

# CME 296: Diffusion & Large Vision Models

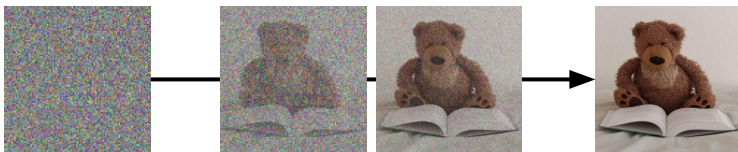


**Afshine Amidi & Shervine Amidi**



# Recap of last episodes...

## Lectures 1, 2, 3 Generation paradigms



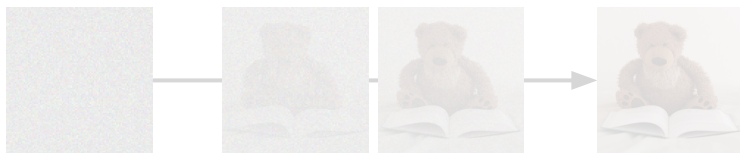
**Diffusion**     $\epsilon_{\theta}(x_t, t)$      $\|\epsilon_{\theta} - \epsilon\|^2$

**Score matching**     $s_{\theta}(x_t, t)$      $\|s_{\theta} - s\|^2$

**Flow matching**     $u_{\theta}(x_t, t)$      $\|u_{\theta} - u\|^2$

# Recap of last episodes...

## Lectures 1, 2, 3 Generation paradigms



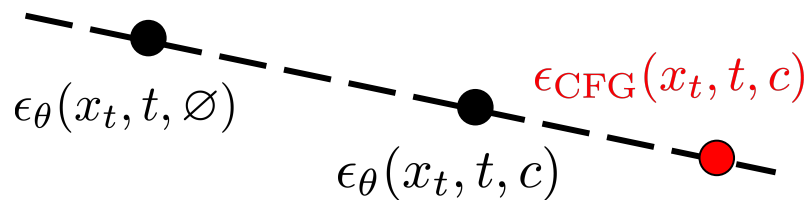
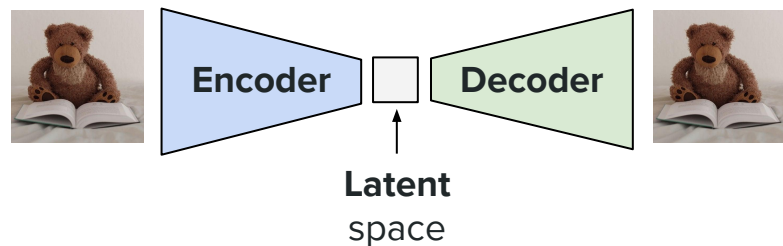
Diffusion  $\epsilon_{\theta}(x_t, t) \quad \|\epsilon_{\theta} - \epsilon\|^2$

Score matching  $s_{\theta}(x_t, t) \quad \|s_{\theta} - s\|^2$

Flow matching  $u_{\theta}(x_t, t) \quad \|u_{\theta} - u\|^2$

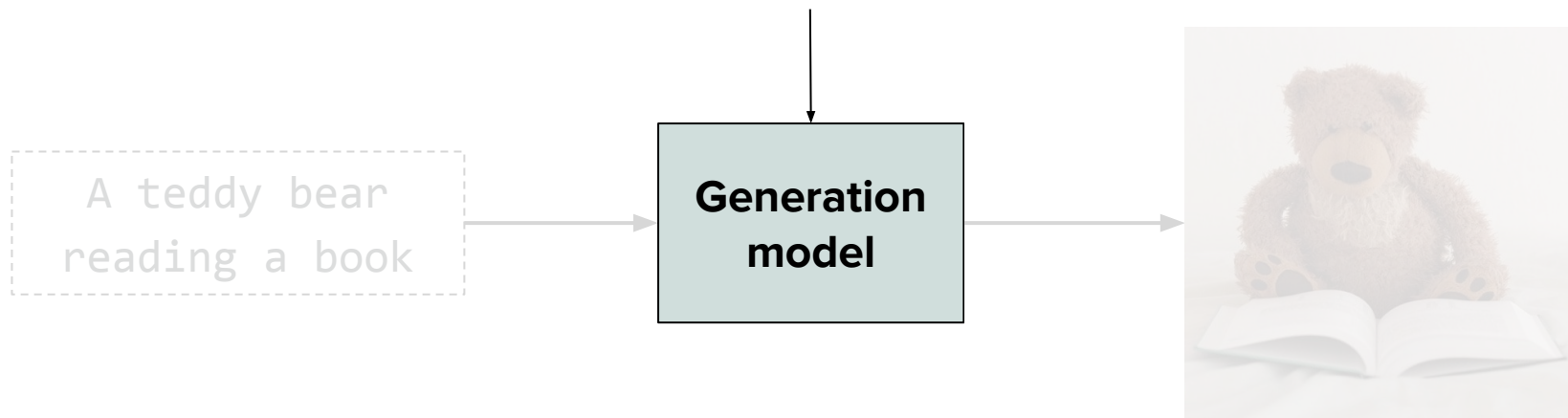
## Lecture 4 Latent space, guidance

### Variational AutoEncoder (VAE)



# Today's lecture

## Model architecture



Today's lecture: **Image generation architectures**



# Diffusion & Large Vision Models

## Motivation

U-Net

Diffusion Transformer

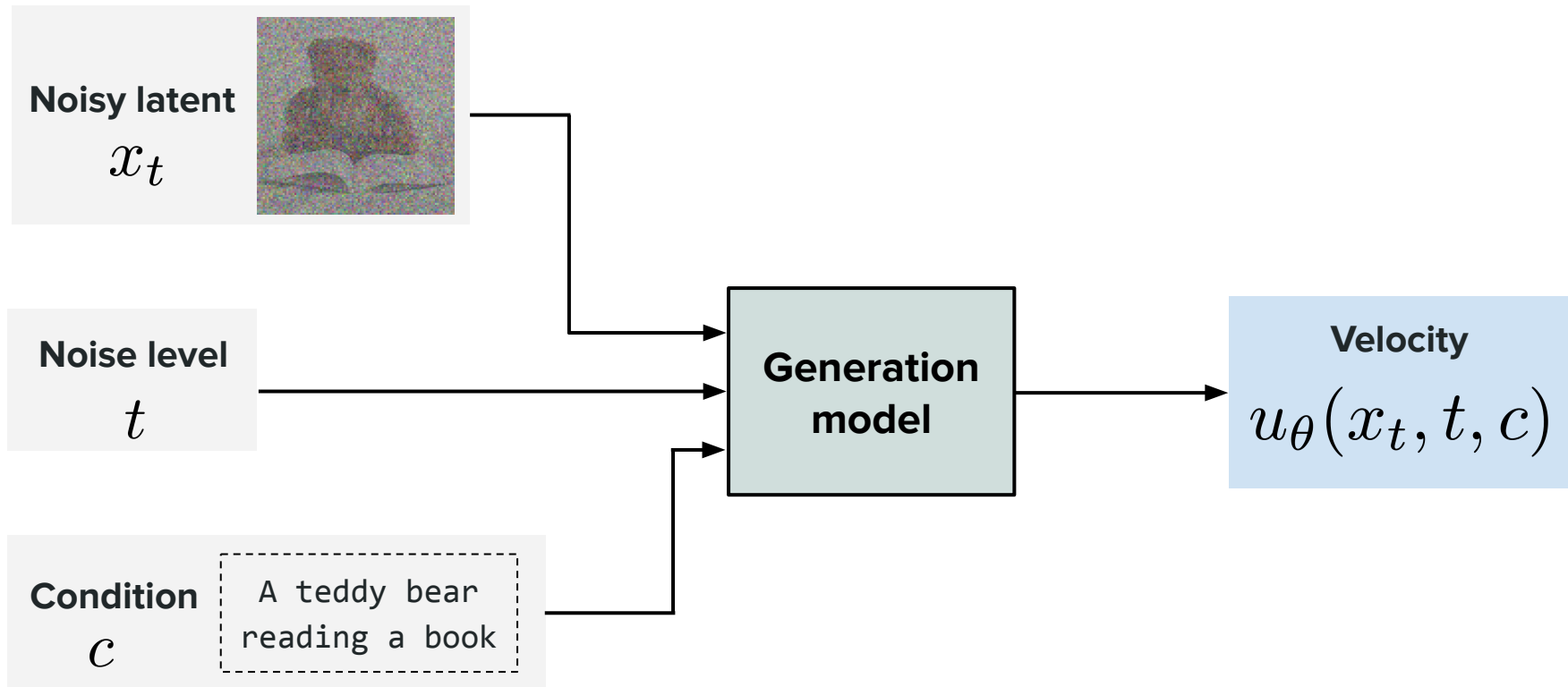
End-to-end example

Multimodal DiT

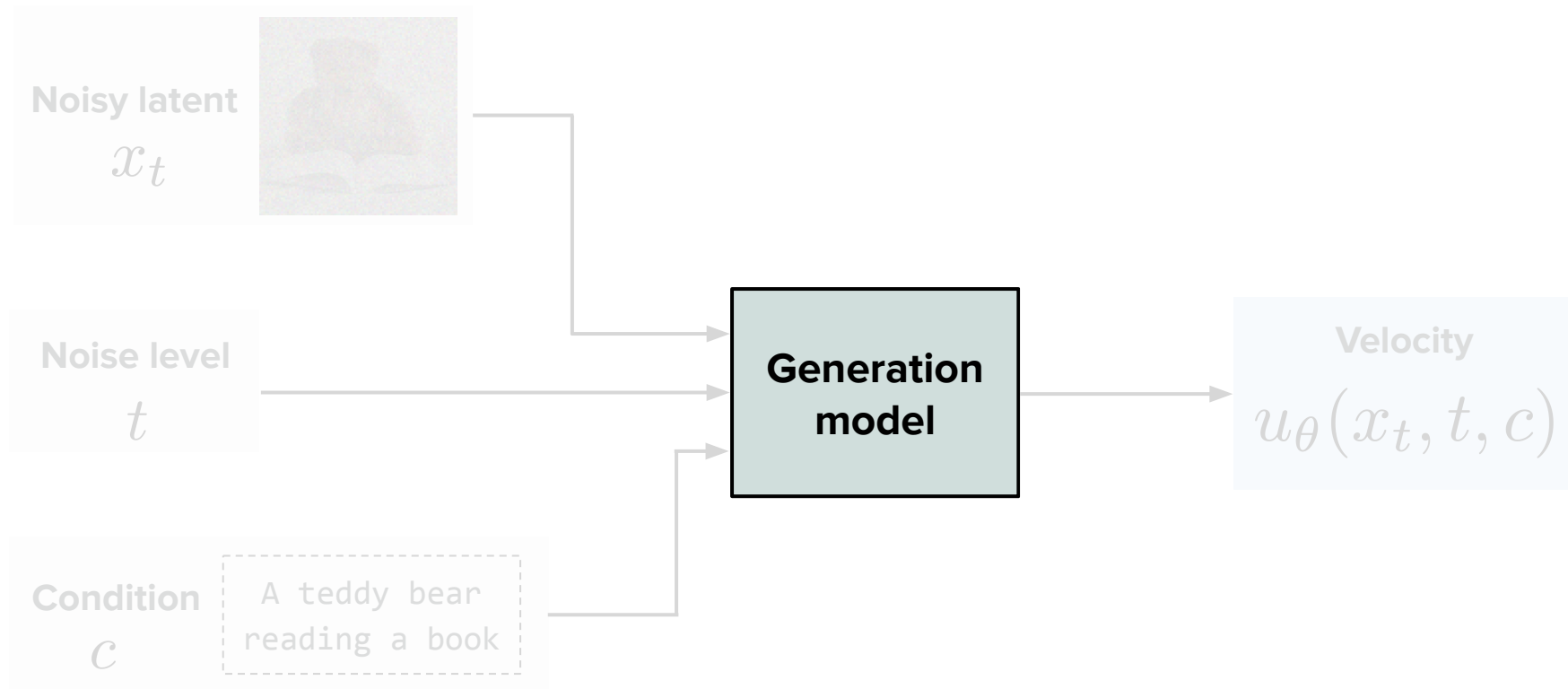
Optimizations

---

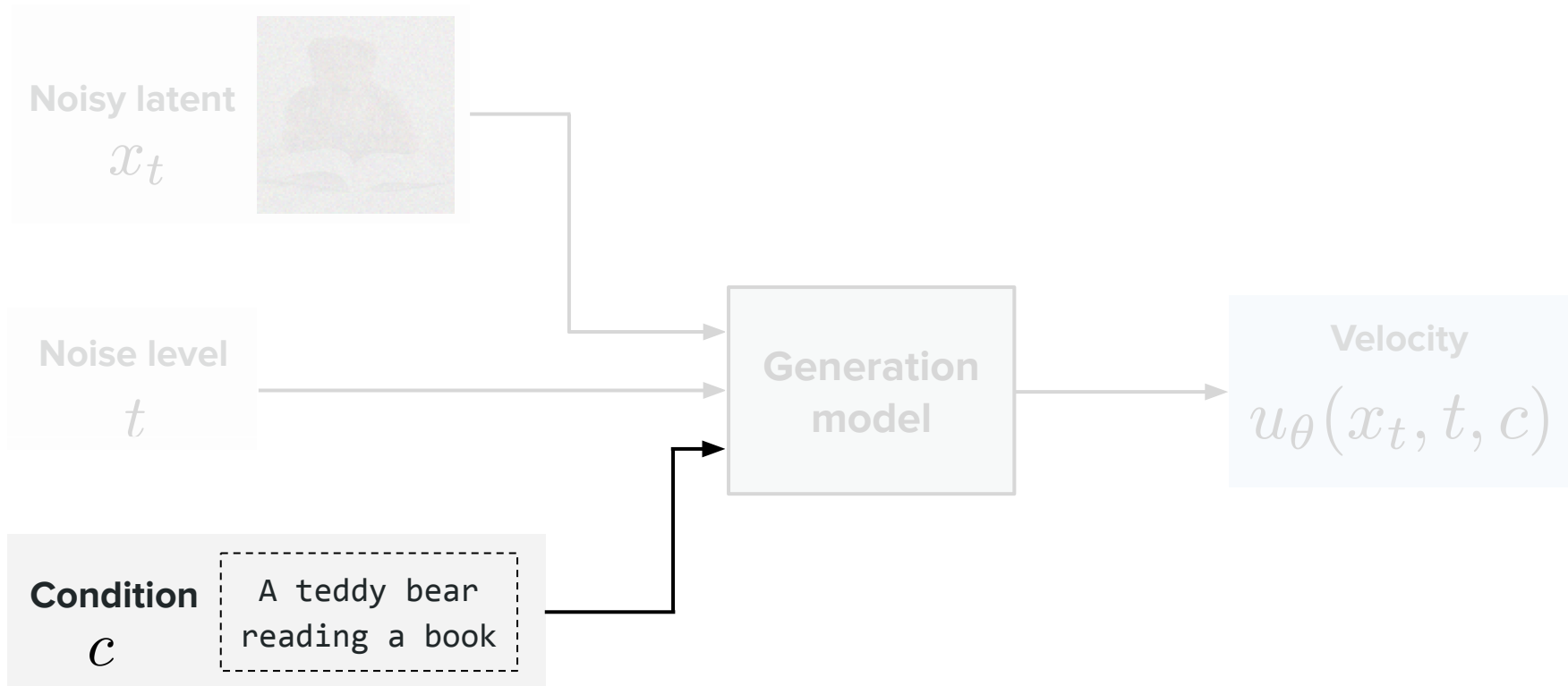
# Problem statement



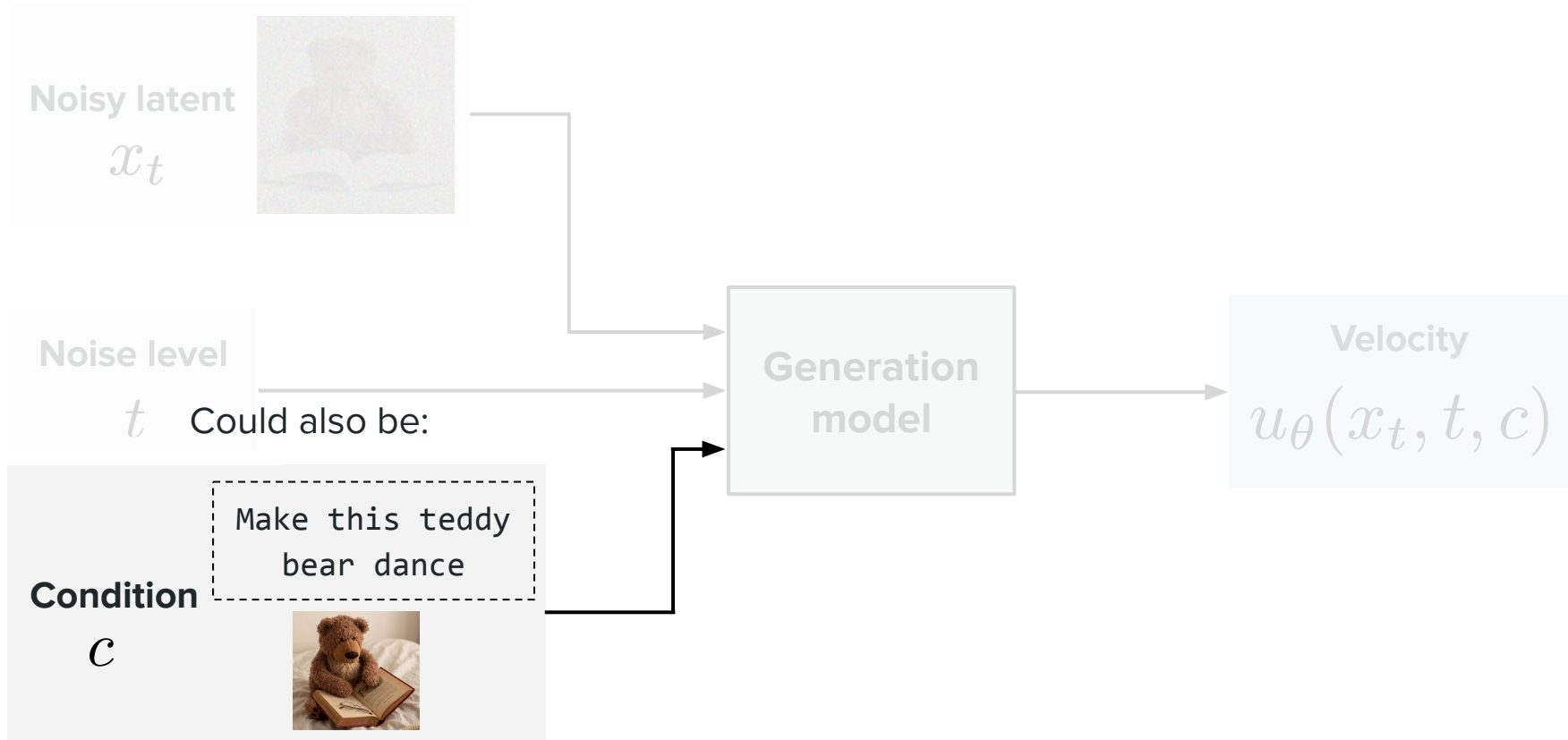
# Objective of this lecture



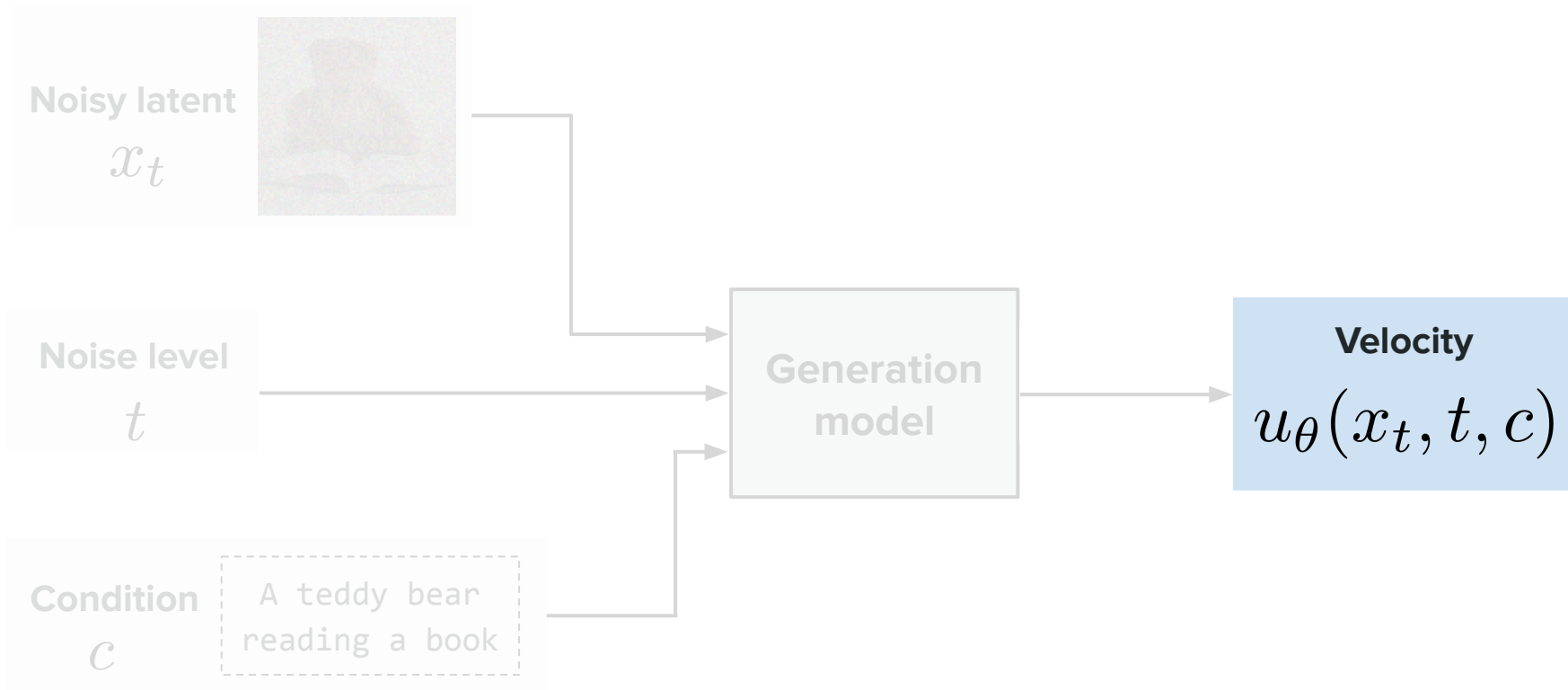
# Side notes



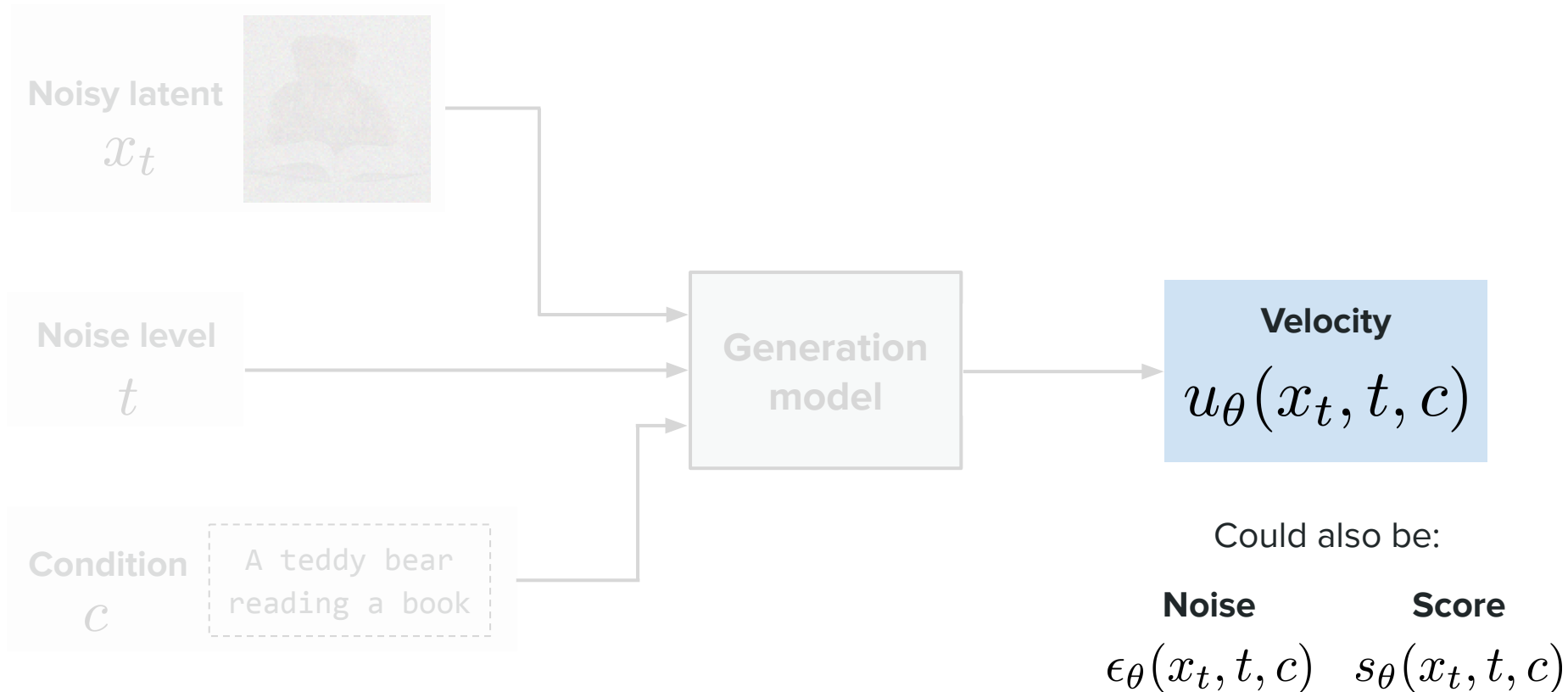
# Side notes



# Side notes



# Side notes



**Objective.** Choose architecture for image generation with criteria below:

- Understand **global structure**
- Preserve **local details**
- Responsive to **external signals** (timestep, condition)
- **Scalable** computationally



# Diffusion & Large Vision Models

Motivation

**U-Net**

Diffusion Transformer

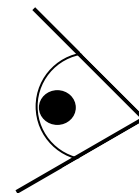
End-to-end example

Multimodal DiT

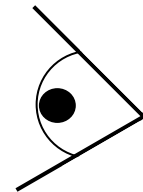
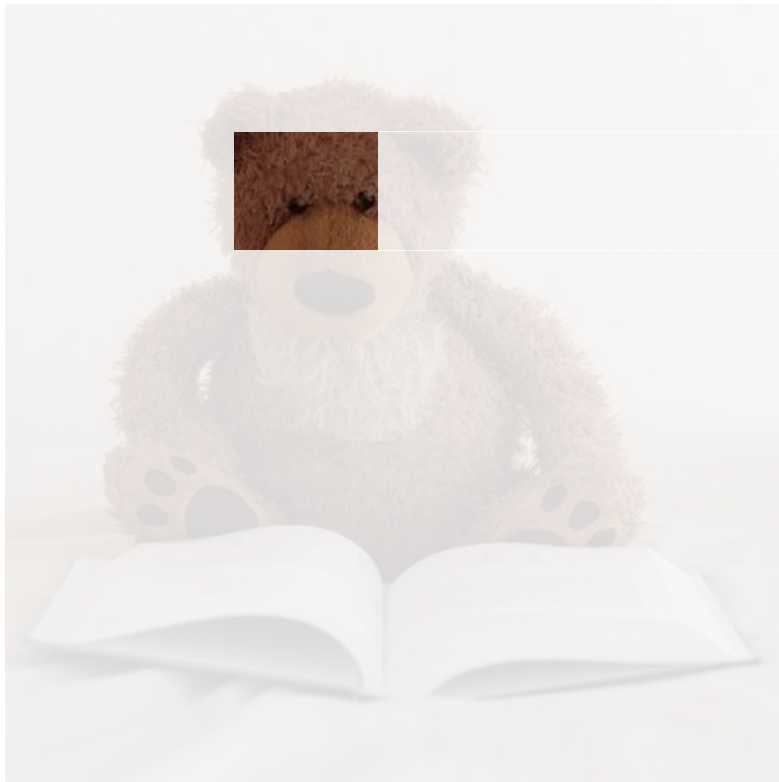
Optimizations

---

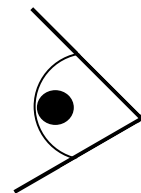
# Intuition



# Intuition

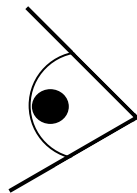
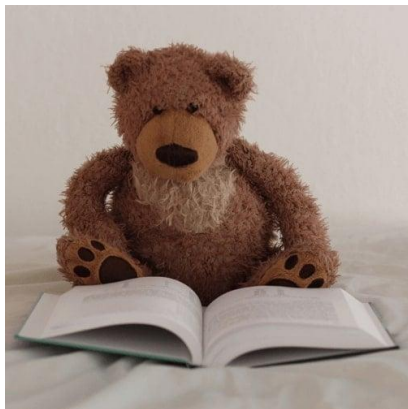


# Intuition



# Intuition behind convolution-based models

**Idea.** Mimic human vision via "inductive" bias baked into the model

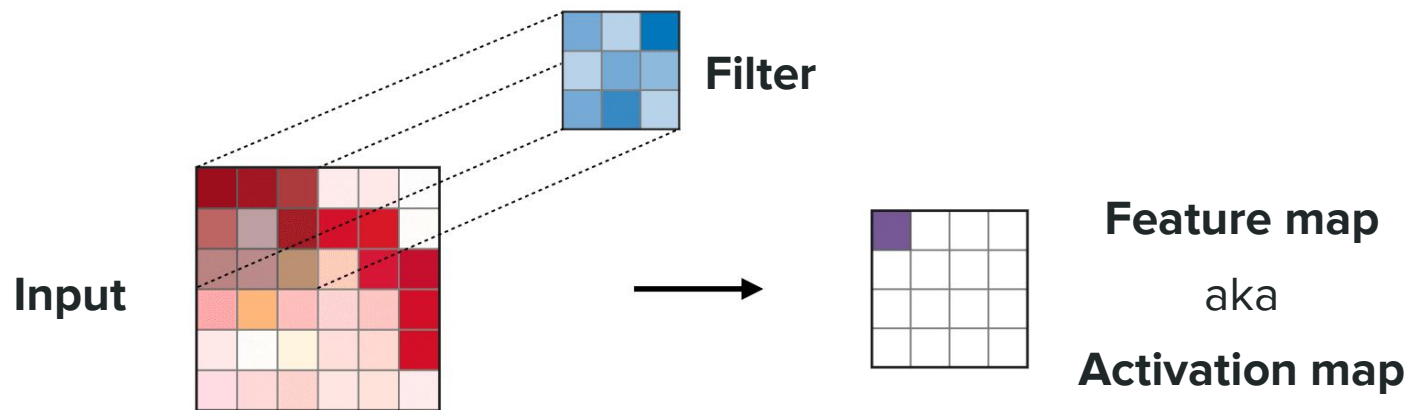


# Convolution operation with filters

**Convolution filter.** Feature detector

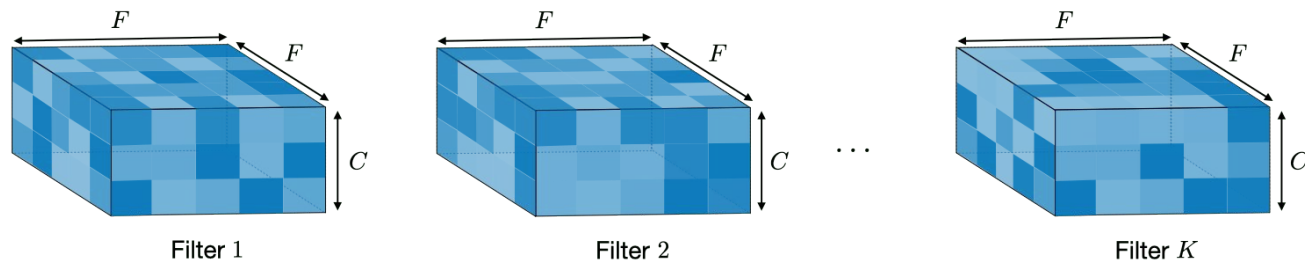


e.g. edges, corners, textures, patterns, shapes



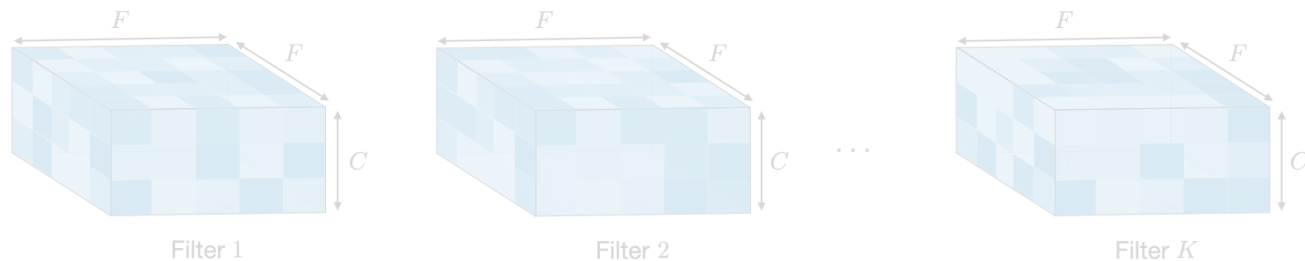
# Convolution operation with filters

**Filter.** Performs convolution operations on the input.

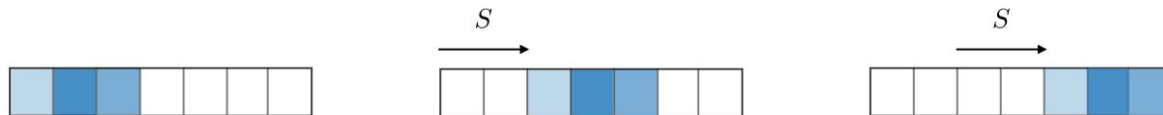


# Convolution operation with filters

**Filter.** Performs convolution operations on the input.

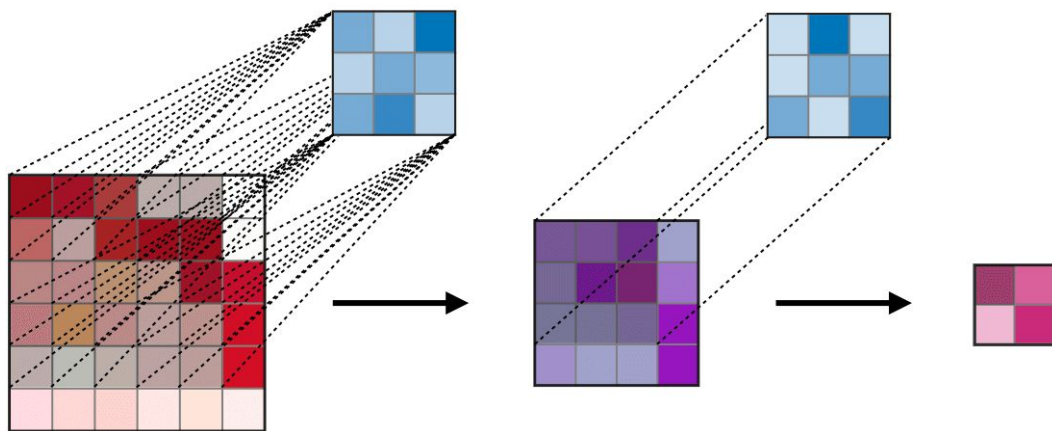


**Stride.** Amount by which the filter moves.



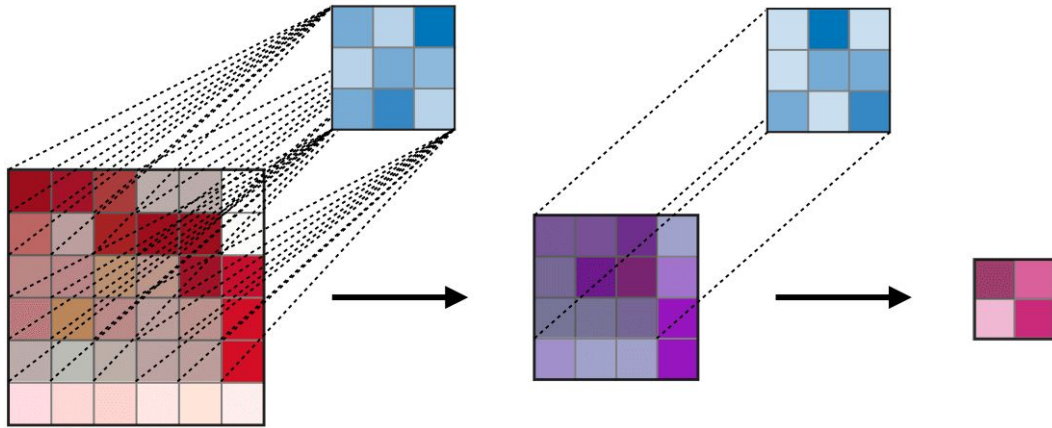
# Receptive field

**Receptive field.** Area that the activation map can "see"



# Receptive field

**Receptive field.** Area that the activation map can "see"



A pixel can see an area of  $R_k \times R_k$

where

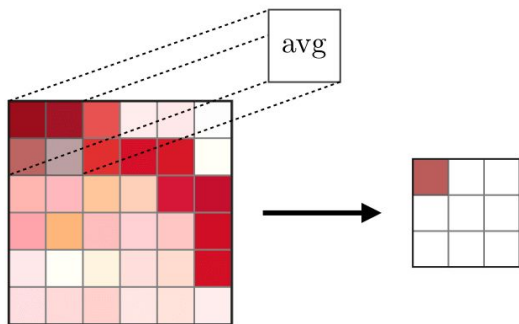
$$R_k = 1 + \sum_{j=1}^k (F_j - 1) \prod_{i=0}^{j-1} S_i$$

Filter size  $\downarrow$   $F_j$

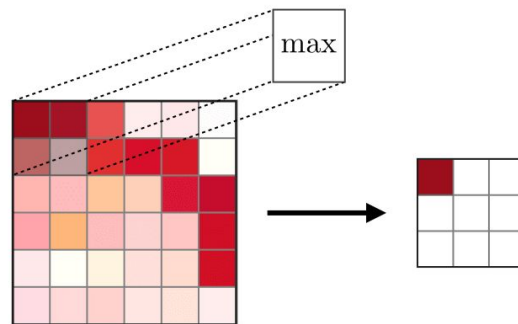
Stride  $\swarrow$   $S_i$

# Downsampling with pooling operation

**Pooling.** Downsampling operation, per channel. Typically after convolution.



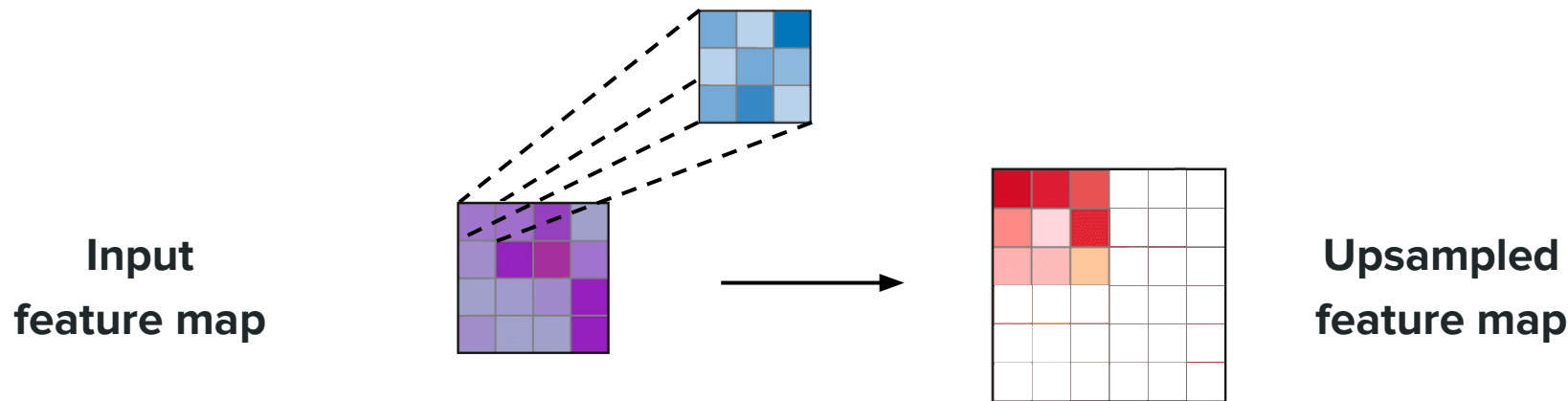
**Average** pooling



**Max** pooling

# Upsampling with transpose convolution

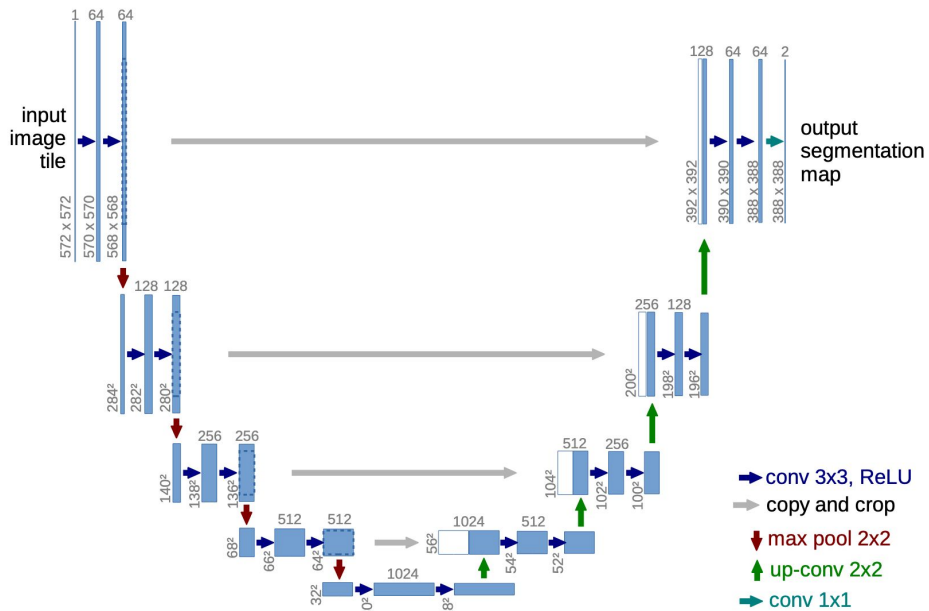
**Transpose convolution.** Upsampling operation, equivalent of "broadcasting"



# Proposal of a convolution-based architecture

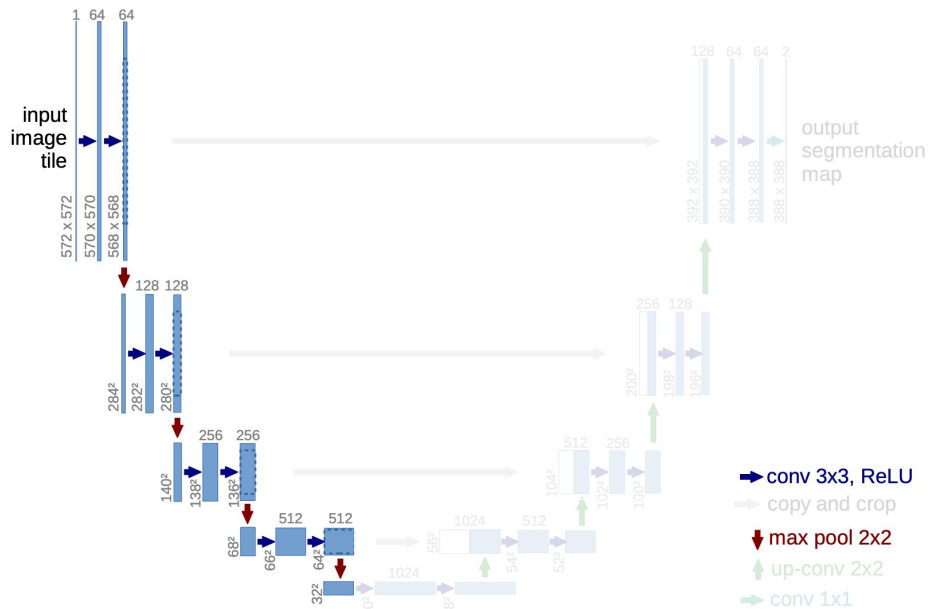
## Overview

- Captures both **local** and **global** features
- **Inductive** bias that is relevant to reverse diffusion process
- Same dimensions for input / output





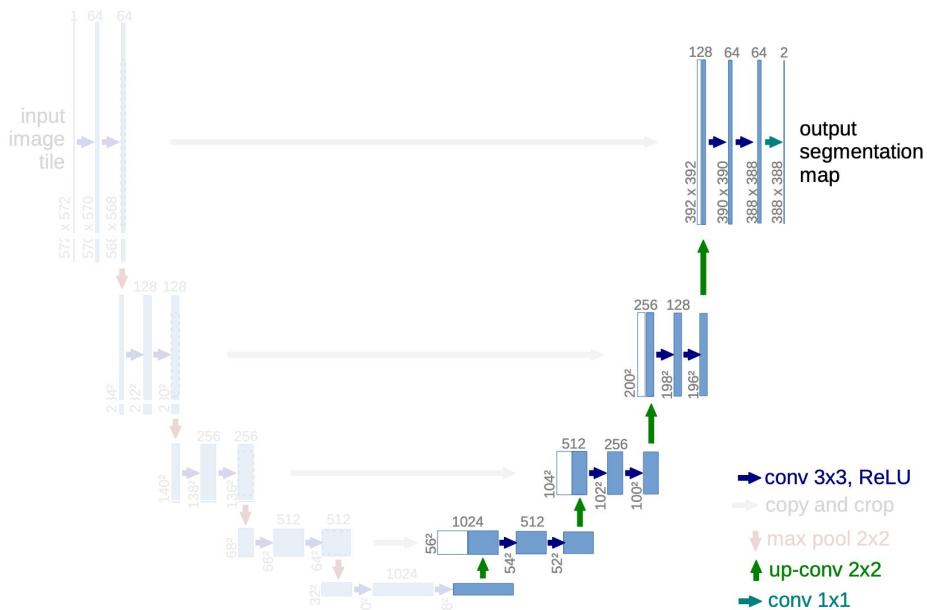
# U-Net: downsampling phase



## Downsampling

- Mix of **convolution** and max **pooling**
- Number of features is **halved** at each step
- **Extraction** of relevant features
- Similar to an "**Encoder**"

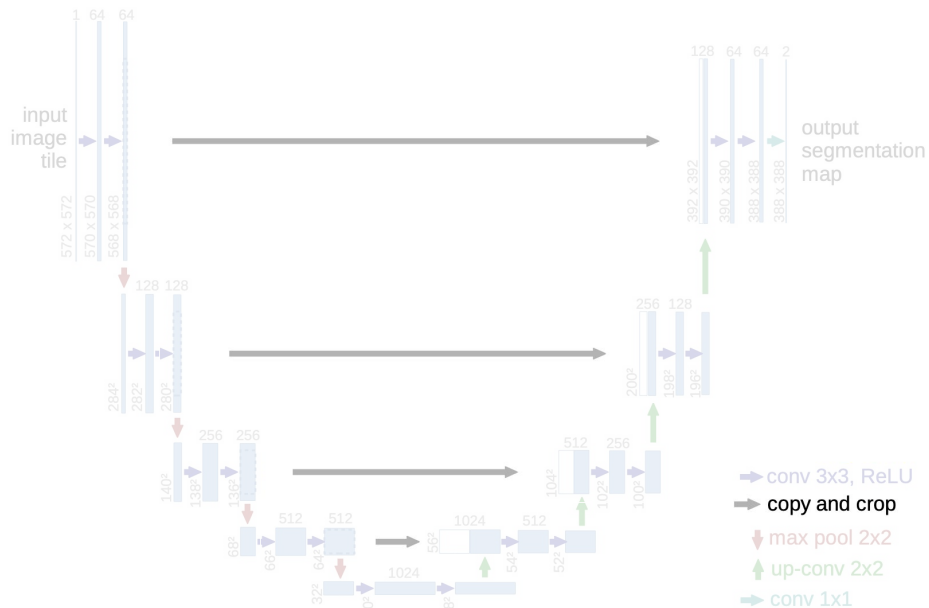
# U-Net: upsampling phase



## Upsampling

- Mix of **up-convolution** and **convolutions**
- **Number of features is doubled** at each step
- **1x1 convolution** at the last step before obtaining output
- Similar to a "**Decoder**"

# U-Net: residual connections

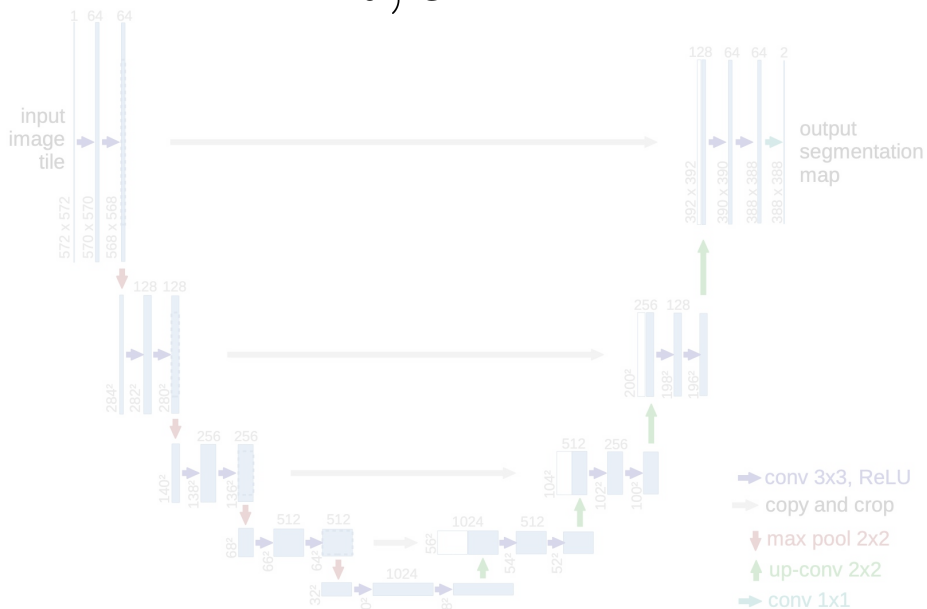


## "Copy and crop" connections

- Avoids losing **local** information
- "**Highway**" of information

# U-Net: adding condition information

$t, c$



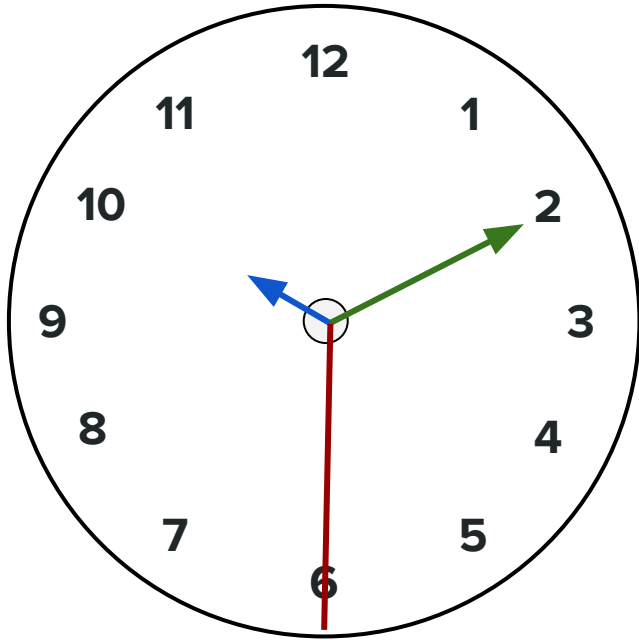
How can we add  
time and class label  
information?

# Representation of timestep

**Hour** (slow)

**Minute** (fast)

**Second** (very fast)

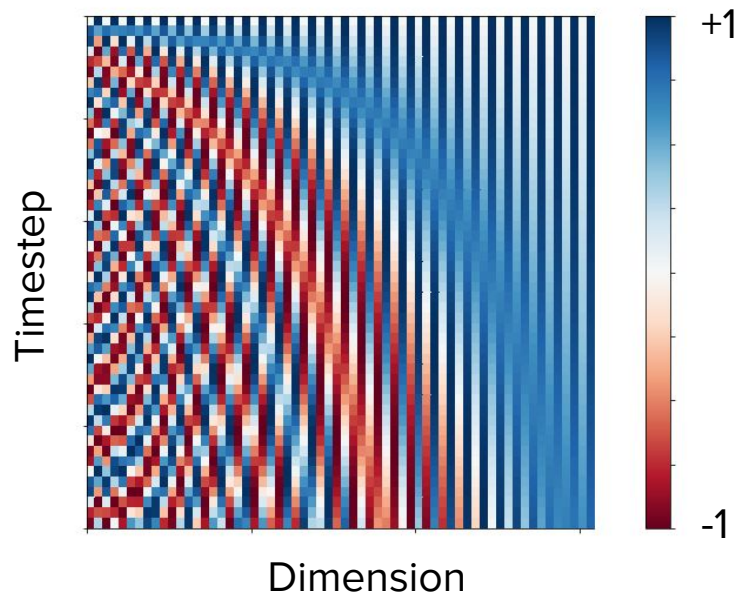


# Representation of timestep

Hour (slow)

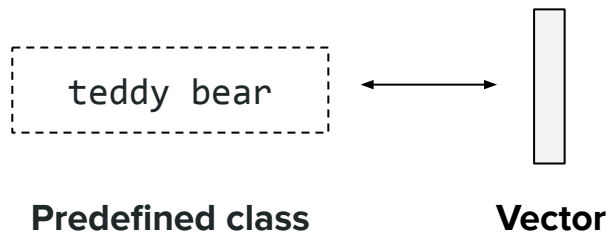
Minute (fast)

Second (very fast)



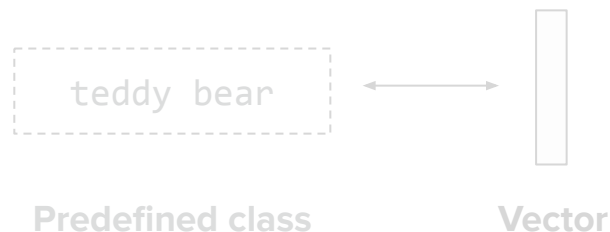
# Representation of class label

- Embedding corresponding to predefined class

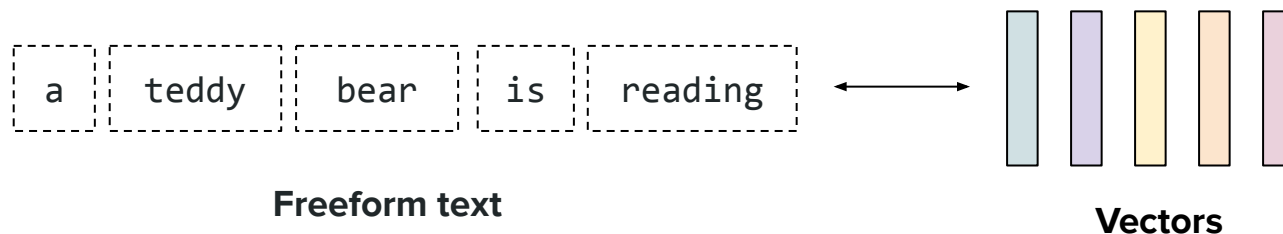


# Representation of class label

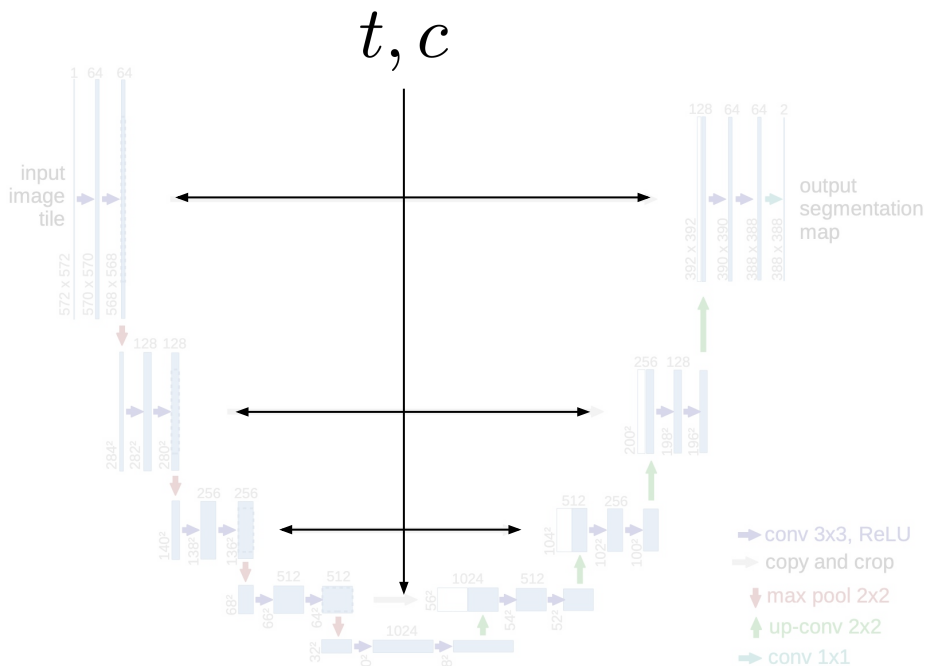
- Embedding corresponding to predefined class



- Rich embeddings from a pretrained LLM



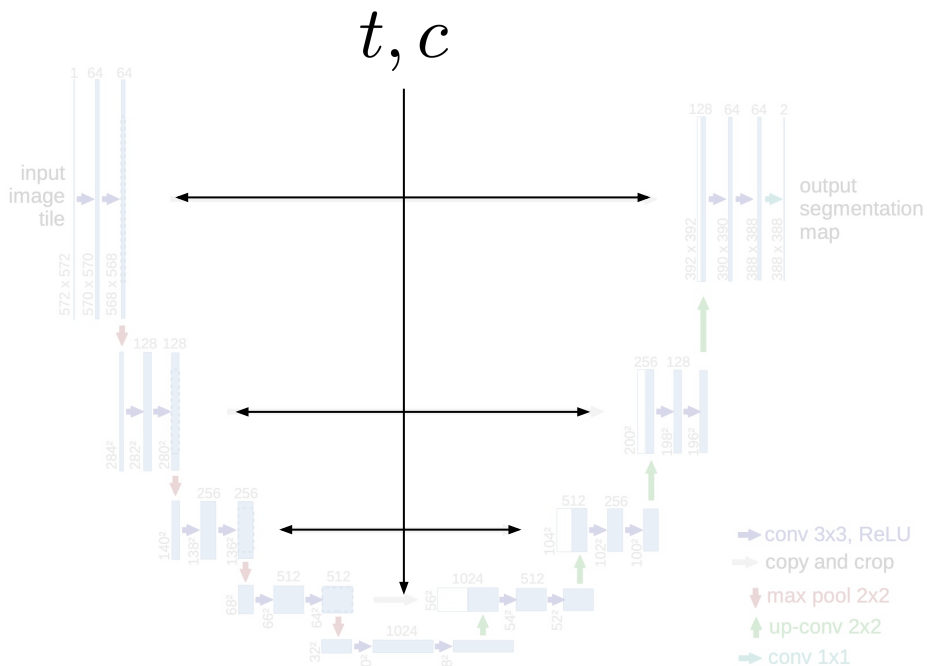
# U-Net: adding condition information



## Method 1

Added to the feature map

# U-Net: adding condition information



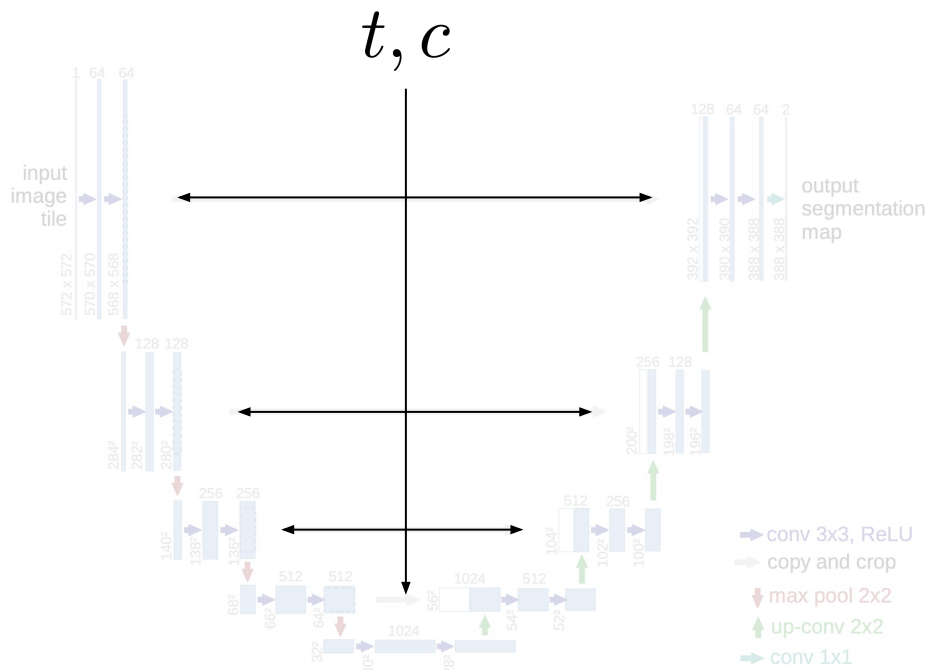
## Method 1

Added to the feature map

## Method 2

Modulate feature map via scaling / shifting

# U-Net: adding condition information



## Method 1

Added to the feature map

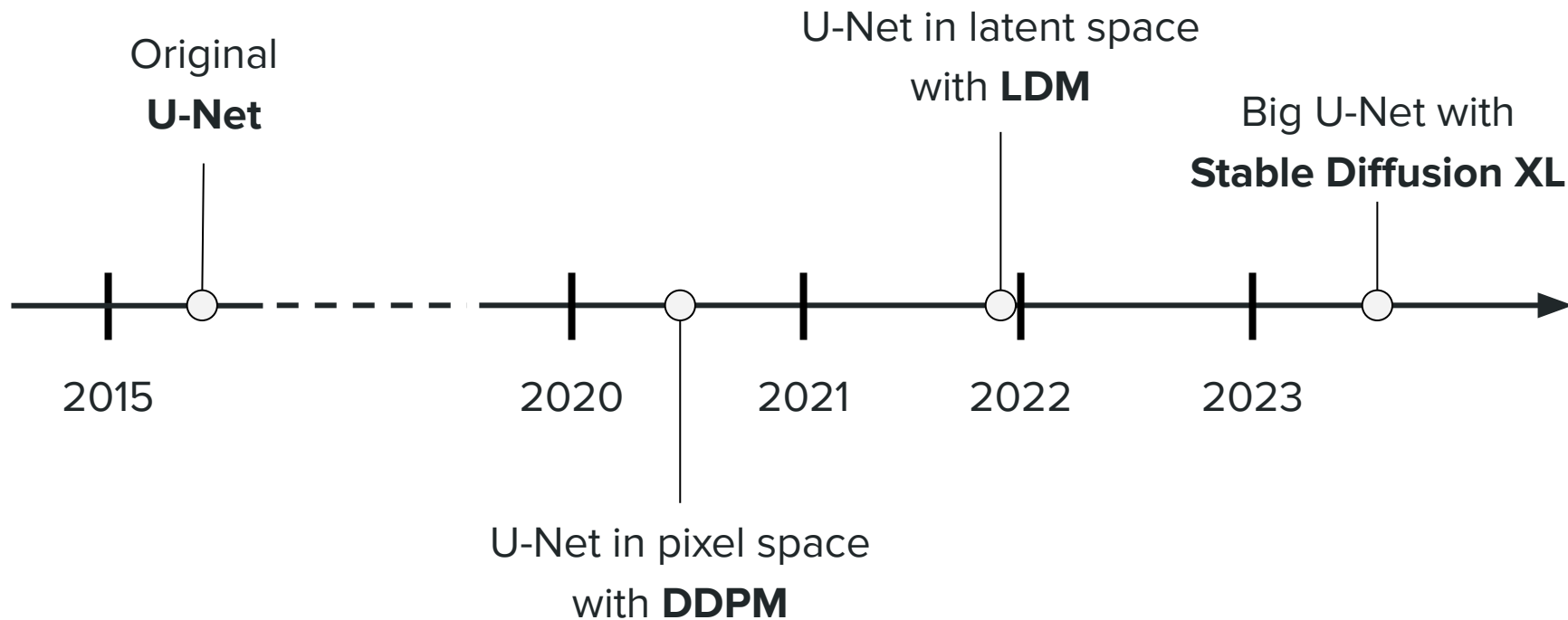
## Method 2

Modulate feature map via scaling / shifting

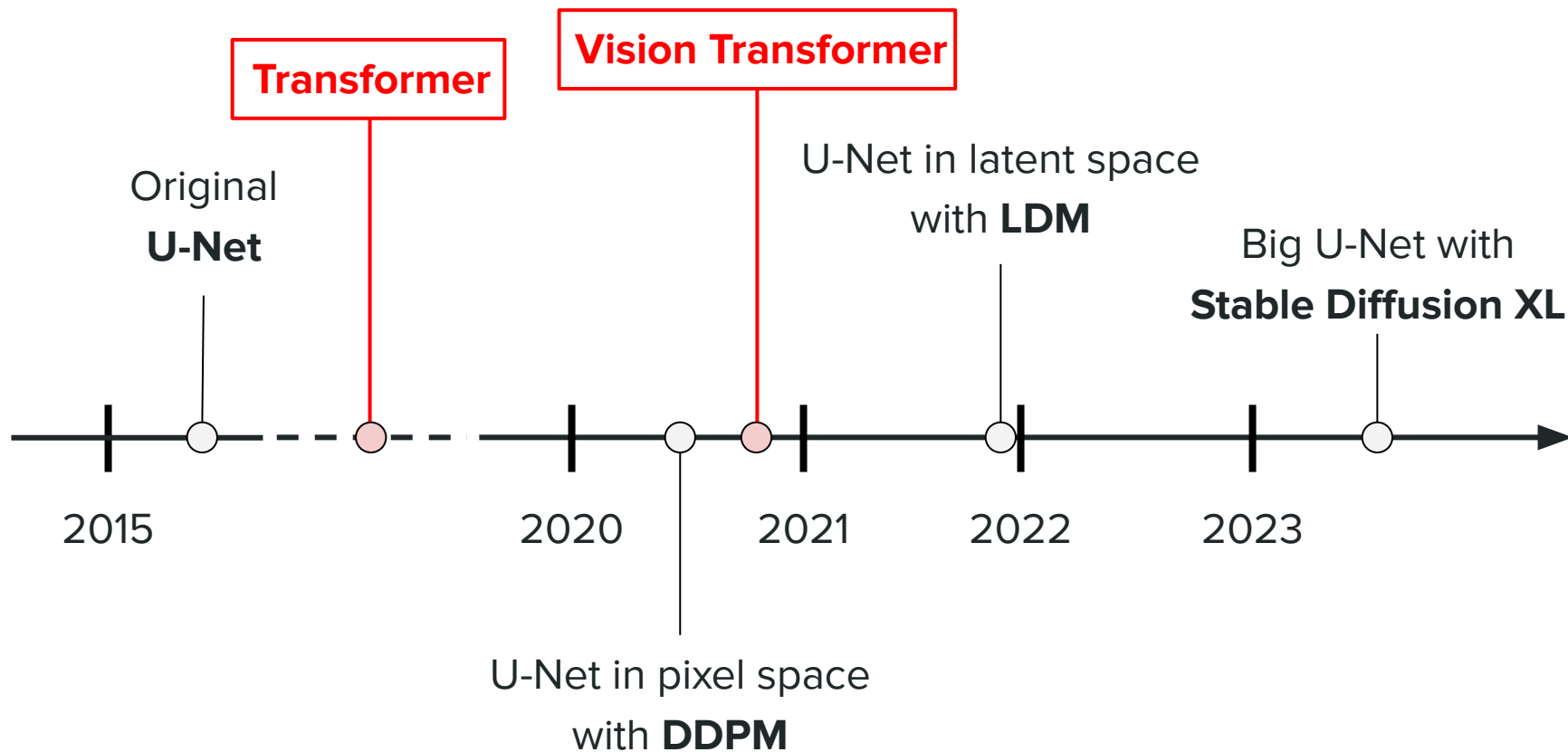
## Method 3

Cross-attention of condition with feature map

# Timeline of U-Net-based models



# Timeline of U-Net-based models





# Diffusion & Large Vision Models

Motivation

U-Net

**Diffusion Transformer**

End-to-end example

Multimodal DiT

Optimizations

---

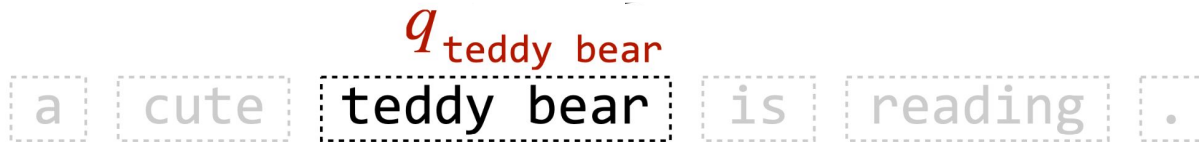
# Motivating example



Need to preserve **local details**  
across **long distances**

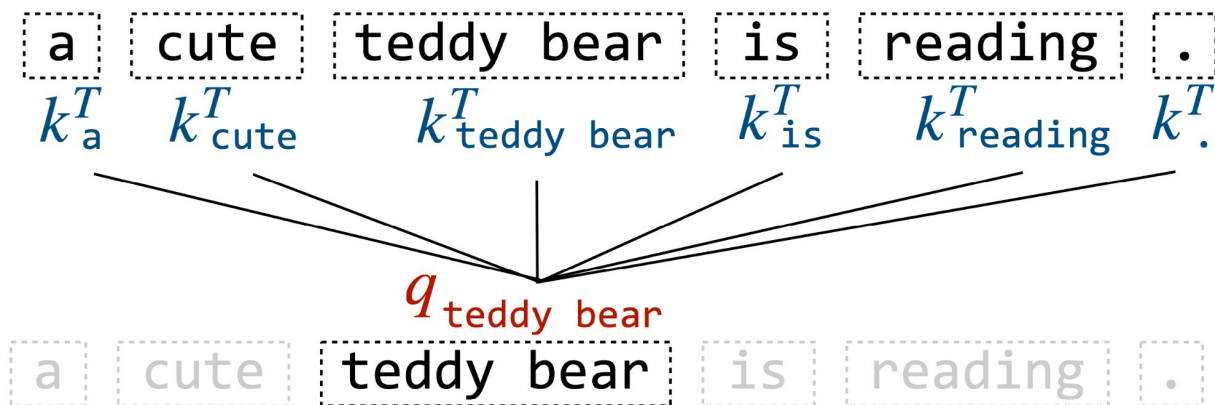
# Attention mechanism

**Idea.** Remove inductive bias



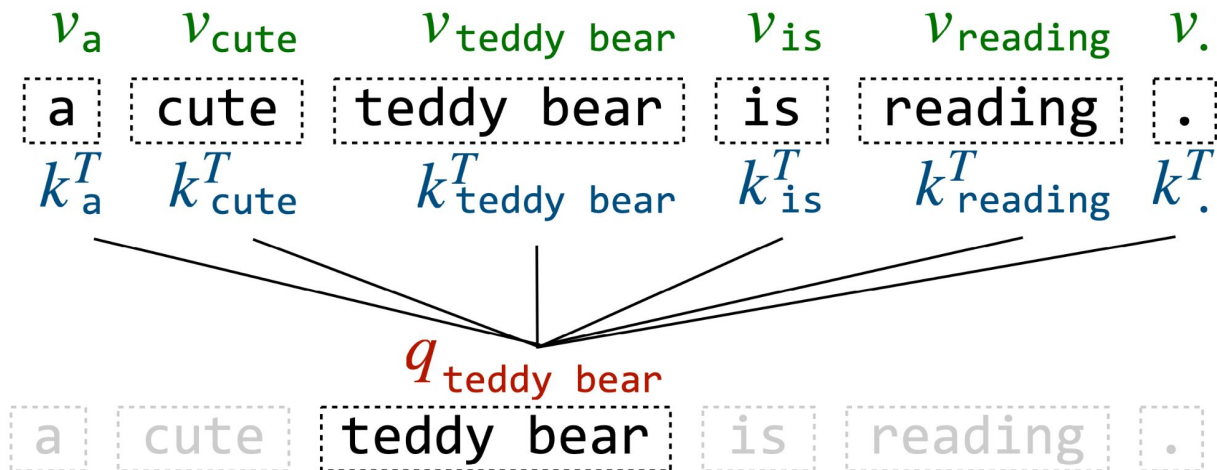
# Attention mechanism

**Idea.** Remove inductive bias

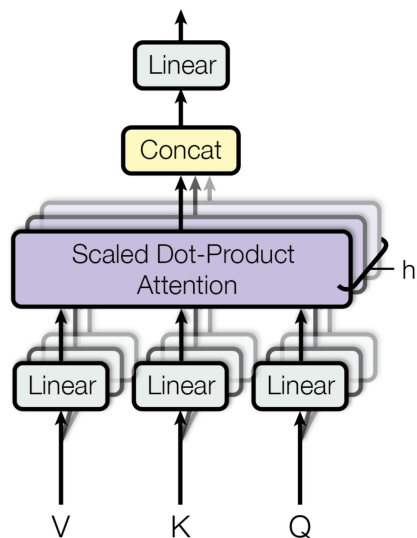


# Attention mechanism

**Idea.** Remove inductive bias



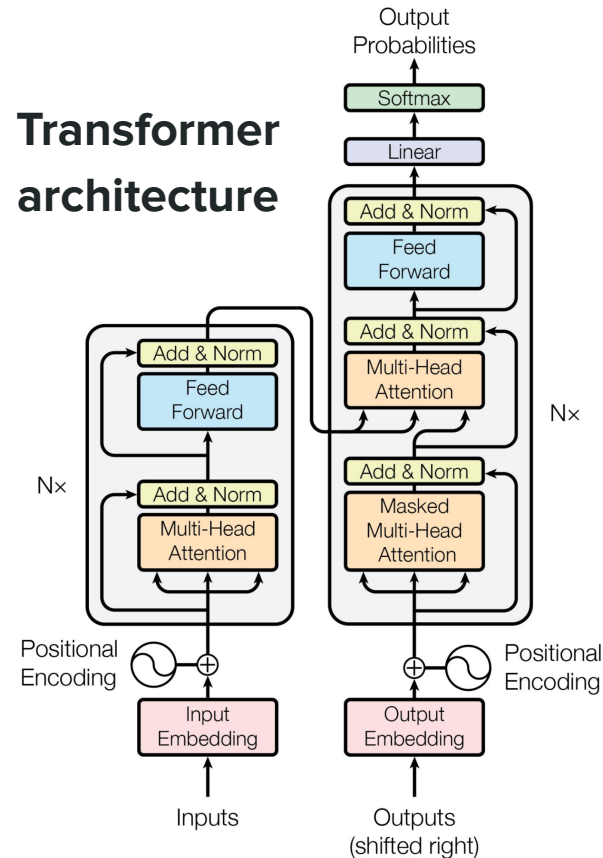
# Transformer



**Multi-head attention layer**

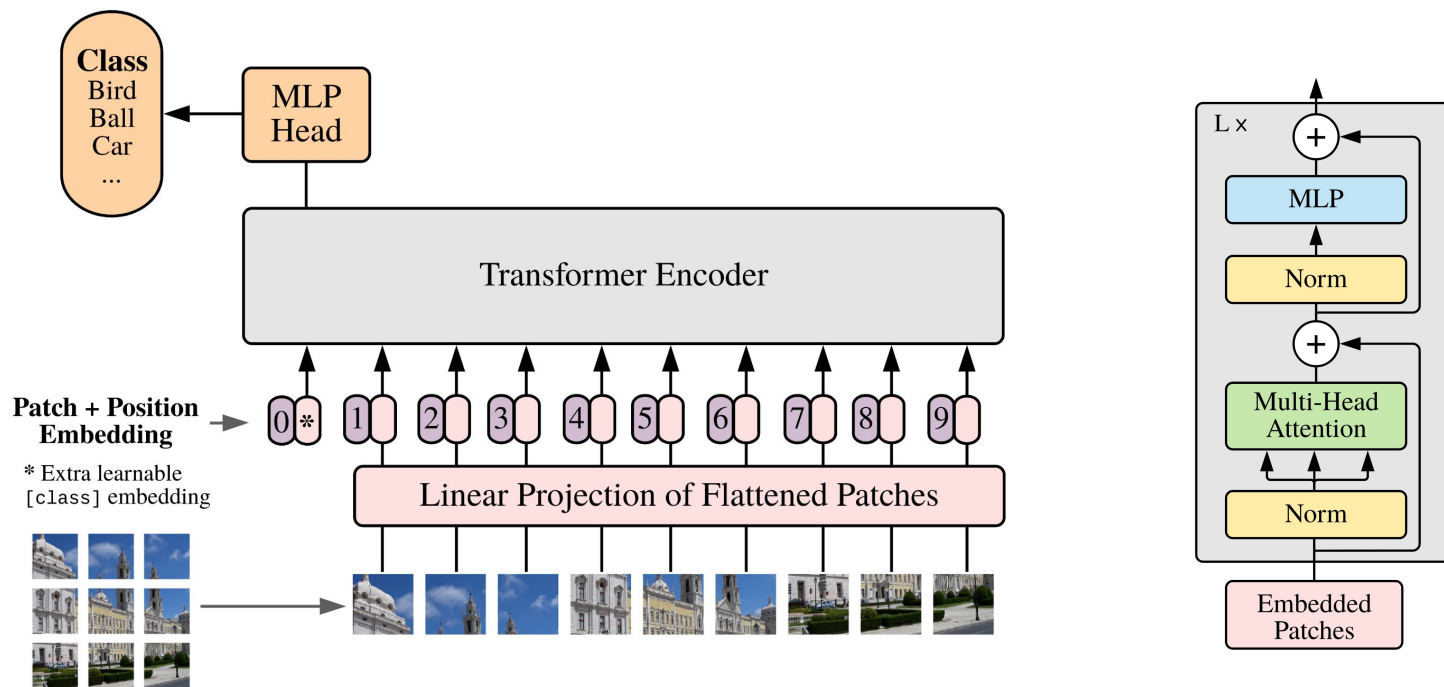
$$\text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

**Transformer architecture**



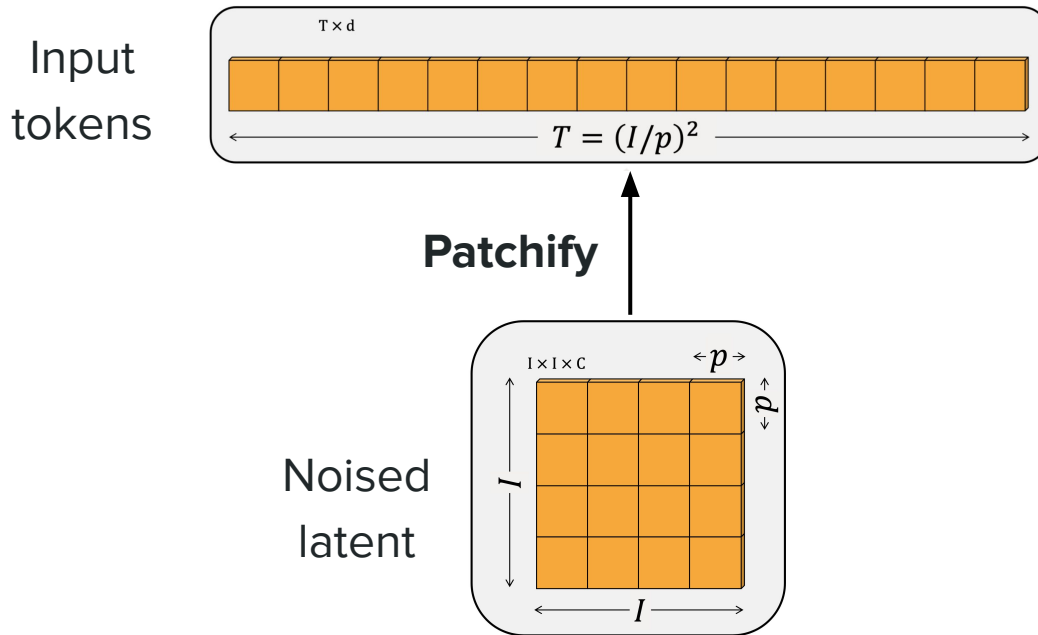
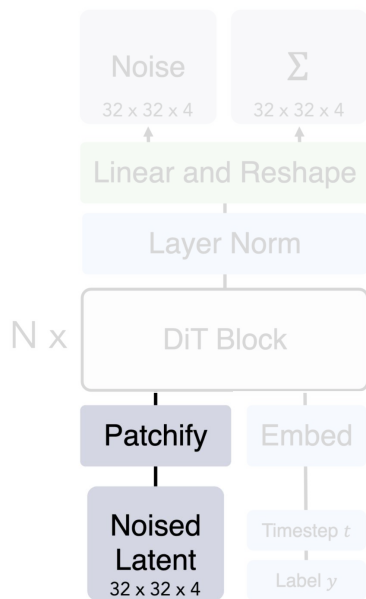
# Transformer for vision understanding...

## ViT = Vision Transformer

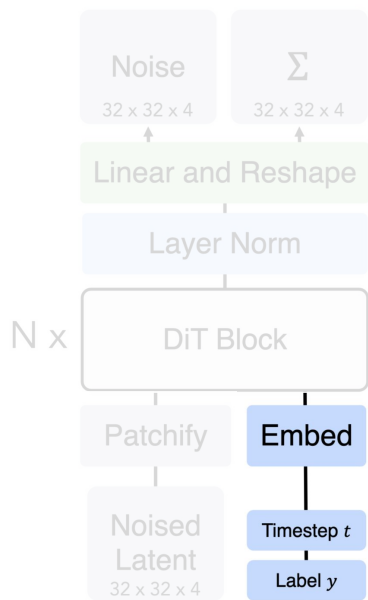




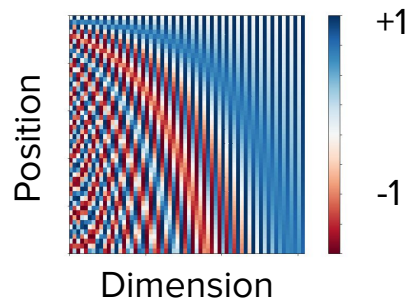
# DiT: tokenization of the input via patches



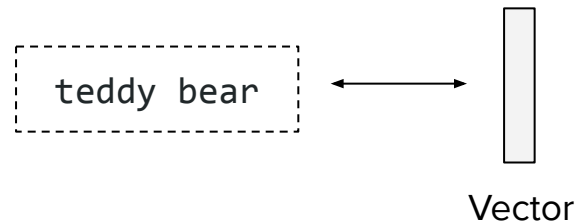
# DiT: embedding of conditions



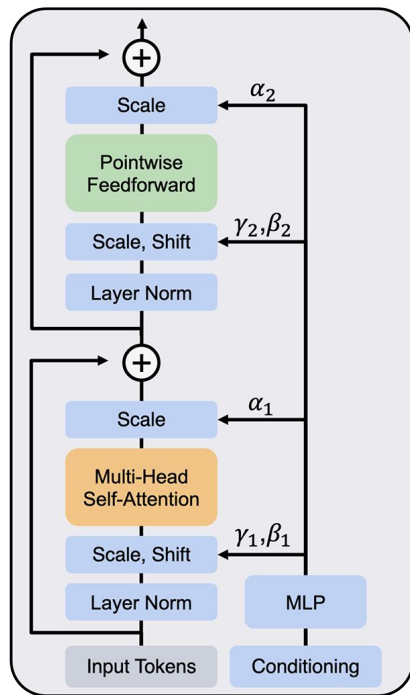
**Timestep**  
embedding



**Label**  
embedding

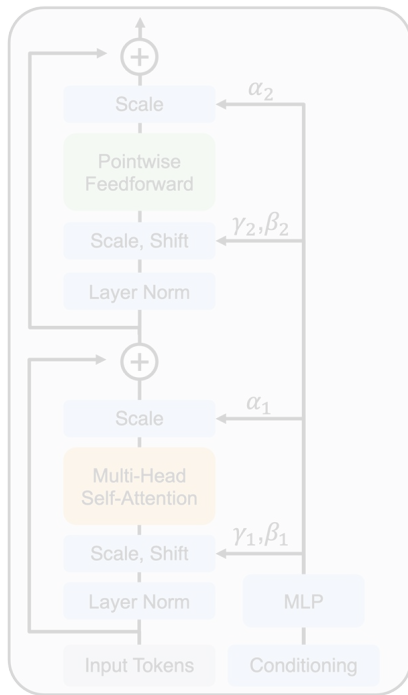


# DiT: injecting conditions with adaptive layer norm

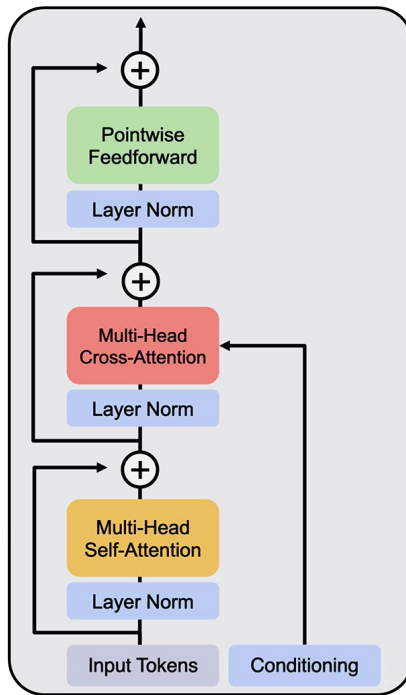


## Adaptive Layer Norm

# DiT: injecting conditions with cross-attention

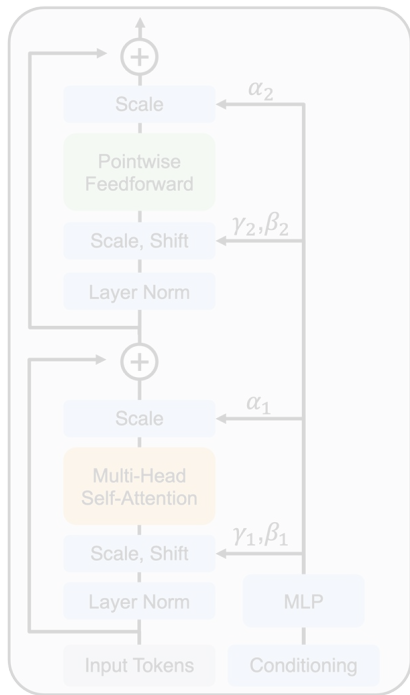


Adaptive Layer Norm

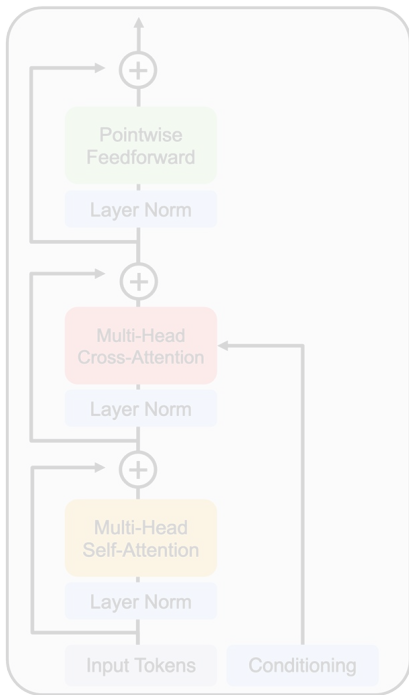


Cross-attention

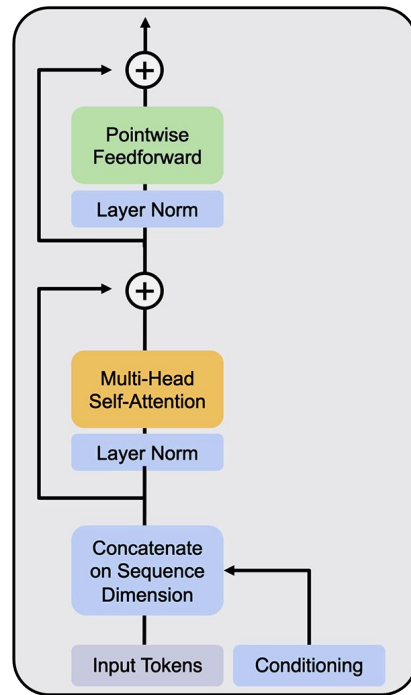
# DiT: injecting conditions with in-context conditioning



Adaptive Layer Norm

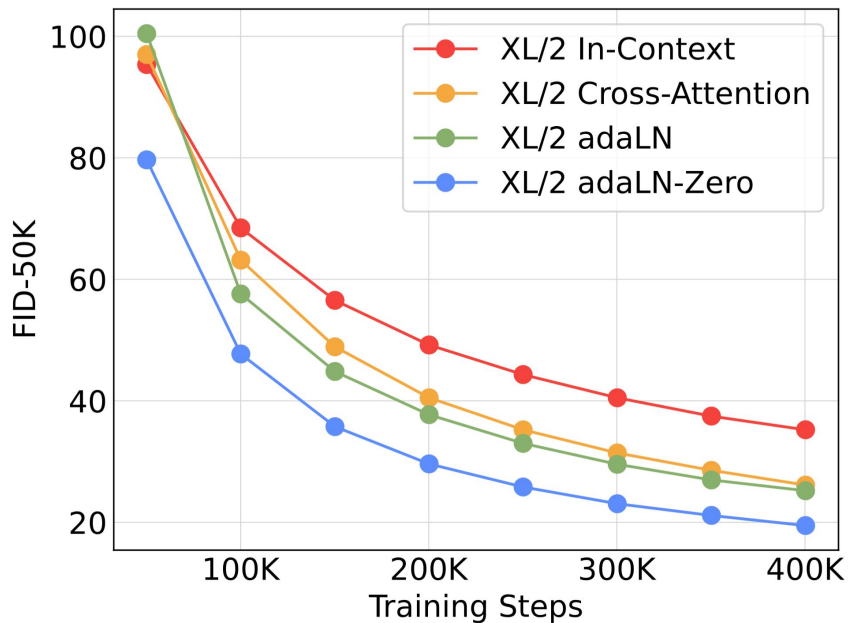


Cross-attention



In-context conditioning

# DiT: comparison of different condition injections



Adaptive Layer Norm

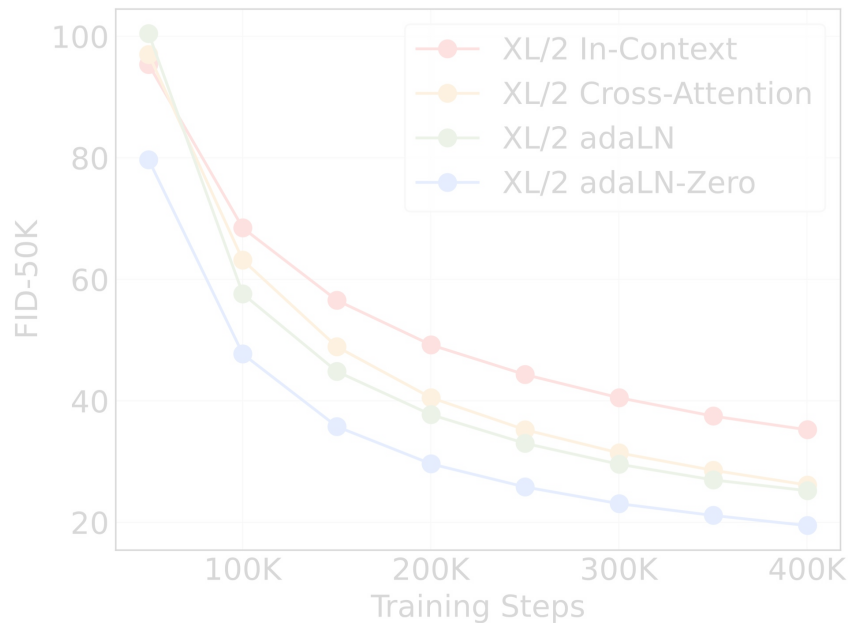


Cross-attention



In-context conditioning

# DiT: comparison of different condition injections



**Adaptive Layer Norm**



Cross-attention



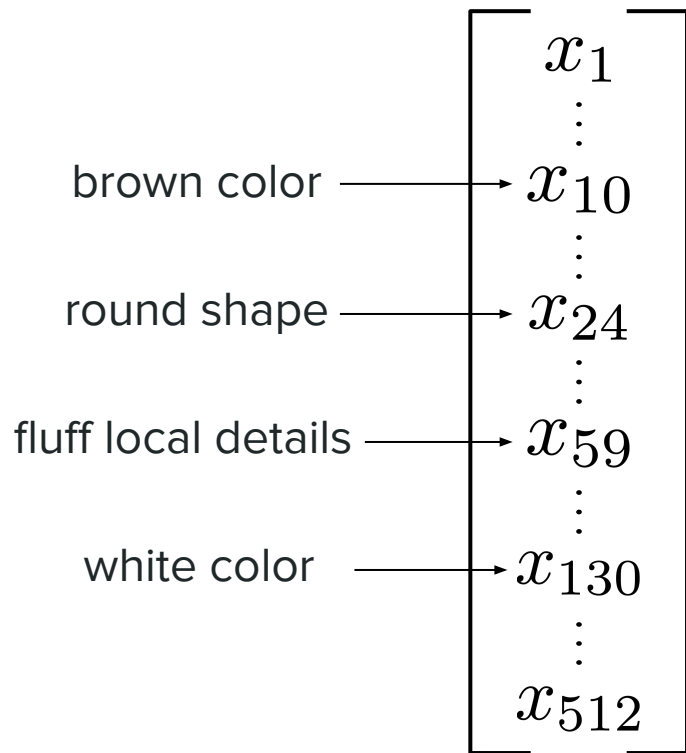
In-context conditioning

# Intuition behind adaptive layer norm

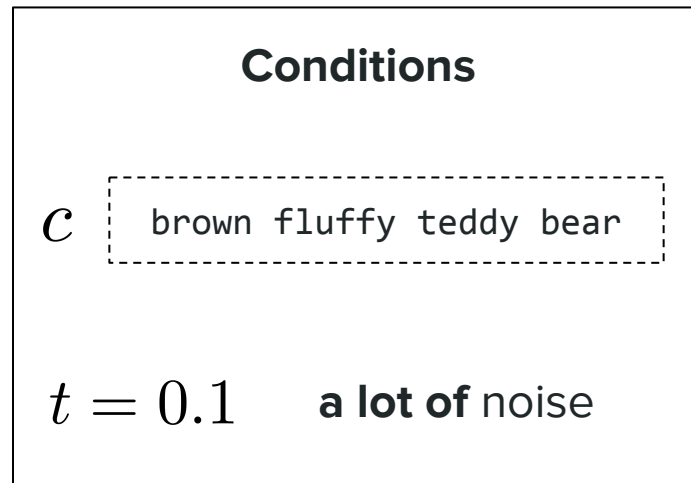
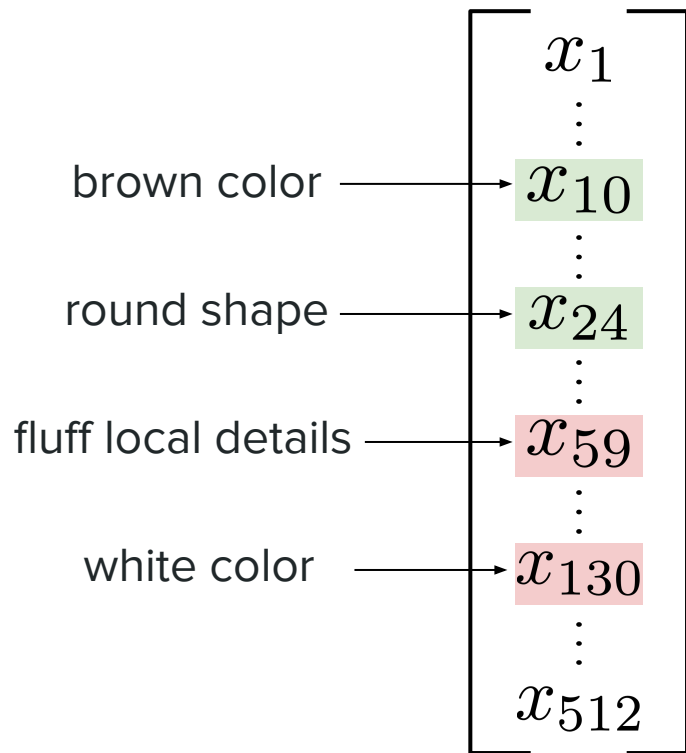
$$\mathcal{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

**Representation of a  
patch embedding**

# Intuition behind adaptive layer norm

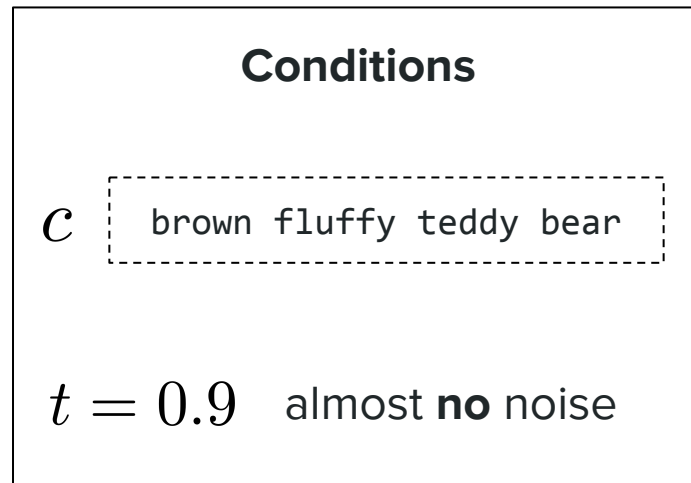
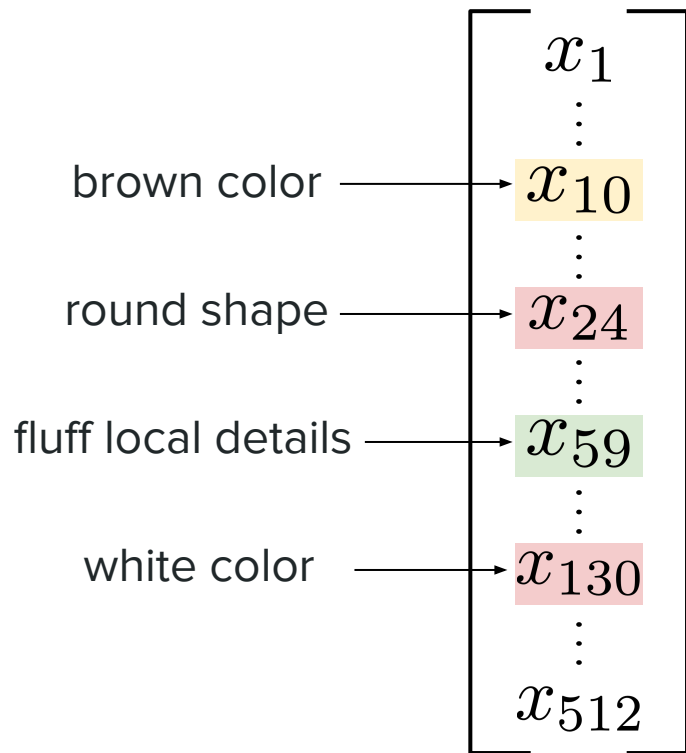


# Intuition behind adaptive layer norm



Need to focus on **global structure** since at the early stage of generation.

# Intuition behind adaptive layer norm



Need to focus on **local details** since almost arriving at the final image.

# Adaptive layer normalization

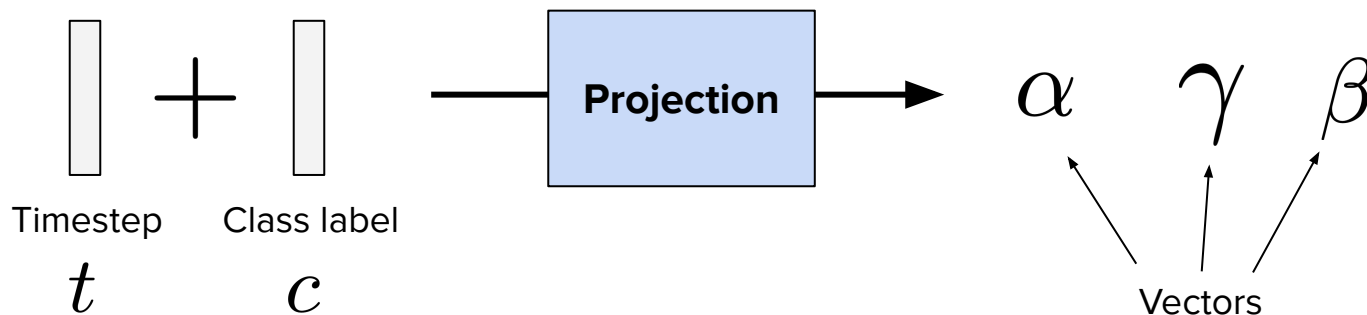
**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

# Adaptive layer normalization

Idea. "Modulate" vector with respect to conditions (timestep and class label)

## Steps.

1. Determine intensity of modulation as a function of conditions



# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

## Steps.

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$x$

# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

## **Steps.**

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$$\text{LN}(x)$$

# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

## Steps.

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$$\text{LN}(x) \star (1 + \gamma)$$

# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

## Steps.

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$$\text{LN}(x) \star (1 + \gamma) + \beta$$

# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

## Steps.

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$$\text{Operation}(\text{LN}(x) \star (1 + \gamma) + \beta)$$

# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

## Steps.

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$$\alpha \star \text{Operation}(\text{LN}(x) \star (1 + \gamma) + \beta)$$

# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

## Steps.

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$$x + \alpha \star \text{Operation}(\text{LN}(x) \star (1 + \gamma) + \beta)$$

# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

## Steps.

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$$x \longleftarrow x + \alpha \star \text{Operation}(\text{LN}(x) \star (1 + \gamma) + \beta)$$

# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

## Steps.

1. Determine intensity of modulation as a function of conditions
- 2. Modulate token with quantities**

$$x \leftarrow x + \alpha \star \text{Operation}(\text{LN}(x)) \star (1 + \gamma) + \beta$$

**element-wise multiplication**

# Adaptive layer normalization

**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

**Steps.**

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$$x \leftarrow x + \underset{\substack{\uparrow \\ \text{Gate}}}{\alpha} \star \text{Operation}(\text{LN}(x)) \star (1 + \underset{\substack{\uparrow \\ \text{Scale}}}{\gamma}) + \underset{\substack{\uparrow \\ \text{Shift}}}{\beta}$$

# Adaptive layer normalization

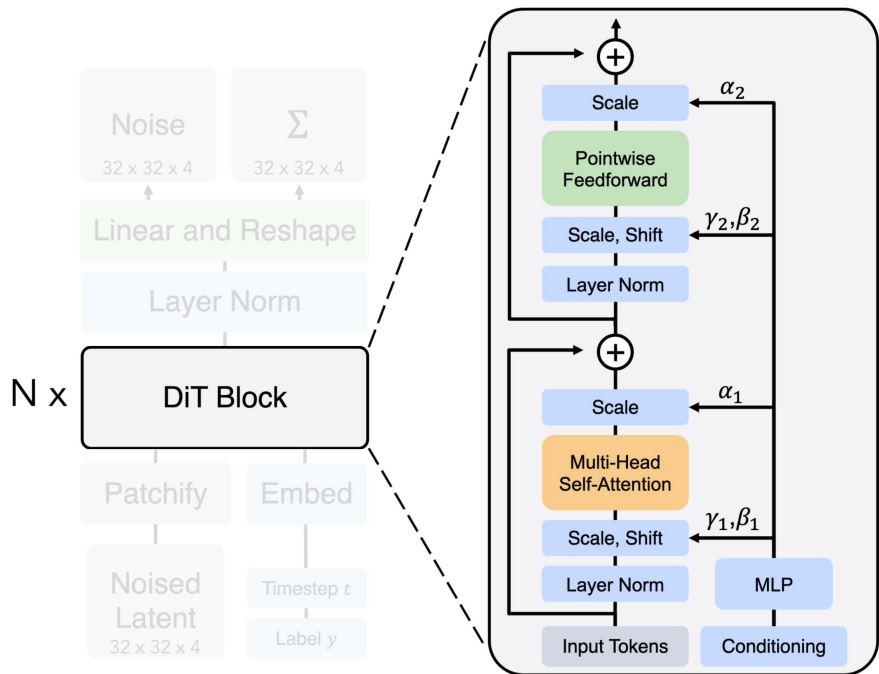
**Idea.** "Modulate" vector with respect to conditions (timestep and class label)

**Steps.**

1. Determine intensity of modulation as a function of conditions
2. Modulate token with quantities

$$x \leftarrow x + \underbrace{\alpha \star \text{Operation}(\text{LN}(x) \star (1 + \gamma) + \beta)}_{\text{At the beginning: } = 0} \quad \text{adaLN-Zero}$$

# DiT: injection of conditions with adaLN-Zero

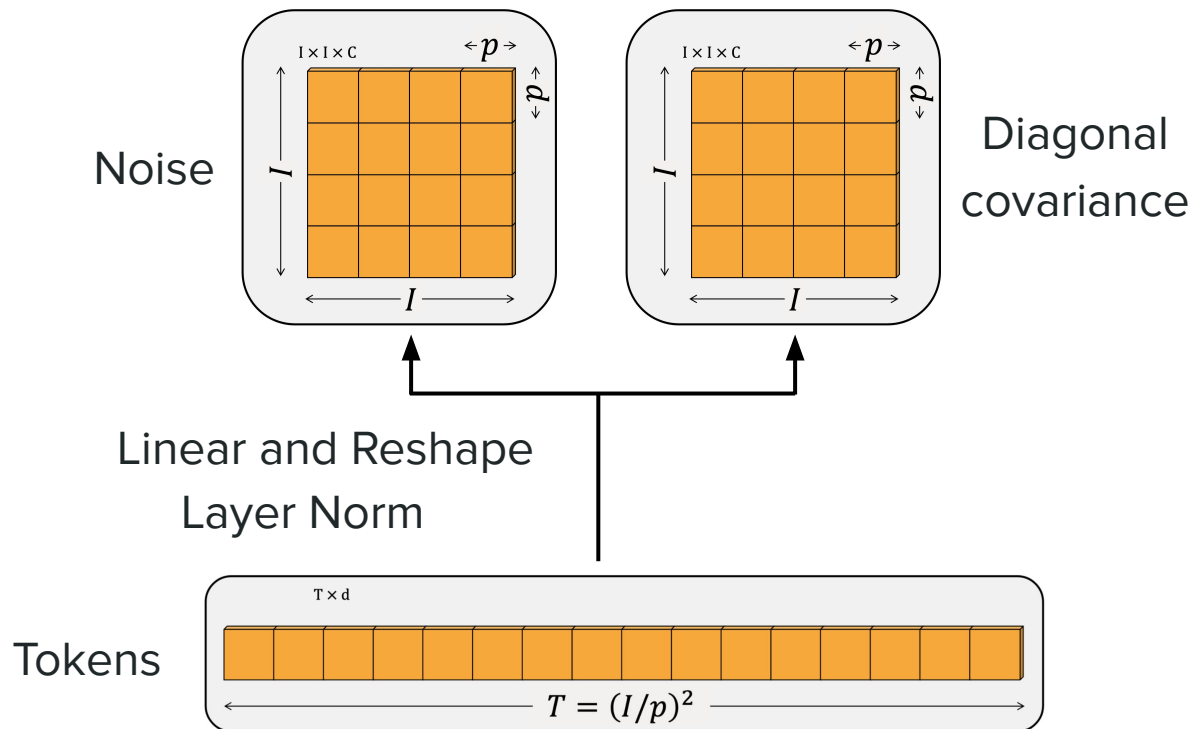
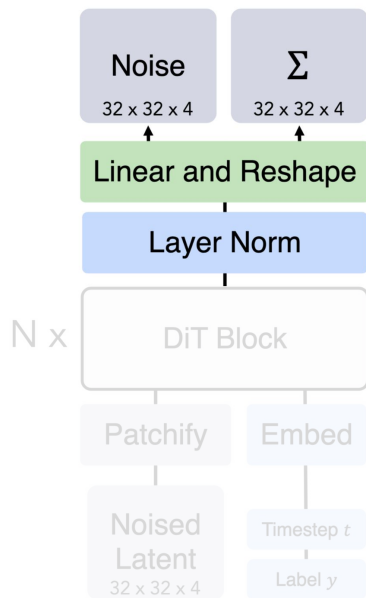


**Modulation** of representation with:

- **Gate**  $\alpha$
- **Scale**  $\gamma$
- **Shift**  $\beta$

as a function of  $t, c$

# DiT: reformatting the output



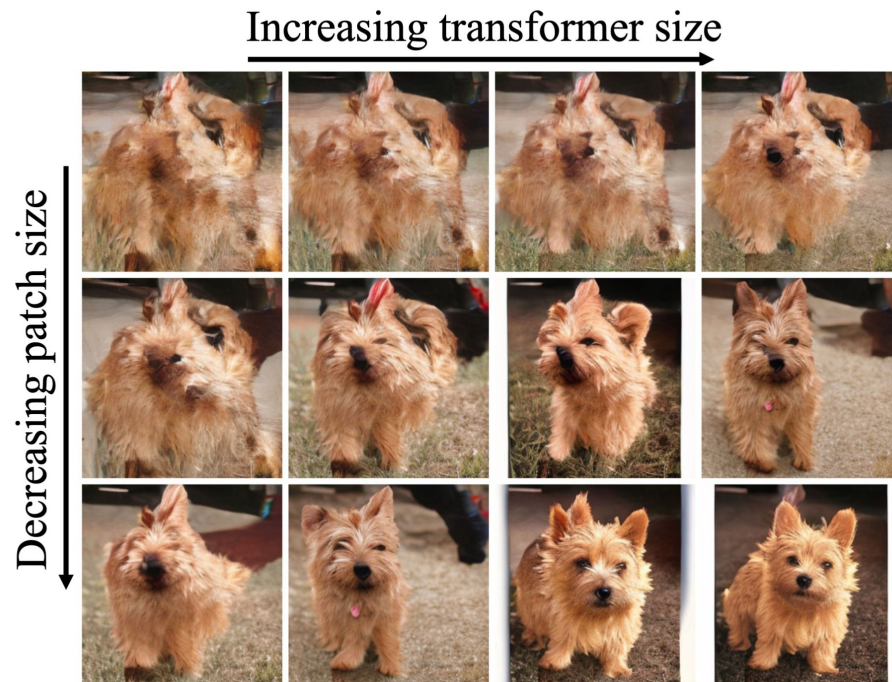
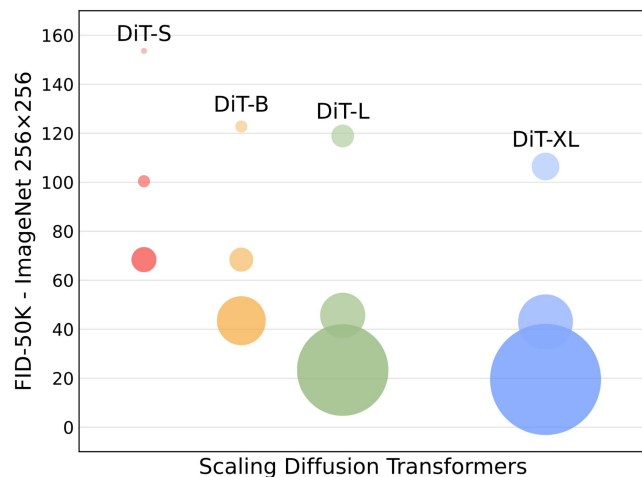
# Scaling up DiT...

Model	Layers $N$	Hidden size $d$	Heads	Gflops ( $L=32, p=4$ )
DiT-S	12	384	6	1.4
DiT-B	12	768	12	5.6
DiT-L	24	1024	16	19.7
DiT-XL	28	1152	16	29.1

- Number of parameters is **not enough** to quantify model complexity
- **FL**loating-point **OP**erations (FLOPs) = # operations in forward pass
- Smaller patch size corresponds to **higher FLOPs**

# Scaling up DiT... improves results!

Model	Layers $N$	Hidden size $d$	Heads	Gflops ( $T=32, p=4$ )
DiT-S	12	384	6	1.4
DiT-B	12	768	12	5.6
DiT-L	24	1024	16	19.7
DiT-XL	28	1152	16	29.1





# Diffusion & Large Vision Models

Motivation

U-Net

Diffusion Transformer

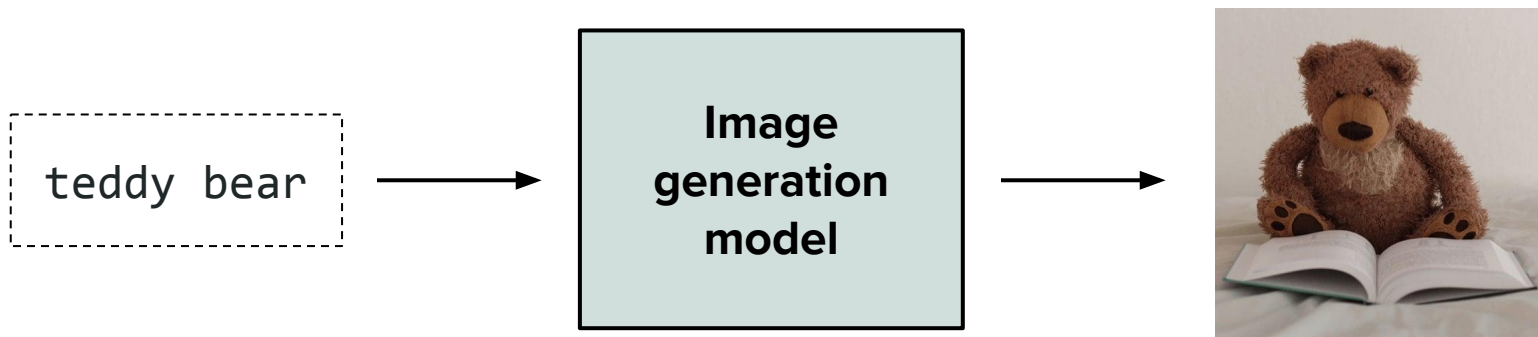
**End-to-end example**

Multimodal DiT

Optimizations

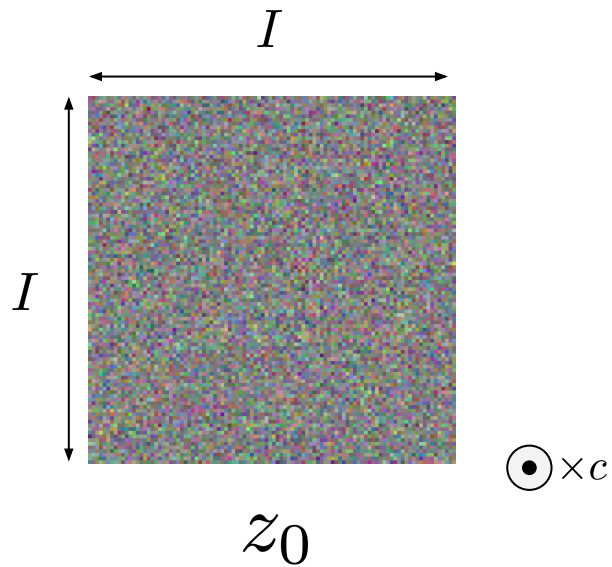
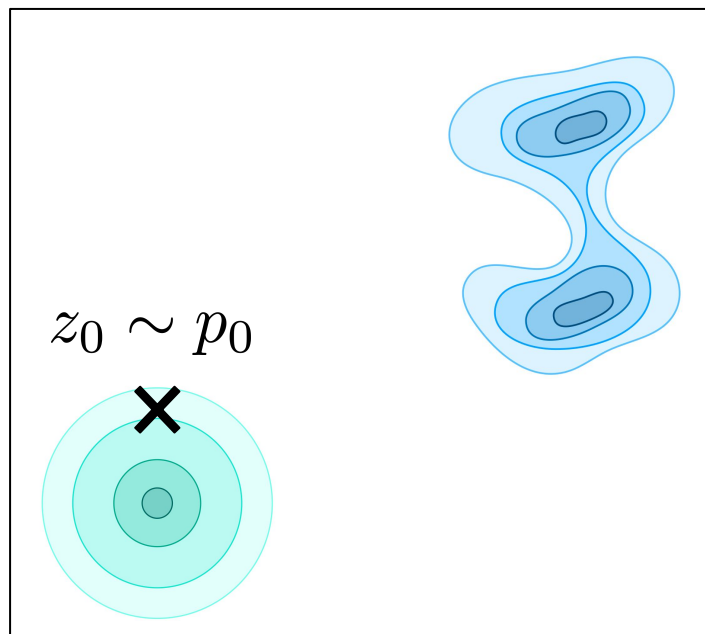
---

# Let's go through an end-to-end example together!

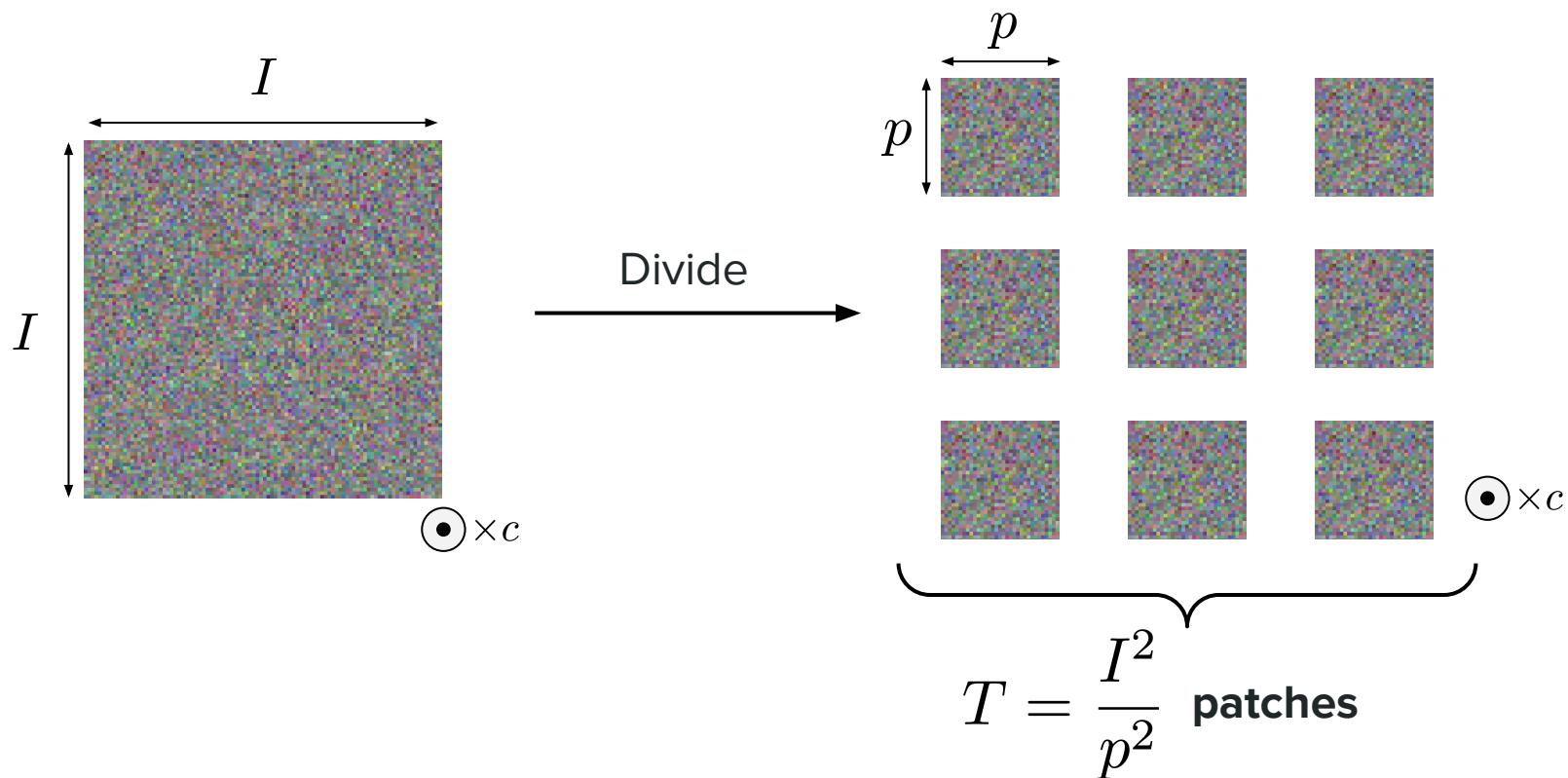


# Sample a noisy latent

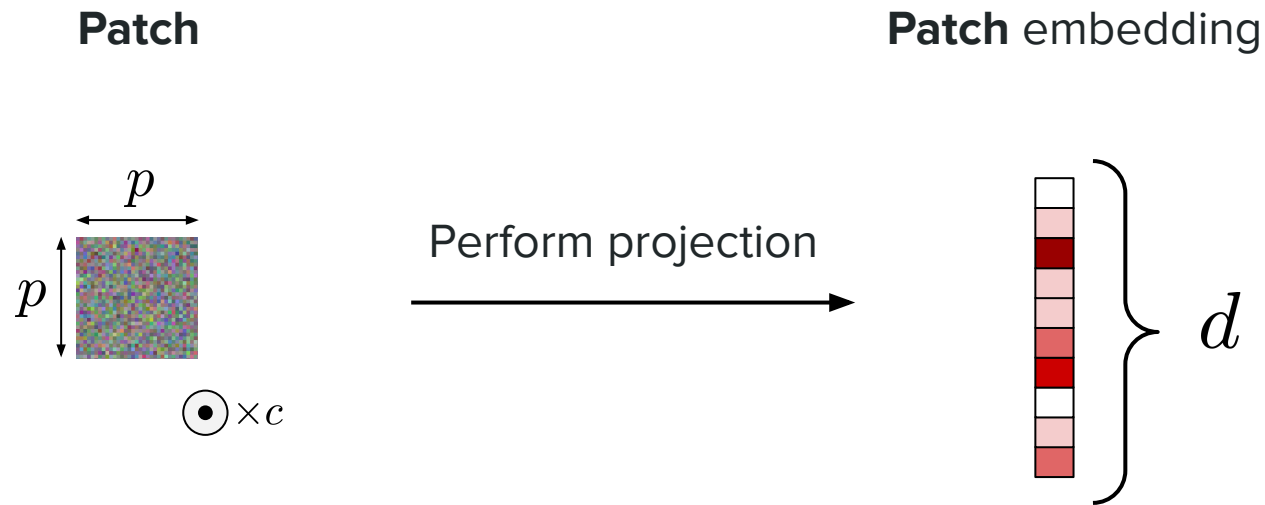
**Latent** space  
learned by the VAE



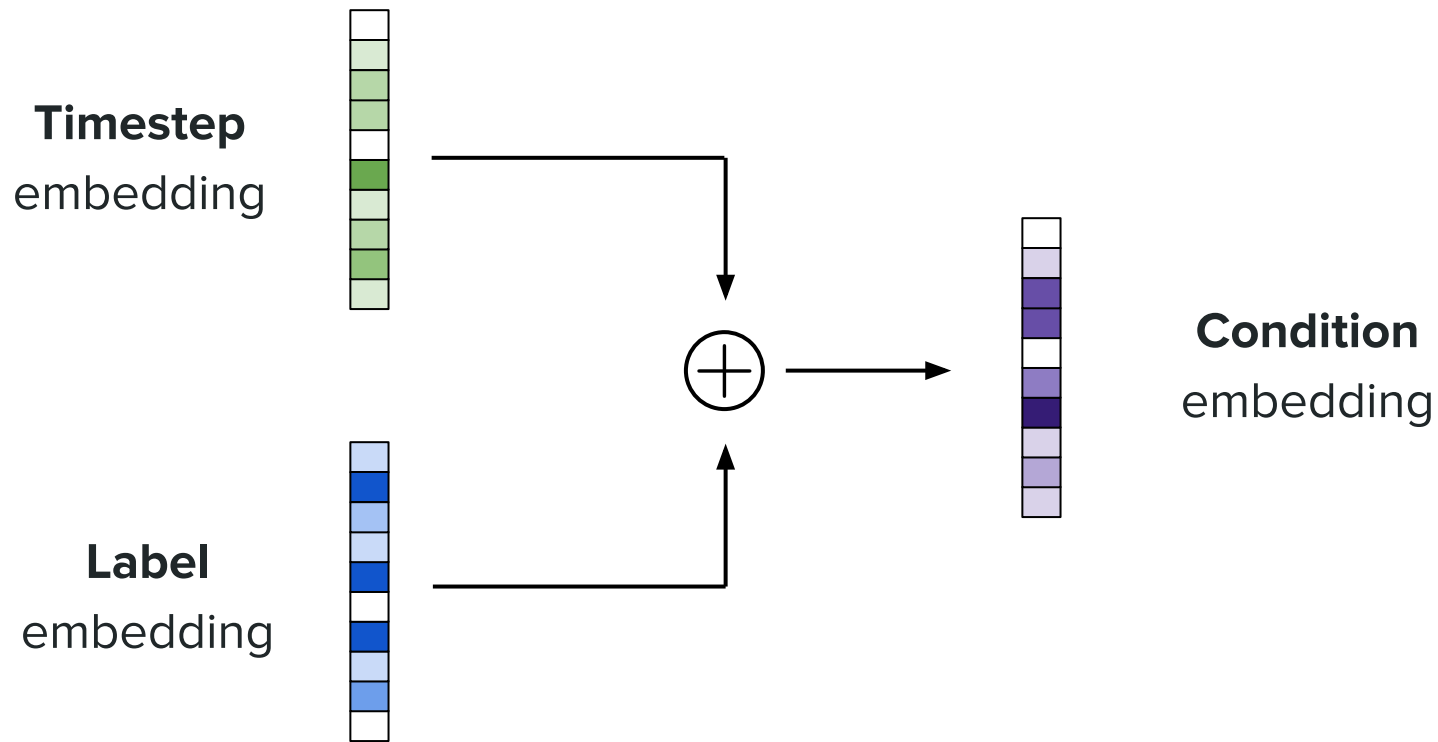
# Patchify noisy latent



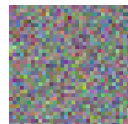
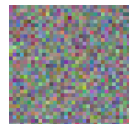
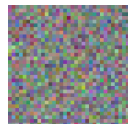
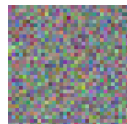
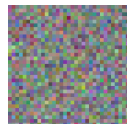
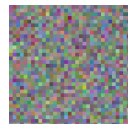
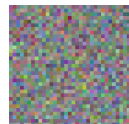
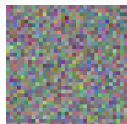
# Patchify noisy latent



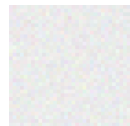
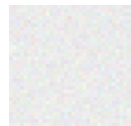
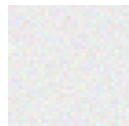
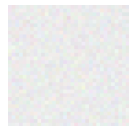
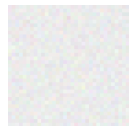
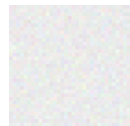
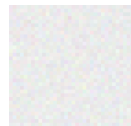
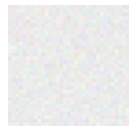
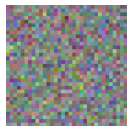
# Embed conditions



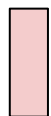
# Attention layers with noisy latent



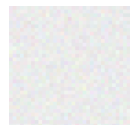
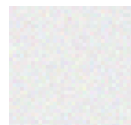
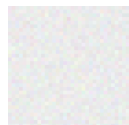
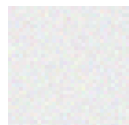
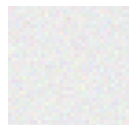
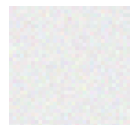
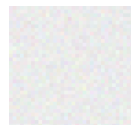
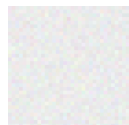
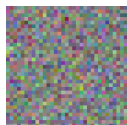
# Attention layers with noisy latent



# Attention layers with noisy latent



**Patch embedding**



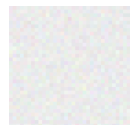
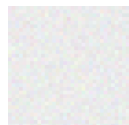
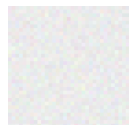
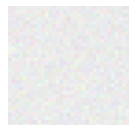
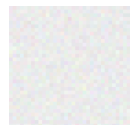
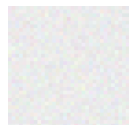
# Attention layers with noisy latent



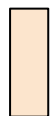
**Position embedding**



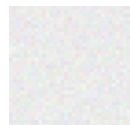
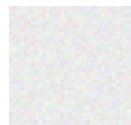
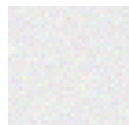
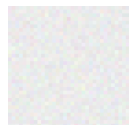
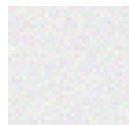
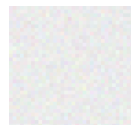
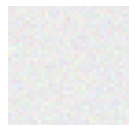
**Patch embedding**



# Attention layers with noisy latent

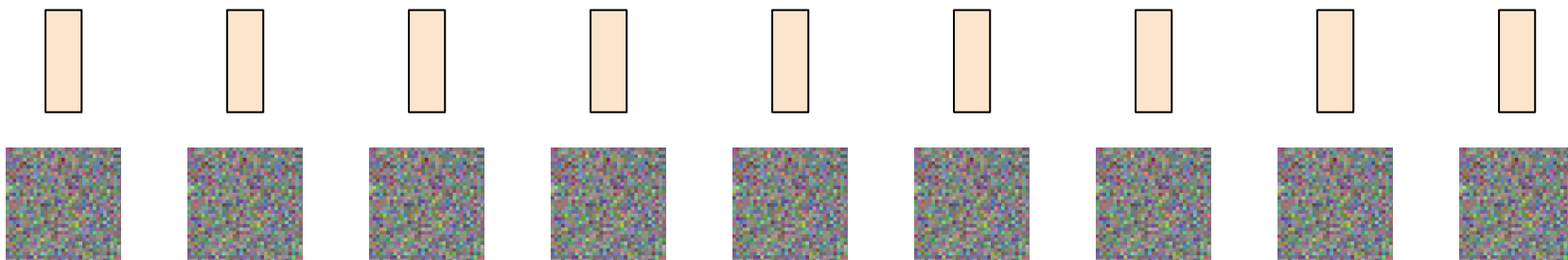


**Position-aware patch embedding**

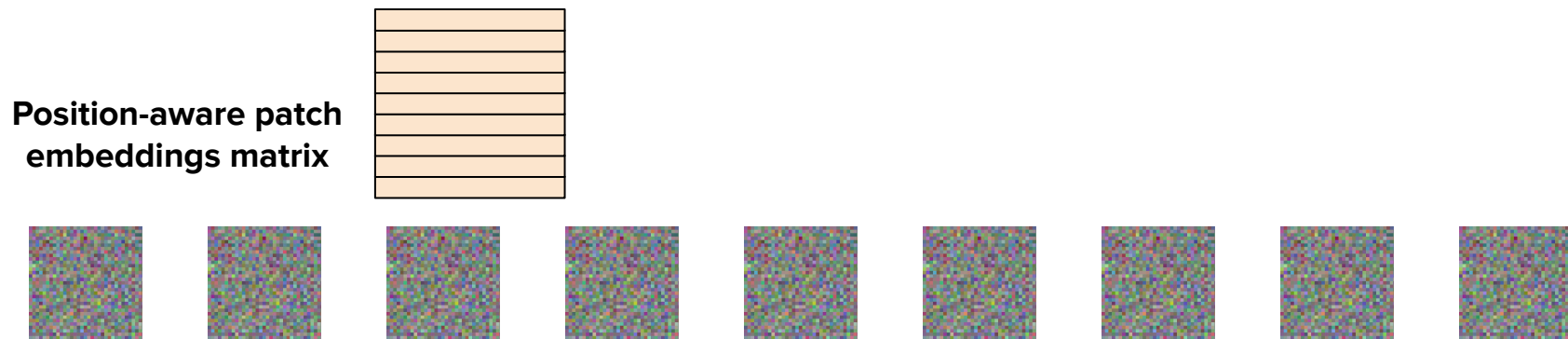


# Attention layers with noisy latent

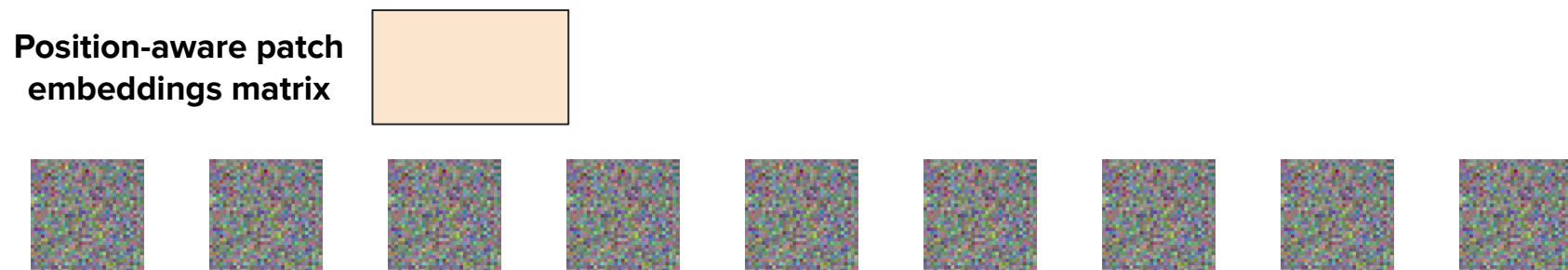
**Position-aware patch embeddings**



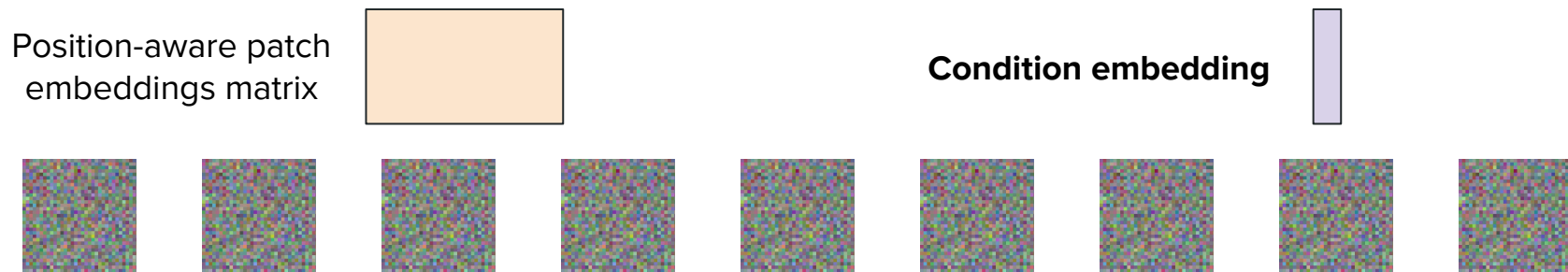
# Attention layers with noisy latent



# Attention layers with noisy latent



# Attention layers with noisy latent



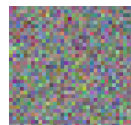
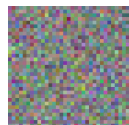
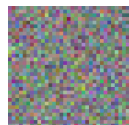
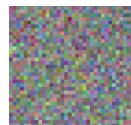
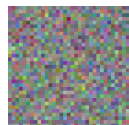
# Attention layers with noisy latent



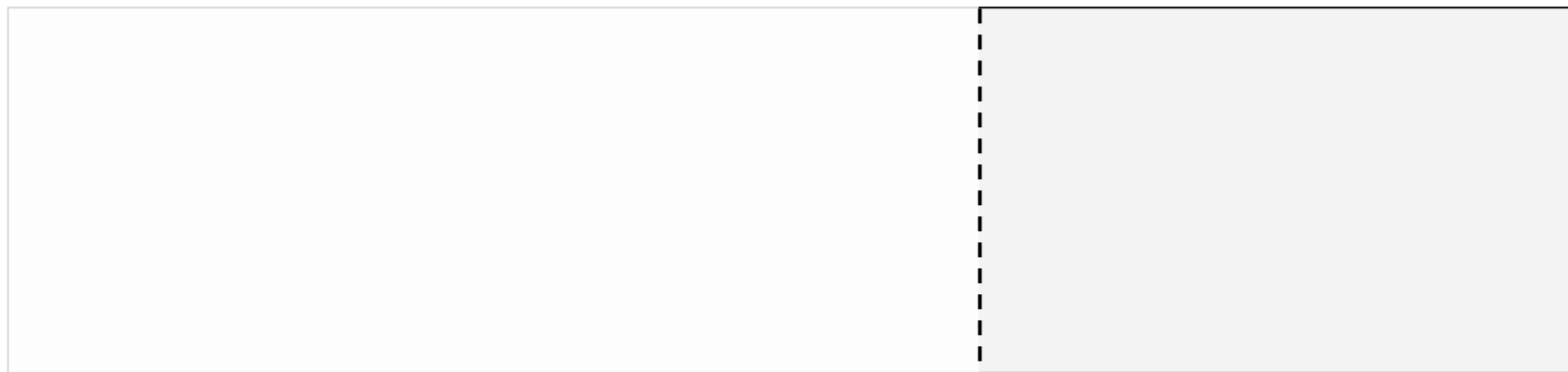
Position-aware patch embeddings matrix



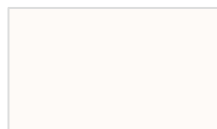
Condition embedding



# Process condition embeddings



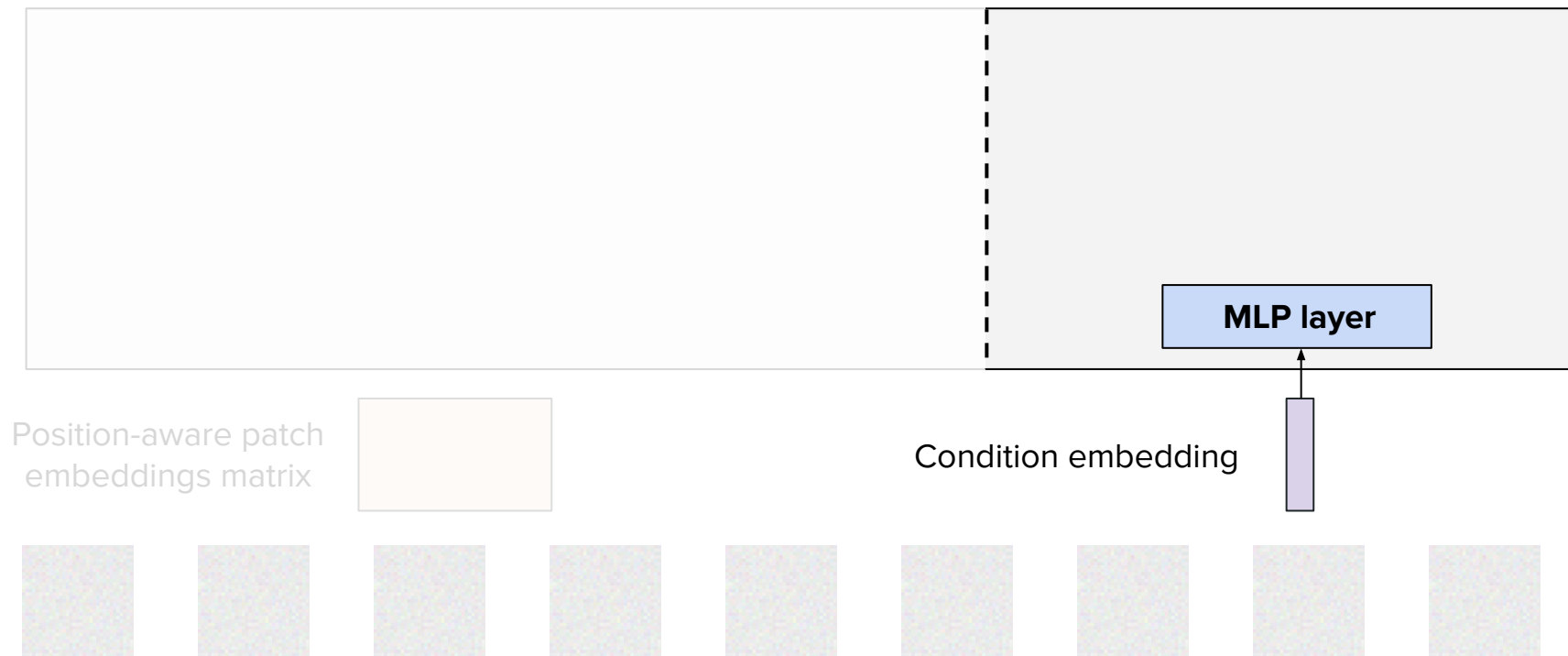
Position-aware patch embeddings matrix



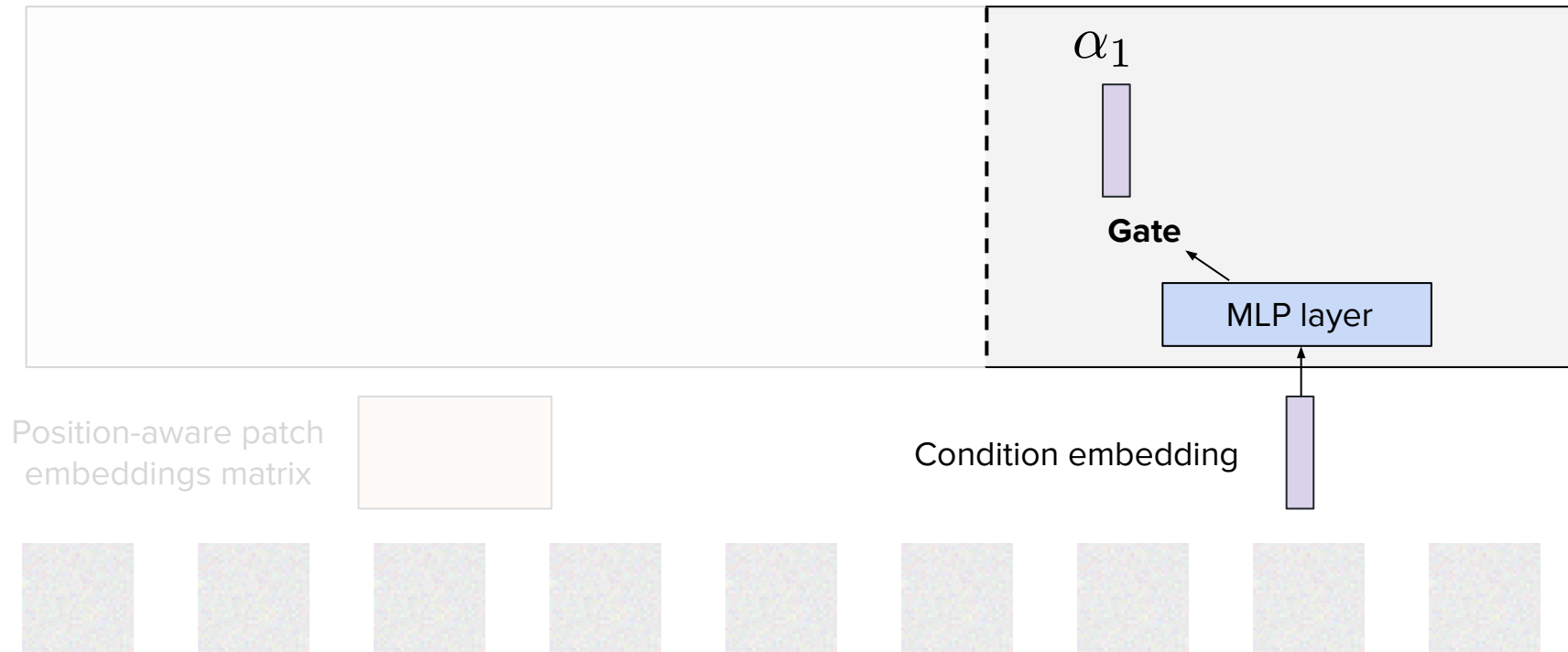
**Condition embedding**



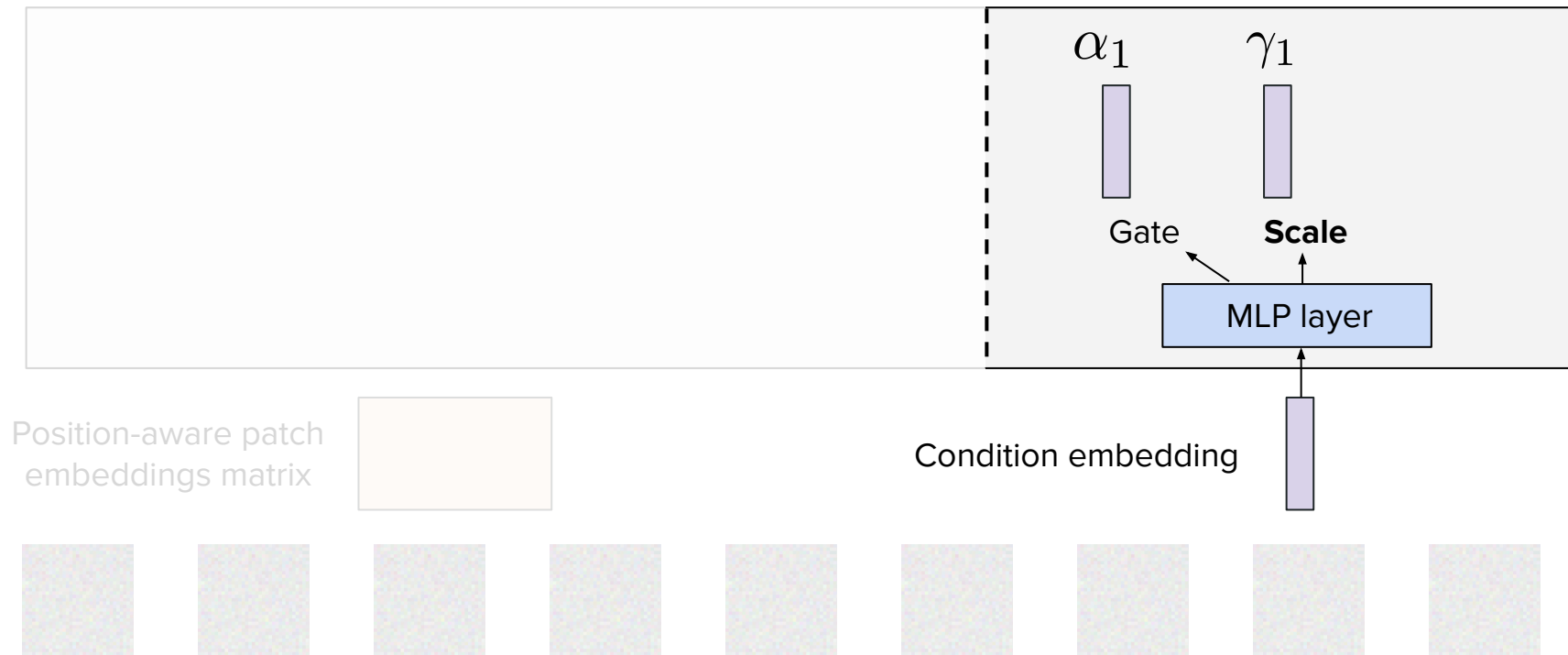
# Process condition embeddings



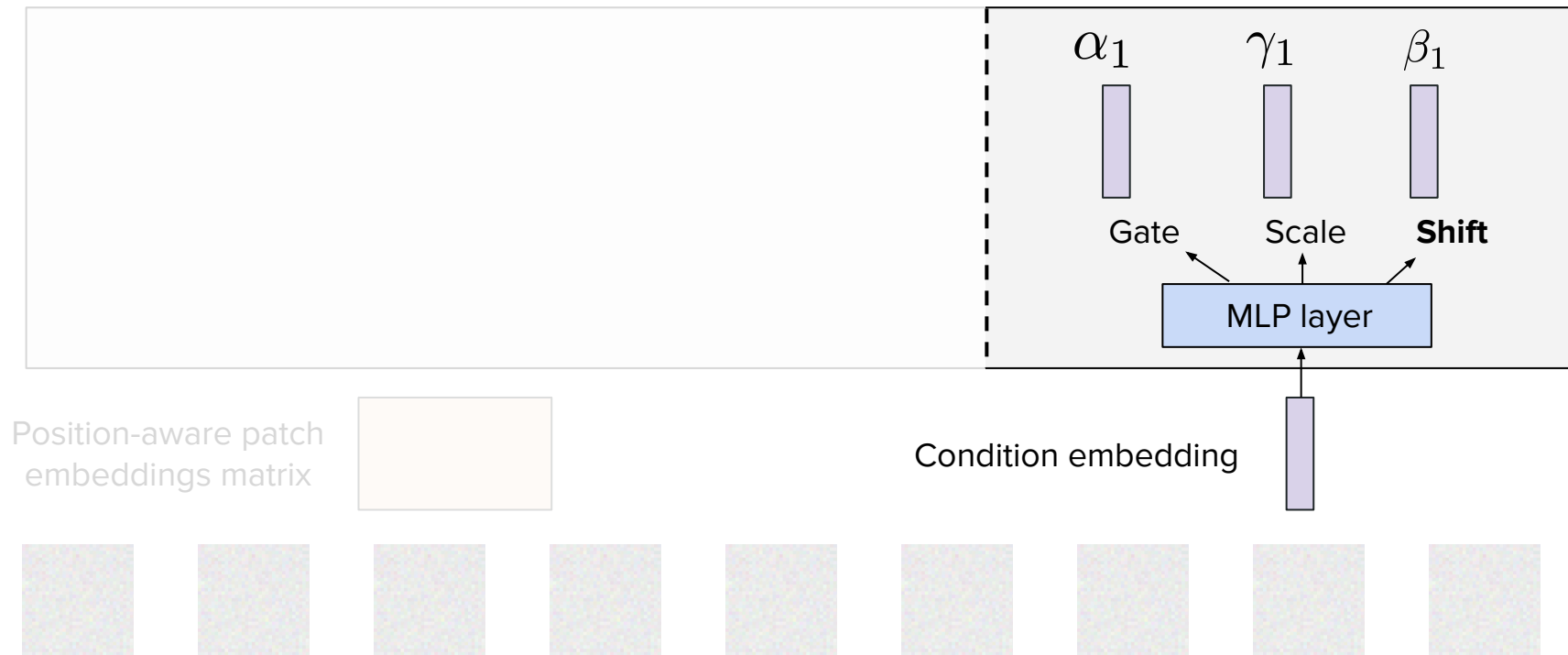
# Obtain gate from conditions



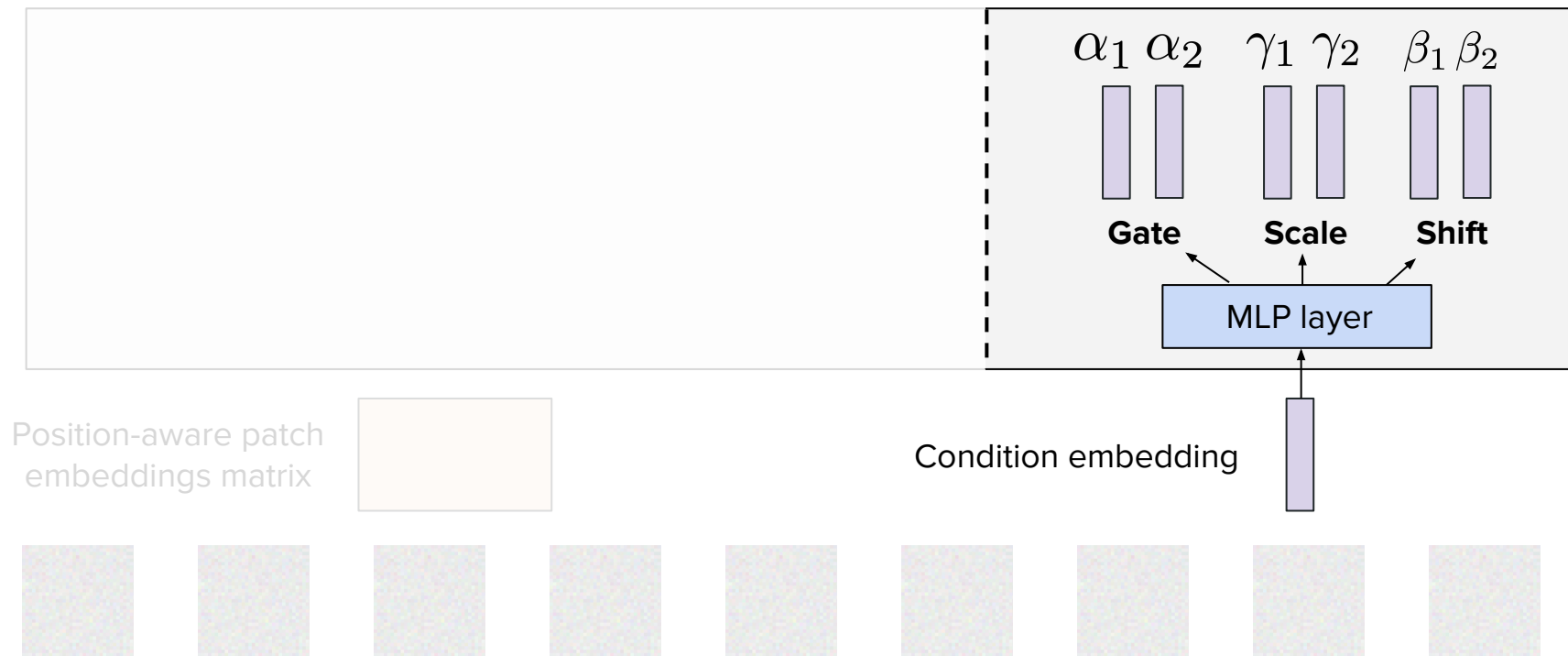
# Obtain scale from conditions



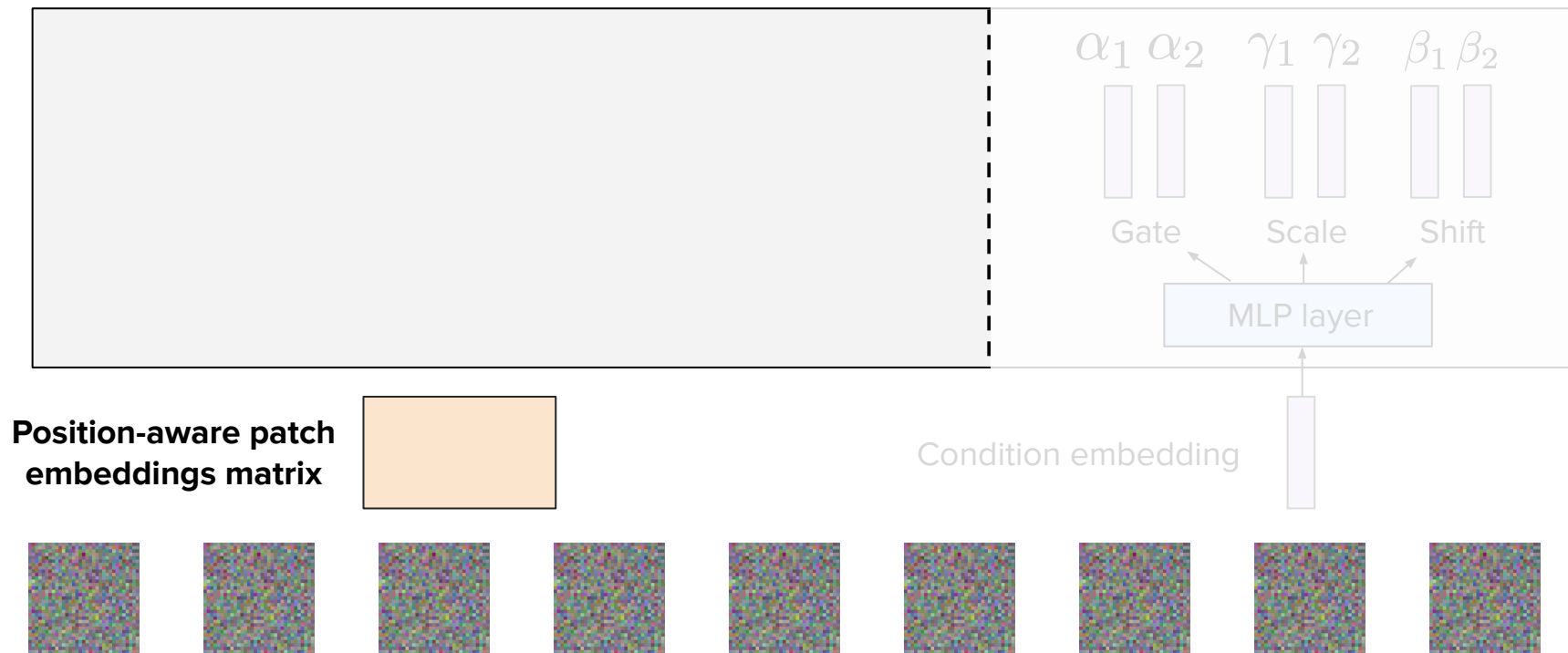
# Obtain shift from conditions



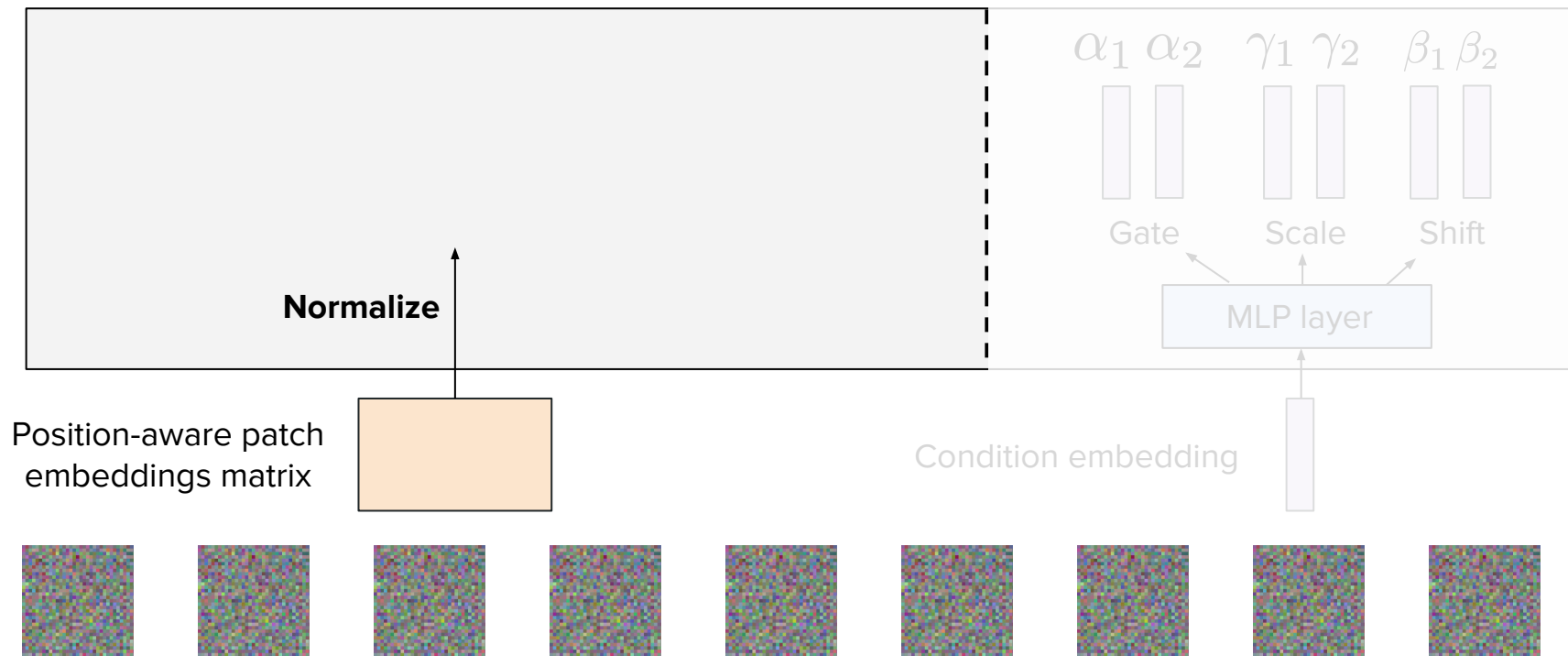
# Obtain quantities from condition



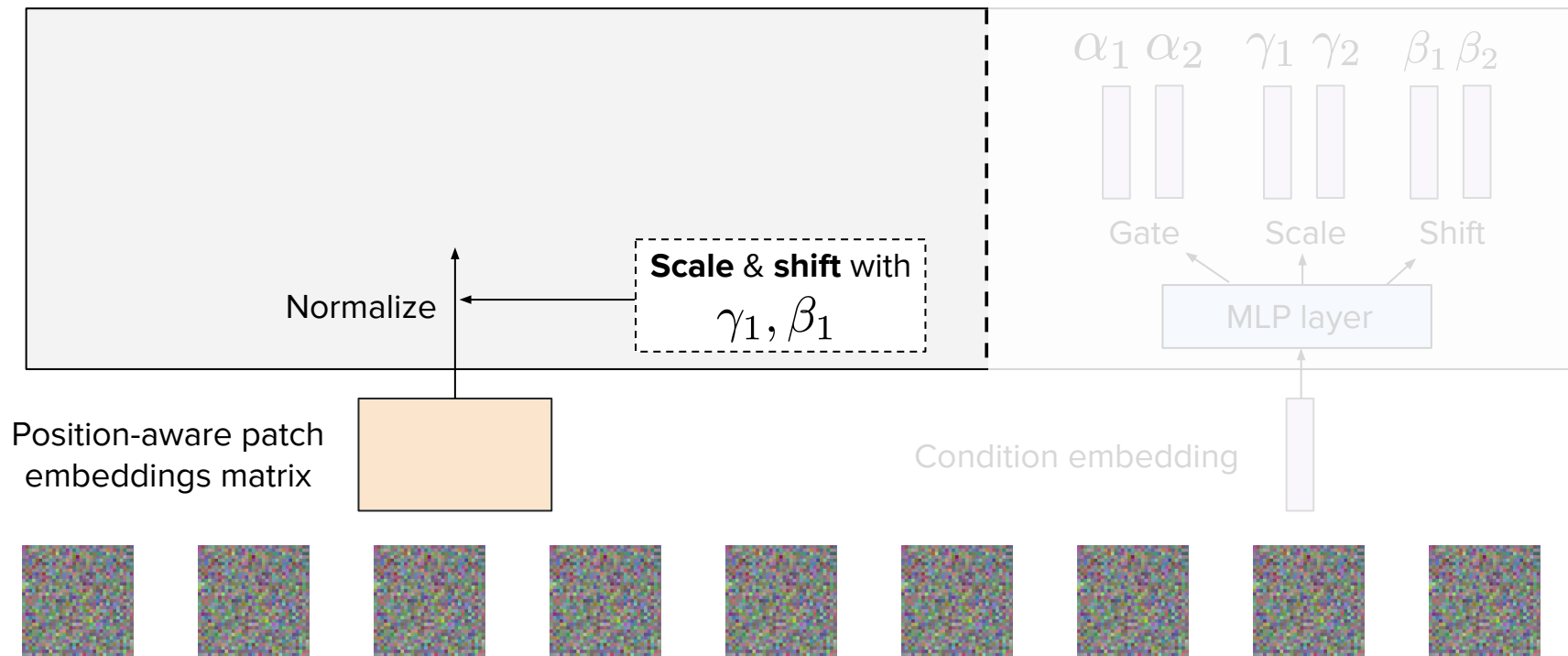
# Process patch embeddings



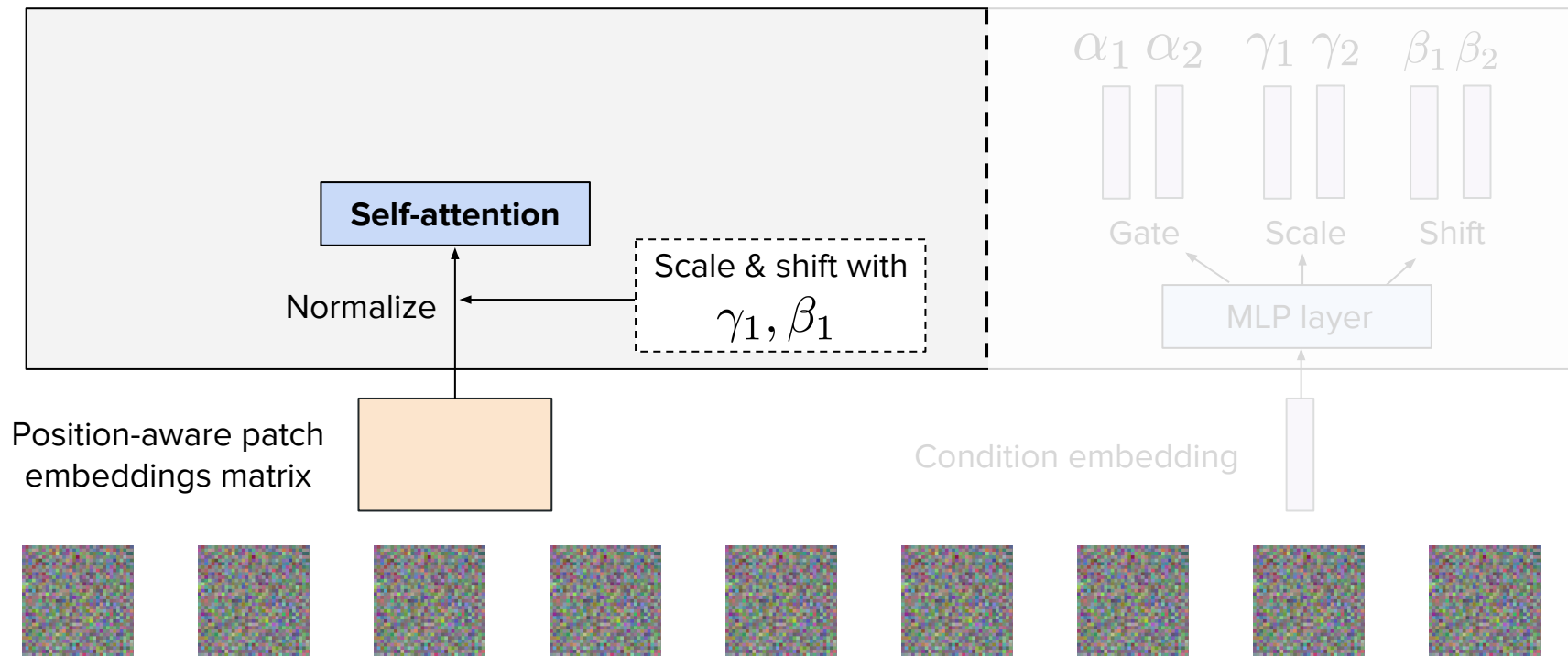
# Process patch embeddings through attention



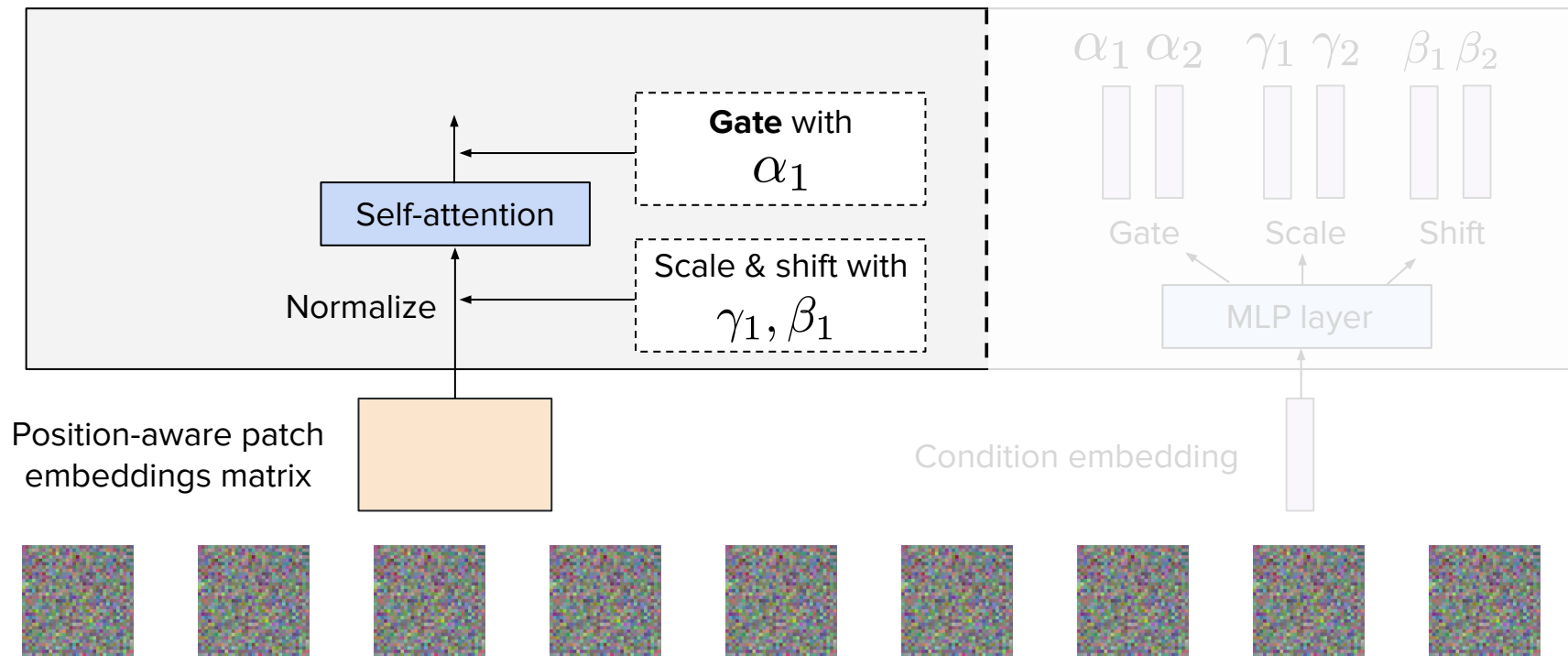
# Process patch embeddings through attention



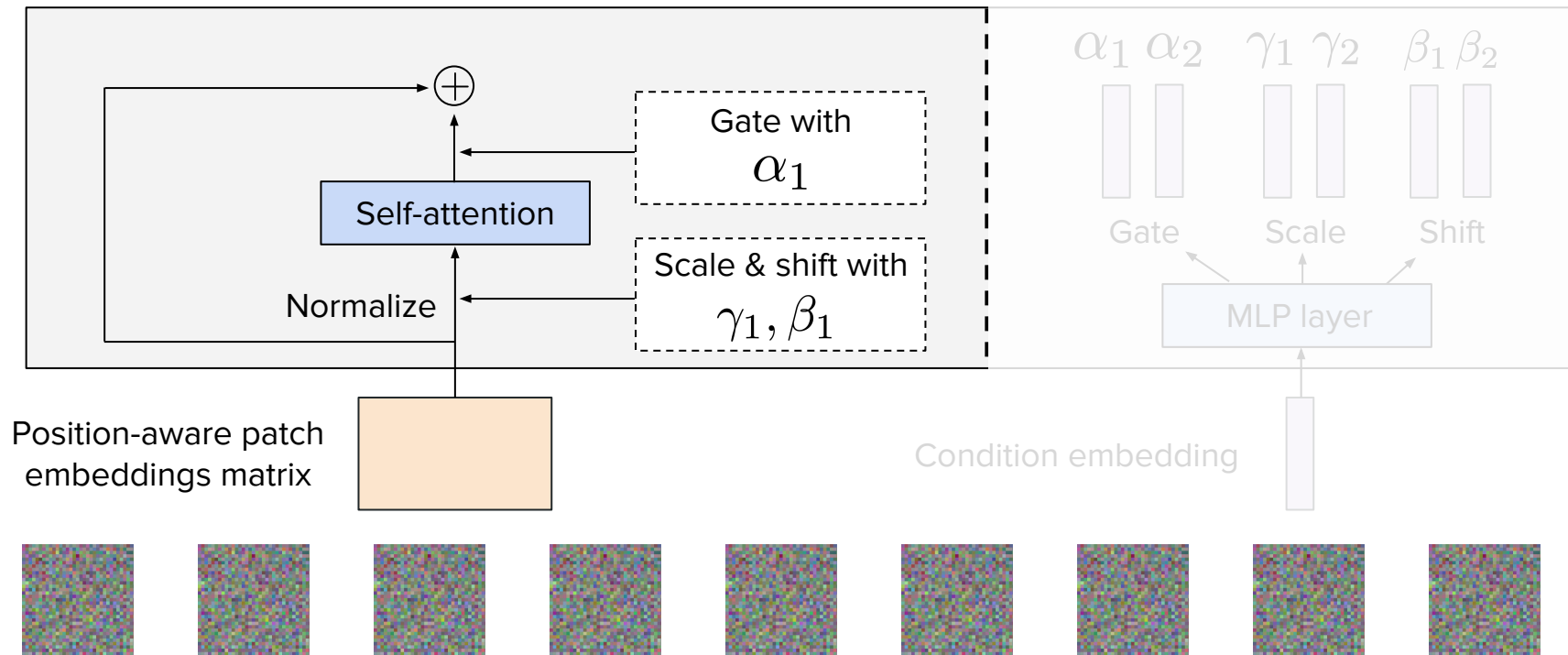
# Process patch embeddings through attention



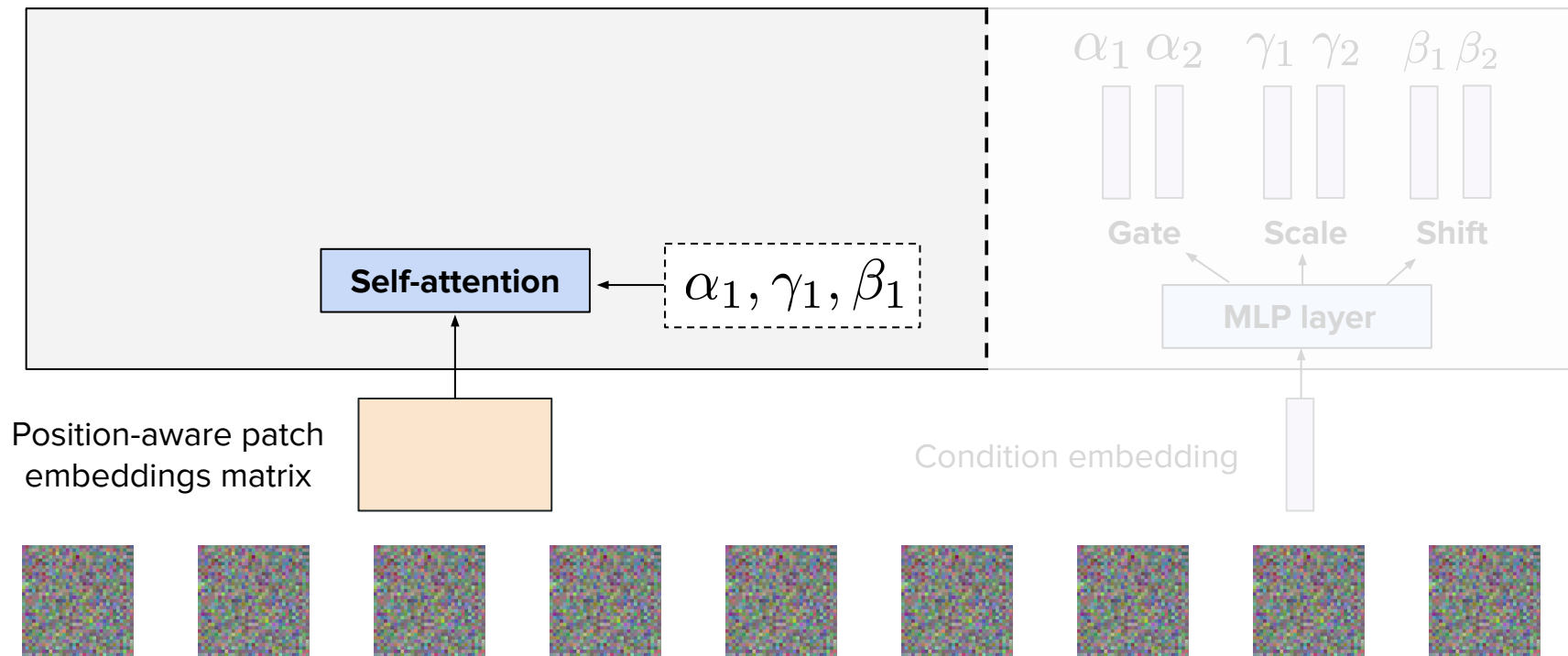
# Process patch embeddings through attention



# Process patch embeddings through attention



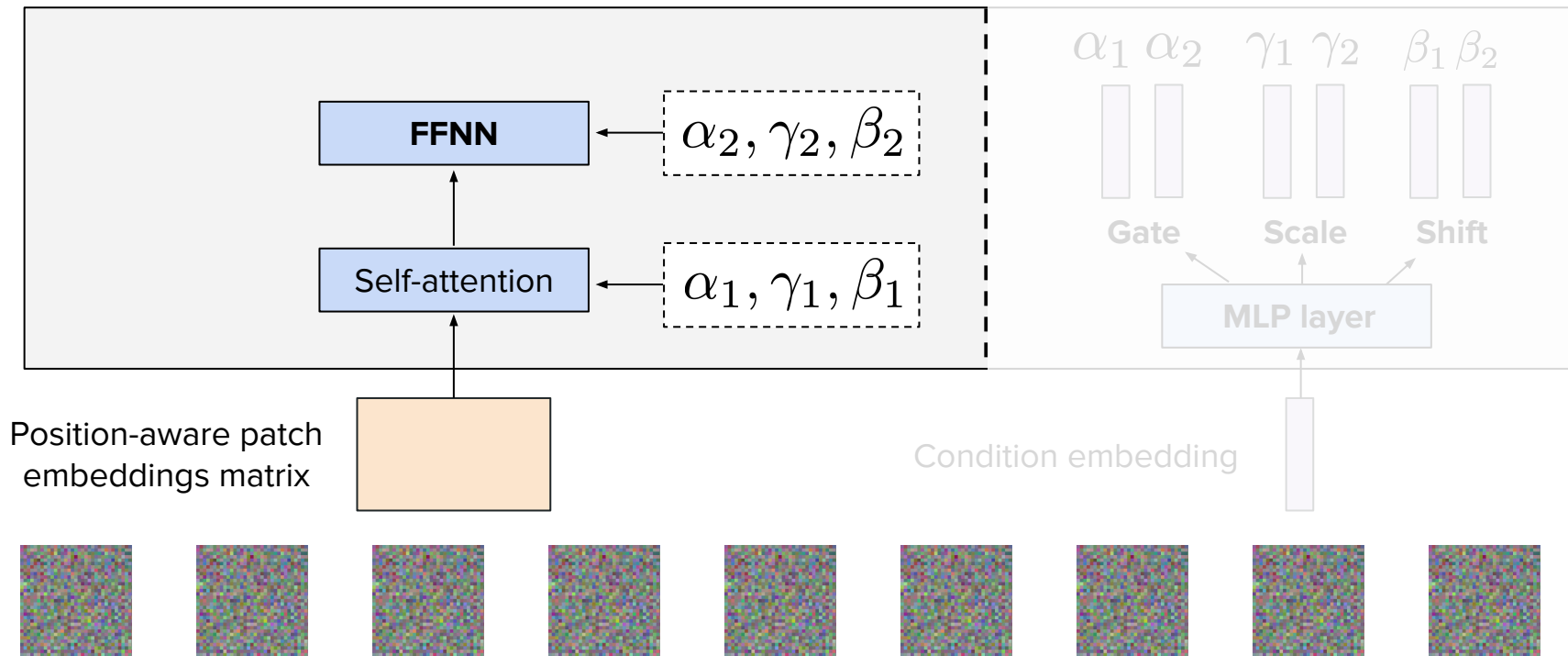
# Process patch embeddings through attention



# Process patch embeddings through attention

$$x \longleftarrow x + \alpha_1 \star \text{Attention}(\text{LN}(x) \star (1 + \gamma_1) + \beta_1)$$

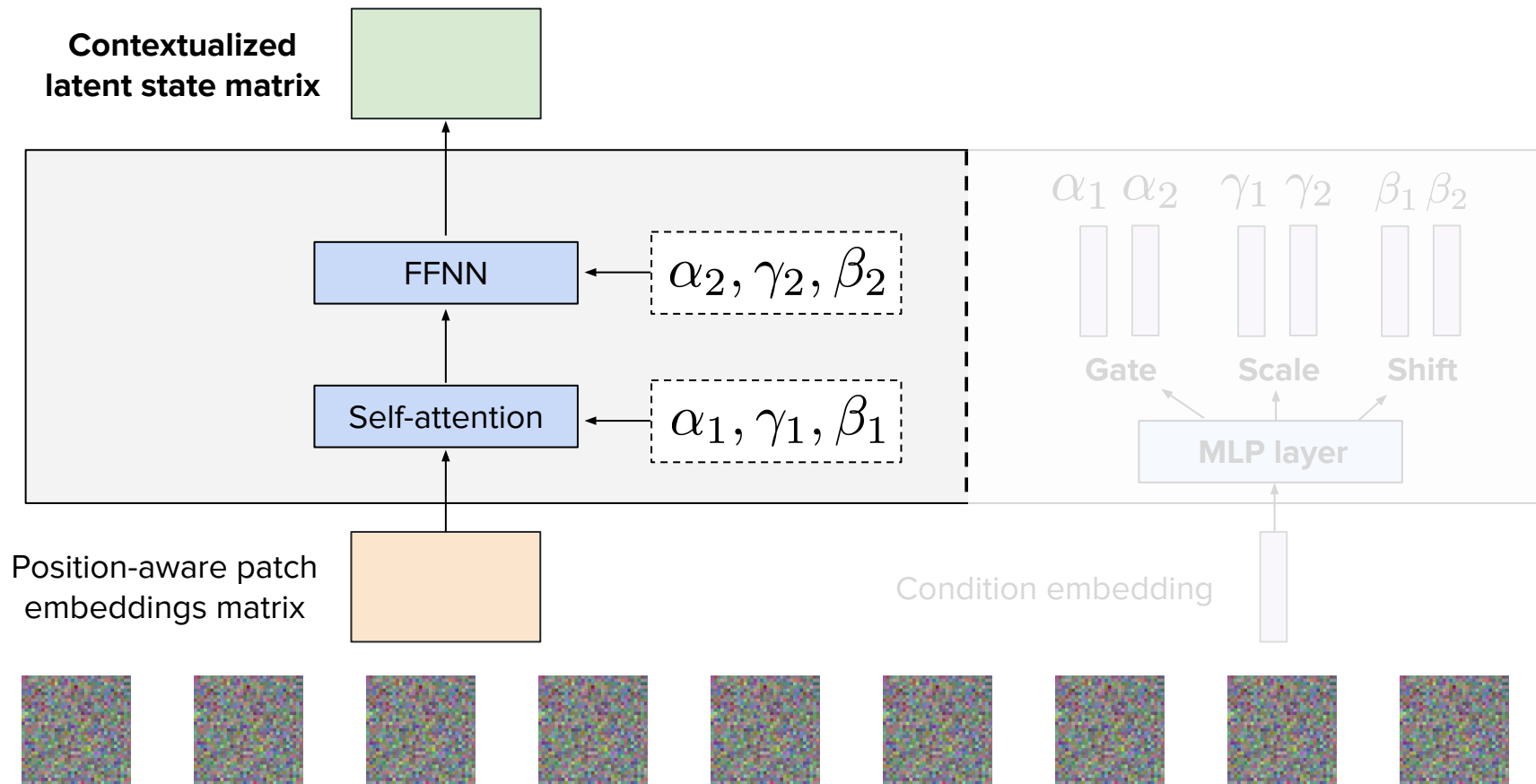
# Process patch embeddings through FFNN



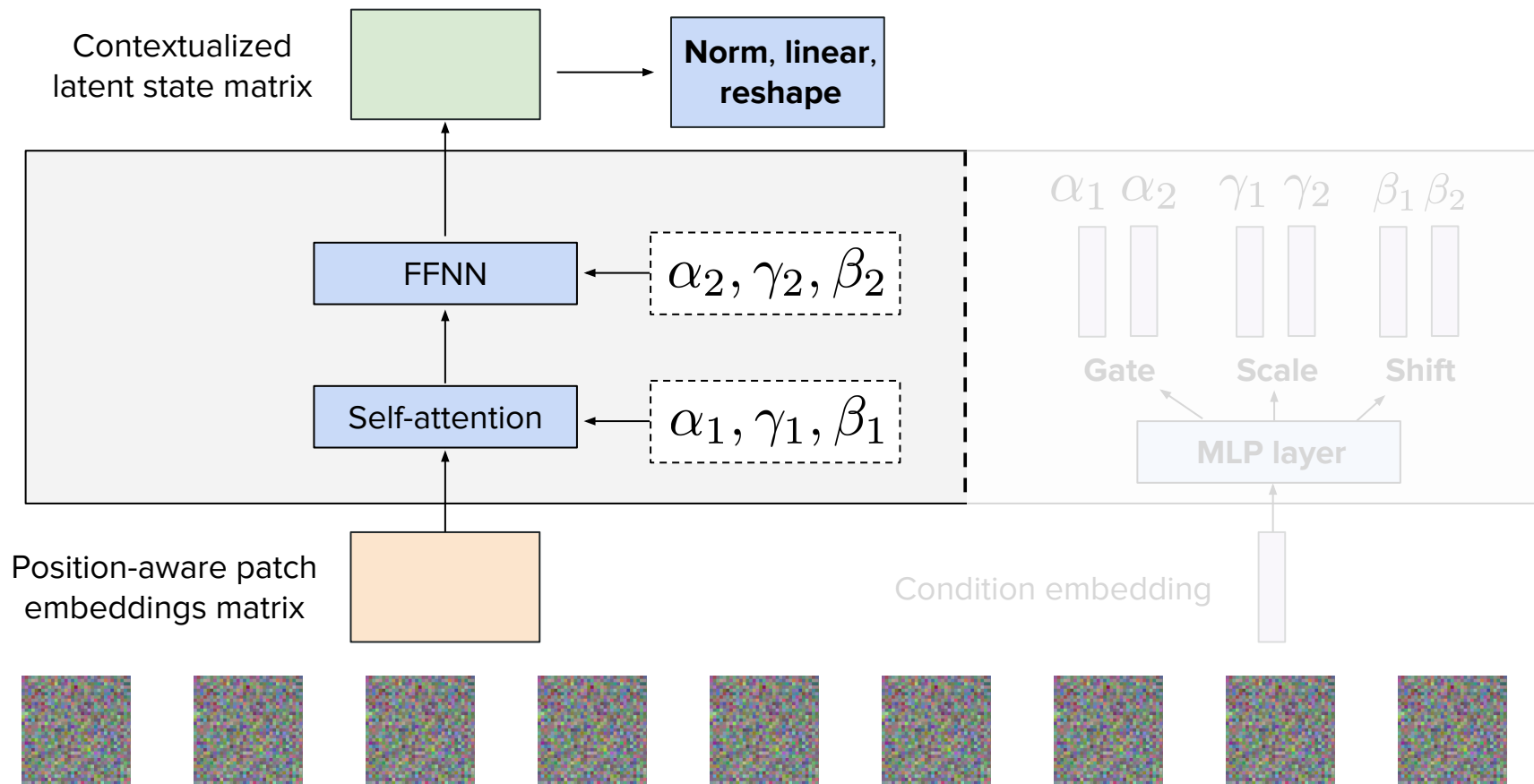
# Process patch embeddings through FFNN

$$x \longleftarrow x + \alpha_2 \star \text{FFNN}(\text{LN}(x) \star (1 + \gamma_2) + \beta_2)$$

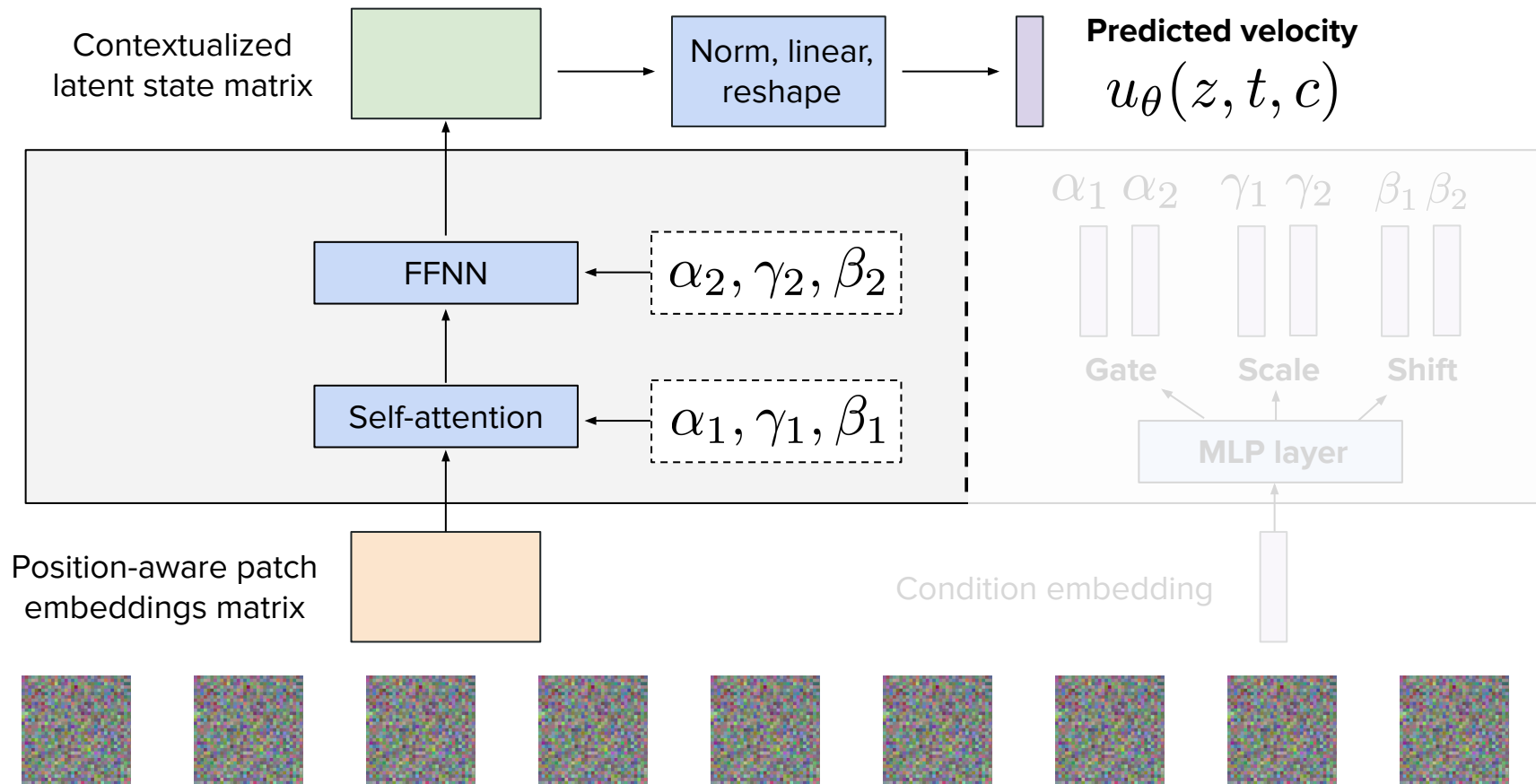
# Obtain output from DiT block



# Project and reshape to desired quantity



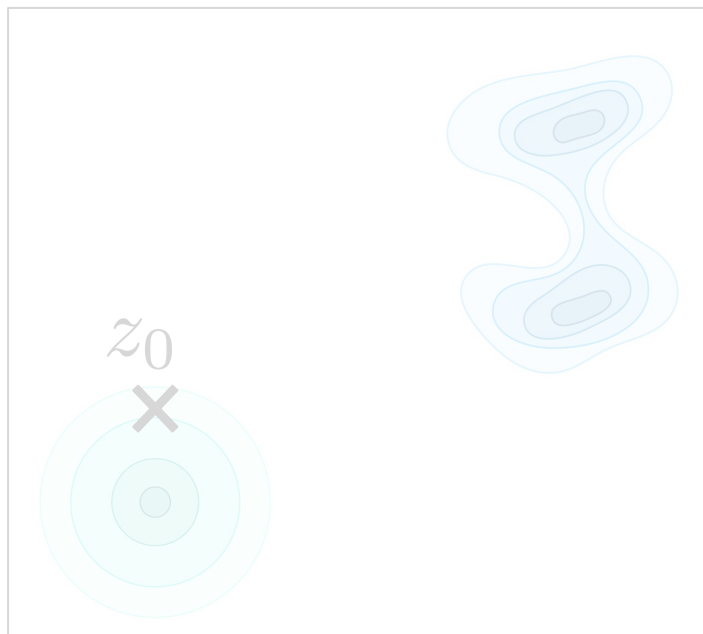
# Project and reshape to desired quantity





# Deduce resulting latent

Latent space  
learned by the VAE

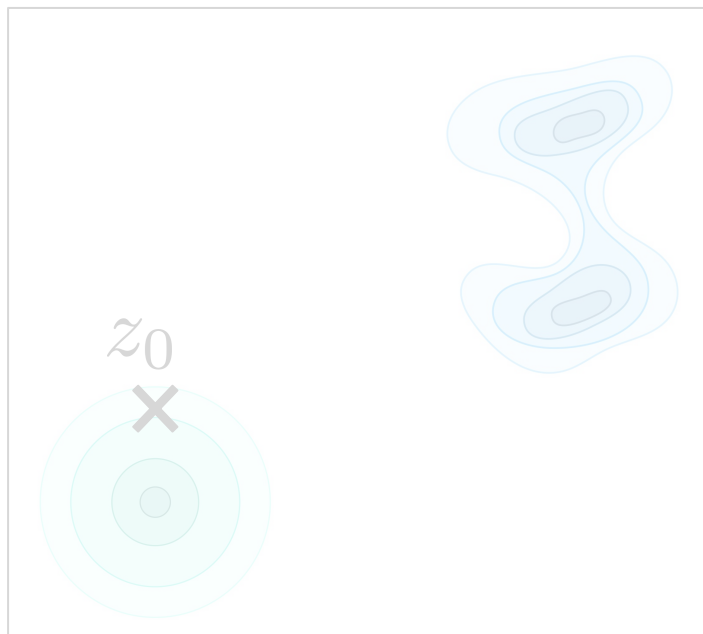


$$\underbrace{u_{\theta}(z_0, 0, c)}$$

Prediction from the DiT

# Deduce resulting latent

Latent space  
learned by the VAE

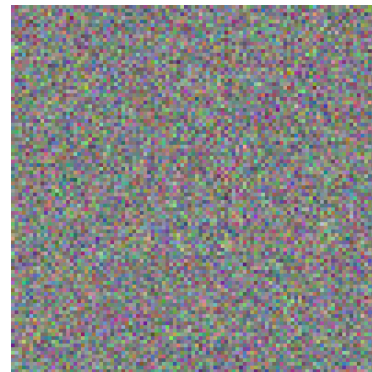
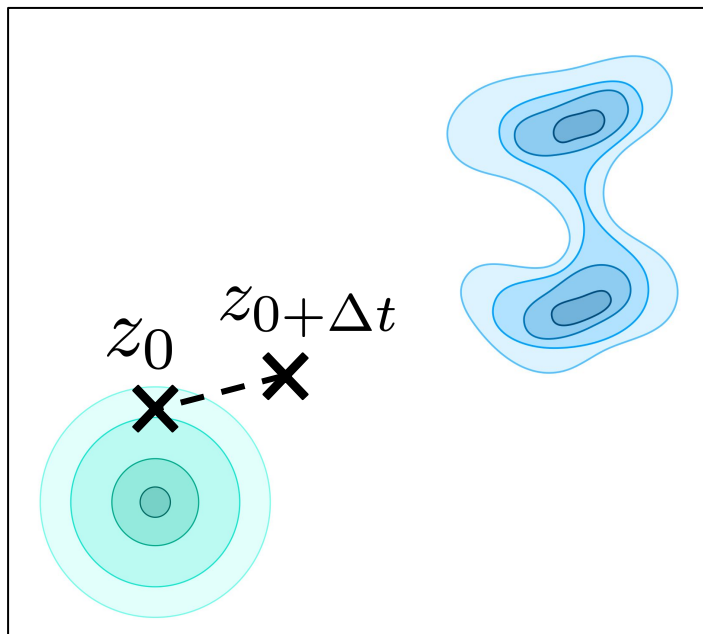


$$z_{0+\Delta t} = z_0 + \underbrace{u_{\theta}(z_0, 0, c)}_{\text{Prediction from the DiT}} \Delta t$$

Prediction from the DiT

# Deduce resulting latent

**Latent** space  
learned by the VAE



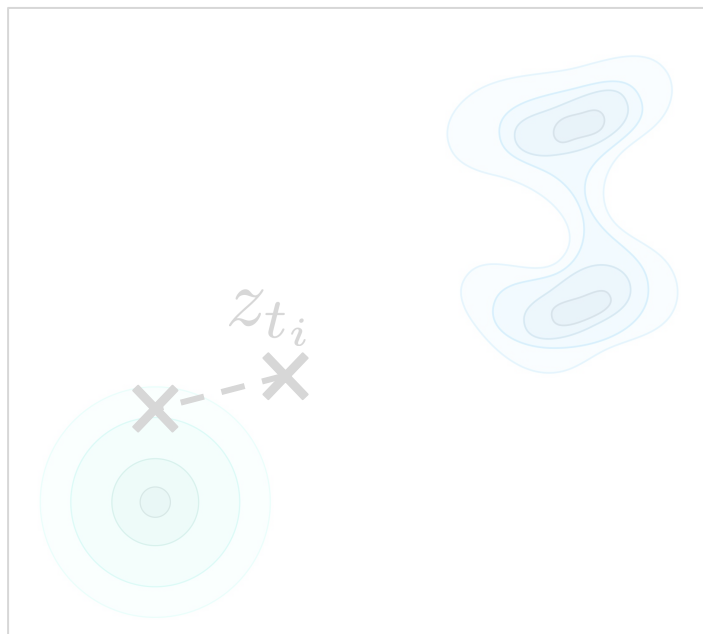
$$z_{0+\Delta t} = z_0 + \underbrace{u_{\theta}(z_0, 0, c)}_{\text{Prediction from the DiT}} \Delta t$$

Prediction from the DiT



# Perform iteration over latents

Latent space  
learned by the VAE

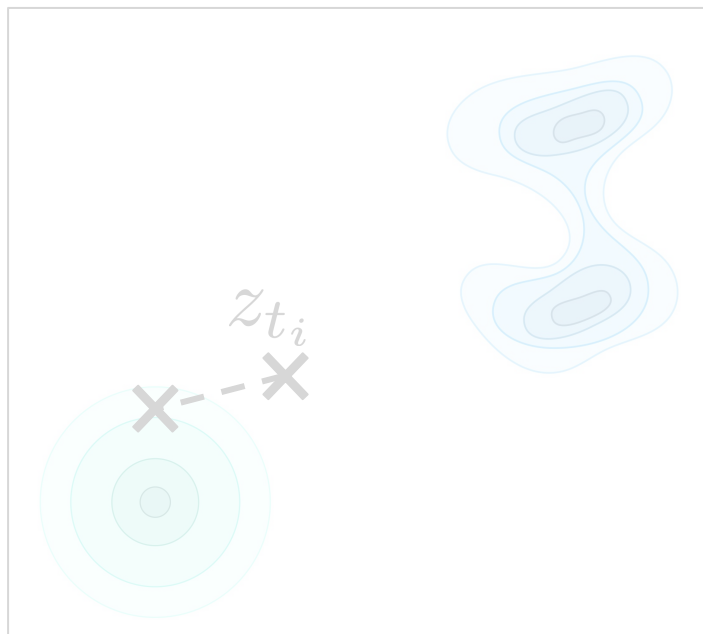


$$\underbrace{u_{\theta}(z_{t_i}, t_i, c)}$$

Prediction from the DiT

# Perform iteration over latents

Latent space  
learned by the VAE

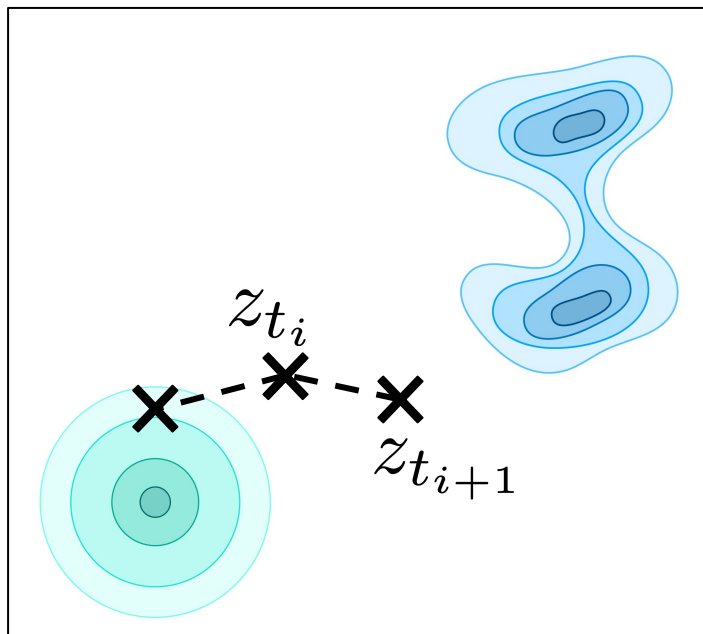


$$z_{t_{i+1}} = z_{t_i} + \underbrace{u_{\theta}(z_{t_i}, t_i, c)}_{\text{Prediction from the DiT}}(t_{i+1} - t_i)$$

Prediction from the DiT

# Perform iteration over latents

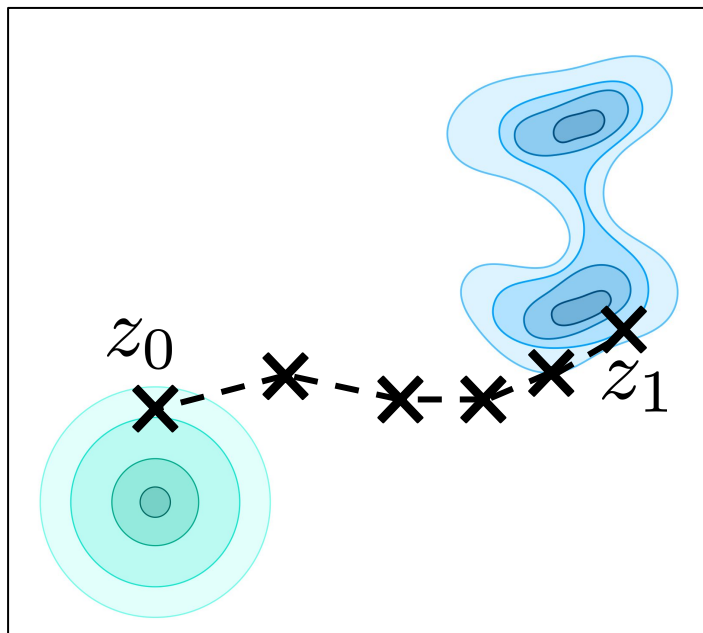
**Latent space**  
learned by the VAE



$$z_{t_{i+1}} = z_{t_i} + \underbrace{u_{\theta}(z_{t_i}, t_i, c)}_{\text{Prediction from the DiT}}(t_{i+1} - t_i)$$

# Obtain final latent

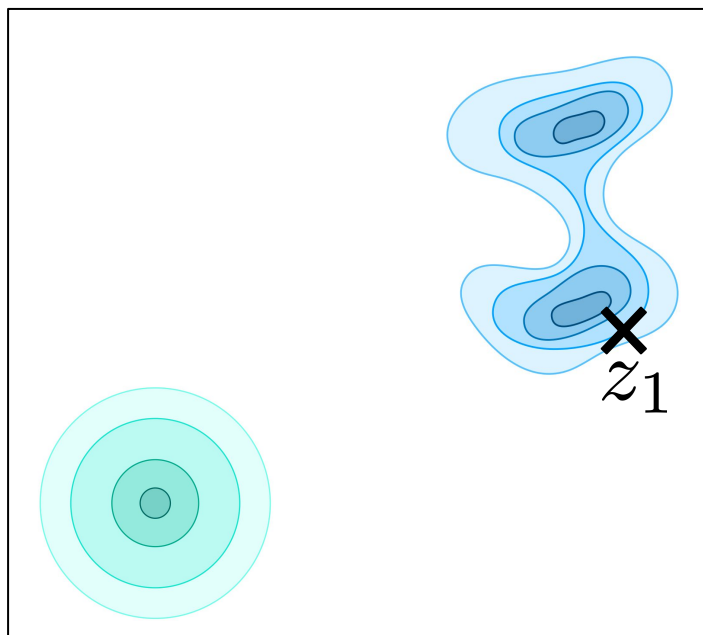
**Latent space**  
learned by the VAE



$z_1$

# Decode latent back into pixel space

**Latent space**  
learned by the VAE



$z_1$



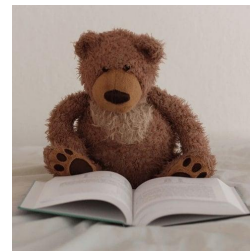
**Decoder**

$p_\theta$



**Pixel space**

$x_1 \sim p_\theta(x_1 | z_1)$





# Diffusion & Large Vision Models

Motivation

U-Net

Diffusion Transformer

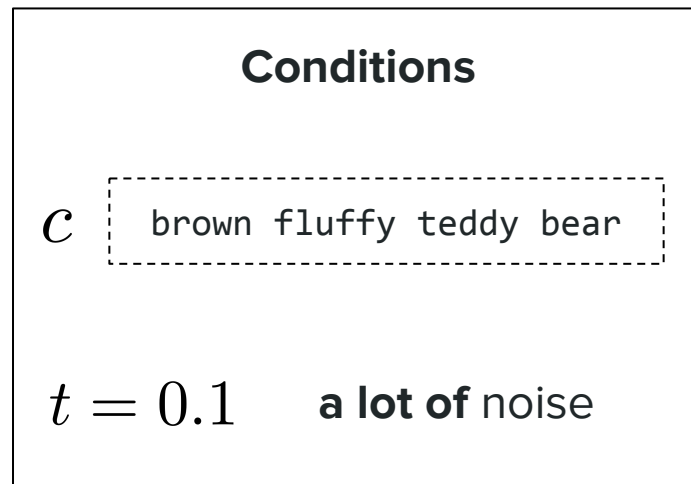
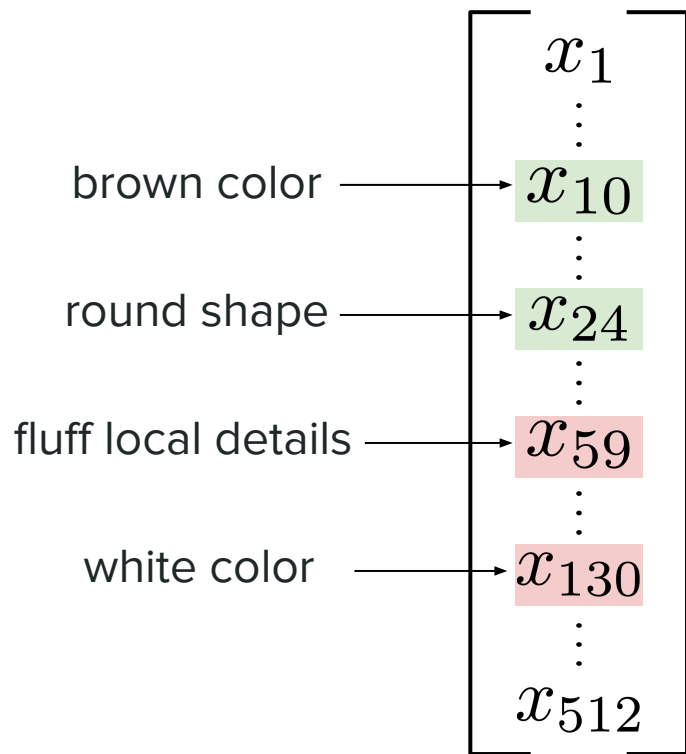
End-to-end example

**Multimodal DiT**

Optimizations

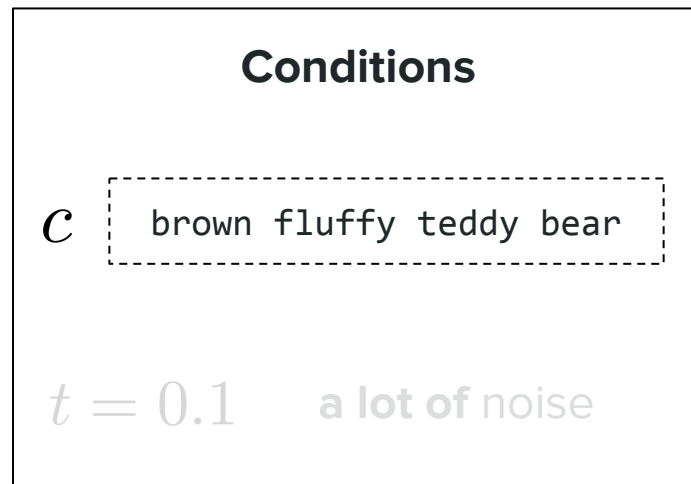
---

# Intuition: back to the example we had



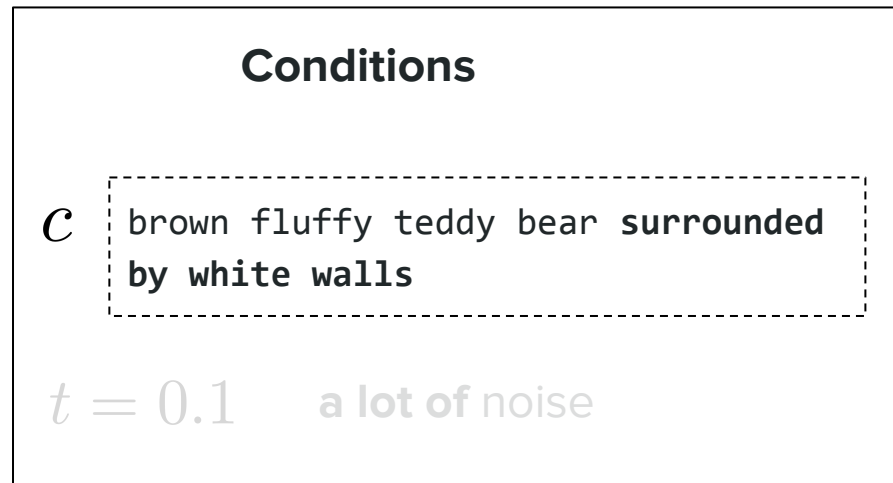
Need to focus on **global structure** since at the early stage of generation.

# Intuition: back to the example we had



Need to focus on **global structure** since at the early stage of generation.

# Intuition: more subtle text prompt



Need to focus on **global structure** since at the early stage of generation.

# Intuition: more subtle text prompt



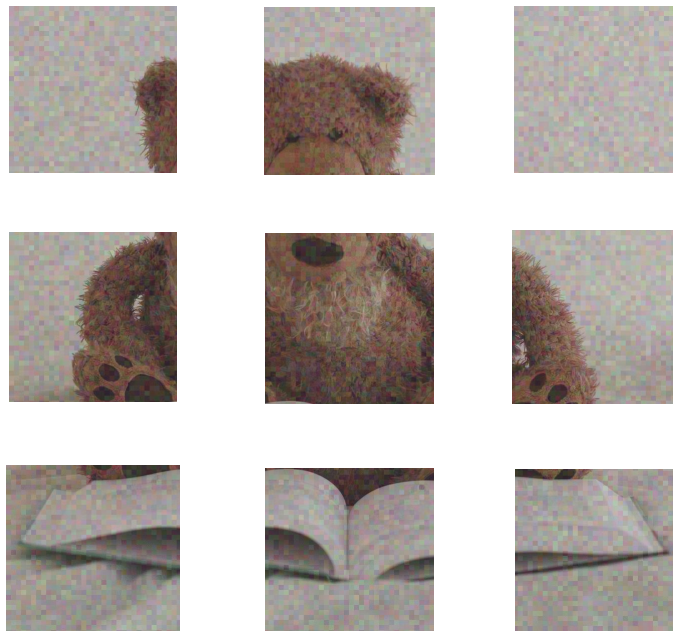
**Limitation.** All patch latents are subject to same "modulation"



# Intuition: more subtle text prompt



**Limitation.** All patch latents are subject to same "modulation"



# Mitigate lack of nuance of previous approach

## Strategy.

- 1 Keep "modulating" using the **timestep** embedding

# Mitigate lack of nuance of previous approach

## Strategy.

- 1 Keep "modulating" using the **timestep** embedding
- 2 **Do something else** to allow for more complex interaction for the **condition**

Cross-attention

Joint attention

# Mitigate lack of nuance of previous approach

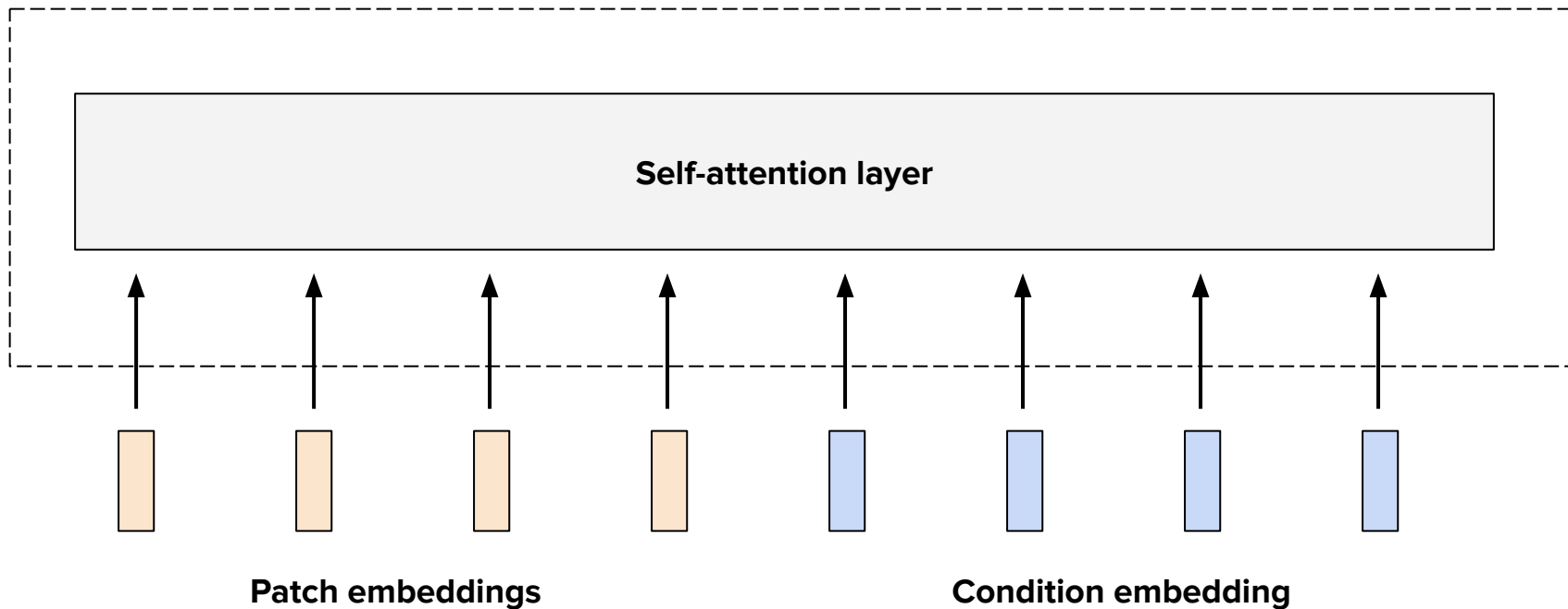
## Strategy.

- 1 Keep "modulating" using the **timestep** embedding
- 2 **Do something else** to allow for more complex interaction for the **condition**

Cross-attention

**Joint attention**

# Intuition behind joint attention



# MultiModal-Diffusion Transformer

## MM-DiT = MultiModal-Diffusion Transformer

- Term coined in the **Stable Diffusion 3** paper published in 2024
- Relies on **joint attention** of different modalities

**Single-stream**

Everyone treated equally

# MultiModal-Diffusion Transformer

## MM-DiT = MultiModal-Diffusion Transformer

- Term coined in the **Stable Diffusion 3** paper published in 2024
- Relies on **joint attention** of different modalities

### Single-stream

Everyone treated equally

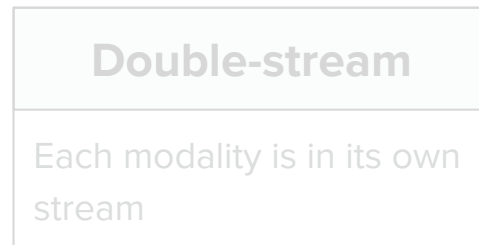
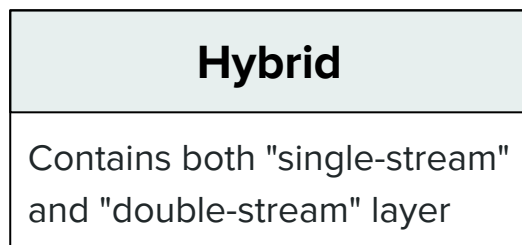
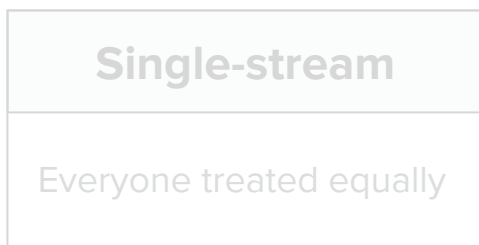
### Double-stream

Each modality is in its own stream

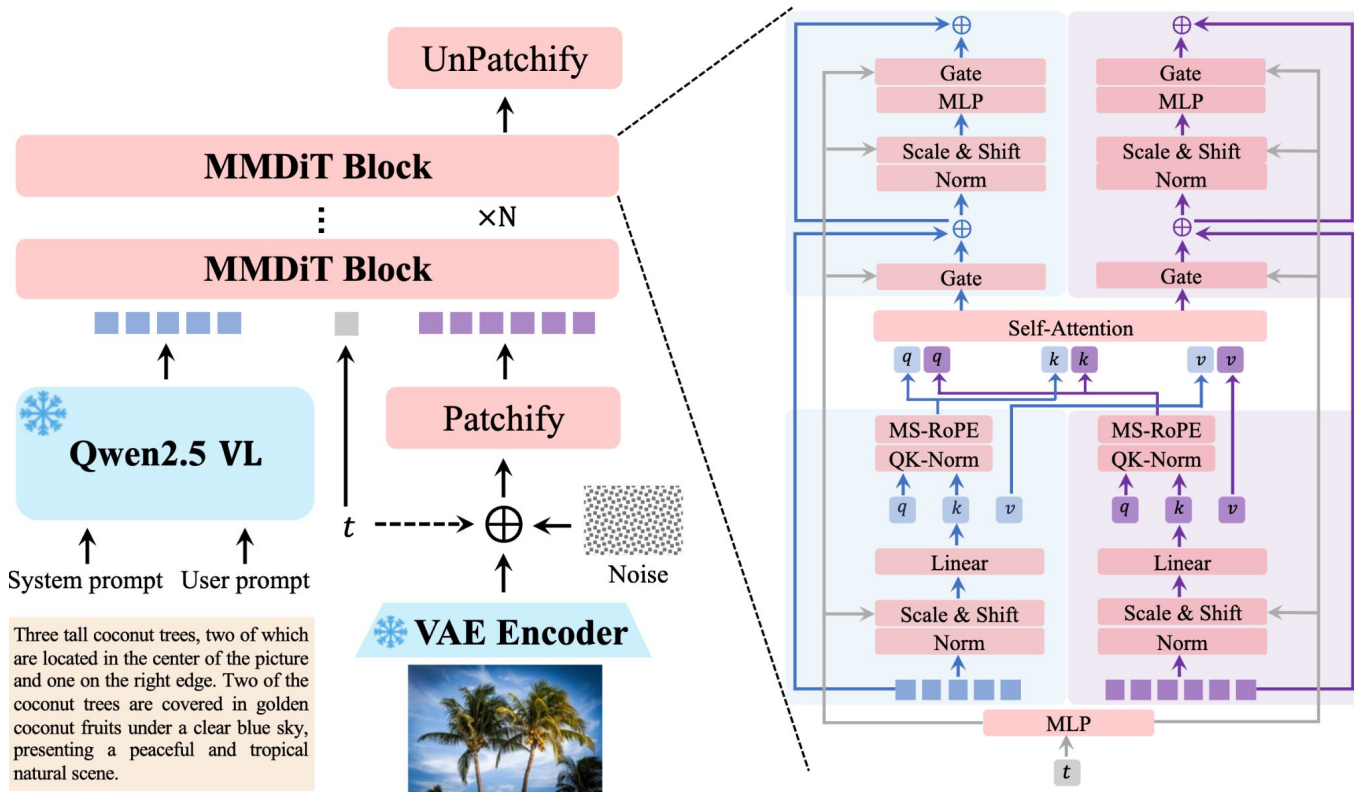
# MultiModal-Diffusion Transformer

## MM-DiT = MultiModal-Diffusion Transformer

- Term coined in the **Stable Diffusion 3** paper published in 2024
- Relies on **joint attention** of different modalities



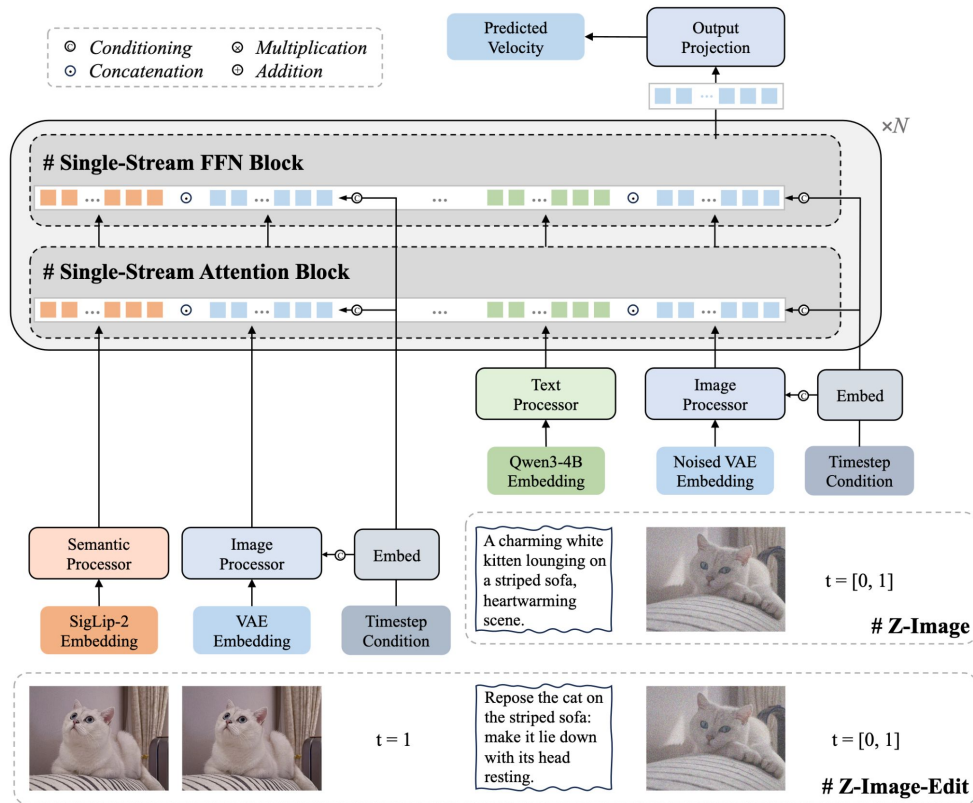
# "Double-stream" DiT variant: SD3, Qwen-Image



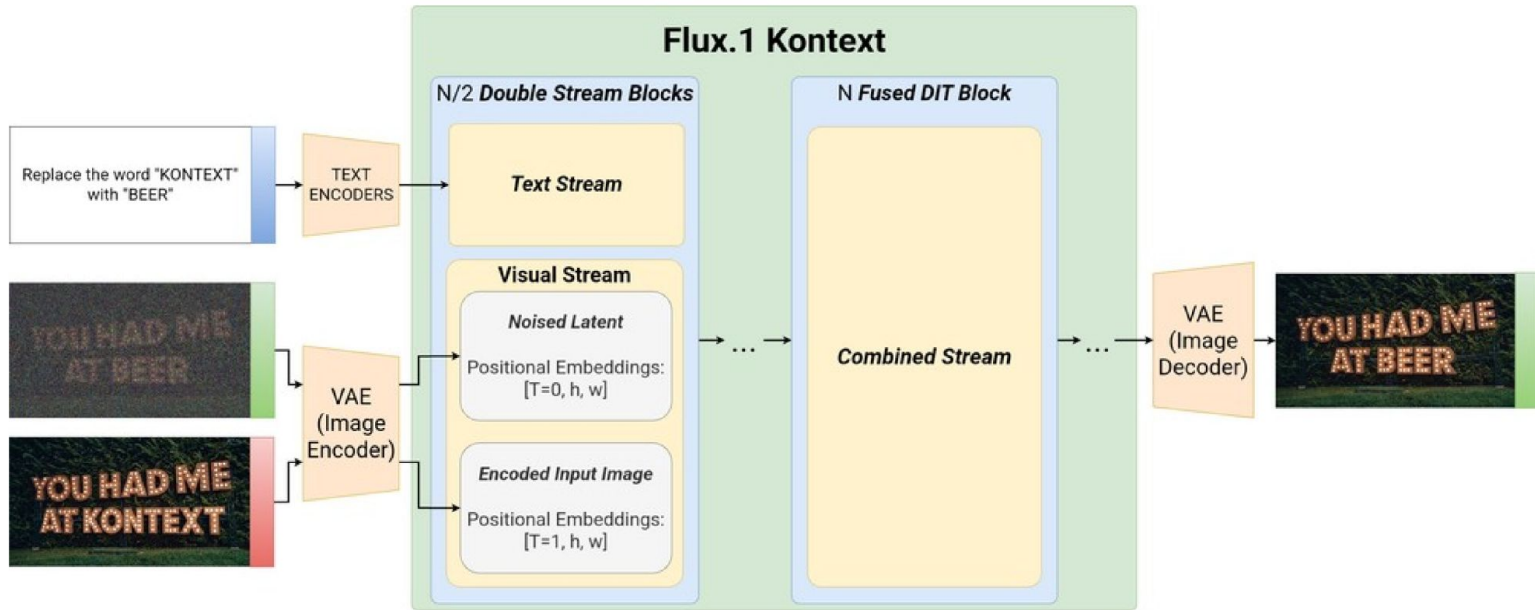
Three tall coconut trees, two of which are located in the center of the picture and one on the right edge. Two of the coconut trees are covered in golden coconut fruits under a clear blue sky, presenting a peaceful and tropical natural scene.



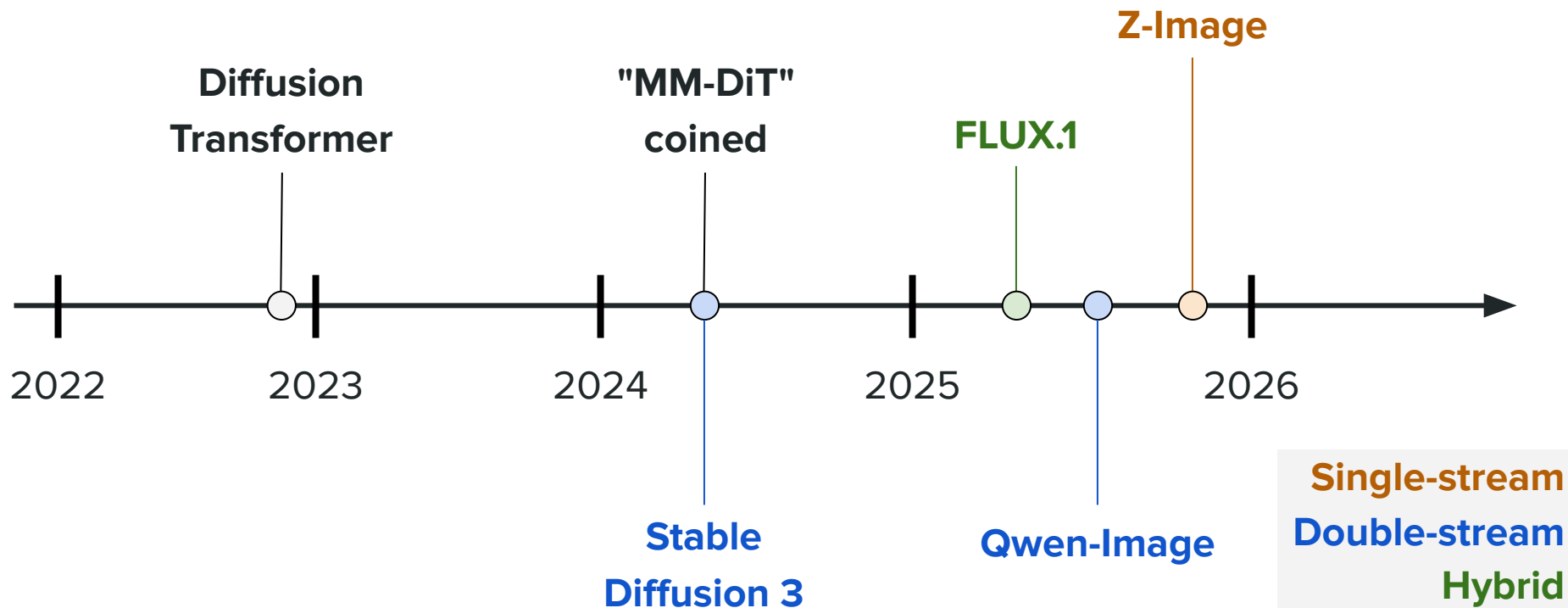
# "Single-stream" DiT variant: Z-Image



# Hybrid approach: FLUX.1 Kontext



# Timeline





# Diffusion & Large Vision Models

Motivation

U-Net

Diffusion Transformer

End-to-end example

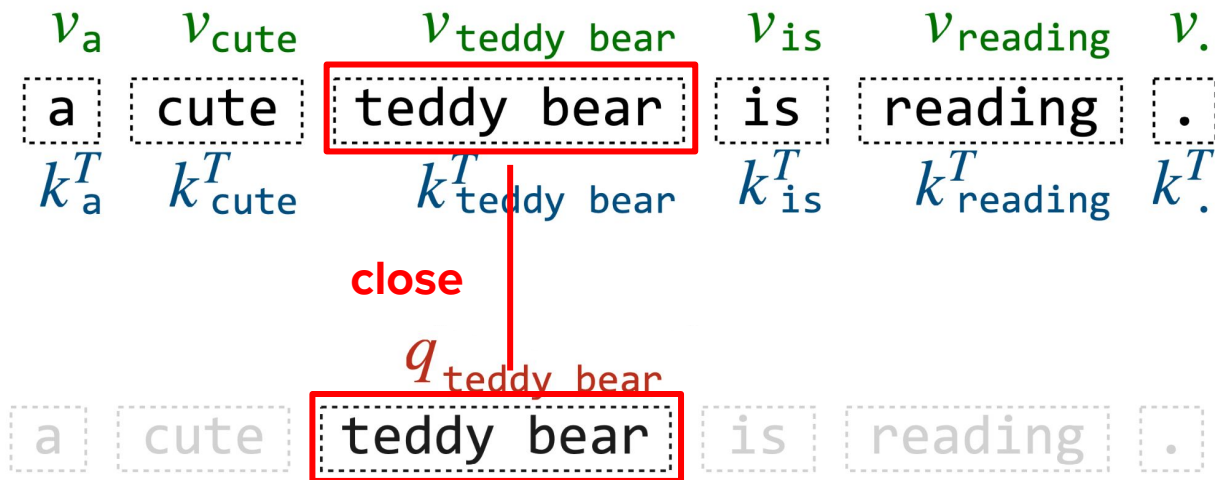
Multimodal DiT

**Optimizations**

---

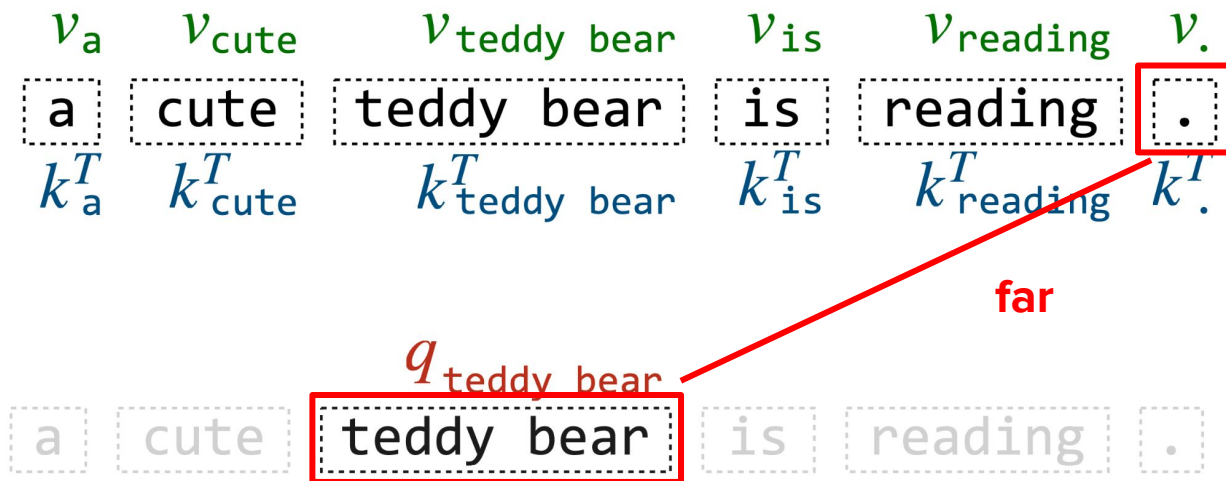
# Need for position information

**Motivation.** Direct links "lose" position info



# Need for position information

**Motivation.** Direct links "lose" position info



# Need for position information

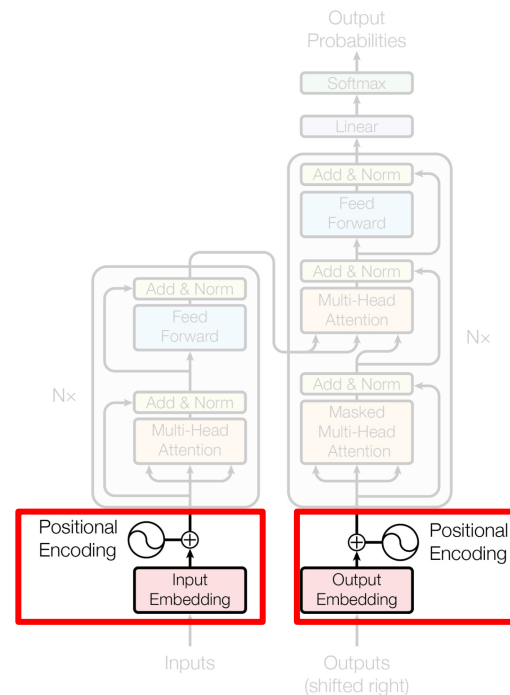
**Hope.** Dot products convey distance between token representations

$$\langle \mathbf{E}_m, \mathbf{E}_n \rangle = f(m - n)$$

# Absolute positional embeddings

**Idea.** Add position-specific embedding to token vector

$$\mathbf{E}_m = \mathbf{IE}_m + \mathbf{PE}_m$$



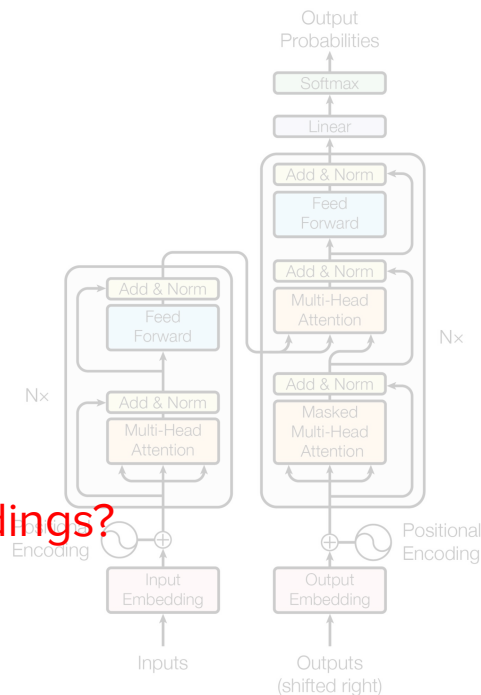
# Absolute positional embeddings

**Idea.** Add position-specific embedding to token vector

$$\mathbf{E}_m = \text{IE}_m + \text{PE}_m$$

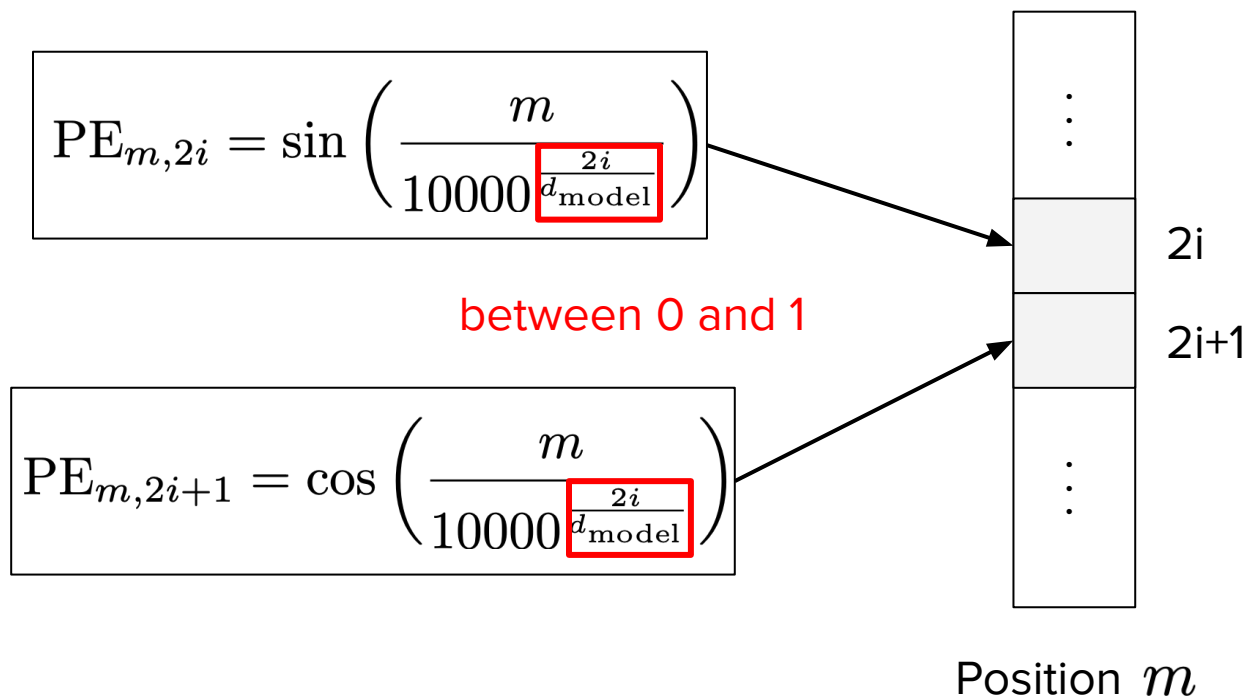
$$\langle \mathbf{E}_m, \mathbf{E}_n \rangle = \text{some terms} + \langle \text{PE}_m, \text{PE}_n \rangle$$

What choice for position embeddings?



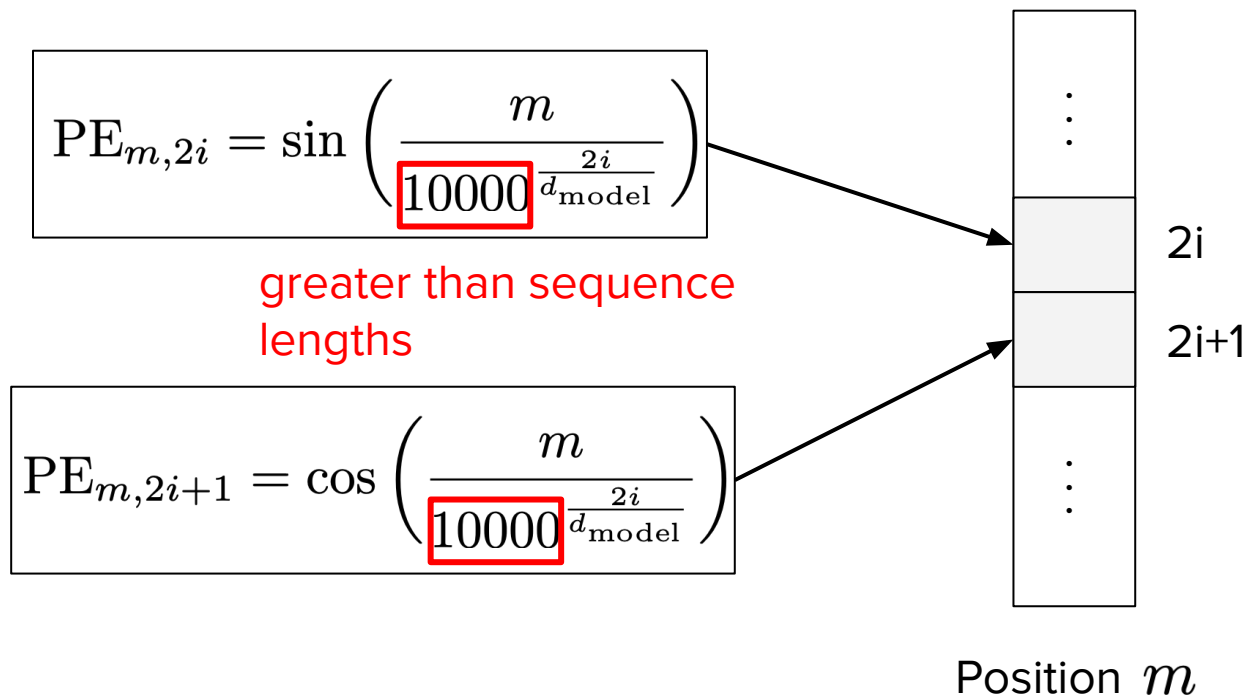
# Hardcoded position embeddings

**Proposal.** A mix of sines and cosines.

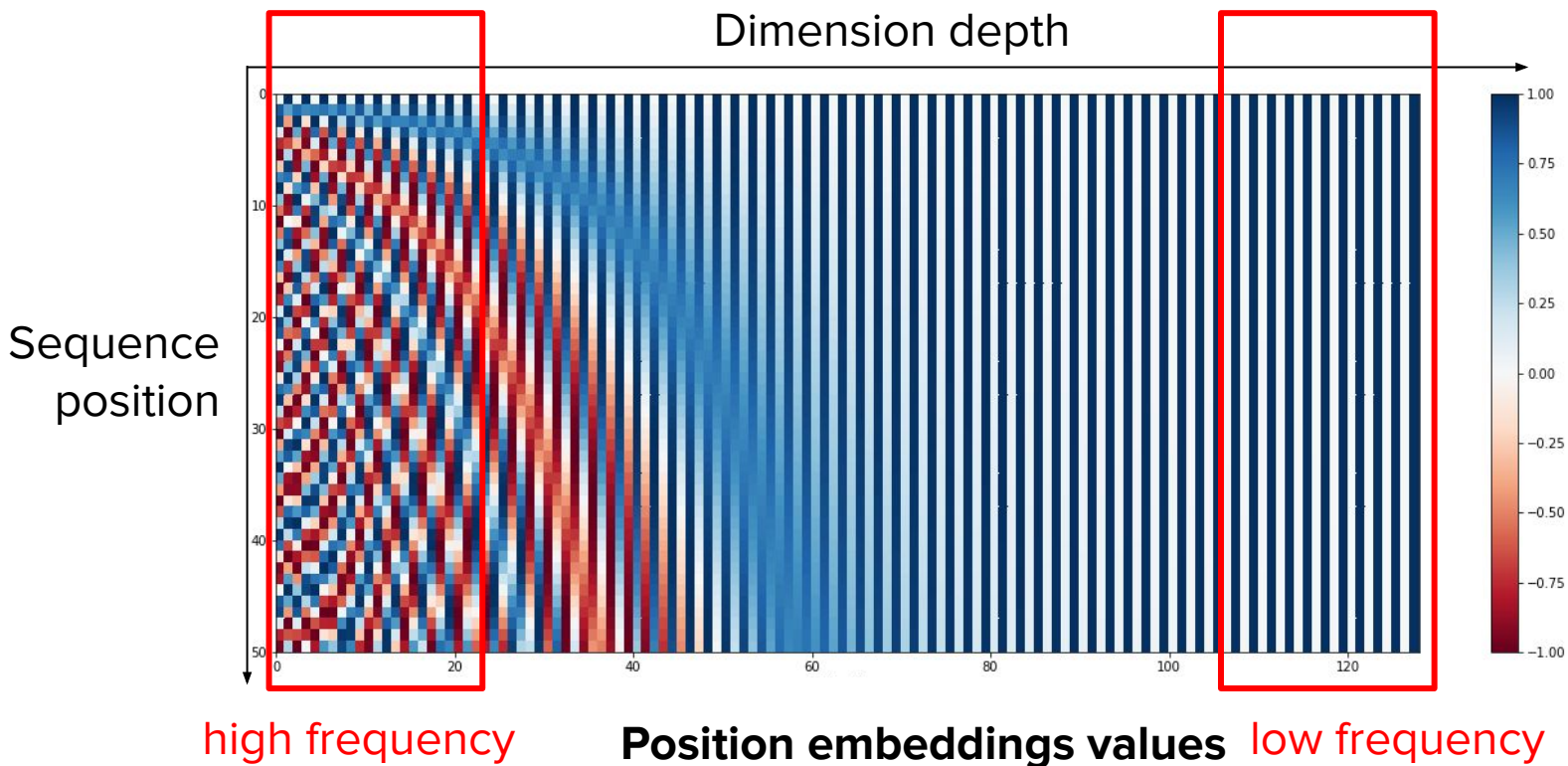


# Hardcoded position embeddings

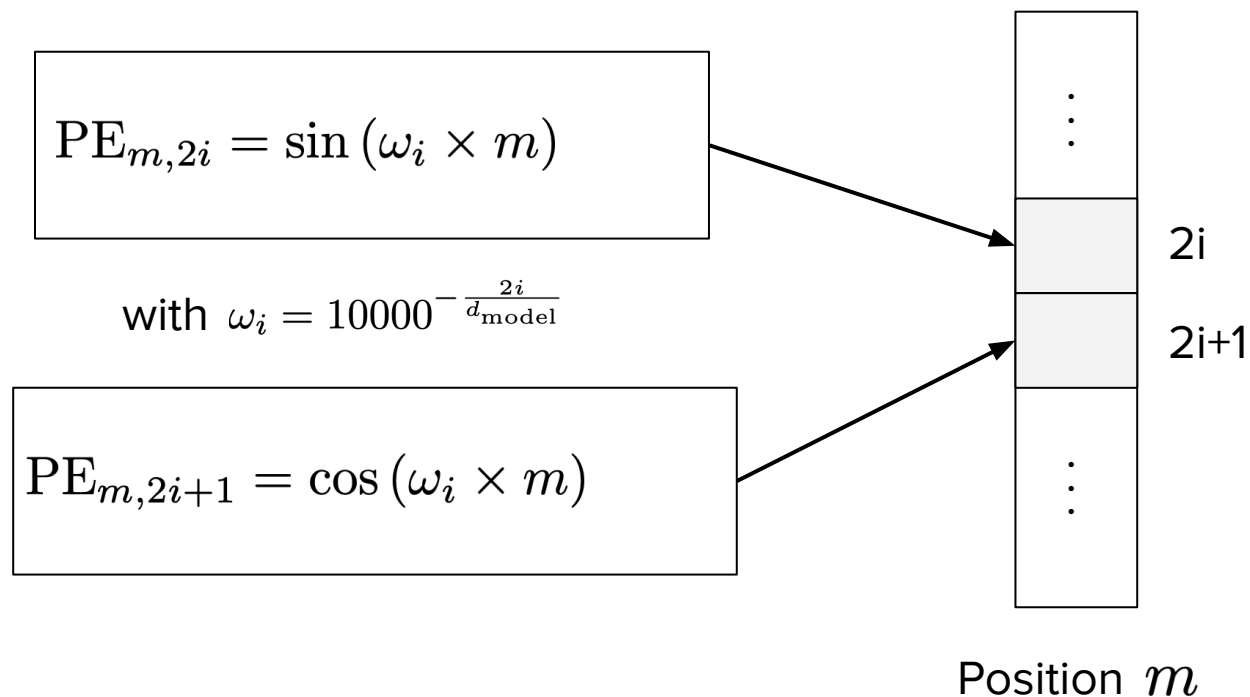
**Proposal.** A mix of sines and cosines.



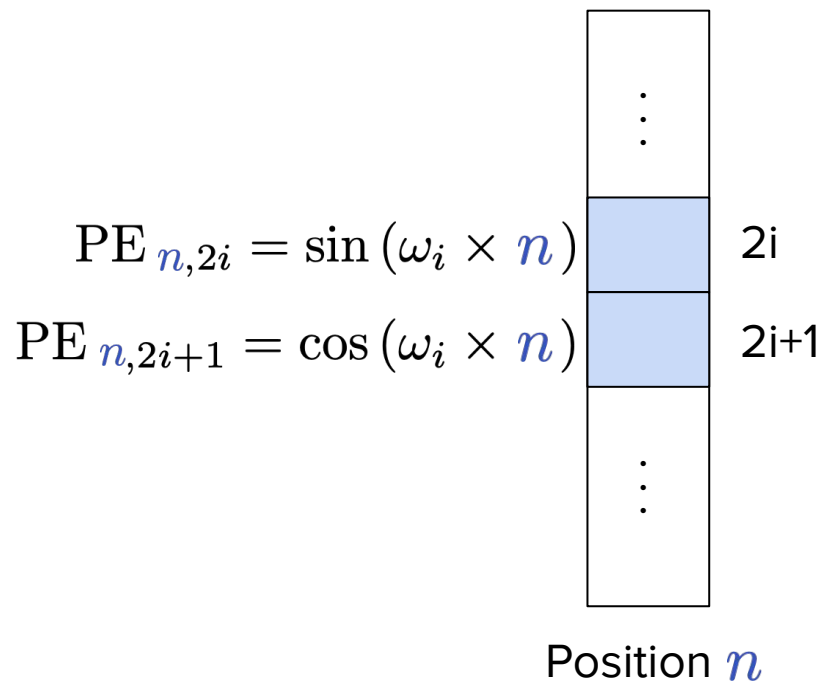
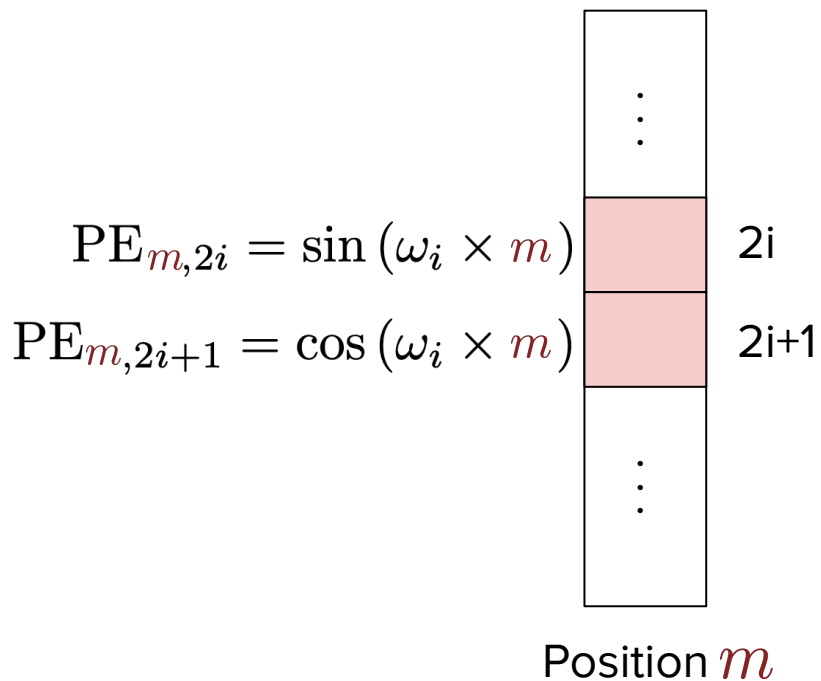
# Hardcoded position embeddings



# Hardcoded position embeddings



# Hardcoded position embeddings



# Hardcoded position embeddings

$$\langle \text{PE}_m, \text{PE}_n \rangle = \dots + \sin(\omega_i m) \sin(\omega_i n) + \cos(\omega_i m) \cos(\omega_i n) + \dots$$

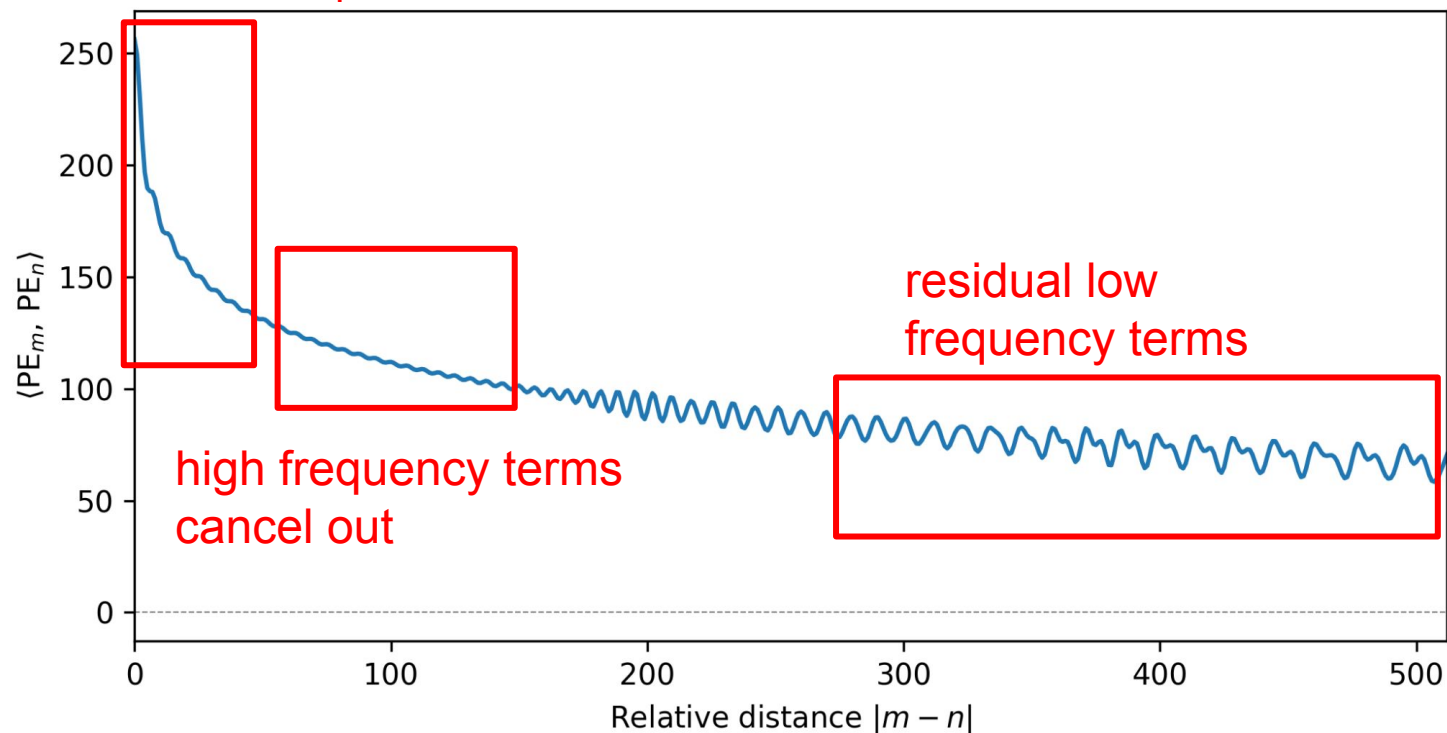


$$\text{Property: } \cos(a - b) = \sin(a) \sin(b) + \cos(a) \cos(b)$$

$$\langle \text{PE}_m, \text{PE}_n \rangle = \dots + \cos(\omega_i(m - n)) + \dots$$

# Hardcoded position embeddings

~all cosines drop



# Discussion

## **Benefits.**

- Simple to implement
- Does the job well as a baseline
- Gives every position a direct representational identity

# Discussion

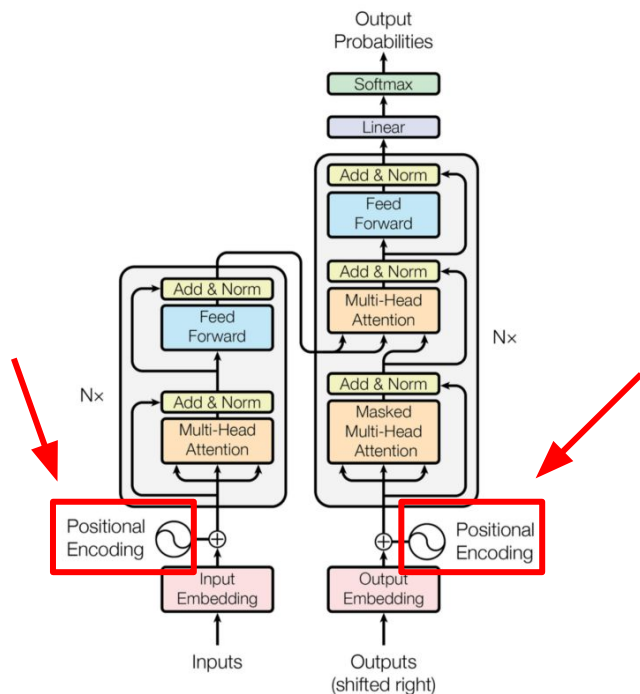
## Benefits.

- Simple to implement
- Does the job well as a baseline
- Gives every position a direct representational identity

## Limitations.

- Doesn't seem to be injected at the appropriate place
- Side effect: extra interaction terms
- Position embedding also makes it to value embeddings

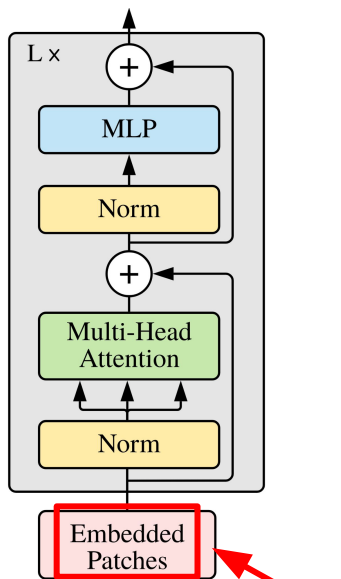
# Absolute position embeddings



Original Transformer

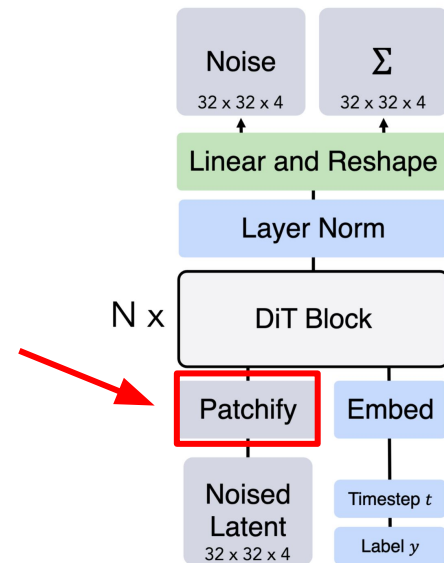
2017

## Who used it?



Original ViT

2020

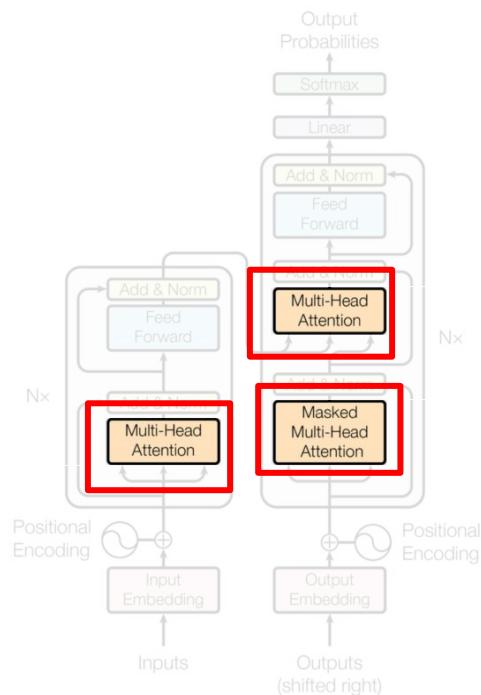


Original DiT

2022

# From absolute to relative position info

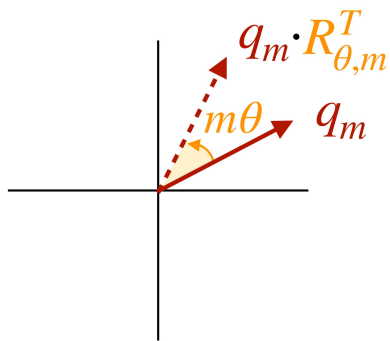
**Refinement.** Any way to move position information to where it's intuitively needed?



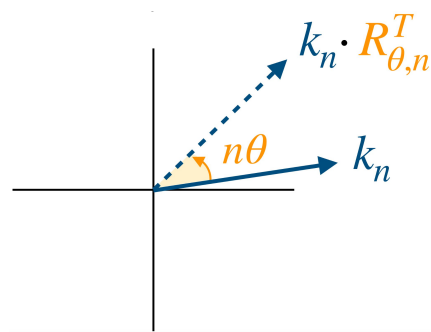
# Default choice nowadays: rotations in attention layer

## RoPE = Rotary Position Embeddings

**Idea.** Rotate query and key vectors with rotation matrix  $R_{\theta,m} = \begin{pmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{pmatrix}$



queries



keys

# Default choice nowadays: rotations in attention layer

**Benefits.** Relative distance nicely captured:

$$q_m k_n^T = x_m W_q \boxed{R_{\theta, n-m}} W_k^T x_n^T$$

effective rotation by angle  
 $\theta(n - m)$

# Generalization to 2D?

✓ Text.

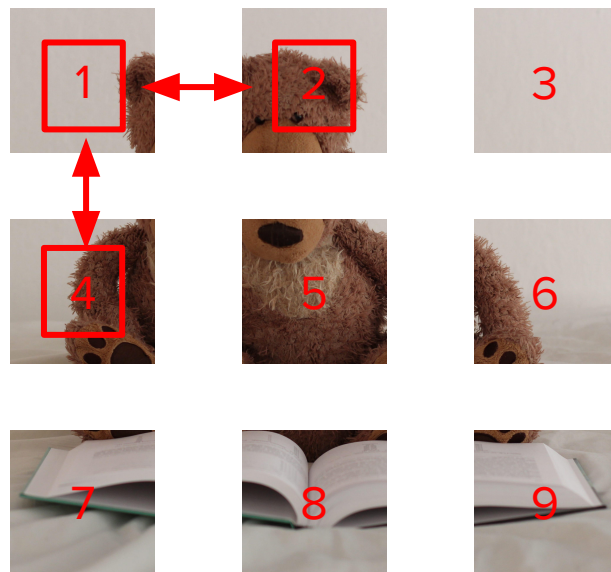
a   cute   teddy bear   is   reading   .  
1   2   3   4   5   6   ...

# Generalization to 2D?

✓ Text.

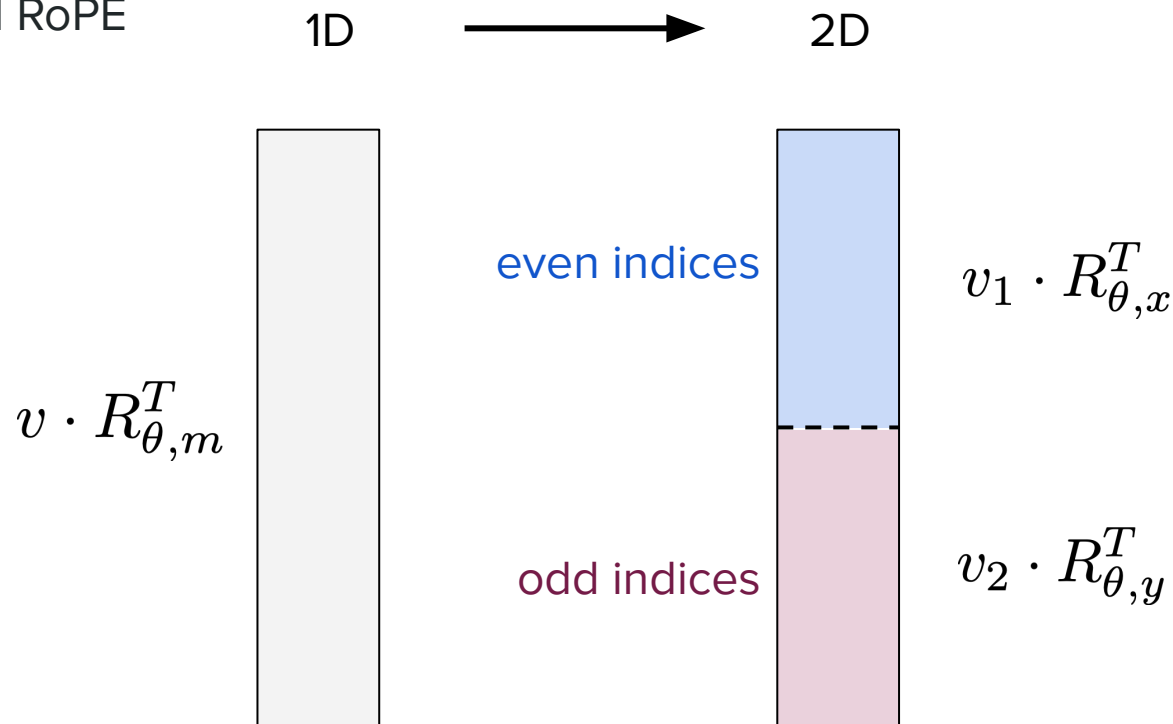


? Images.



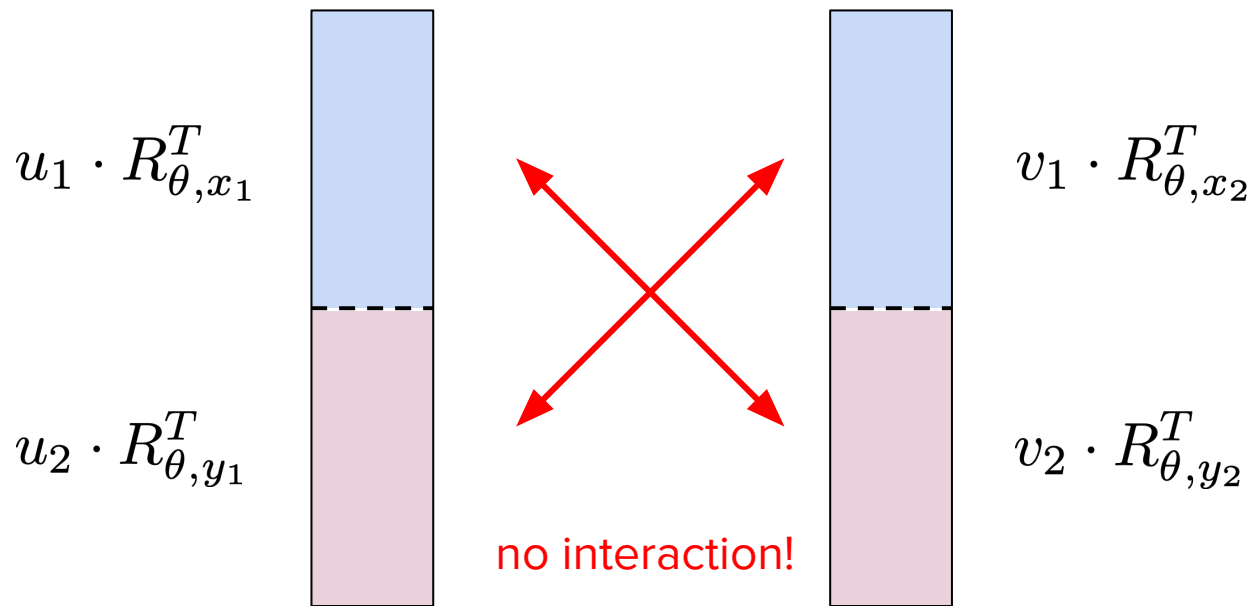
# Generalization to 2D?

First idea. Axial RoPE



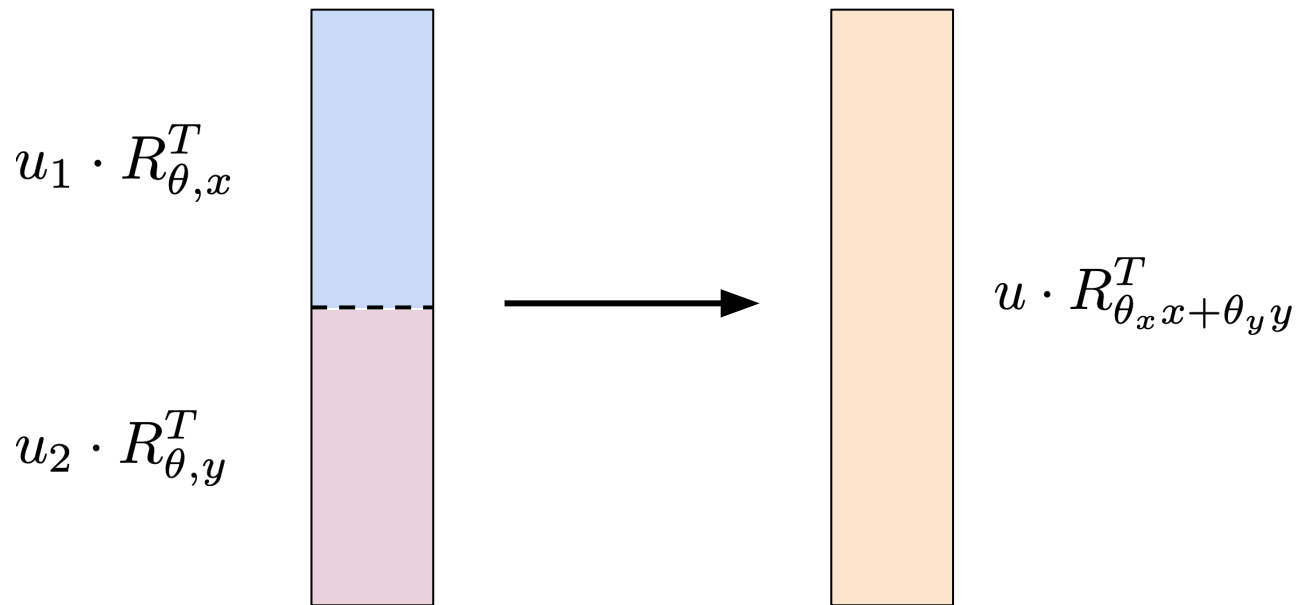
# Generalization to 2D?

**Problems.** No cross axes interactions + segregated information is arbitrary



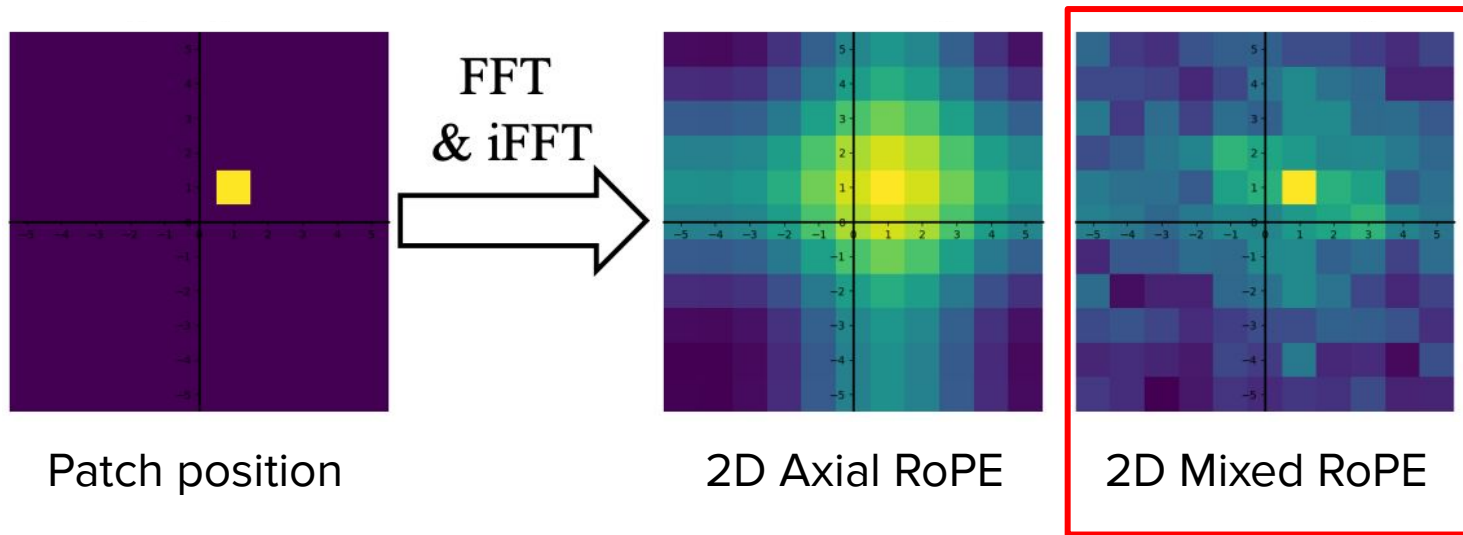
# 2D RoPE

**Proposal.** Mix both axes as part of the same rotation matrix



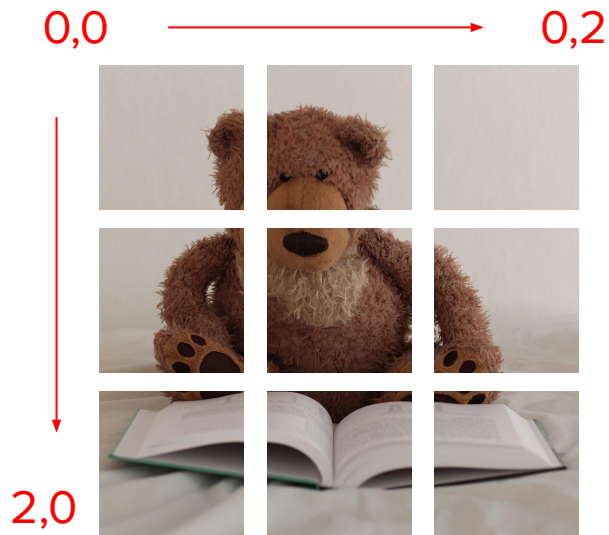
# 2D RoPE

**Consequence.** Avoid axes artifacts.

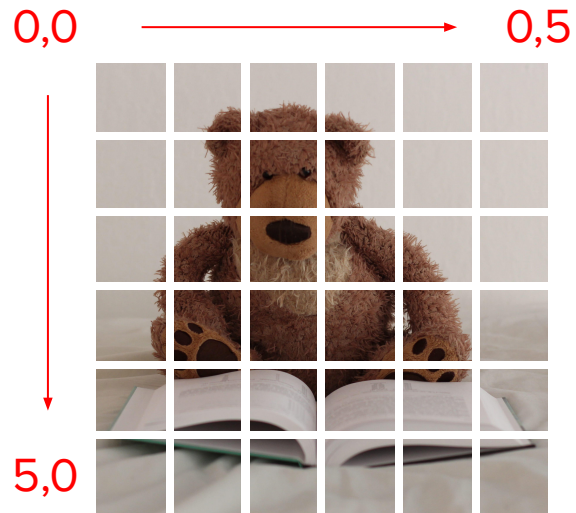


# Scaling RoPE

**Common case.** Resolution variations?

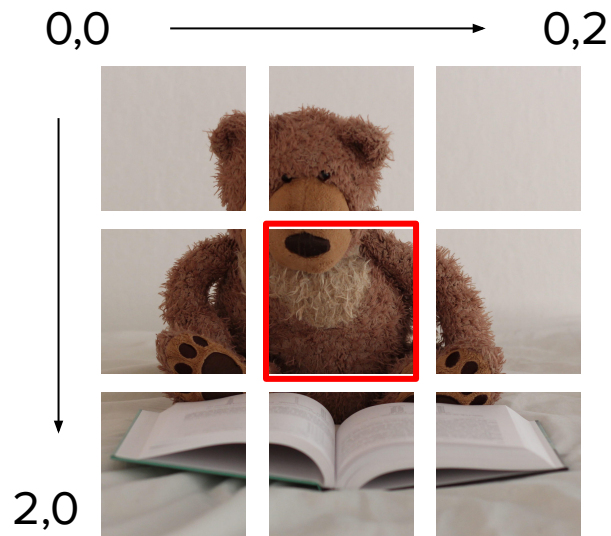


VS



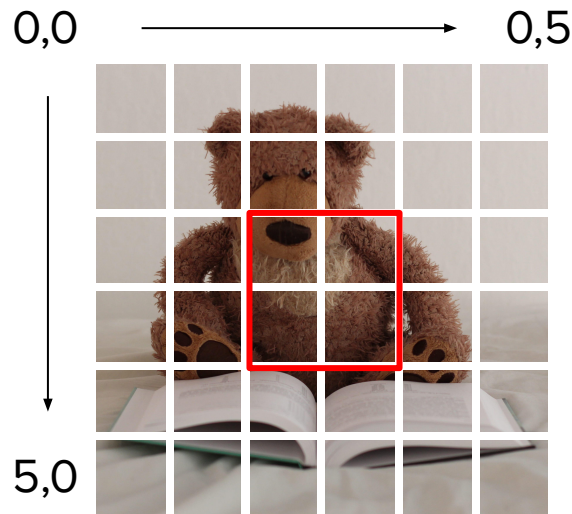
# Scaling RoPE

**Problem.** Spatial meaning also changes.



"center"

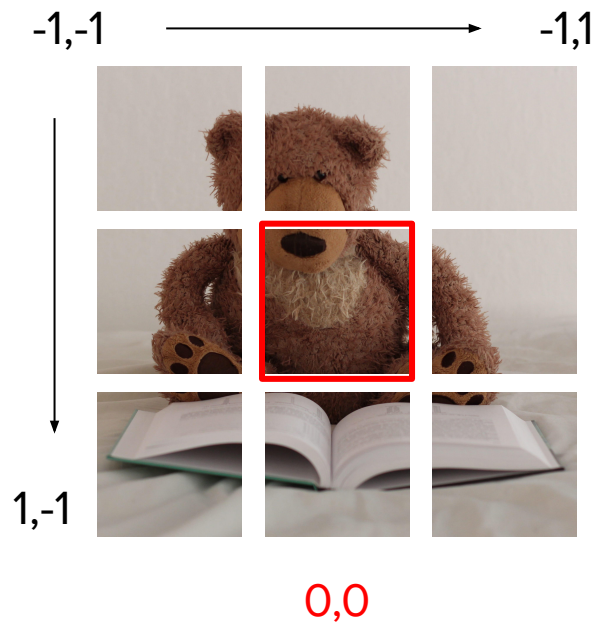
vs



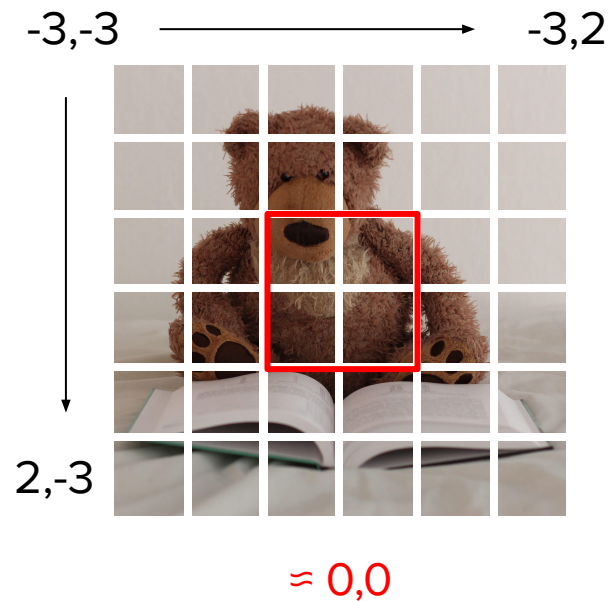
"center"

# Scaling RoPE

**Remedy.** Scale to canonical coordinates.



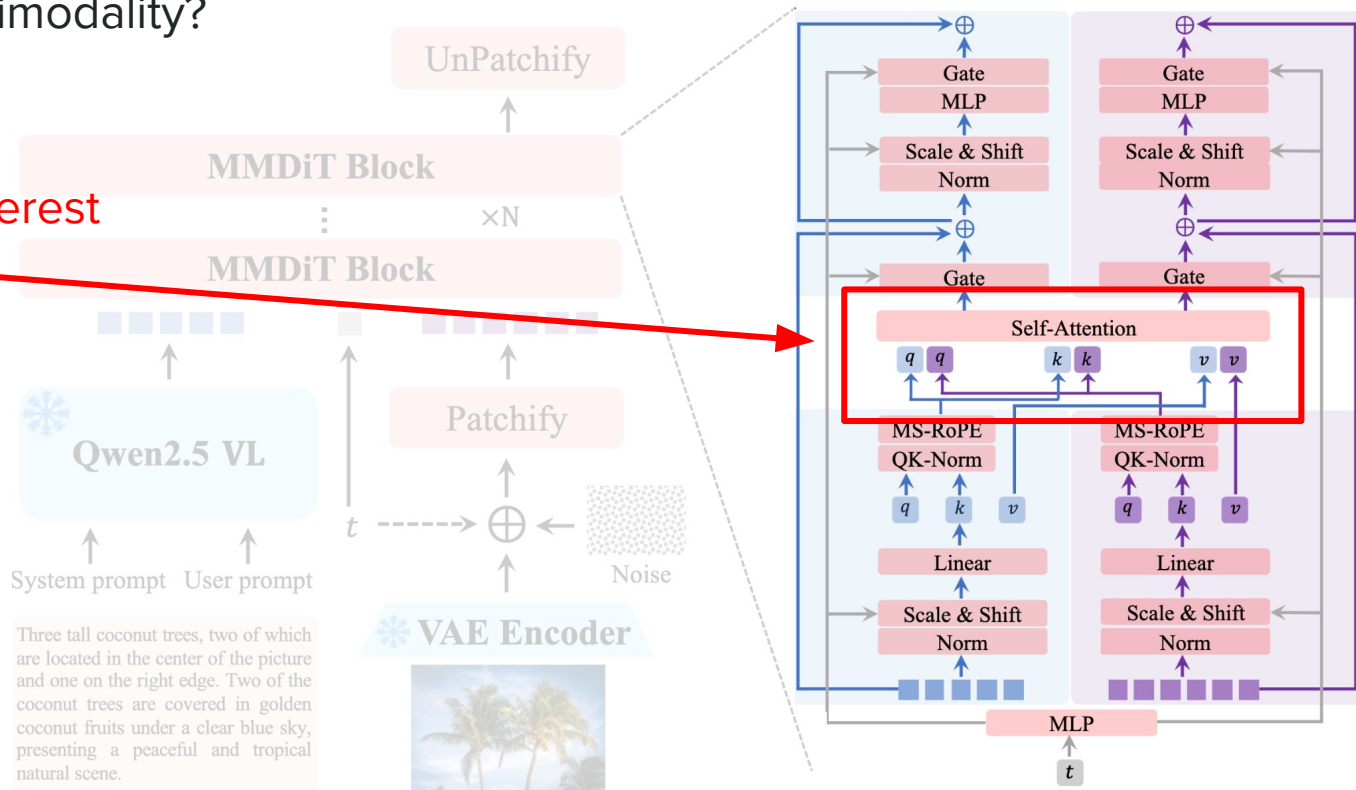
VS



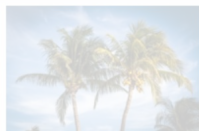
# Multimodal RoPE

Common case. Multimodality?

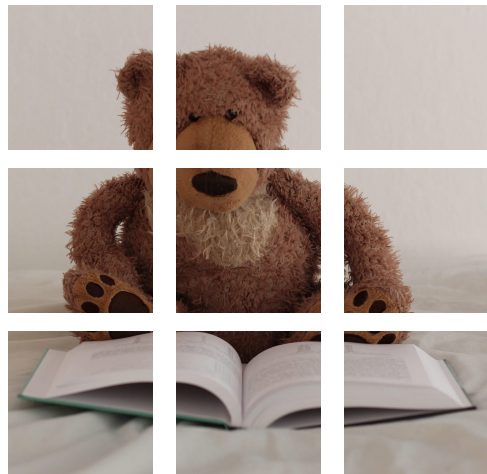
Potential area of interest



Three tall coconut trees, two of which are located in the center of the picture and one on the right edge. Two of the coconut trees are covered in golden coconut fruits under a clear blue sky, presenting a peaceful and tropical natural scene.



# Multimodal RoPE



2D



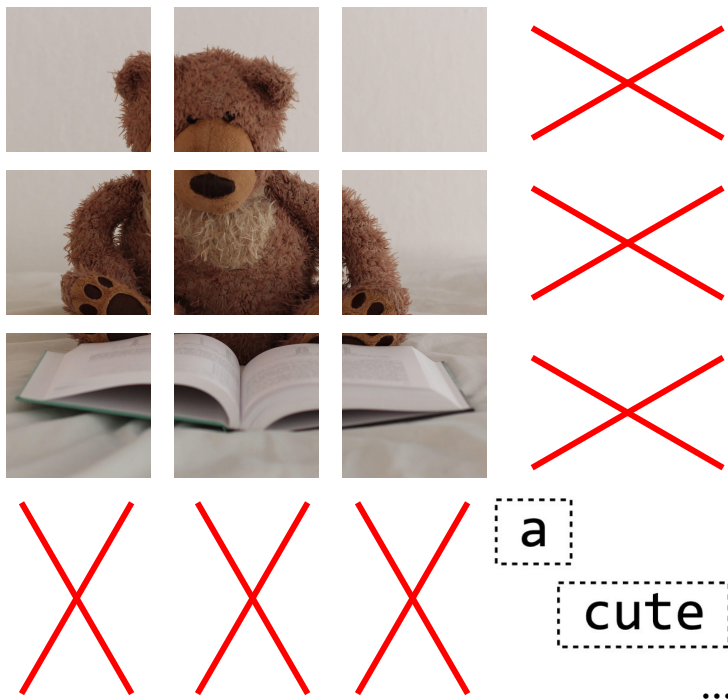
1D

a cute teddy bear is reading .

How to reconcile the two?

# Multimodal RoPE

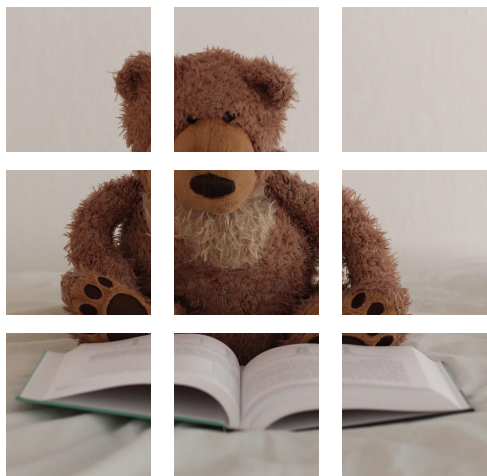
**MSRoPE = Multimodal Scalable RoPE**



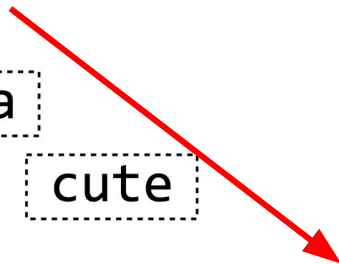
## Benefits.

- Image and text can cohabit without interference

## MSRoPE = Multimodal Scalable RoPE



a  
cute  
...



### Benefits.

- Image and text can cohabit without interference
- Text is functionally equivalent to 1D RoPE

# Conclusion

**Embedding positions = open problem.**

- Many variations out there
- Dust hasn't settled yet
- Trade-offs

Thank you for your attention!

