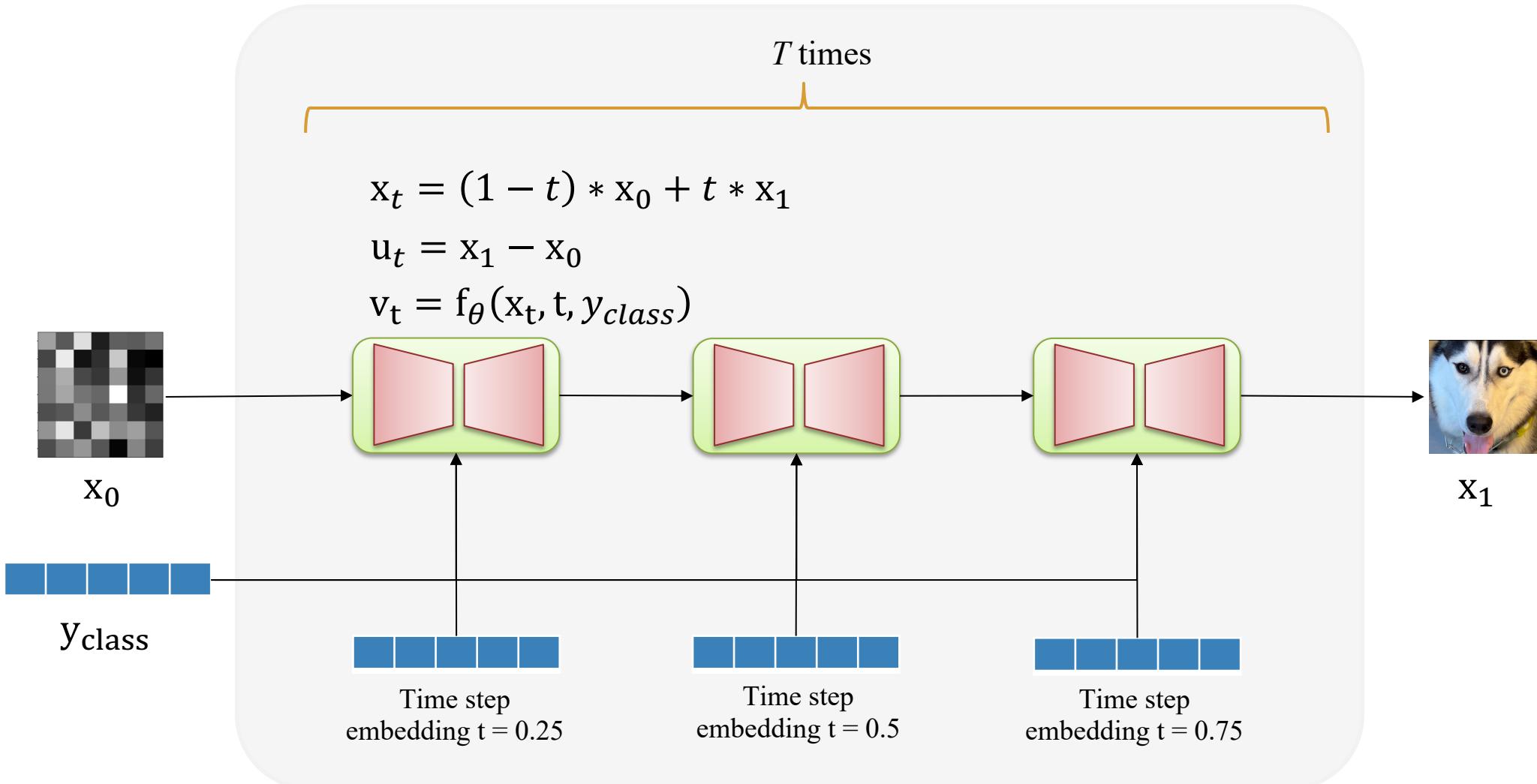
## Training in Conditional Flow Matching

$$\mathbb{E}_{t, X_0, X_1} \|u_t^\theta(X_t) - (X_1 - X_0)\| \longrightarrow \mathbb{E}_{t, X_0, X_1, Y} \|u_t^\theta(X_t, Y) - (X_1 - X_0)\|$$

$$X_t = (1 - t) * X_0 + t * X_1$$



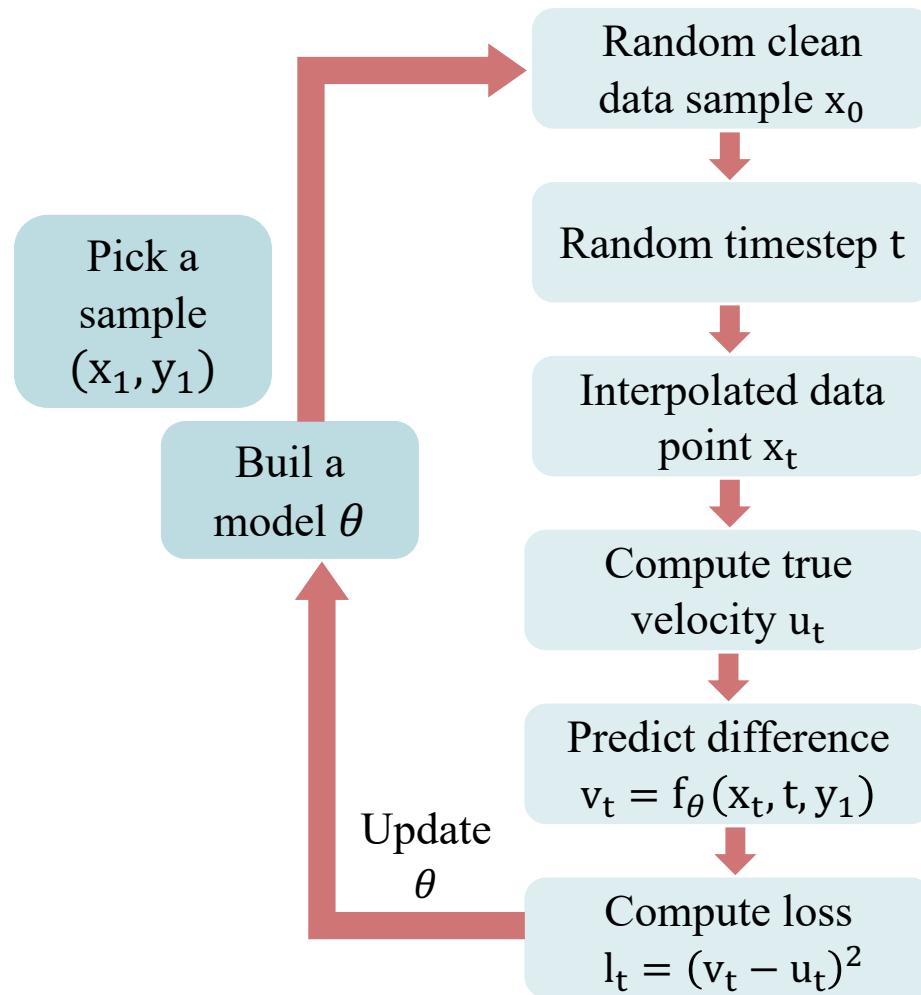
# Conditional Flow Matching



# Conditional Flow Matching

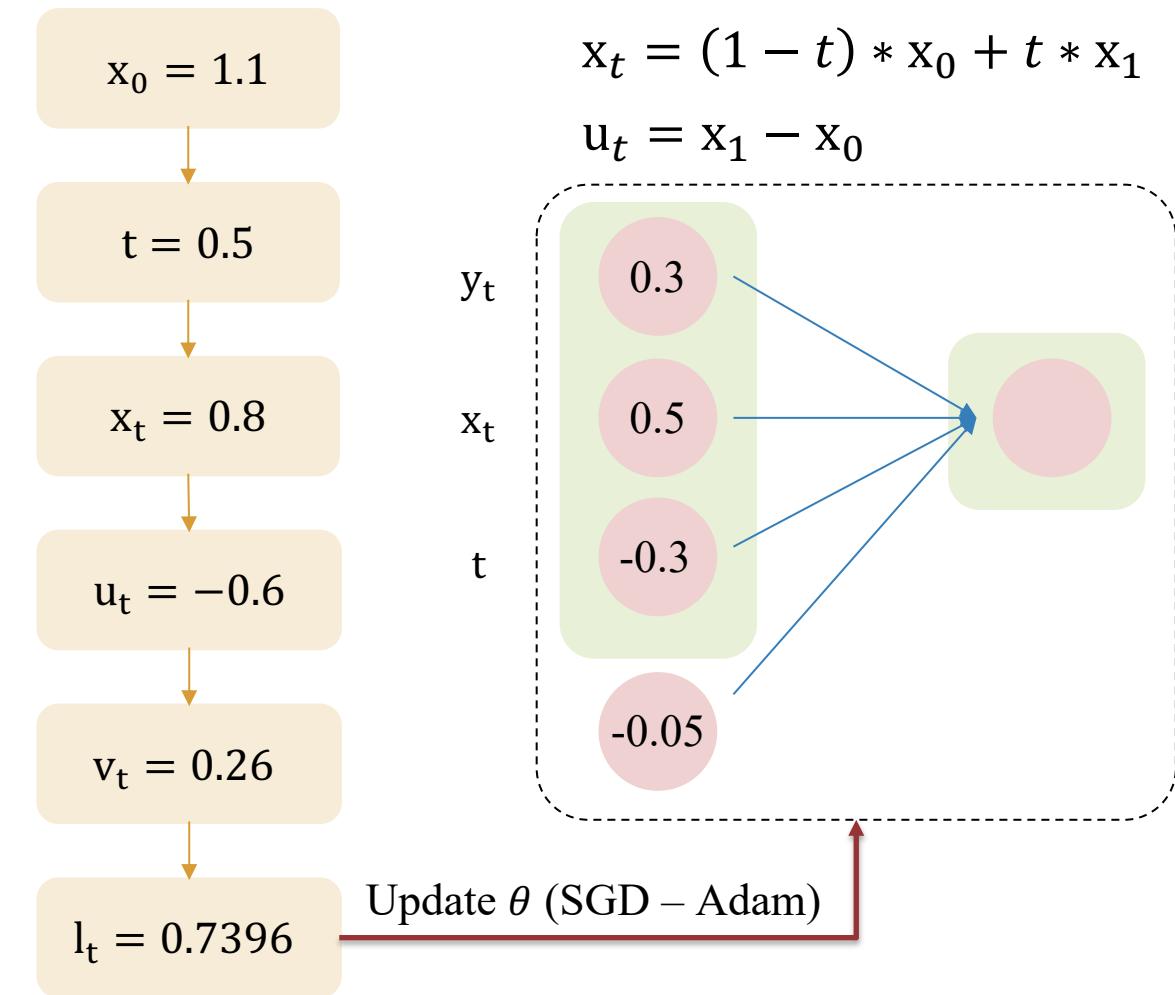


## Training in Conditional Flow Matching



$$5 \text{ classes} \quad y_1 = \frac{1}{5} = 0.2$$

$$x_1 = 0.5$$



# Conditional Flow Matching



## Sampling in Conditional Flow Matching

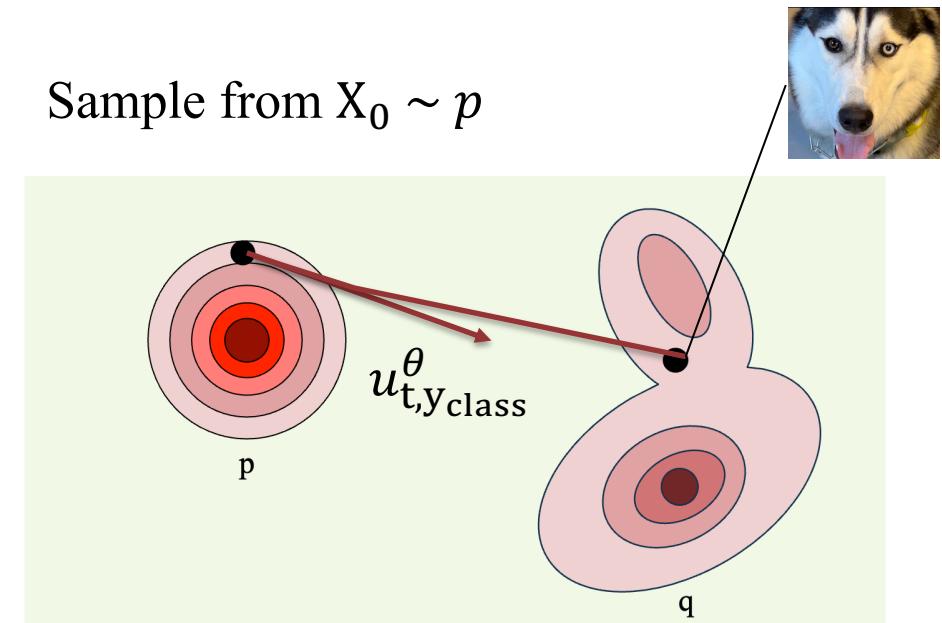
- Euler's method for solving ordinary differential equations (ODEs)

$$\frac{dy}{dt} = f_{\theta}(y, t, y_{\text{class}})$$

$$y(t + \Delta t) \approx y(t) + \left( \frac{dy(t)}{dt} \right) \Delta t$$

$$y(t + \Delta t) \approx y(t) + f_{\theta}(y, t, y_{\text{class}}) \Delta t$$

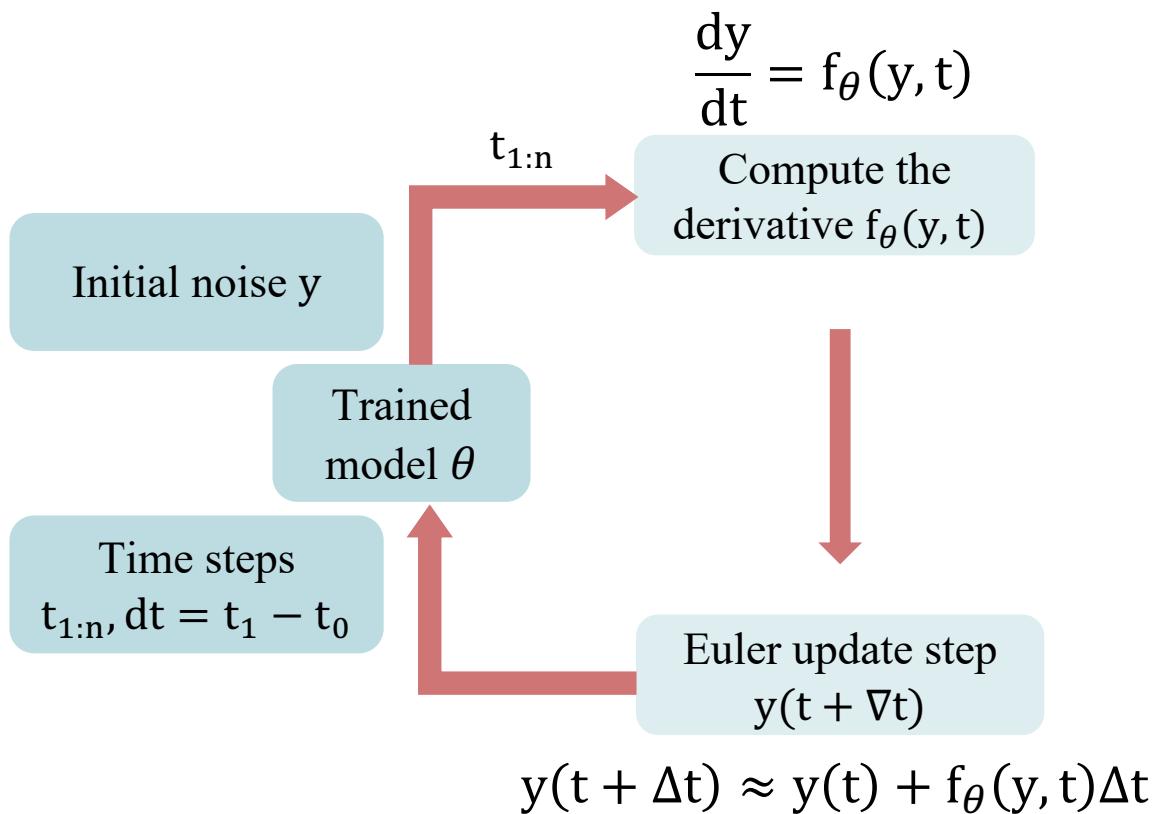
Sample from  $X_0 \sim p$



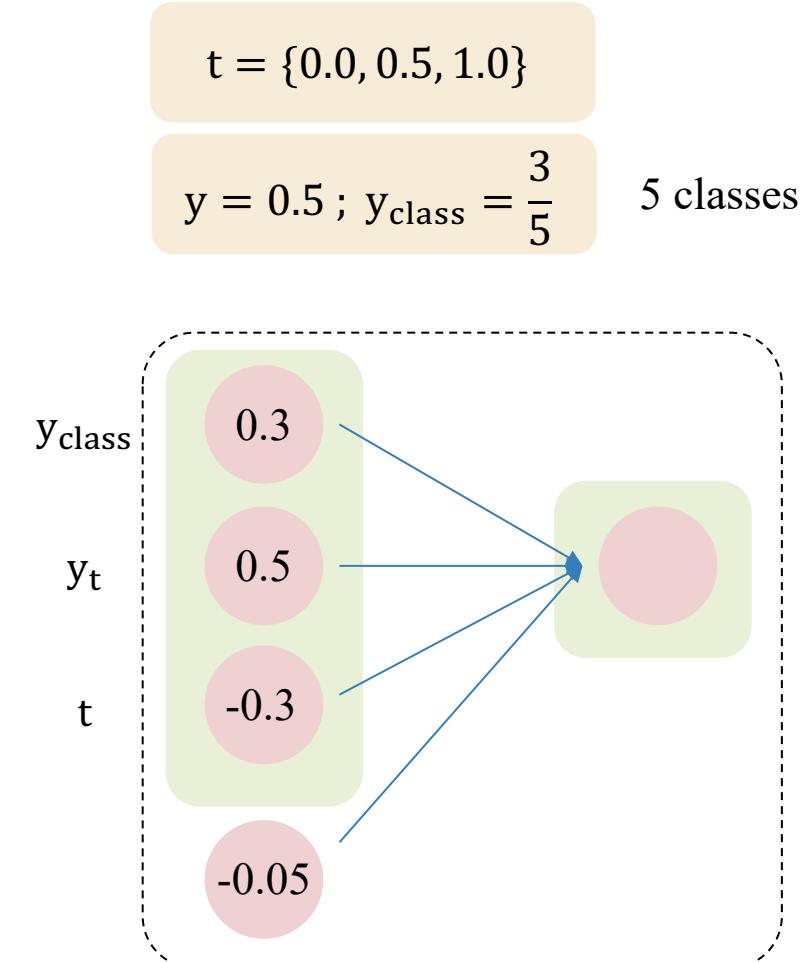
# Conditional Flow Matching



## Sampling



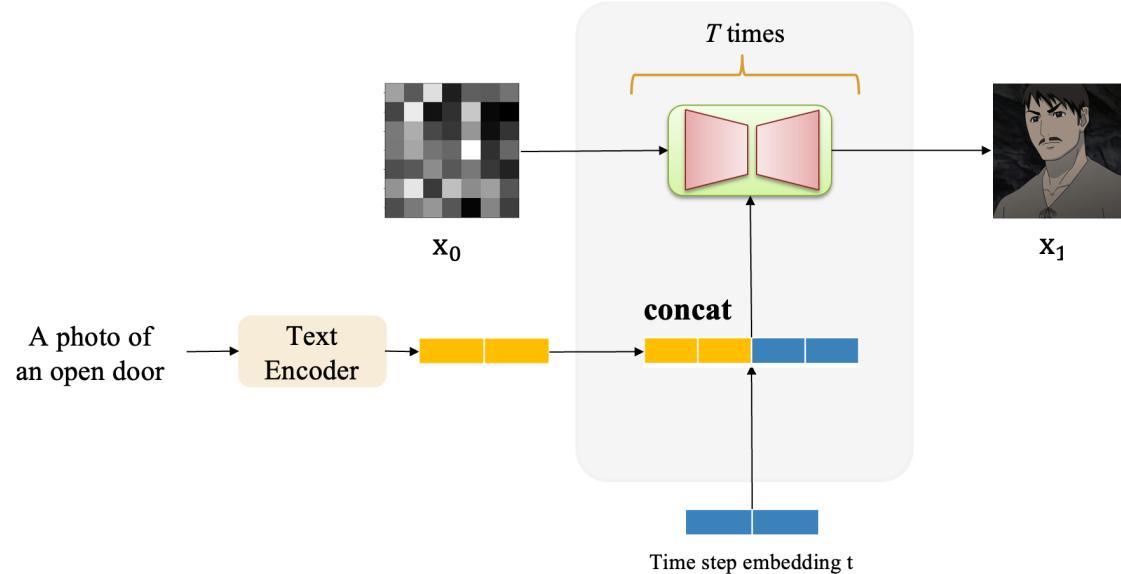
$$\begin{aligned} t &= 0.5 \\ f_{\theta}(y, t, y_{\text{class}}) &= 0.28 \\ y_{t+\Delta t} &= 0.64 \end{aligned}$$



# Outline

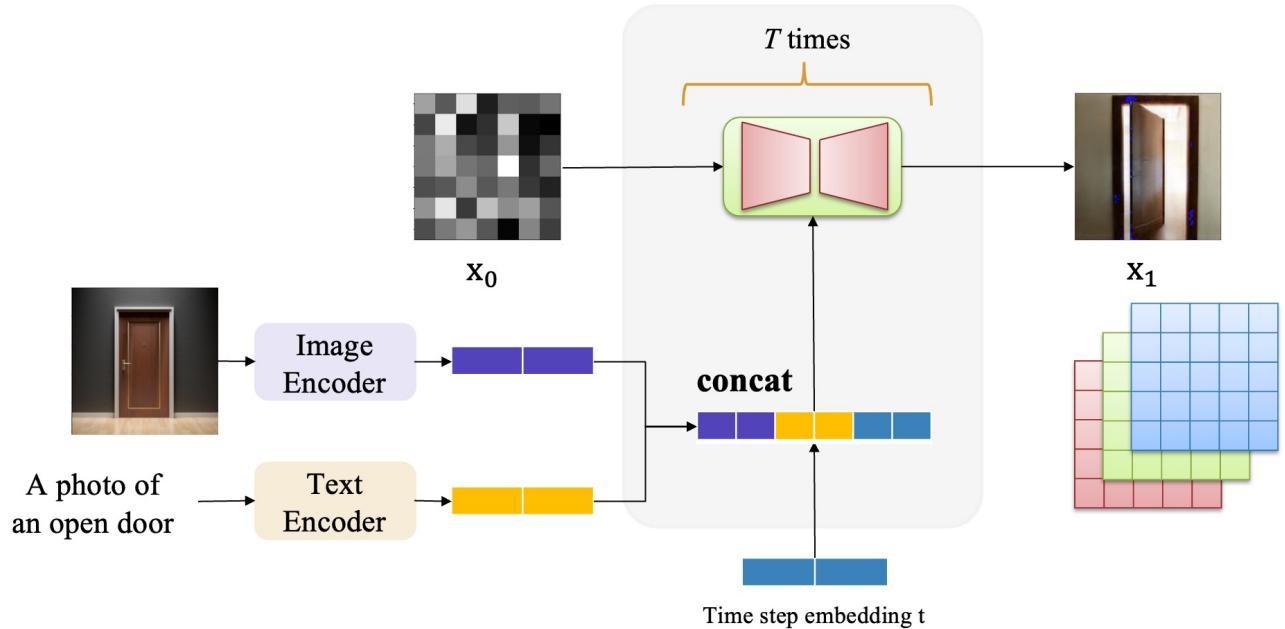
SECTION 1

## Conditional Flow Matching



SECTION 2

## Text-Guided Image Generation



# Text-to-Image Generation

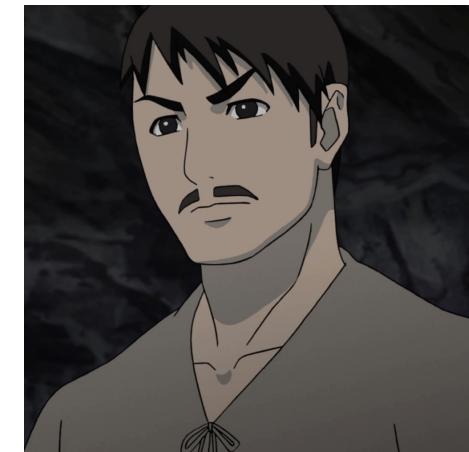


## Description

- **Text-to-Image Generation** is a process that combines text embeddings with image generation techniques to create images from text description
- Datasets: Oxford-102 Flower, CUB-200, Naruto-Captions,...



This pink and yellow flower has a beautiful yellow center with many stamens



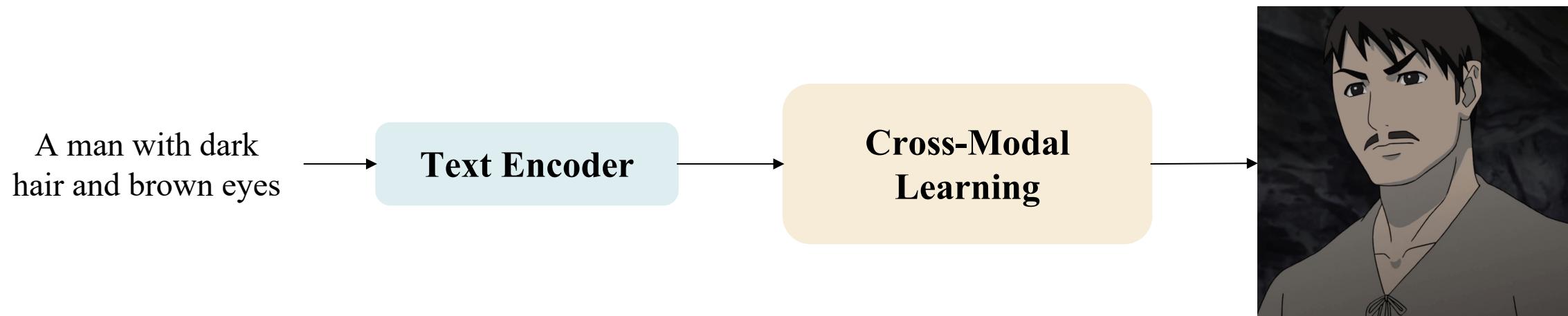
A man with dark hair and brown eyes

# Text-to-Image Generation



## Methods

- **Text-to-Image Generation** is a process that combines text embeddings with image generation techniques to create images from text description
- Datasets: Oxford-102 Flower, CUB-200, Naruto-Captions,...

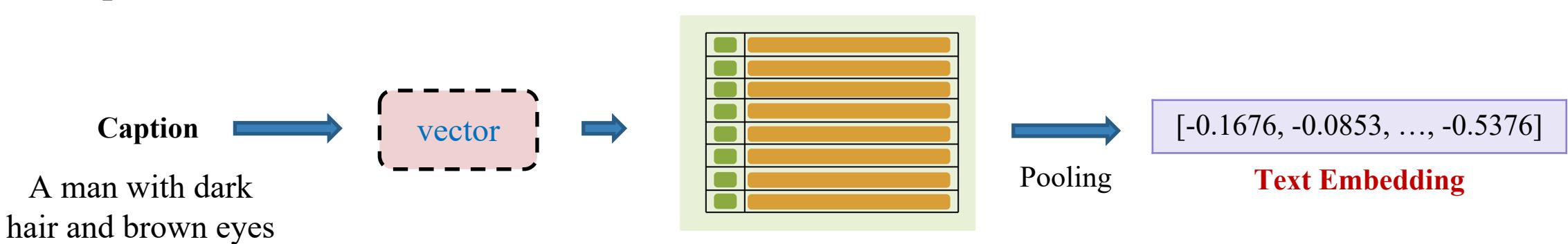


# Text-to-Image Generation



## Preprocessing

### ➤ Text Representation (Text Encoder)



index	token
0	<unk>
1	<pad>
2	a
3	man
4	bright
5	dark
...	...

2
3
0
5
...

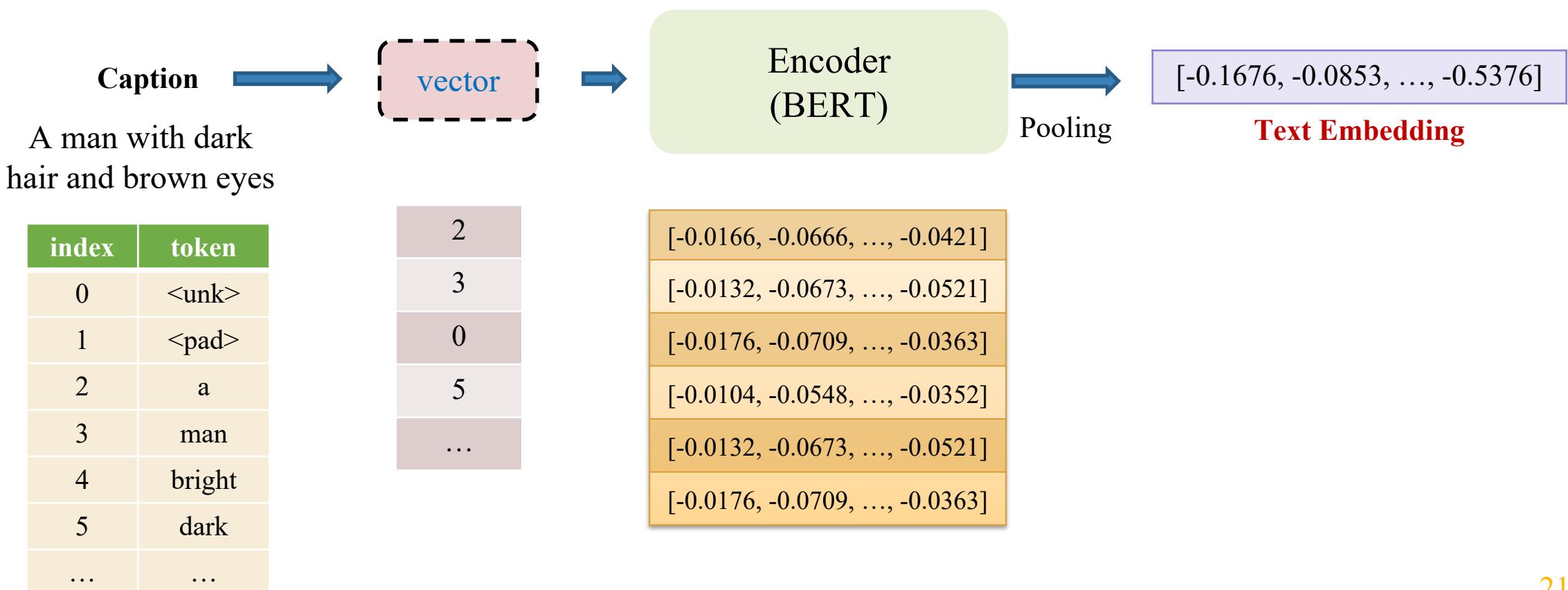
[-0.0166, -0.0666, ..., -0.0421]
[-0.0132, -0.0673, ..., -0.0521]
[-0.0176, -0.0709, ..., -0.0363]
[-0.0104, -0.0548, ..., -0.0352]
[-0.0132, -0.0673, ..., -0.0521]
[-0.0176, -0.0709, ..., -0.0363]

# Text-to-Image Generation



## Preprocessing

- Text Representation (Text Encoder): Pre-trained LMs

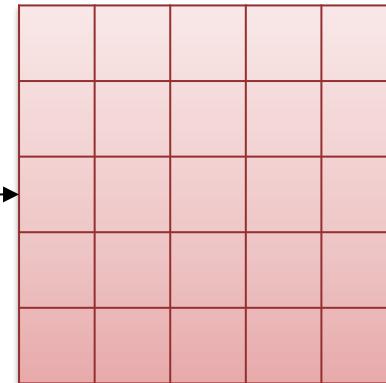


# Text-to-Image Generation



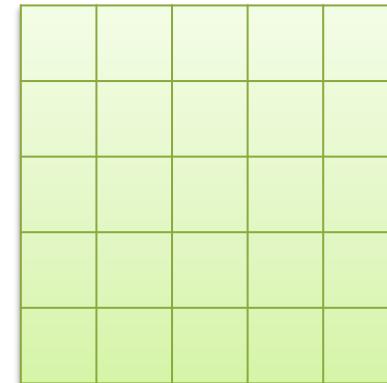
## Preprocessing

- **Image Representation:** The generated image is output as a tensor or array representing pixel values, typically of shape (channels, height, width) for neural network processing

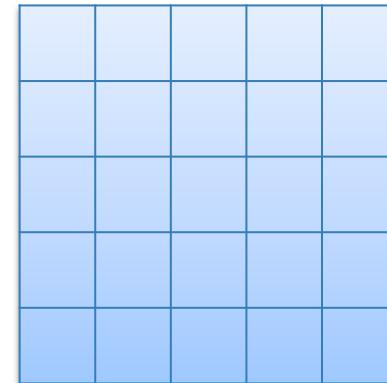


R (Red)

**Image Encoding**



G (Green)

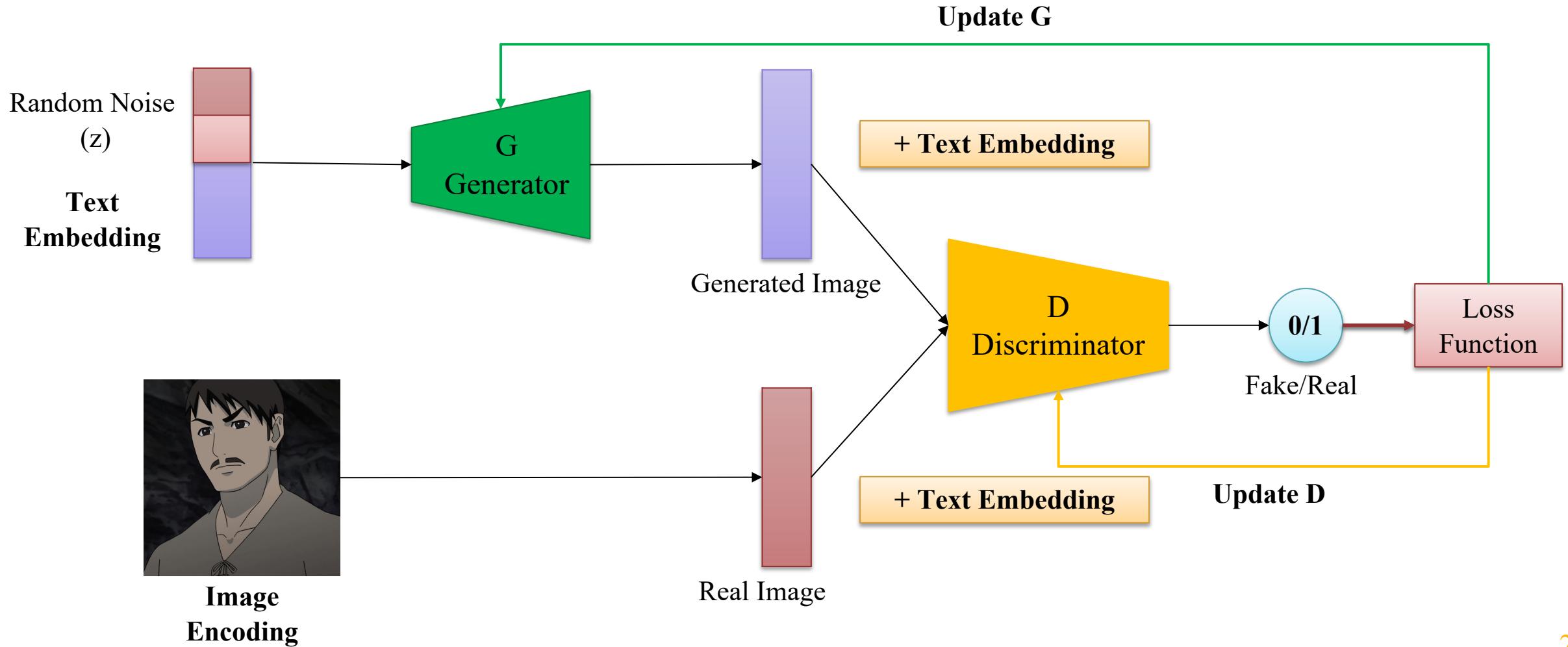


B (Blue)

# Text-to-Image Generation



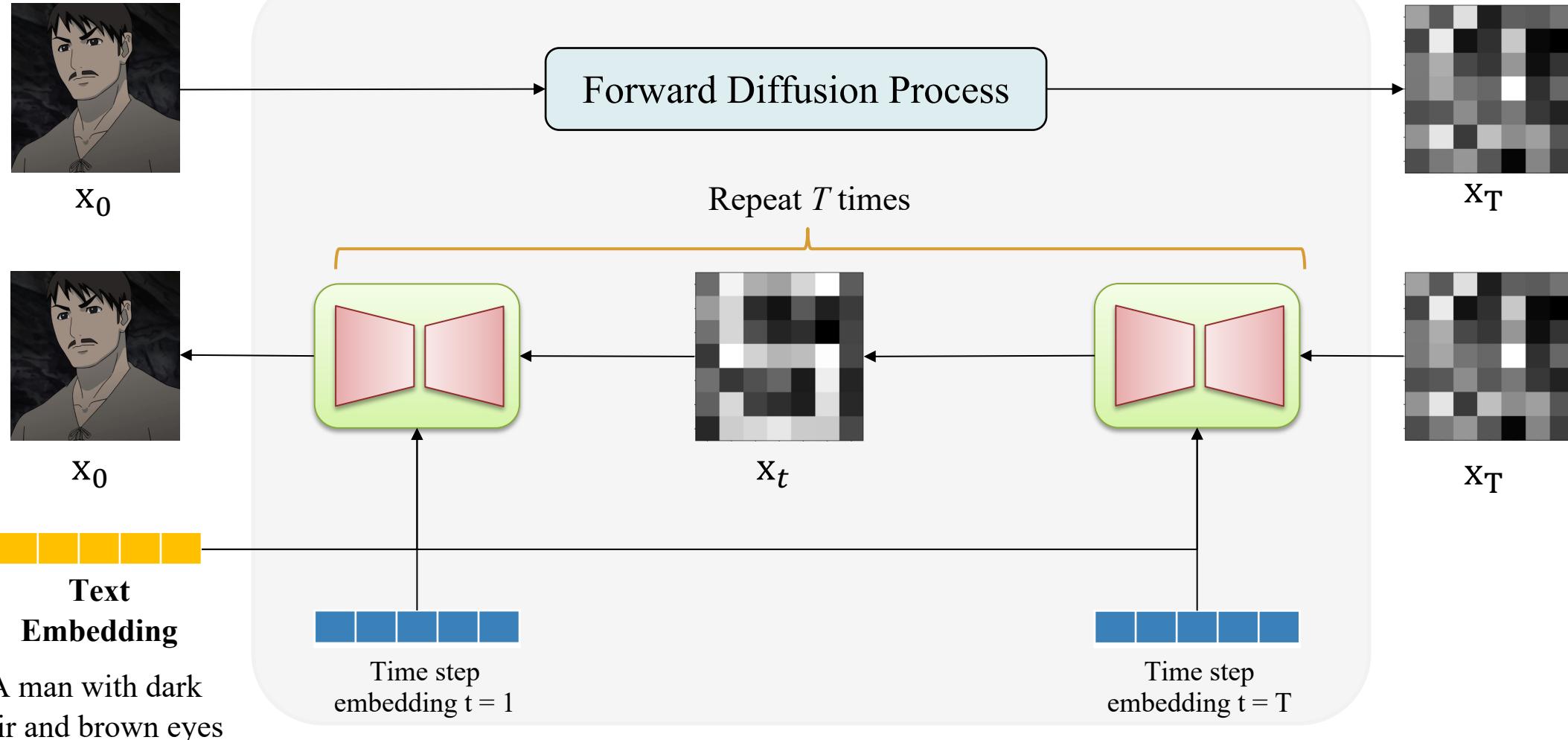
## DCGANs



# Text-to-Image Generation



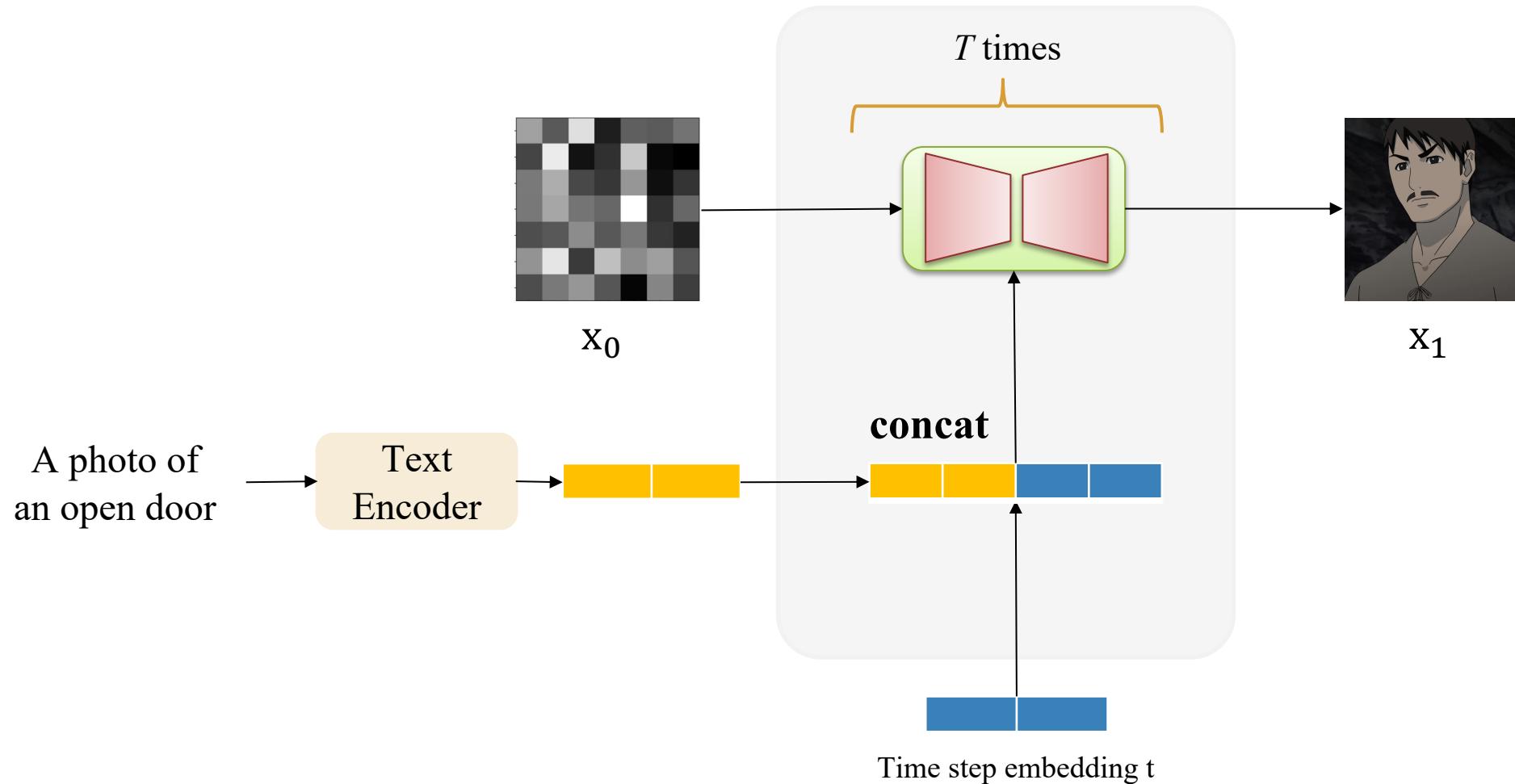
## Diffusion Models



# Text-to-Image Generation



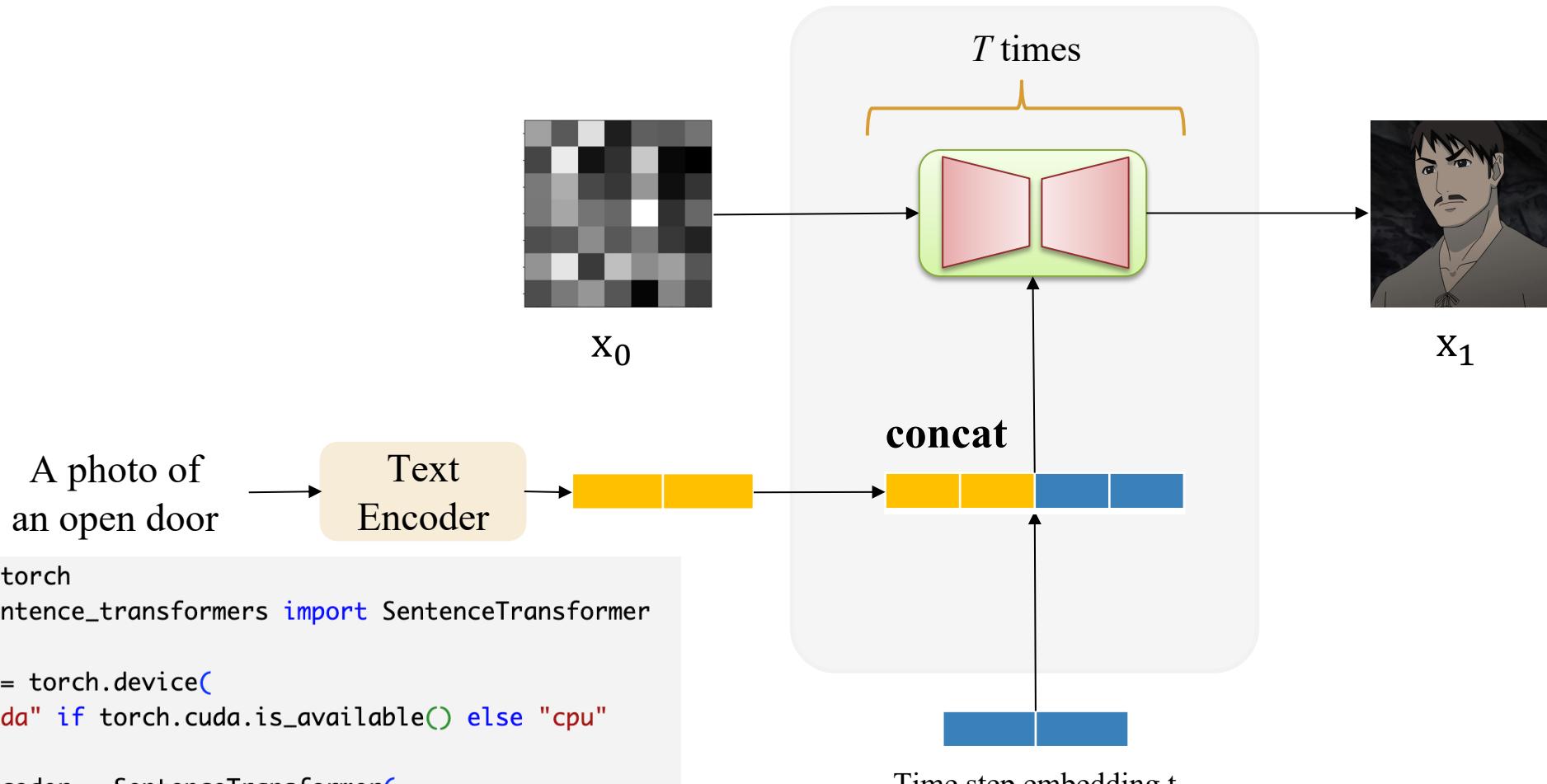
## Conditional Flow Matching



# Text-to-Image Generation



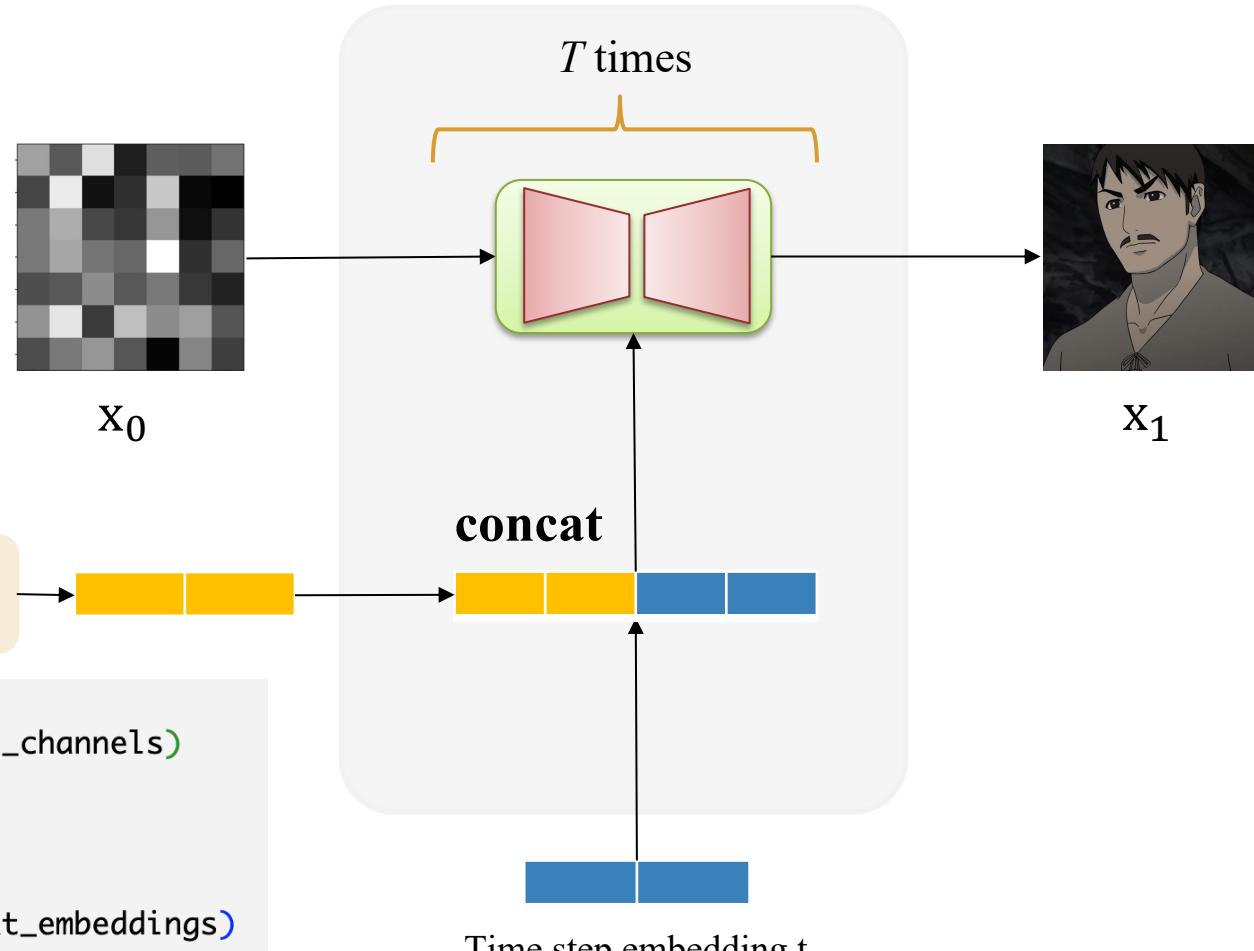
## Conditional Flow Matching



# Text-to-Image Generation



## Conditional Flow Matching



# Text-to-Image Generation



## Conditional Flow Matching

```
x1 = batch["image"].to(device)
text_embeddings = batch["caption_embedding"].to(device)

x0 = torch.randn_like(x1).to(device)

t = torch.rand(x0.shape[0], 1, 1, 1).to(device)

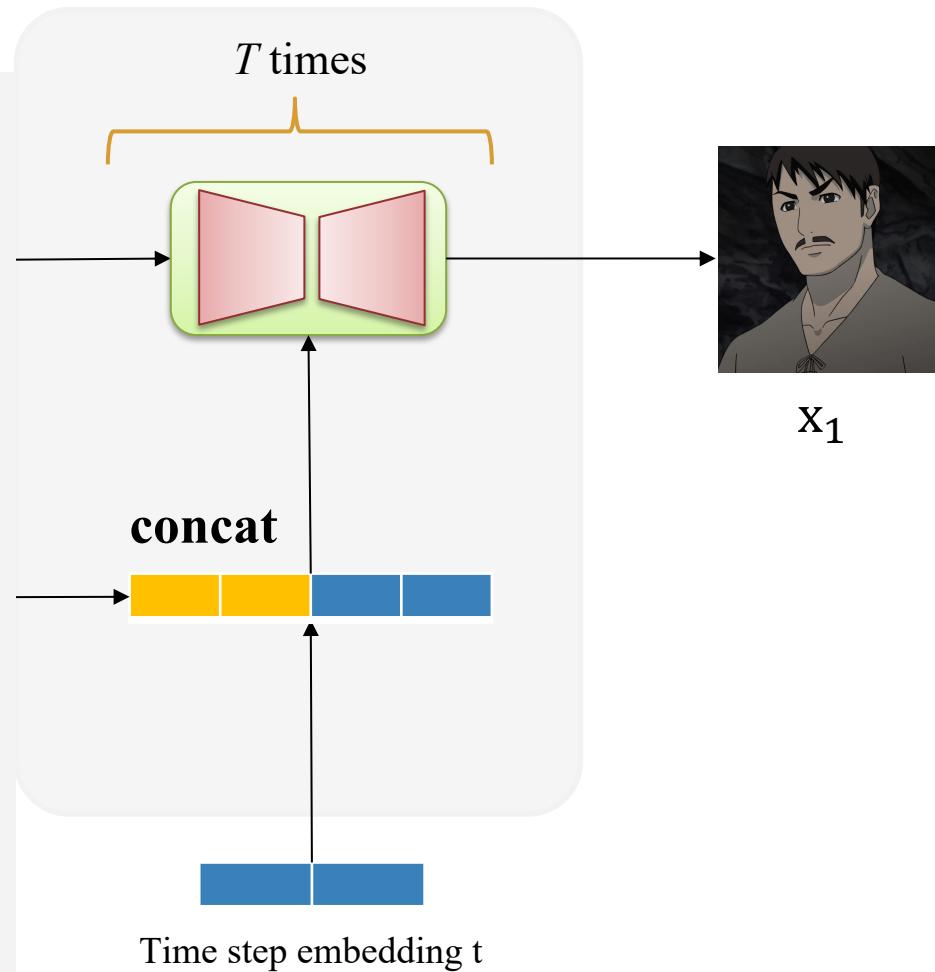
xt = t * x1 + (1 - t) * x0
ut = x1 - x0

t = t.squeeze()

vt = model(t, xt, text_embeddings=text_embeddings)

loss = torch.mean((vt - ut) ** 2)

loss.backward()
optimizer.step()
```



# Text-to-Image Generation



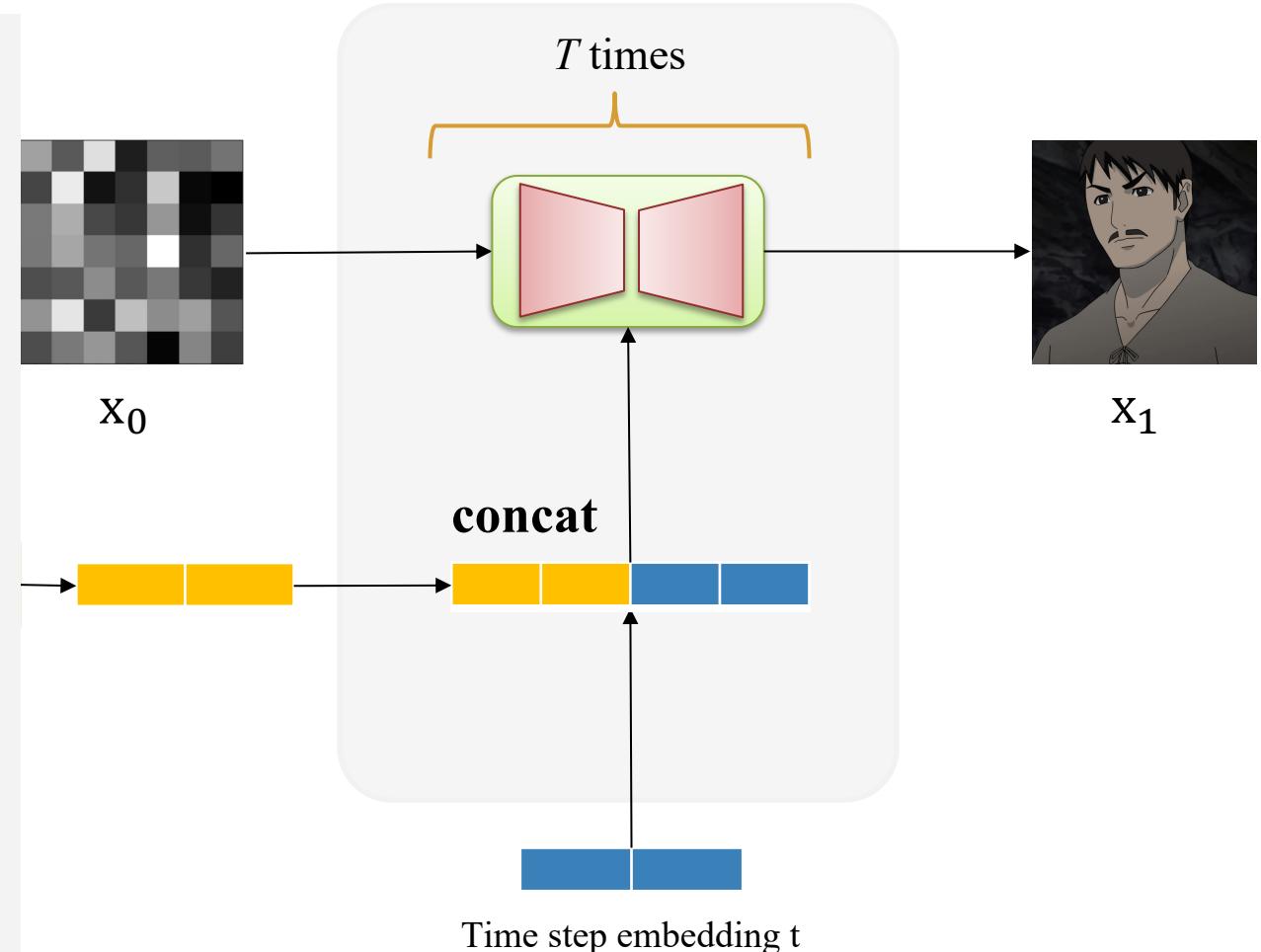
## Inference

```
model.eval()
def euler_method(model, text_embedding, t_steps, dt, noise):
    y = noise
    y_values = [y]
    with torch.no_grad():
        for t in t_steps[1:]:
            t = t.reshape(-1, )
            dy = model(
                t.to(device), y,
                text_embeddings=text_embedding
            )
            y = y + dy * dt
            y_values.append(y)
    return torch.stack(y_values)

# Initial random image and class (optional)
noise = torch.randn(3, 64, 64), device=device).unsqueeze(0)
text_embedding = text_embeddings[1].unsqueeze(0).to(device)

# Time parameters
t_steps = torch.linspace(0, 1, 100, device=device)
dt = t_steps[1] - t_steps[0]

# Solve the ODE using Euler method
results = euler_method(model, text_embedding, t_steps, dt, noise)
```



# Text-Guided Image Generation



## Description

- **Text-Guided Image Generation** aims to edit the given synthetic or real image to meet the specific requirements from the prompt.
- Based on the development of text-to-image models, which generate images according to text prompts



This pink and yellow flower has a beautiful yellow center with many stamens



A photo of an open door

# Text-Guided Image Generation



## Content-Aware Editing: Local Editing

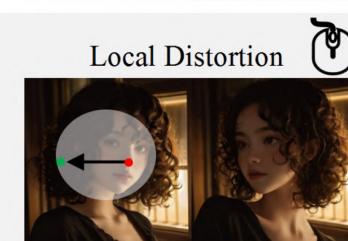
Object Manipulation



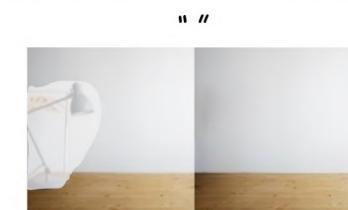
Attribute Manipulation



Spatial Transformation



Inpainting



# Text-Guided Image Generation



## Content-Aware Editing: Global Editing

Style Change

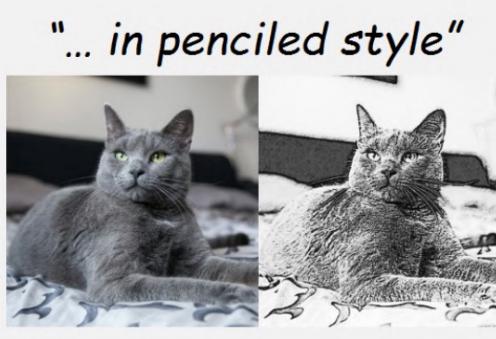


Image Translation



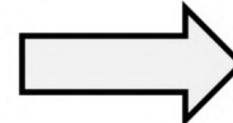
Style Transfer



# Text-Guided Image Generation



## Content-Free Editing



&lt;style&gt;



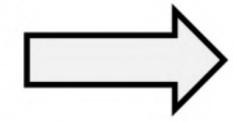
&lt;color&gt;



&lt;layout&gt;



&lt;object&gt;



“A dog &lt;action&gt; basket”



# Text-Guided Image Generation



## Dataset

- Oxford 120 Flowers, CUB-200, Recipe1M, ...
- TedBench Dataset



A photo of a chair swed in half



A photo of an open door

# Text-Guided Image Generation

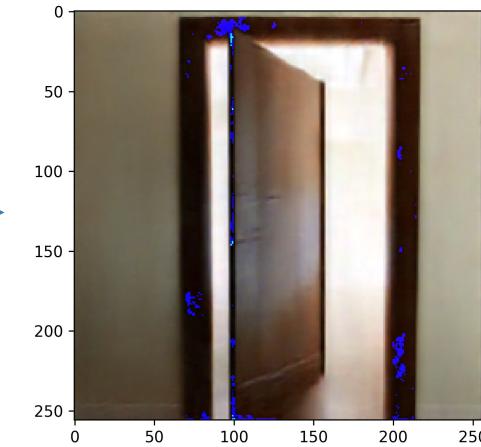
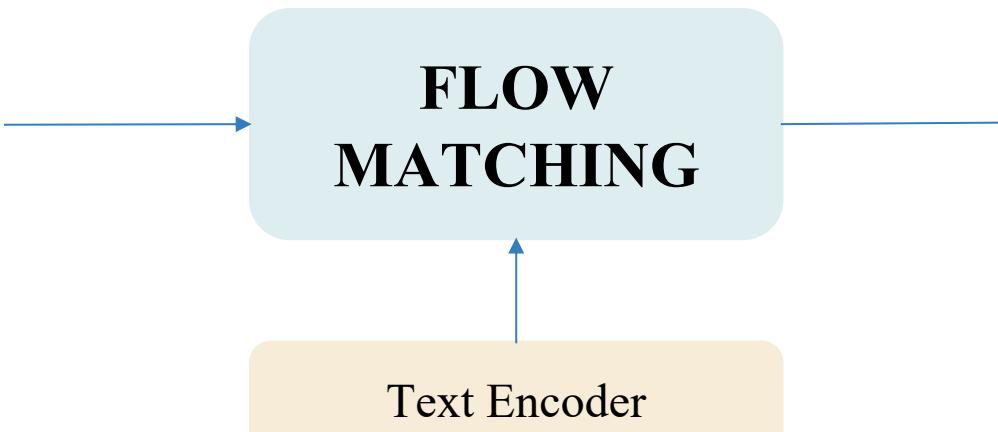


## Methods

- **Text Embedding:** Convert text into a vector representation
- **Image Encoder:** The original image is encoded to extract its feature
- **Text-Image Alignment:** the image aligns with the text description requires robust cross-modal learning
- **Generator:** The decoder reconstructs the image from the modified features



Original Image



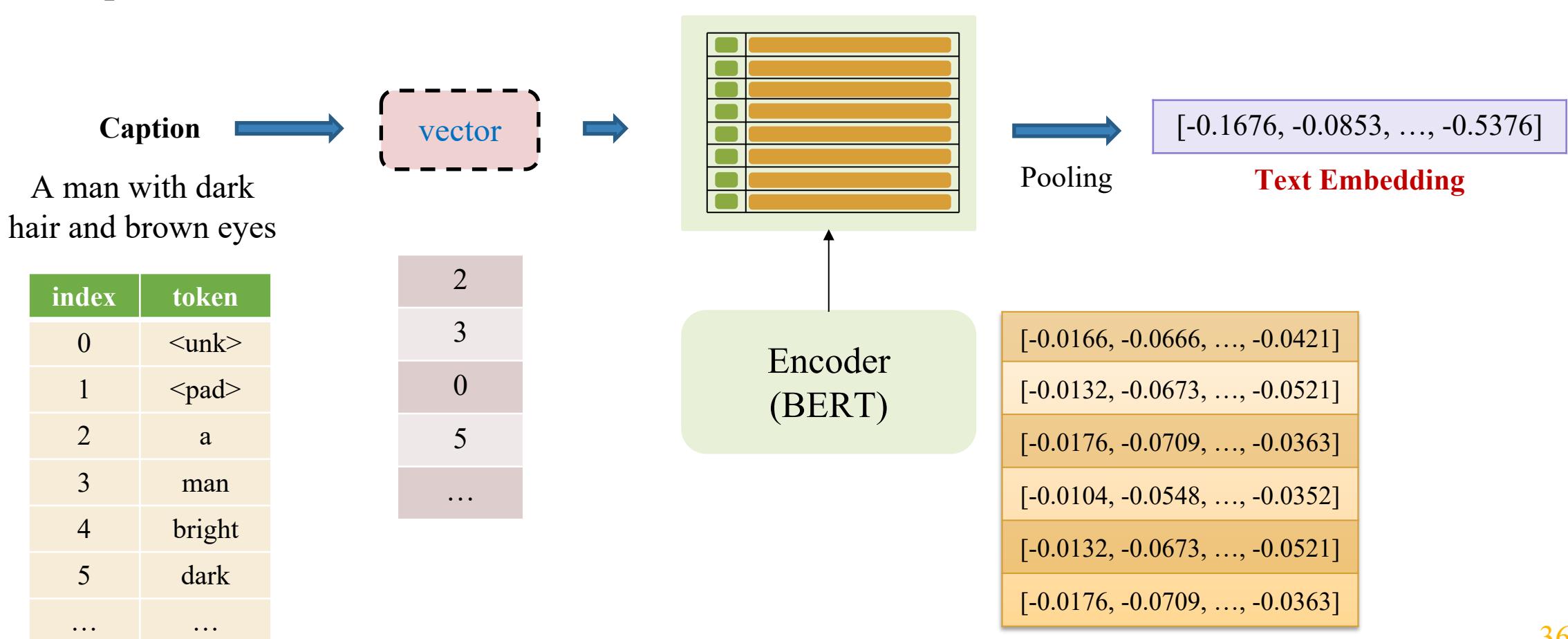
Edited Image

# Text-Guided Image Generation



## Preprocessing

### ➤ Text Representation (Text Encoder)



# Text-Guided Image Generation

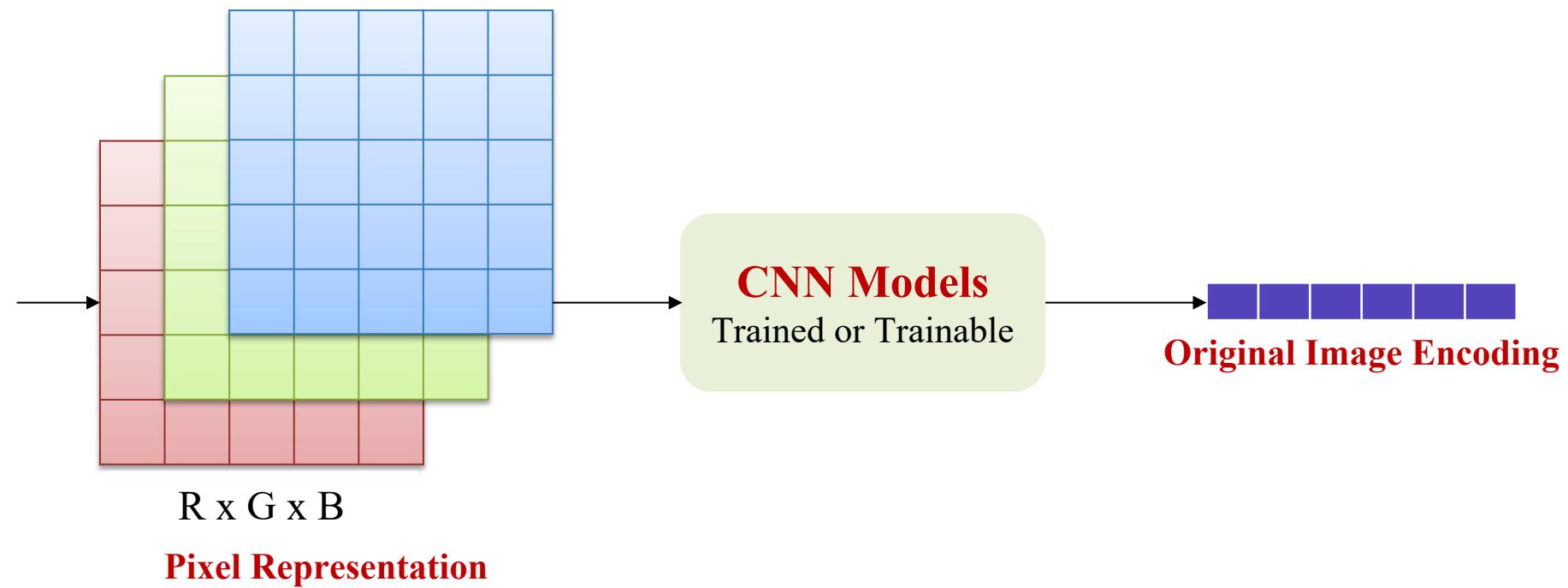


## Preprocessing

- **Original Image Encoding:** Encoded to extract high-level features that can guide the generation progress through a CNN or a pre-trained model as a image encoder.



Original Image

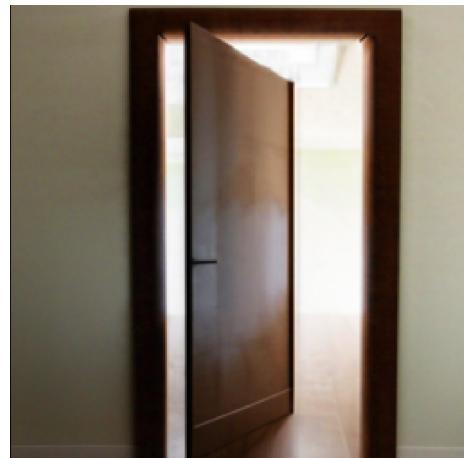


# Text-Guided Image Generation

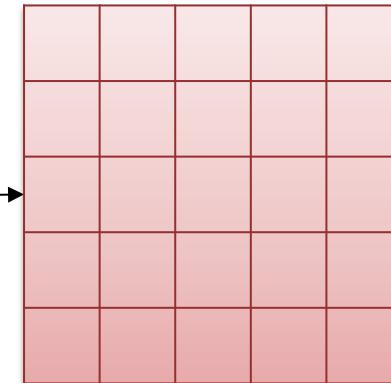


## Preprocessing

- **Edited Image Encoding:** The generated image is output as a tensor or array representing pixel values, typically of shape (channels, height, width) for neural network processing, would be normalized

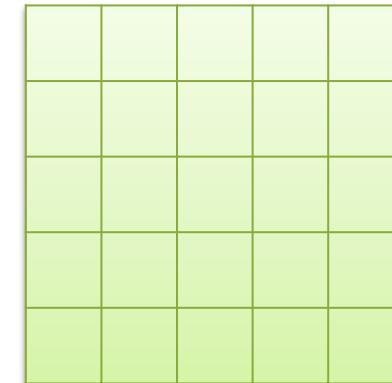


Edited Image

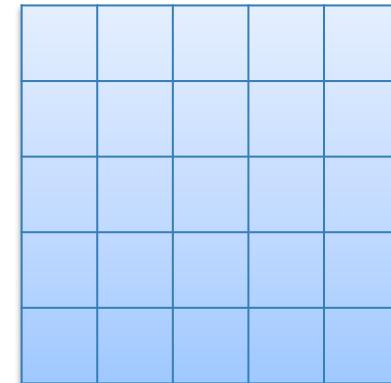


R (Red)

Edited Image Encoding



G (Green)

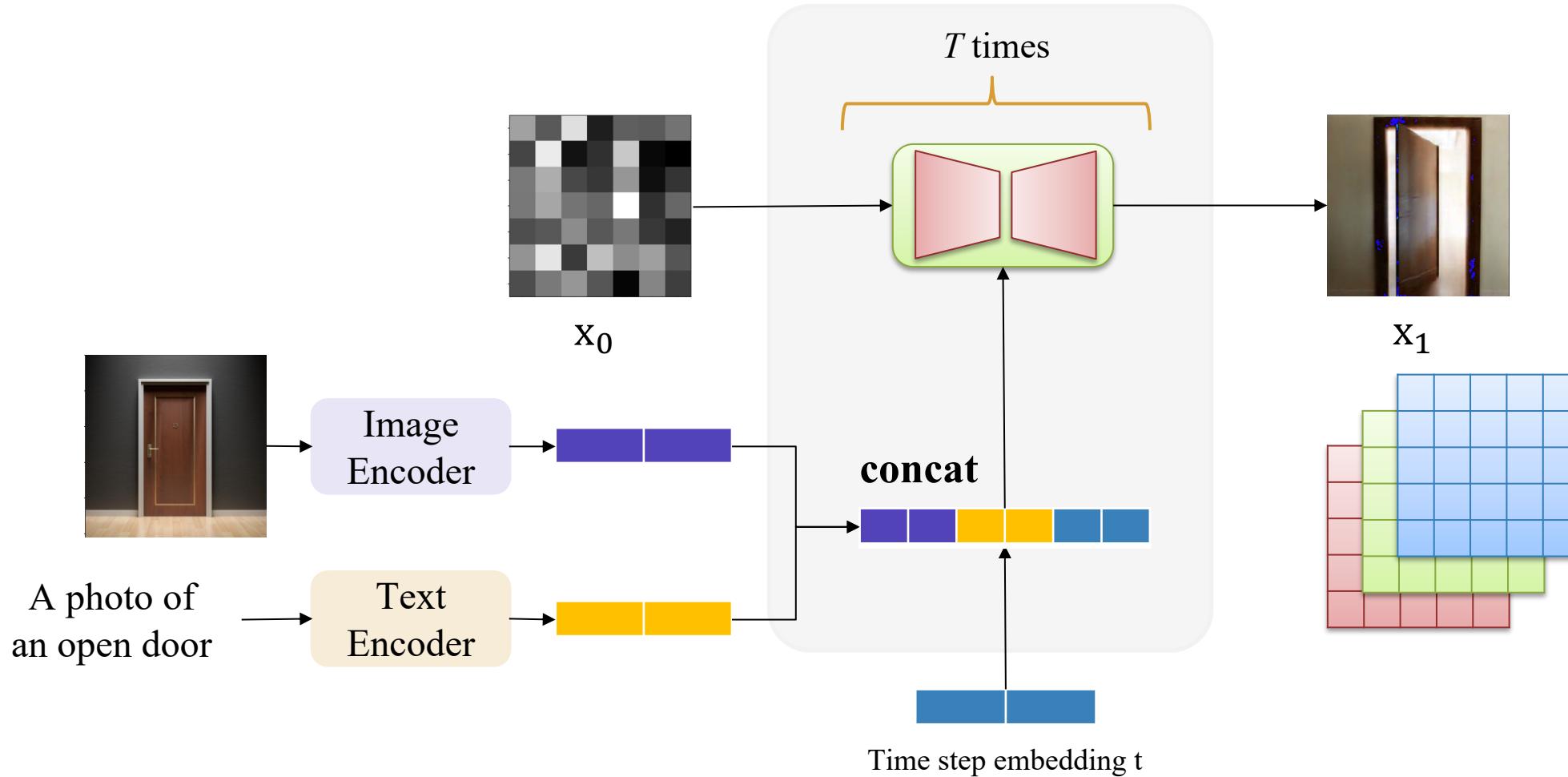


B (Blue)

# Text-Guided Image Generation



## Conditional Flow Matching

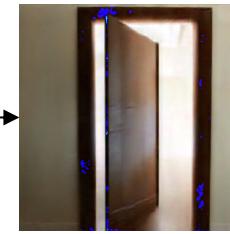
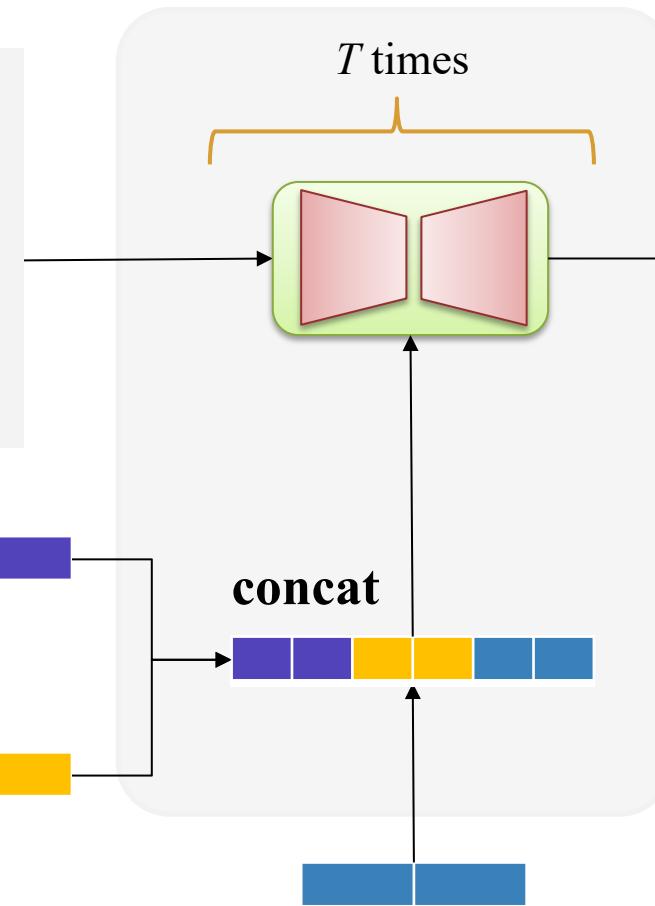
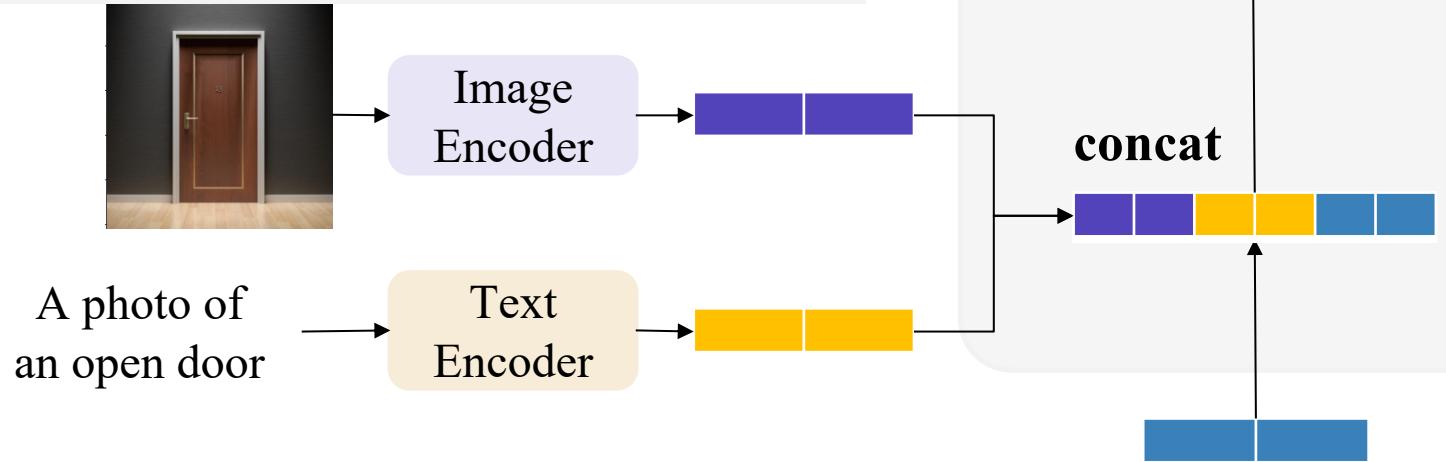


# Text-Guided Image Generation

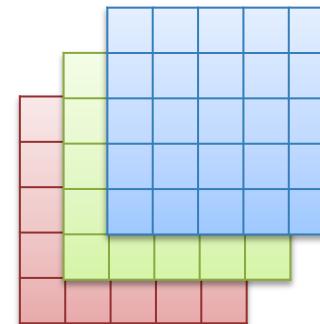


## Conditional Flow Matching

```
1 import torch
2 from sentence_transformers import SentenceTransformer
3
4 device = torch.device(
5     "cuda" if torch.cuda.is_available() else "cpu"
6 )
7 text_encoder = SentenceTransformer(
8     "all-mpnet-base-v2"
9 ).to(device)
```



$X_1$

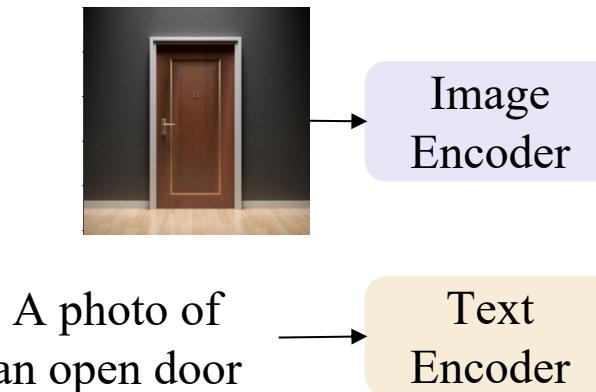


Time step embedding  $t$

# Text-Guided Image Generation



## Conditional Flow Matching

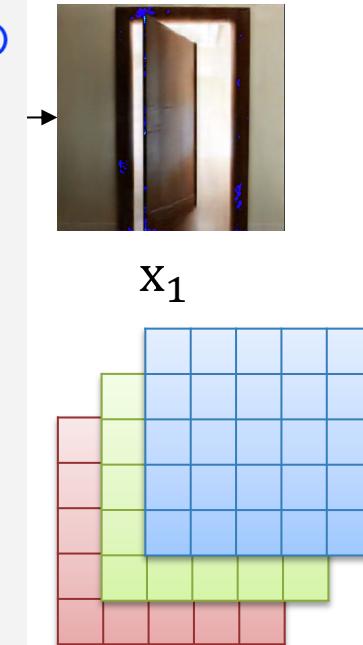


```
# get an original image
original_image = self.original_images[idx]
original_image = self.transform(original_image)

# get a output image
edited_image = self.edited_images[idx]
edited_image = self.transform(edited_image)

# get a text
caption = self.captions[idx]
caption_embedding = self.embed_captions[idx]

return {
    "original_image": original_image,
    "edited_image": edited_image,
    "caption": caption,
    "caption_embedding": caption_embedding,
}
```

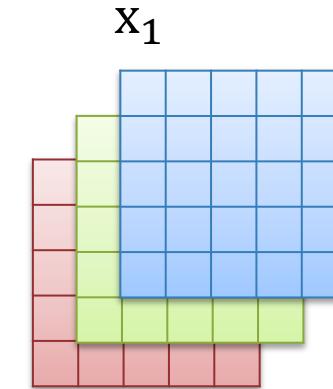
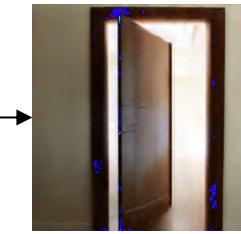
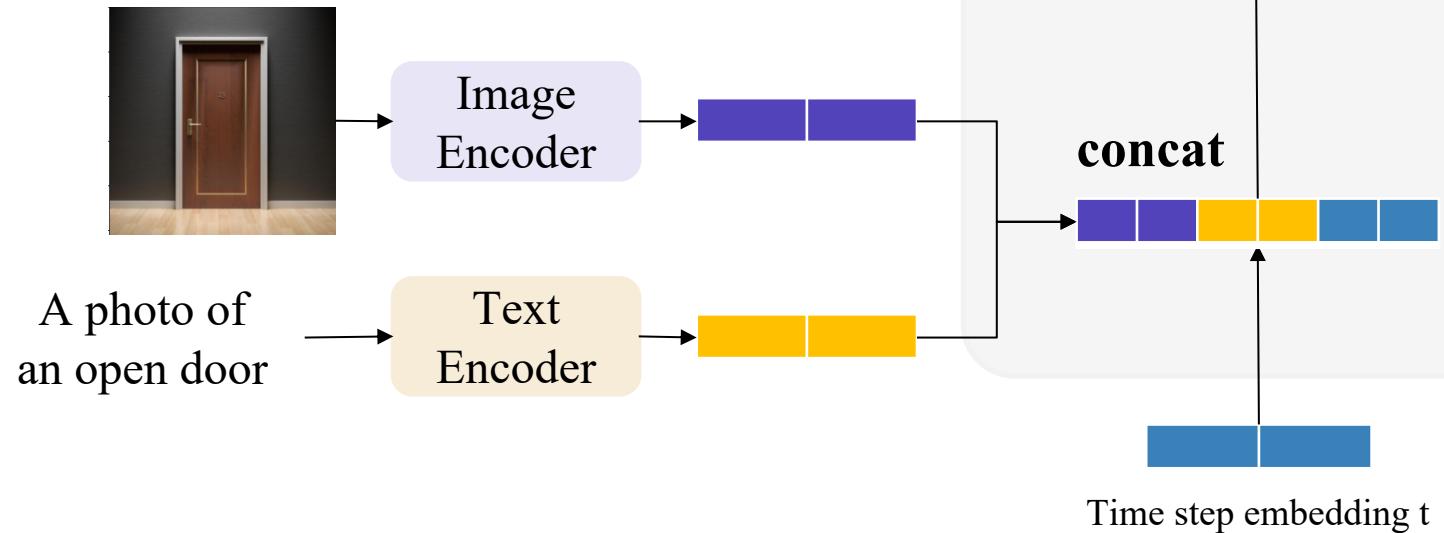


# Text-Guided Image Generation



## Conditional Flow Matching

```
if (text_embeddings is not None) and (original_image is not None):  
    text_embedded = self.embedding_layer(text_embeddings)  
    image_embedded = self.image_encoder(original_image).squeeze(2, 3)  
    emb = torch.cat([emb, text_embedded, image_embedded], dim=1)  
    emb = self.fc(emb)
```



# Text-Guided Image Generation



## Conditional Flow Matching

```
optimizer.zero_grad()
x1 = edited_image

x0 = torch.randn_like(x1).to(device)

t = torch.rand(x0.shape[0], 1, 1, 1).to(device)

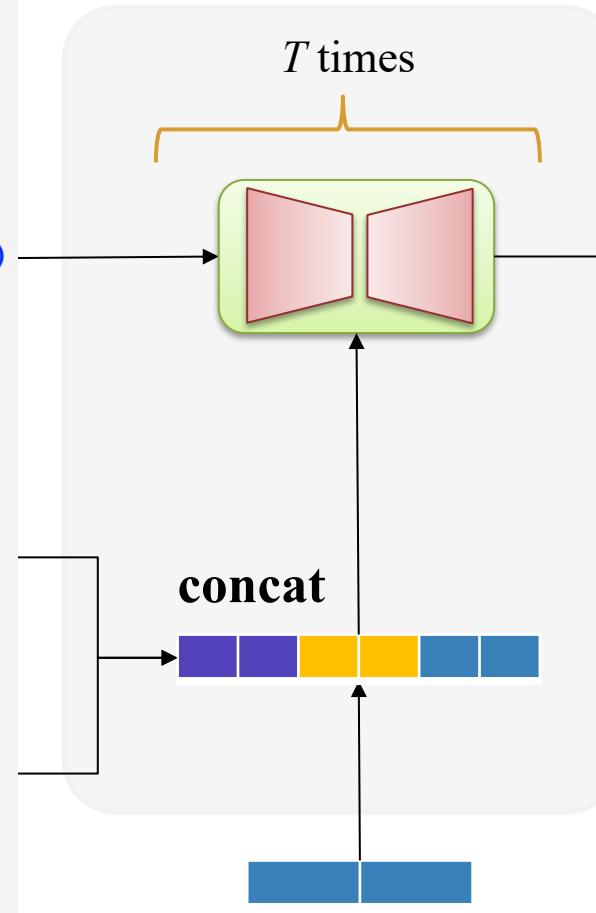
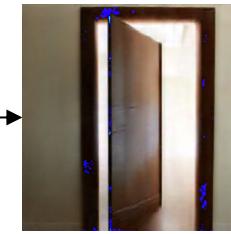
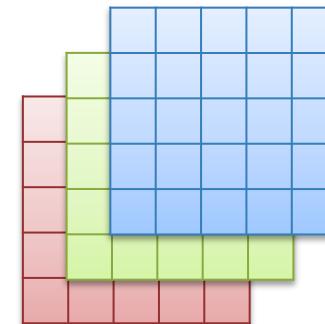
xt = t * x1 + (1 - t) * x0
ut = x1 - x0

t = t.squeeze()

vt = model(
    t, xt,
    text_embeddings=caption_embedding,
    original_image=original_image
)

loss = torch.mean((vt - ut) ** 2)

loss.backward()
optimizer.step()
```

Time step embedding  $t$  $x_1$ 

# Text-Guided Image Generation



## Deployment

- ❖ [Github - Demo](#)



A photo of an open door

### Text-Guided Image Generation using Conditional Flow Matching

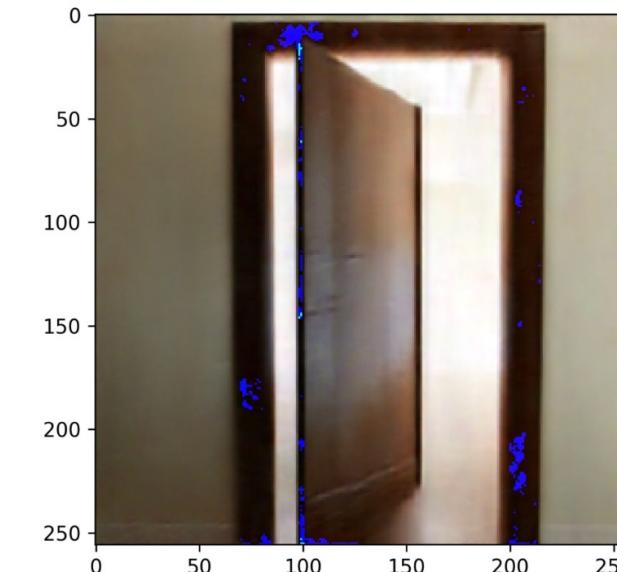
Model: Conditional Flow Matching. Dataset: Tedbench

Instruction:

A photo of an open door.

How would you like to give the input?

Run Example Image



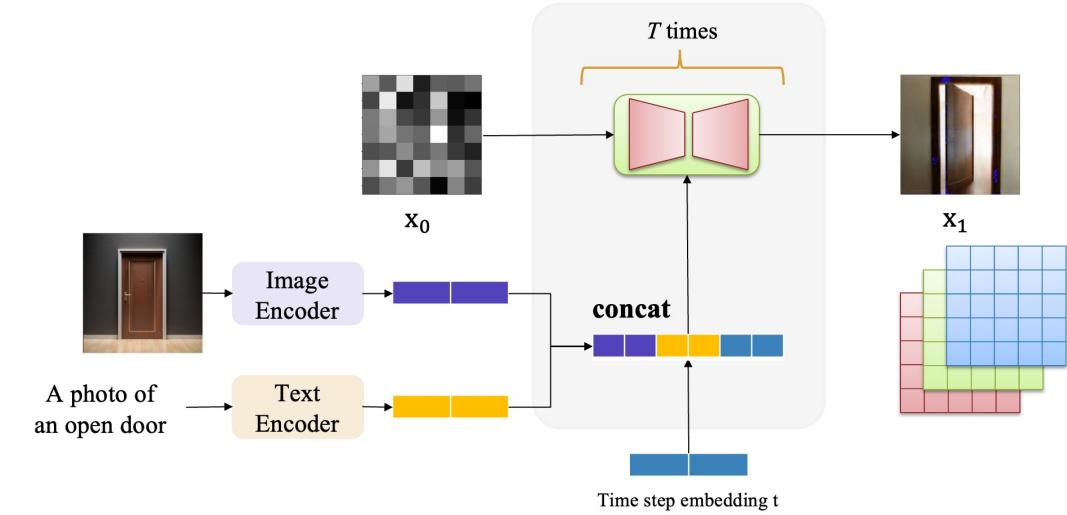
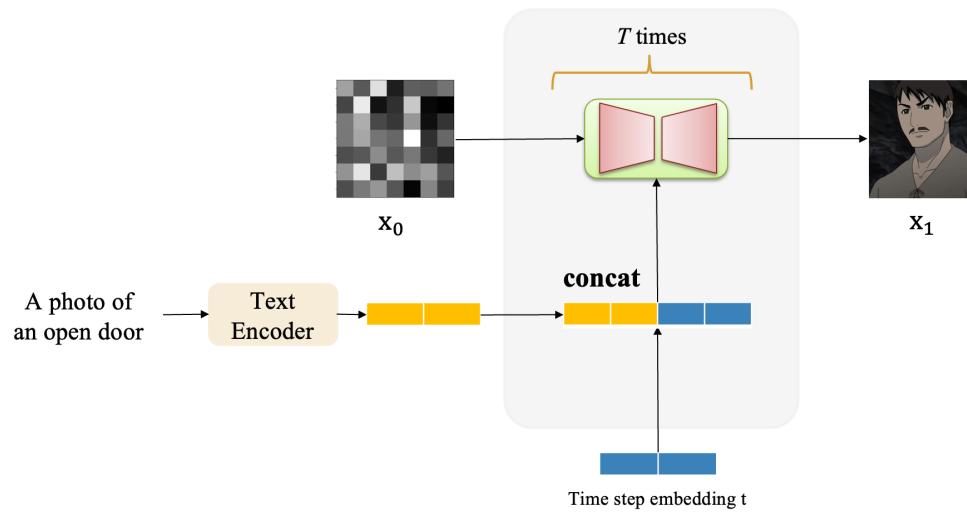
# Objectives

## Conditional Flow Matching (Review)

- ❖ Flow Matching Model
- ❖ Training & Sampling
- ❖ Conditional Flow Matching

## Text-Guided Image Generation

- ❖ Text-To-Image Generation
- ❖ Text-Guided Image Generation
- ❖ TedBench Dataset



# Thanks!

Any questions?