

# ViT and Self-Supervised Models

CS280

Spring 2025

Angjoo Kanazawa

# Attention is expensive

For Memory and Computation.

Prevents long context!

## Flash Attention

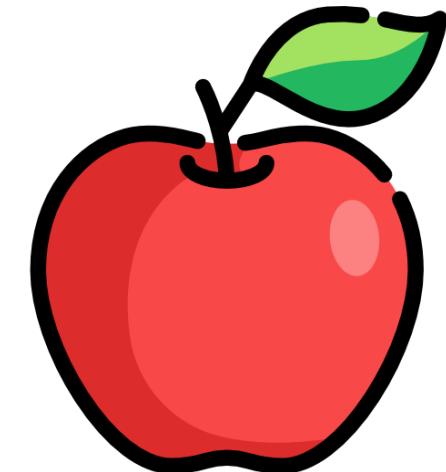
- Does not instantiate the entire  $N \times N$  attention matrix  $A$ , i.e.  $QK^T$
- Dedicated GPU kernel for efficient implementation
- Running sum for softmax...
- Memory reduced from  $O(N^2)$  to  $O(Nd)$
- Need operation to be on  $Q$  and  $K$  independently

# Positional Encoding

I bought an apple watch

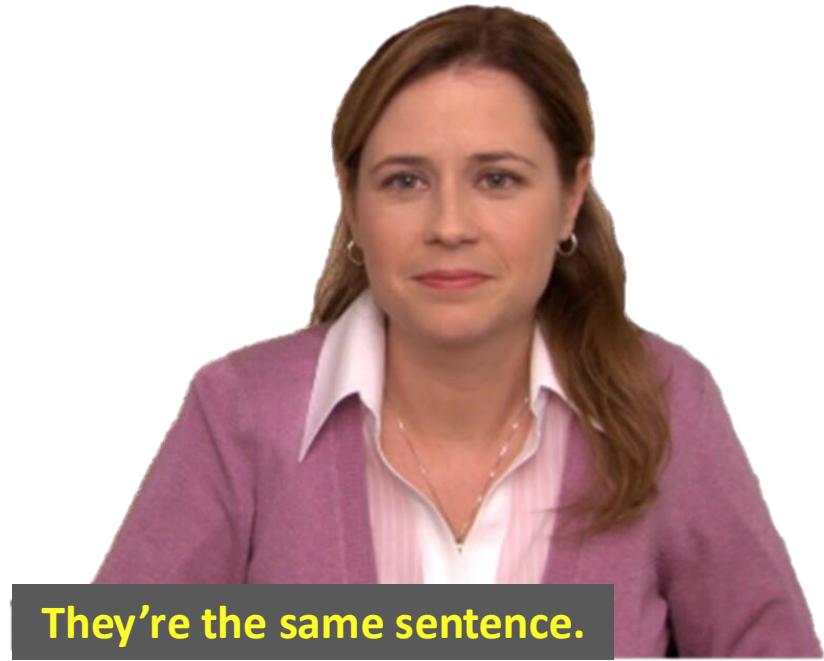


watch an apple I bought



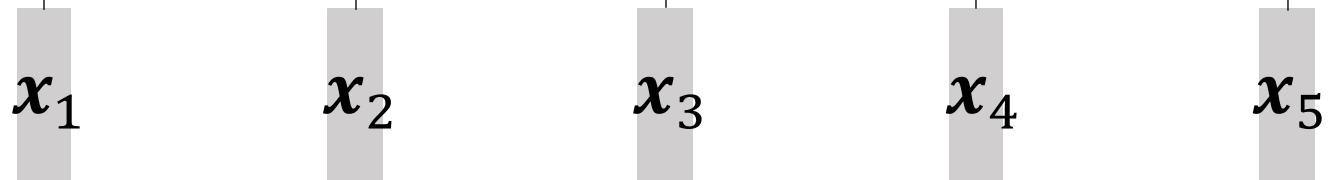
I bought an apple watch

watch an apple I bought



They're the same sentence.

Embedded  
Tokens



Tokens

I      **bought**      an      apple      watch

Position

$m = 1$        $m = 2$        $m = 3$        $m = 4$        $m = 5$

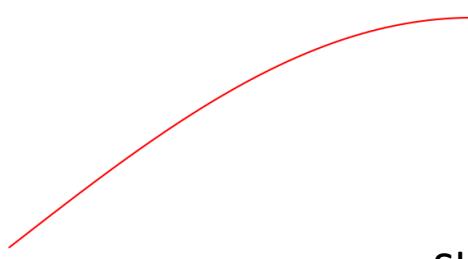
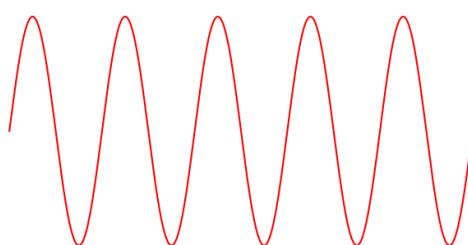
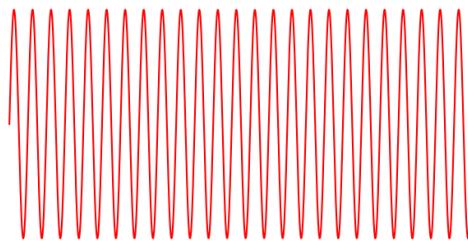
Position  $k$

Angular frequency

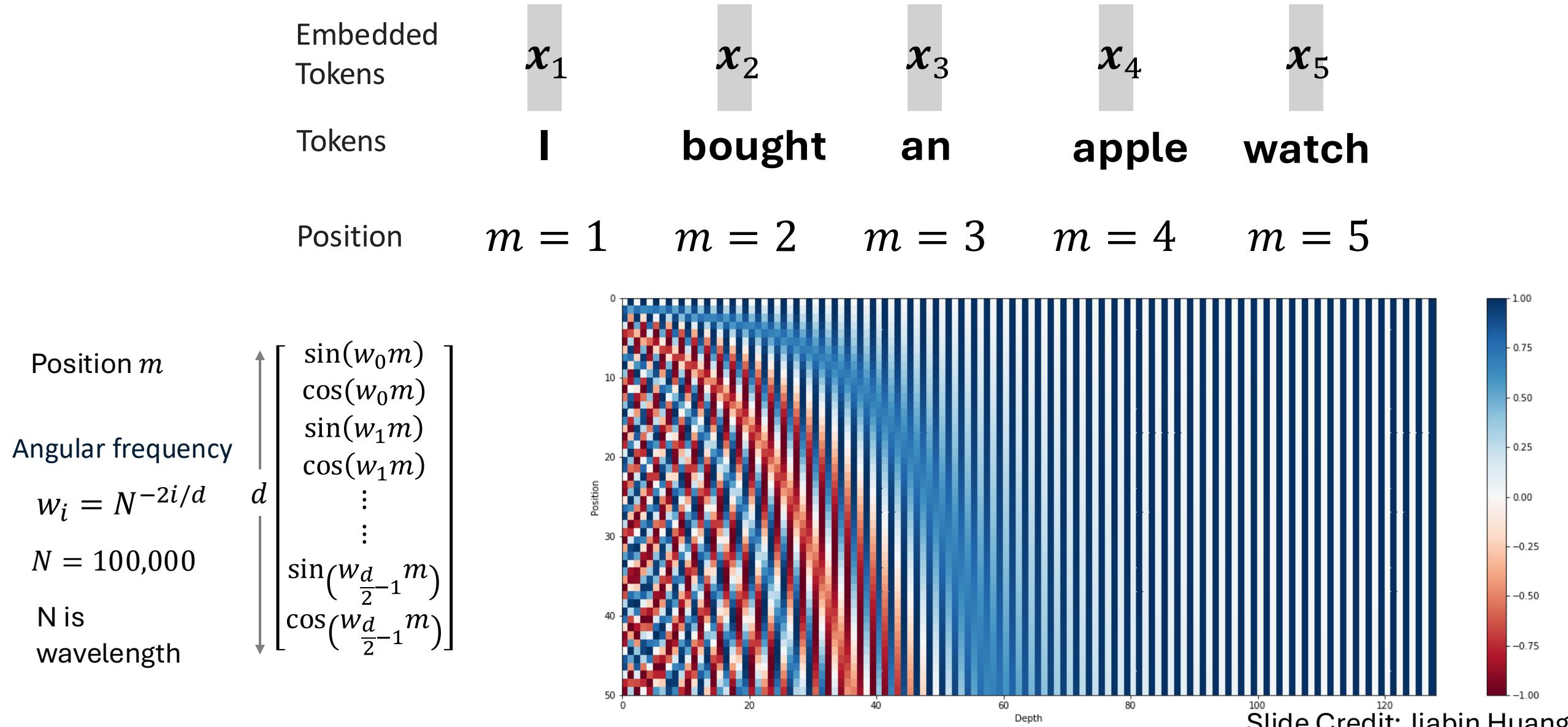
$w_i = N^{-2i/d}$

$N = 100,000$

$$d \begin{bmatrix} \sin(w_0 m) \\ \cos(w_0 m) \\ \sin(w_1 m) \\ \cos(w_1 m) \\ \vdots \\ \vdots \\ \sin(w_{\frac{d}{2}-1} m) \\ \cos(w_{\frac{d}{2}-1} m) \end{bmatrix} \leftarrow \begin{array}{l} \text{Fast oscillating} \\ \text{Slow oscillating} \end{array}$$

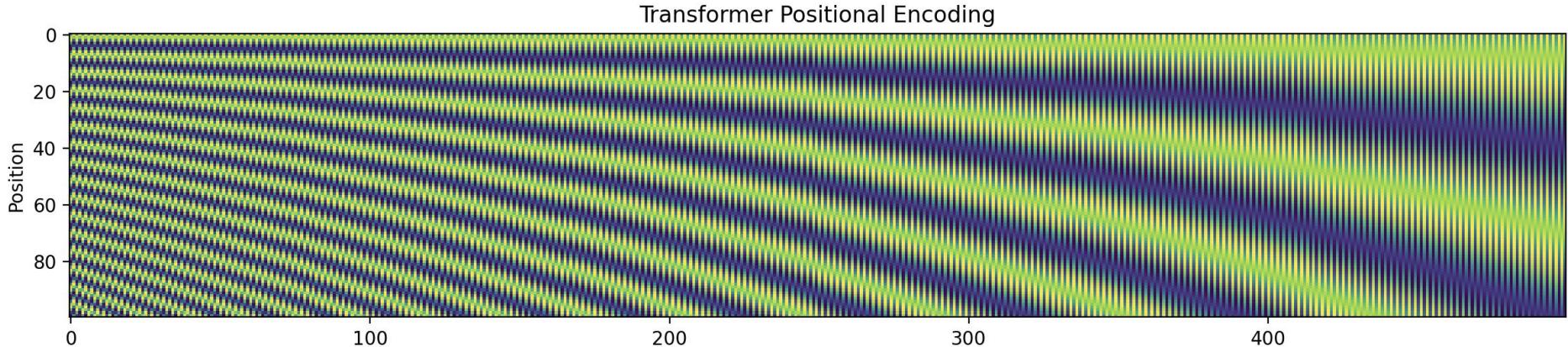


# Positional Encoding



# Changing N

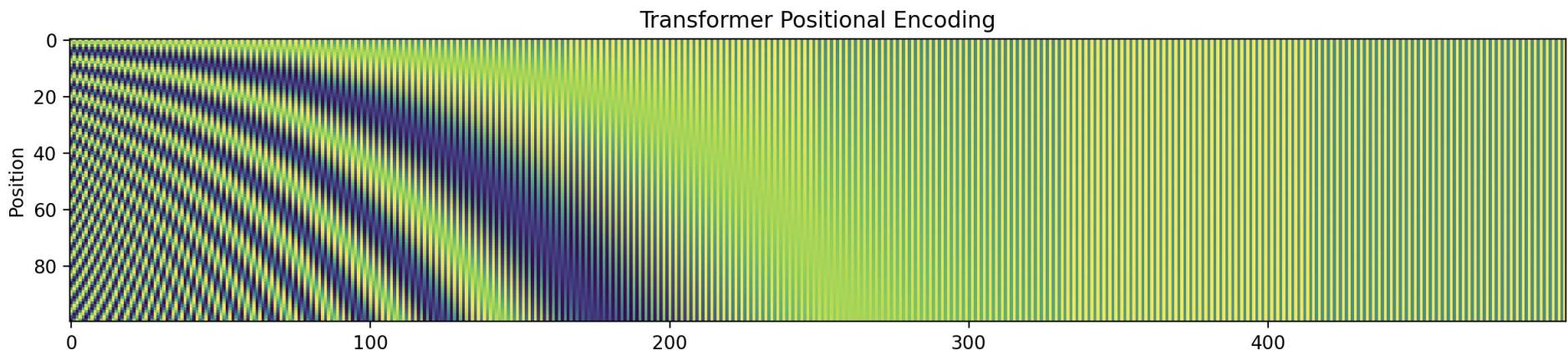
N=10

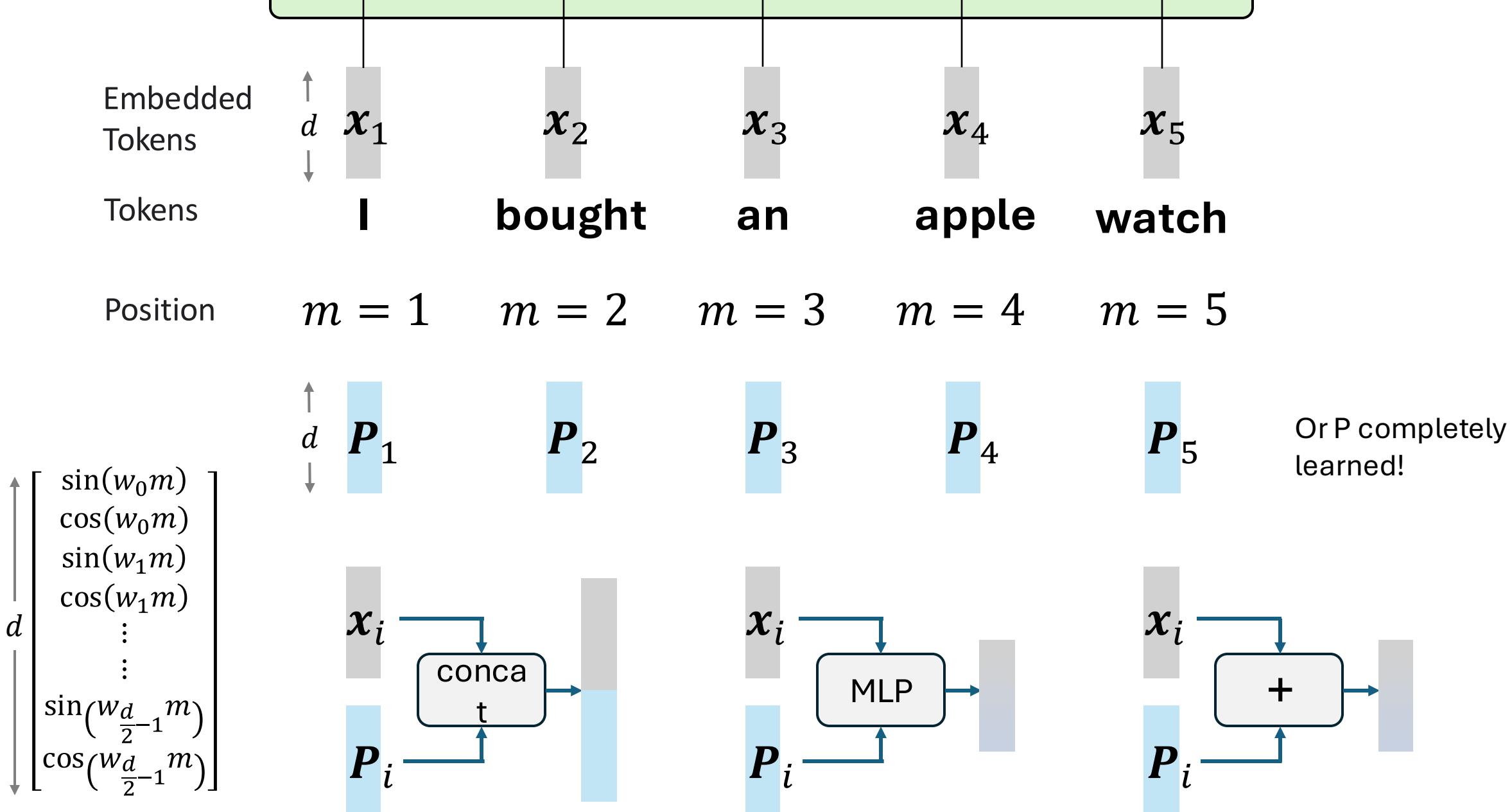


Goal: Unique positional encoding for each of your tokens.

N is the wavelength of the low-frequency band, larger it is, larger number you can represent.

What is used  
N=10,000



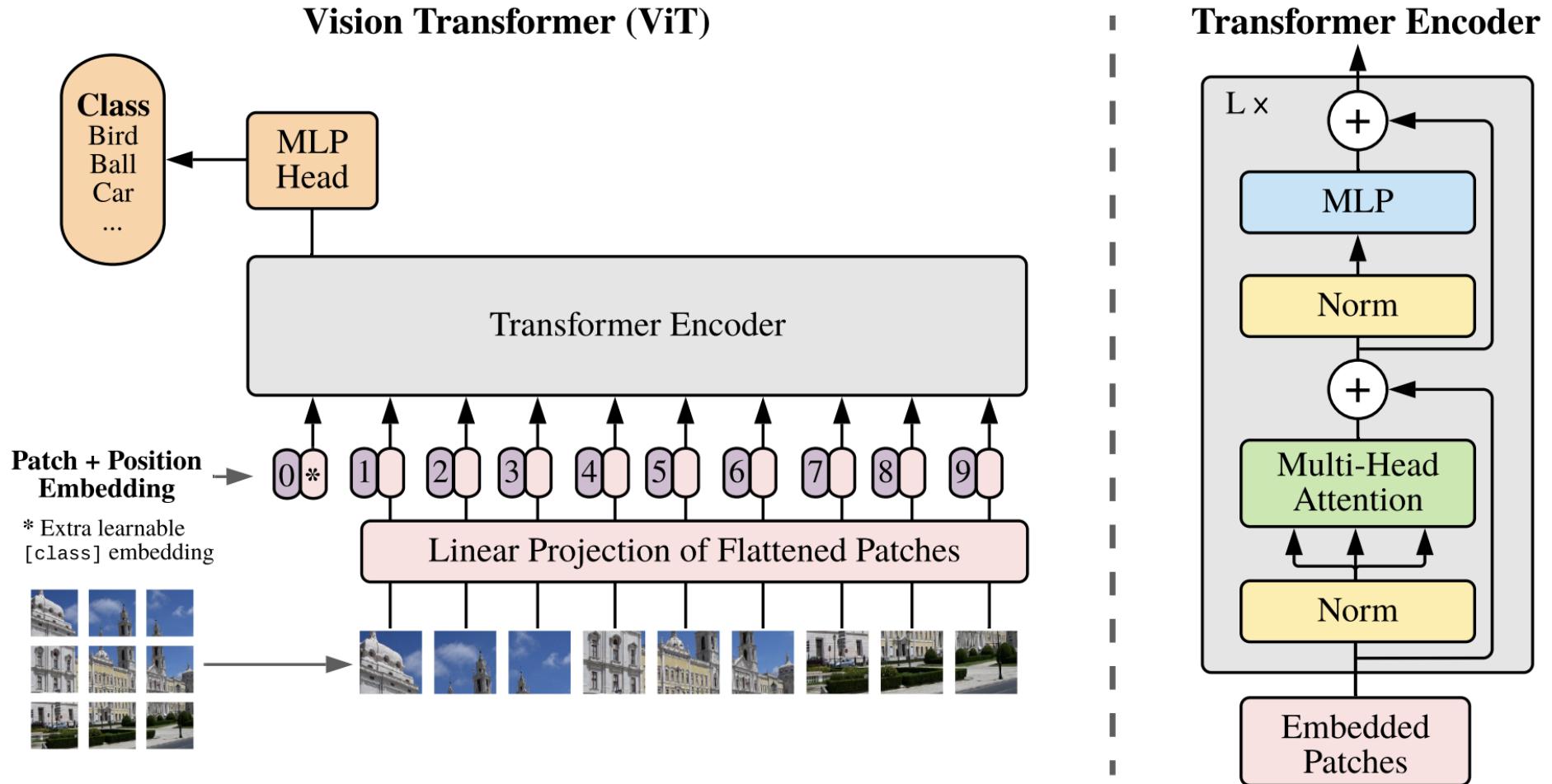


**Back to pixels!**

# How to apply transformers to pixels?

- Naïve tokenization of a pixel for 224x224 image:
  - Over 500k tokens!
- Too much!
- Solution?

# Vision Transformer



# Positional Encoding

- Originally, they just flattened the patches from 1-256, treated it as a 1D signal!
- Or learned positional embedding
- Problem?
- Both solutions assume fixed size!
- You can't change your image size!
- Need to retrain..

# Relative POsitional Encoding: ROPE

# Absolute

# Positional Encoding

Position

1

2

3

4

5

6

I        wal        my        dog    every    day  
                    k

# Absolute

# Positional Encoding



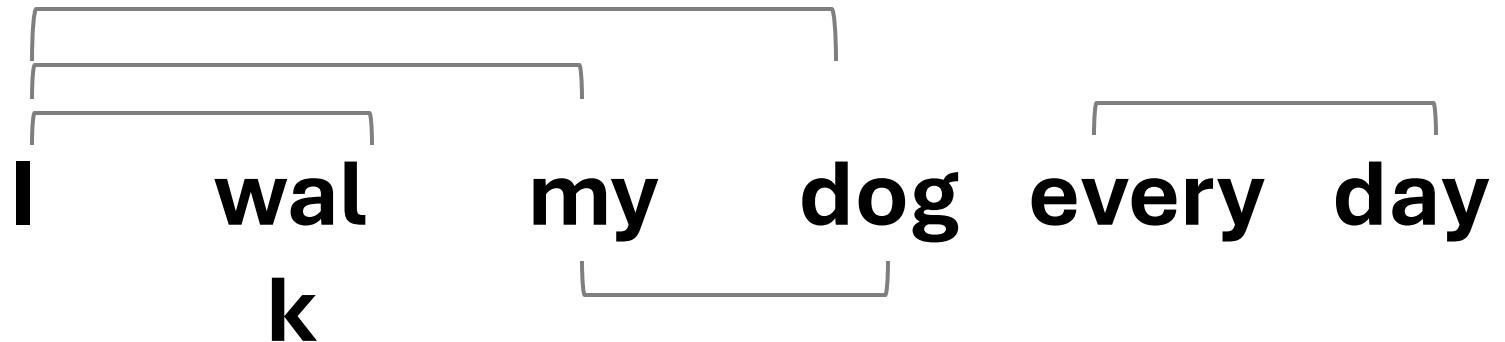
Position	1	2	3	4	5	6
	I	<b>wal</b>	<b>my</b>	<b>dog</b>	<b>every</b>	<b>day</b>
		<b>k</b>				

Position	1	2	3	4	5	6
	<b>every</b>	<b>day</b>	I	<b>wal</b>	<b>my</b>	<b>dog</b>
				<b>k</b>		

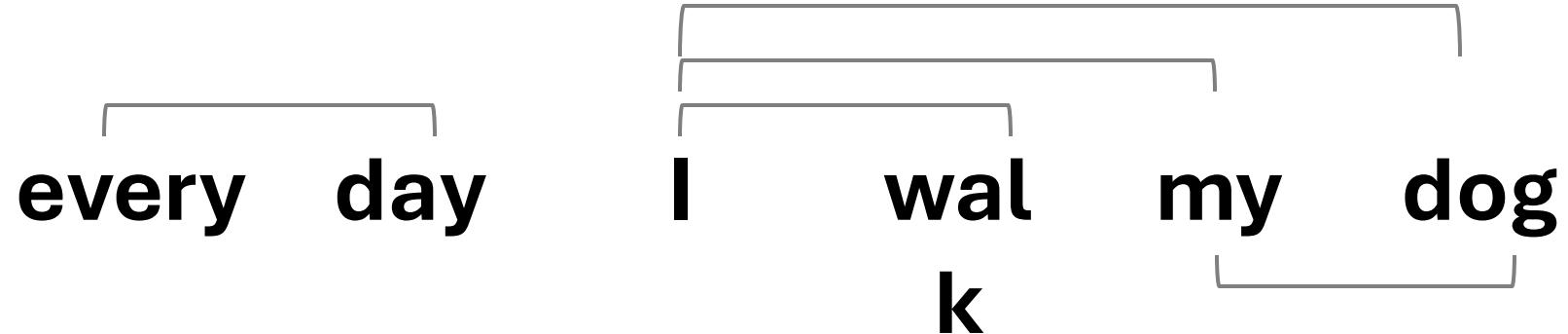
RELATIVe

# Positional Encoding

Position



Position



# Re-write the original posenc

$$\text{PE}(p) = \begin{bmatrix} \sin(w_0 p) \\ \cos(w_0 p) \\ \sin(w_1 p) \\ \cos(w_1 p) \\ \vdots \\ \sin\left(w_{\frac{d}{2}-1} p\right) \\ \cos\left(w_{\frac{d}{2}-1} p\right) \end{bmatrix}^d$$

- Focus on one frequency,  $\theta = w_d$
- Then this is a 2D vector
$$\text{PE}(p) = [\sin(\theta p), \cos(\theta p)].$$
- This is a coordinate on an unit circle!
- You can write this fourier basis in its exponential form (Euler's formula)
$$e^{i\theta p} = \cos(\theta p) + i \sin(\theta p).$$
- Has a nice property (holds across frequencies)!:

$$\langle \text{PE}(p_1), \text{PE}(p_2) \rangle = e^{-i\theta p_1} e^{i\theta p_2} = e^{i\theta(p_2 - p_1)}$$

Multiplying by  $e^{i\theta}$  is same as applying a 2D rotation matrix

# Re-write the original posenc

- Focus on one frequency,  $\theta = w_d$

- Then this is a 2D vector

$$\text{PE}(p) = [\sin(\theta p), \cos(\theta p)].$$

- This is a coordinate on an unit circle!

- You can write this fourier basis in its exponential form  
**(Euler's formula)**

$$e^{i\theta p} = \cos(\theta p) + i \sin(\theta p).$$

Multiplying by  $e^{i\theta}$  is same as applying a 2D rotation matrix

- Has a nice property (holds across frequencies)!:

$$\langle \text{PE}(p_1), \text{PE}(p_2) \rangle = e^{-i\theta p_1} e^{i\theta p_2} = e^{i\theta(p_2 - p_1)}$$

# ROPE — rotate tokens for relative positional encoding

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V. \quad Q = W_q X, \quad K = W_k X.$$

$$QK^T = (W_q X)(X^T W_k^T).$$

$$Q = W_q(Xe^{im}), \quad K = W_k(Xe^{in}).$$

$$QK^T = W_q X e^{i(m-n)} X^T W_k^T.$$

Absolute PosEnc:

$$X = \cancel{X'} + PE(m)$$

Pre-multiply the tokens with rotations

Flash Attention friendly!!

Unfortunately only works in 1D SO(2) bc of  
commutativity

Aurich and Weule 1995  
Tomasi and Manduchi 1998...

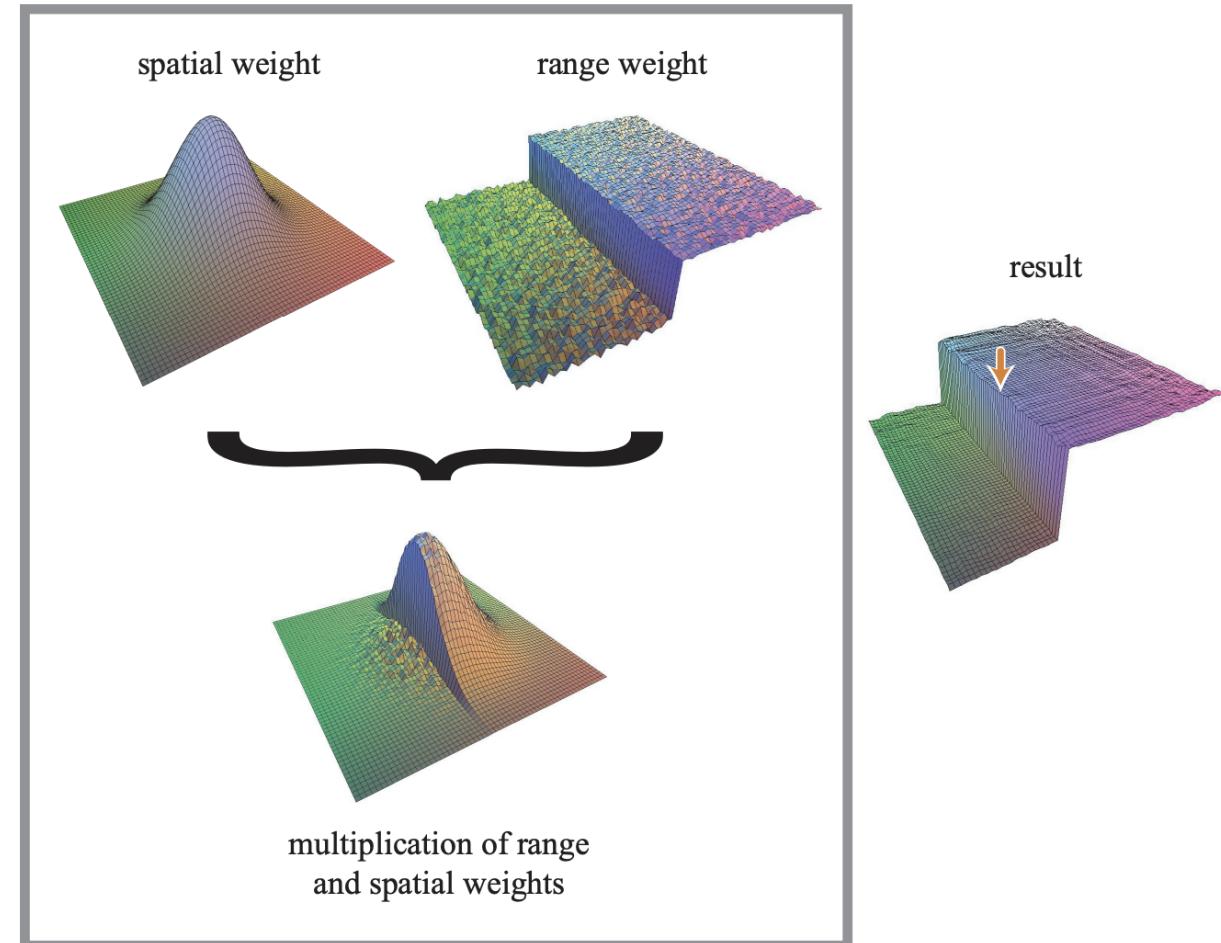
# Bilateral Filter

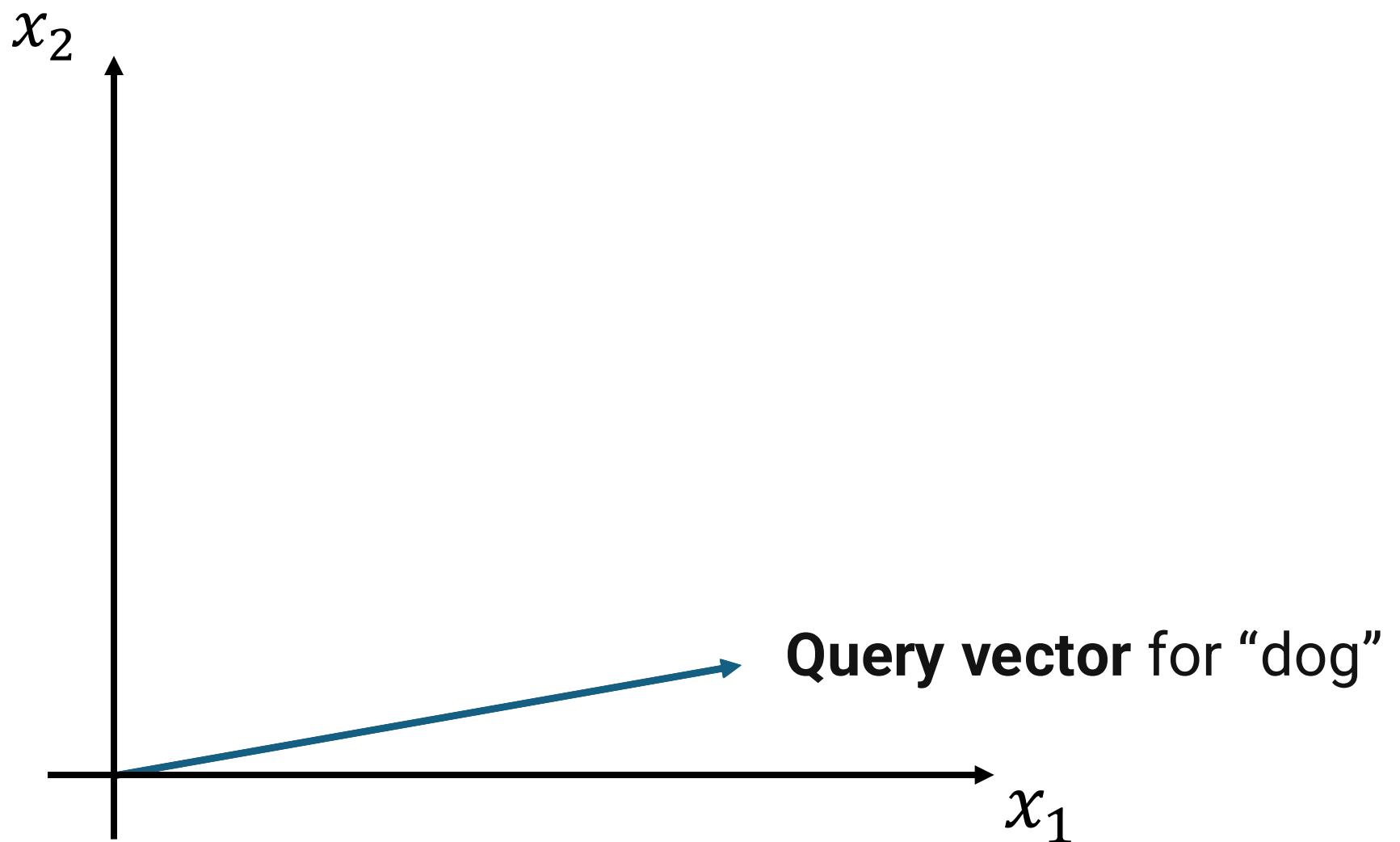
$$BF[I]_p = \frac{1}{W_p} \sum_{q \in S} G_{\sigma_s}(\|p - q\|) G_{\sigma_r}(I_p - I_q) I_q$$

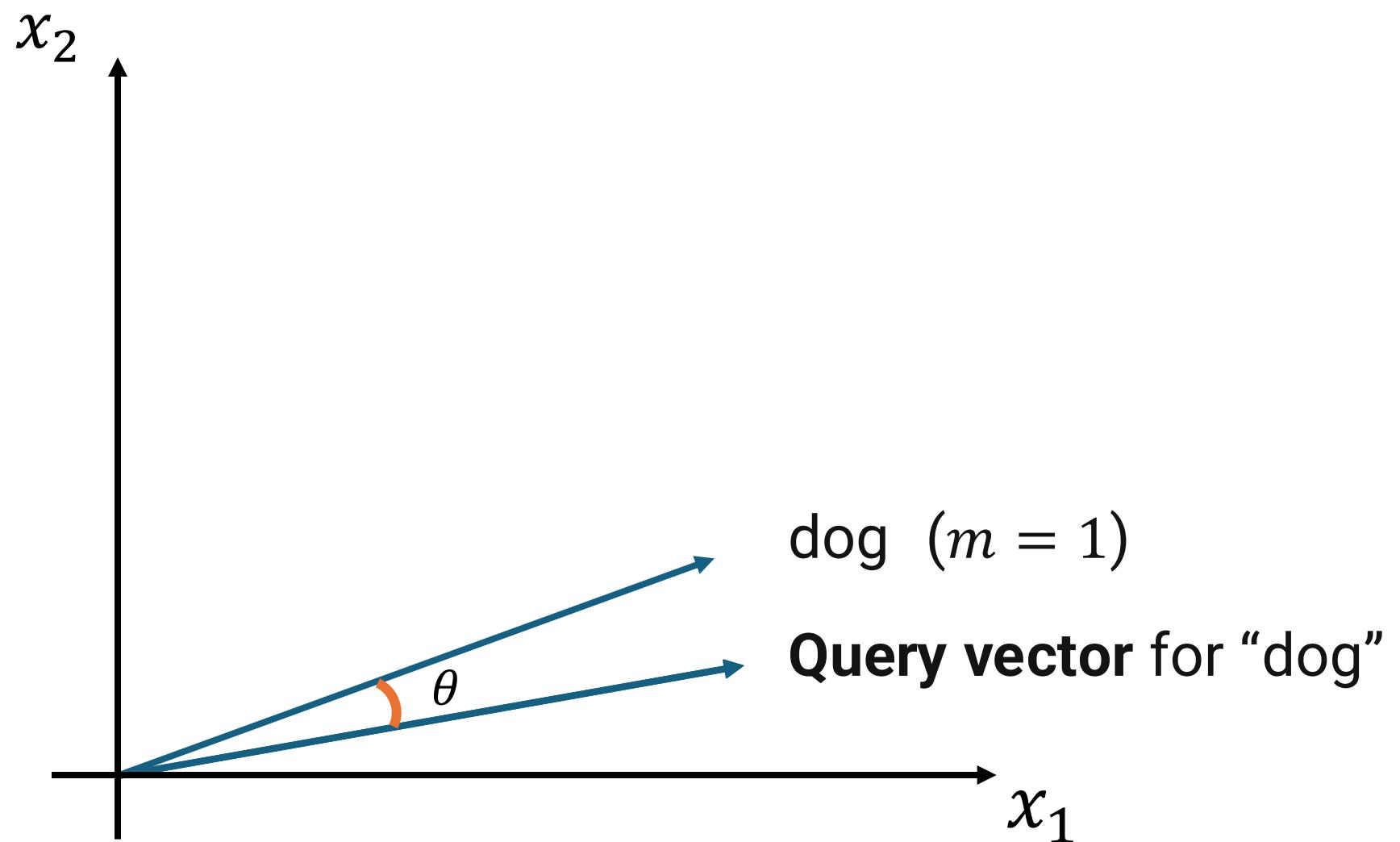
Distance Similarity

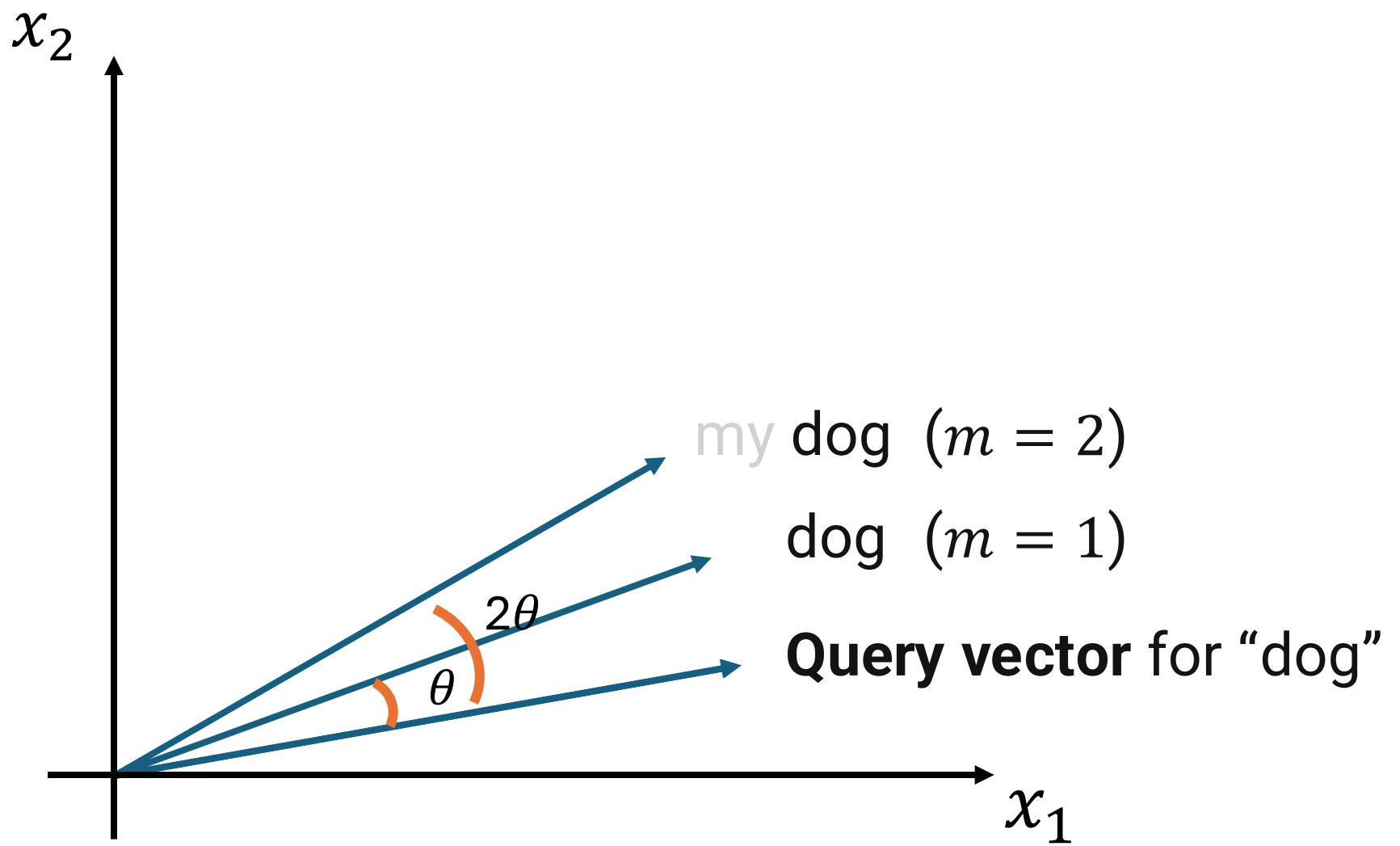
Pixel Similarity

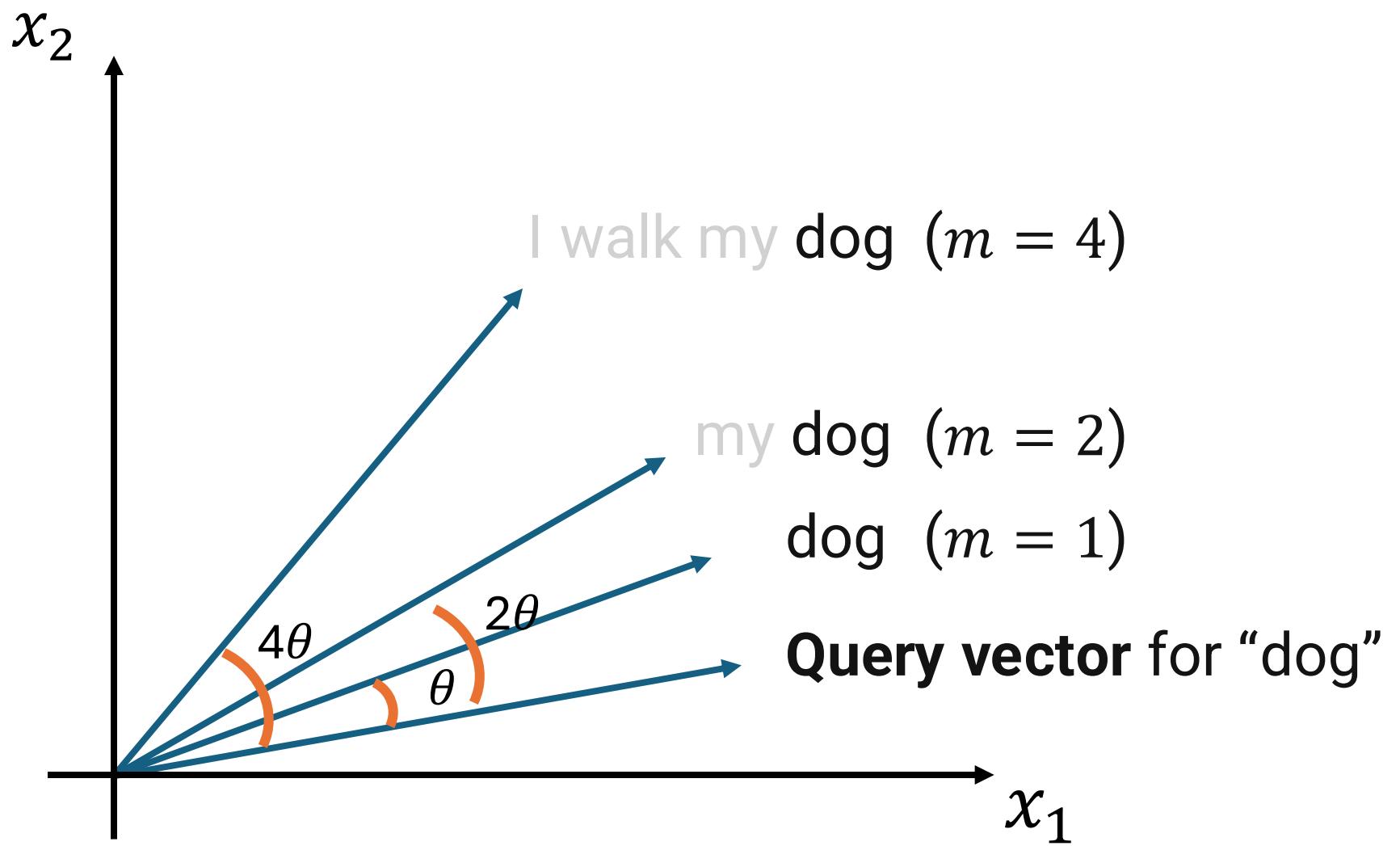
bilateral filter weights of the central pixel

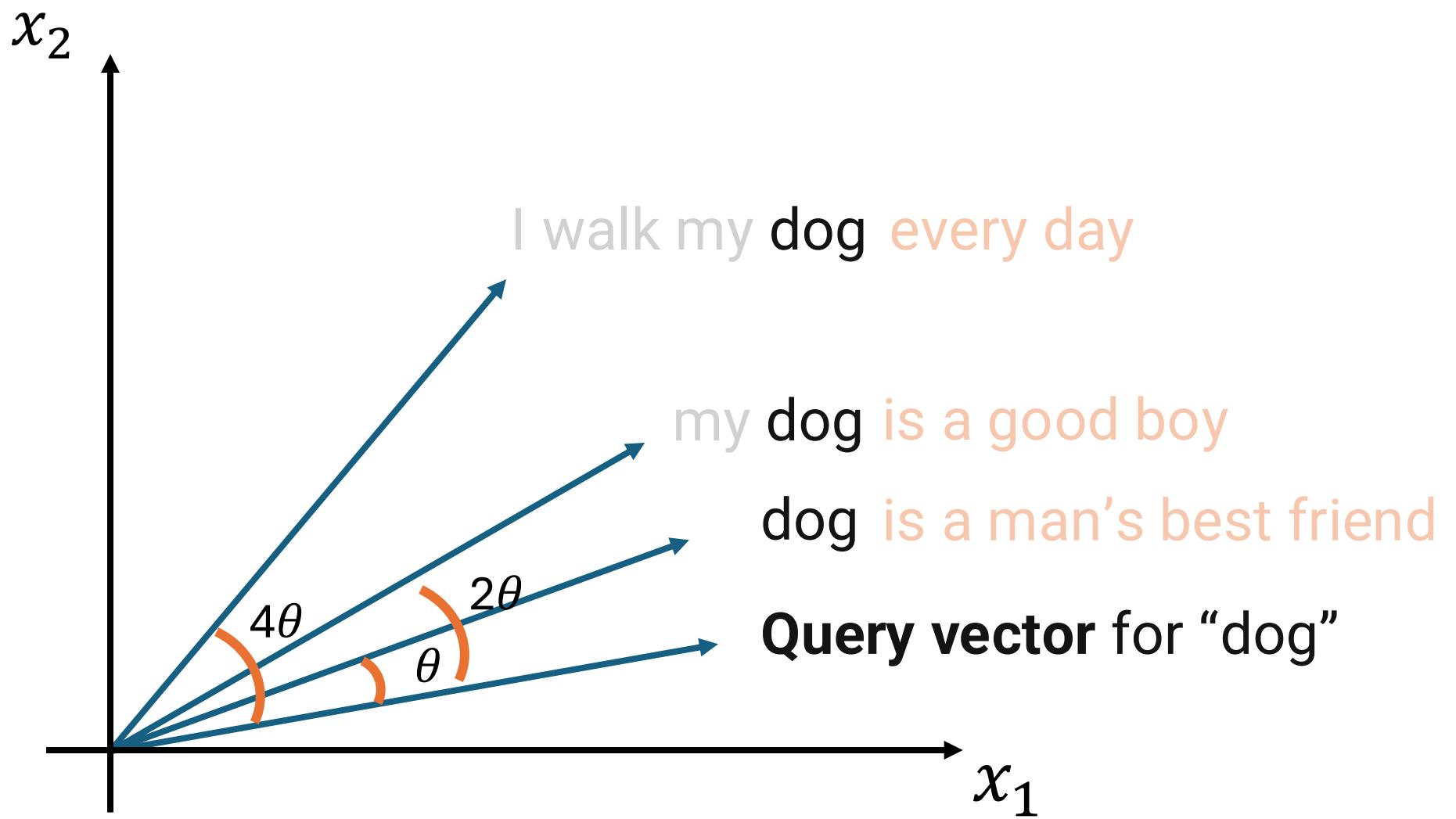


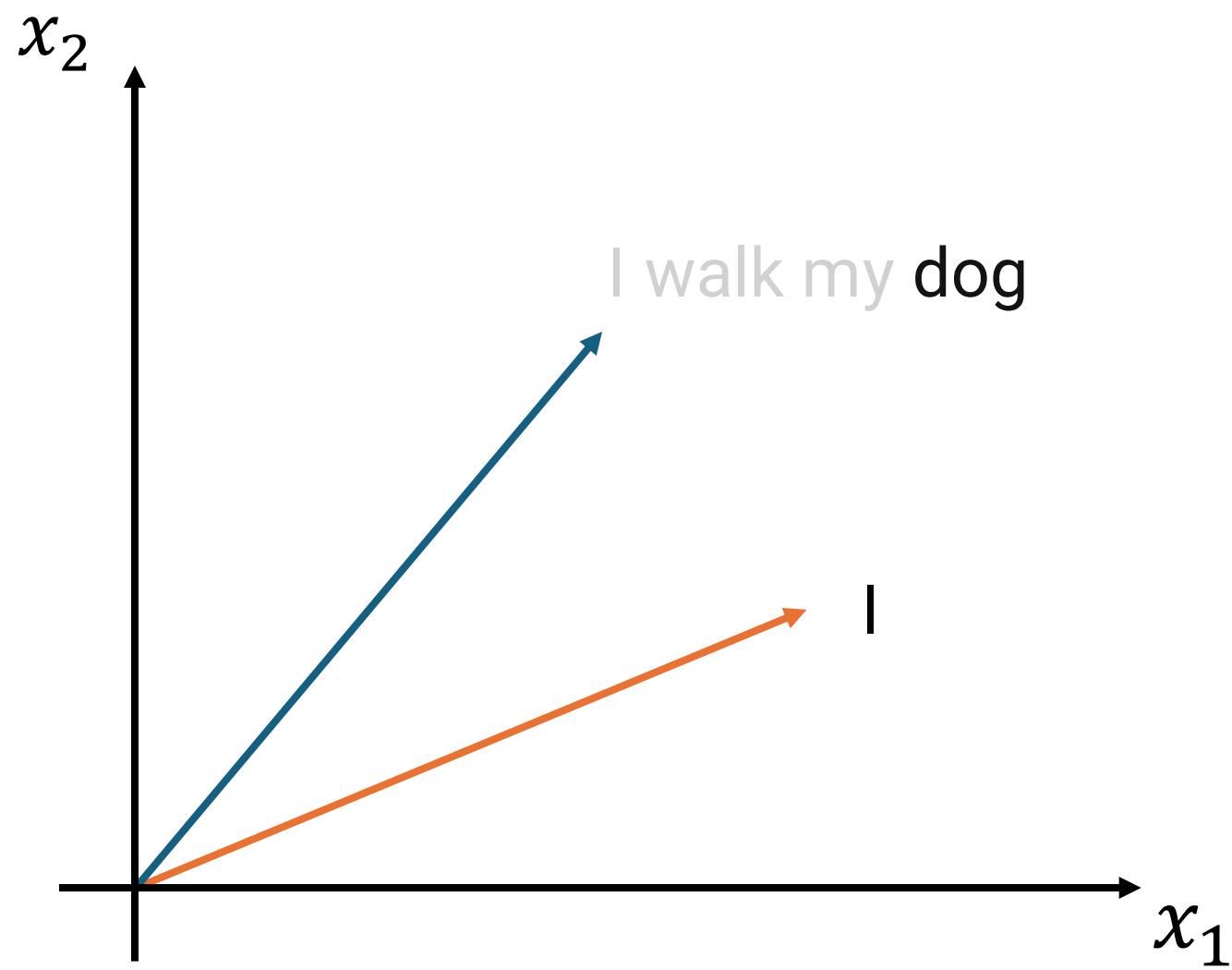


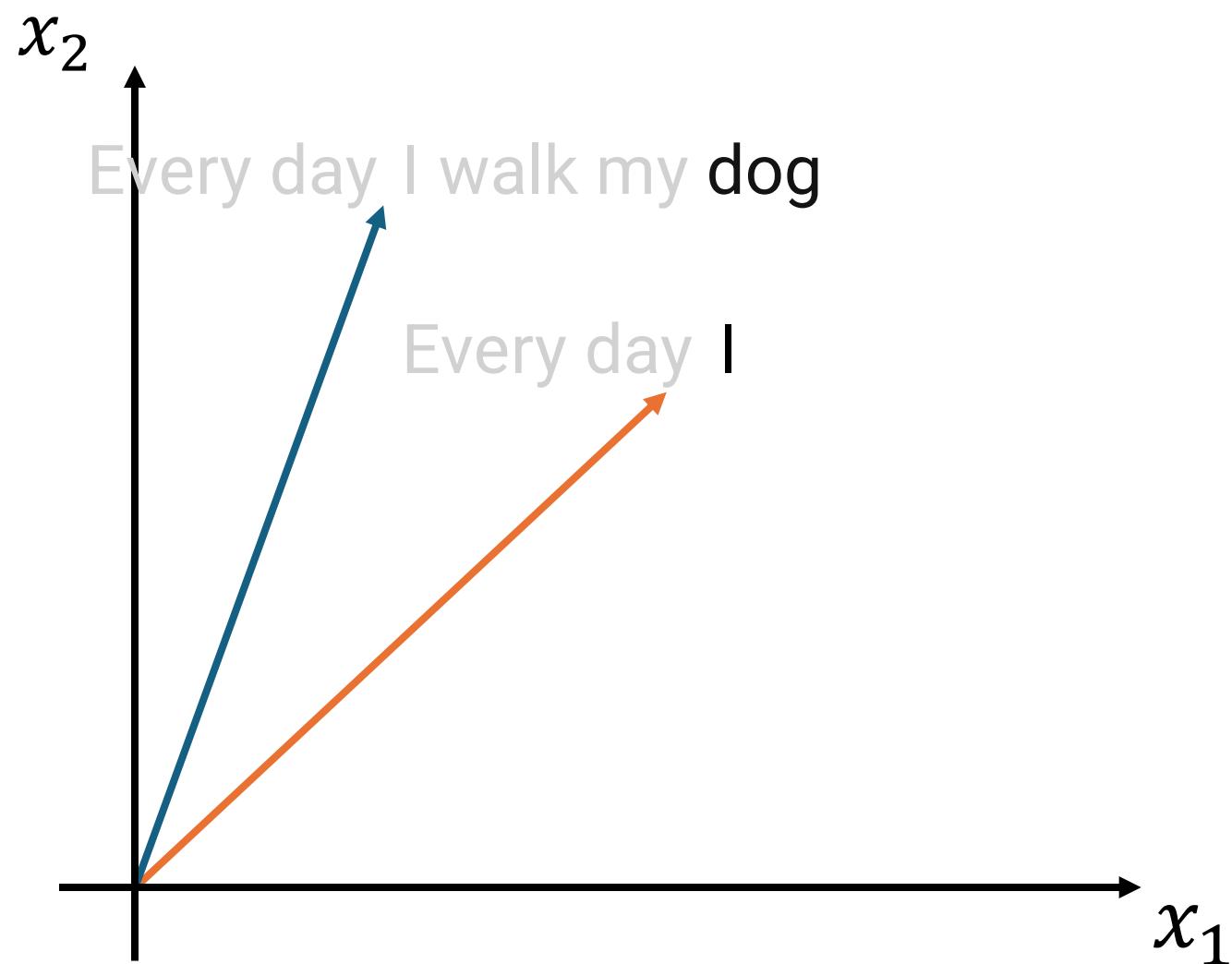












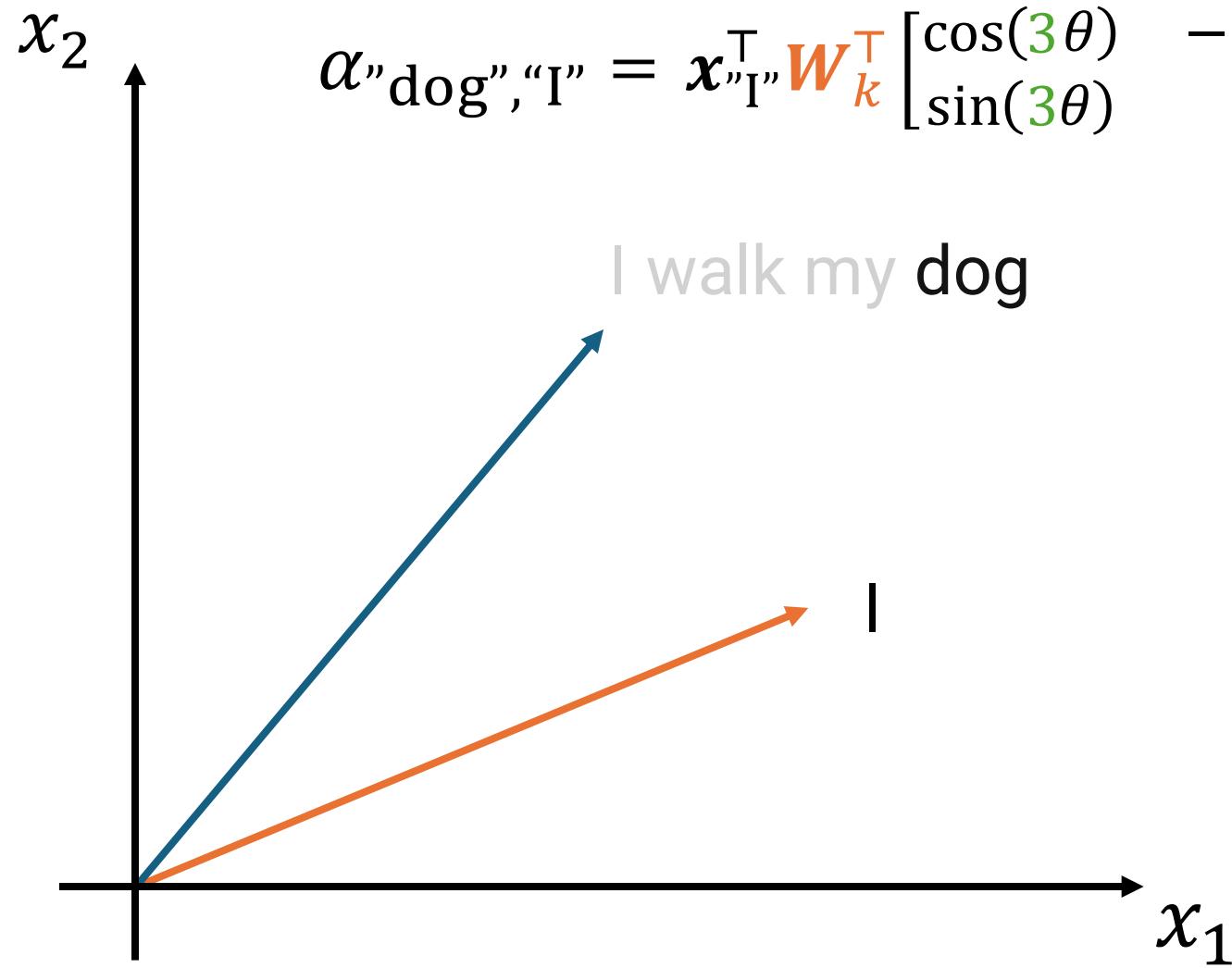
$$q\text{"dog"} = \begin{bmatrix} \cos(4\theta) & -\sin(4\theta) \\ \sin(4\theta) & \cos(4\theta) \end{bmatrix} \mathbf{W}_q \mathbf{x}\text{"dog"}$$

Query vector for "dog"

$$k\text{"I"} = \begin{bmatrix} \cos(1\theta) & -\sin(1\theta) \\ \sin(1\theta) & \cos(1\theta) \end{bmatrix} \mathbf{W}_k \mathbf{x}\text{"I"}$$

Key vector for "I"

$$\alpha\text{"dog", "I"} = \mathbf{x}\text{"I"}^\top \mathbf{W}_k^\top \begin{bmatrix} \cos(3\theta) & -\sin(3\theta) \\ \sin(3\theta) & \cos(3\theta) \end{bmatrix} \mathbf{W}_q \mathbf{x}\text{"dog"}$$



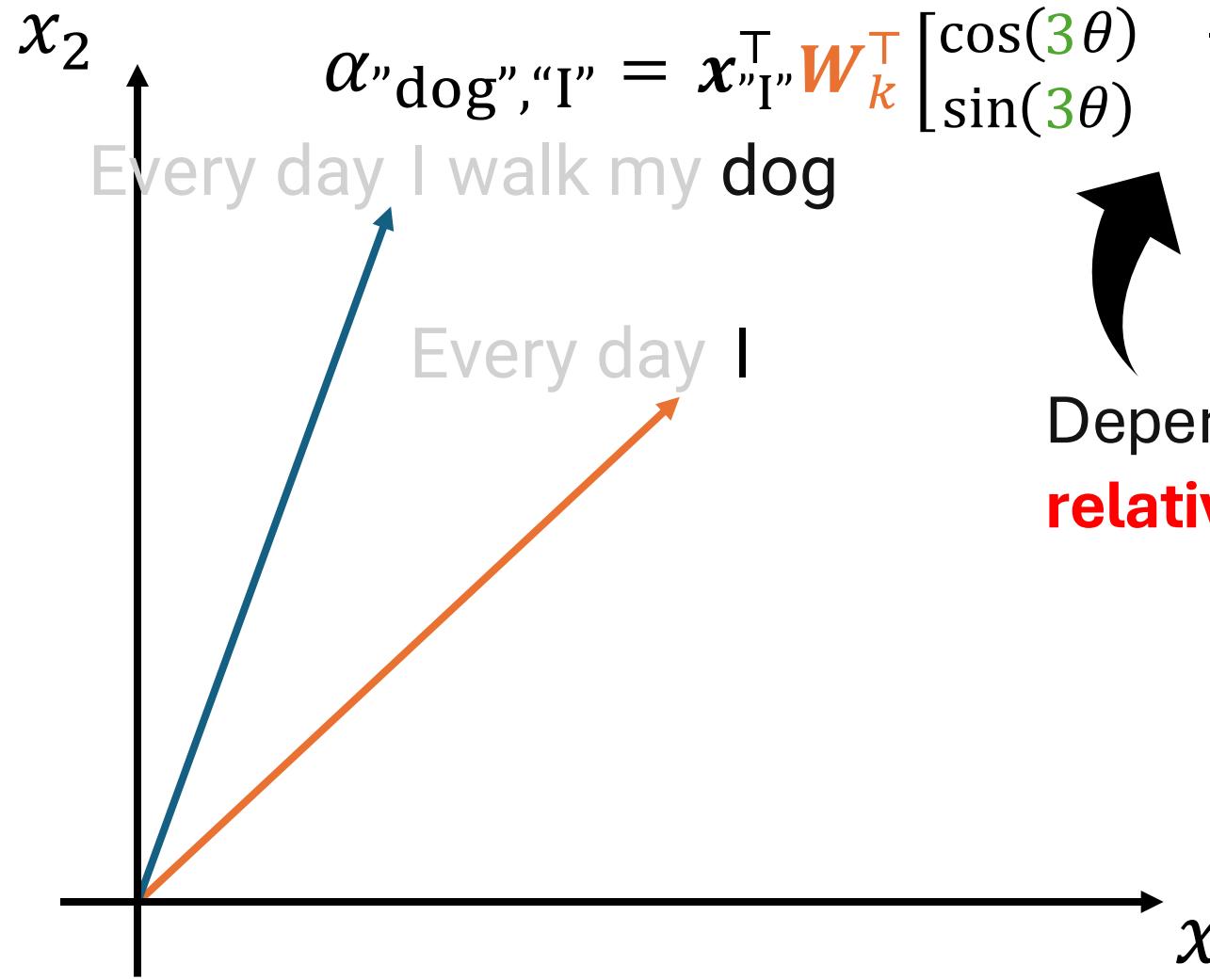
$$q\text{"dog"} = \begin{bmatrix} \cos(6\theta) & -\sin(6\theta) \\ \sin(6\theta) & \cos(6\theta) \end{bmatrix} W_q x\text{"dog"}$$

Query vector for "dog"

$$k\text{"I"} = \begin{bmatrix} \cos(3\theta) & -\sin(3\theta) \\ \sin(3\theta) & \cos(3\theta) \end{bmatrix} W_k x\text{"I"}$$

Key vector for "I"

$$\alpha\text{"dog", "I"} = x\text{"I"}^T W_k^T \begin{bmatrix} \cos(3\theta) & -\sin(3\theta) \\ \sin(3\theta) & \cos(3\theta) \end{bmatrix} W_q x\text{"dog"}$$



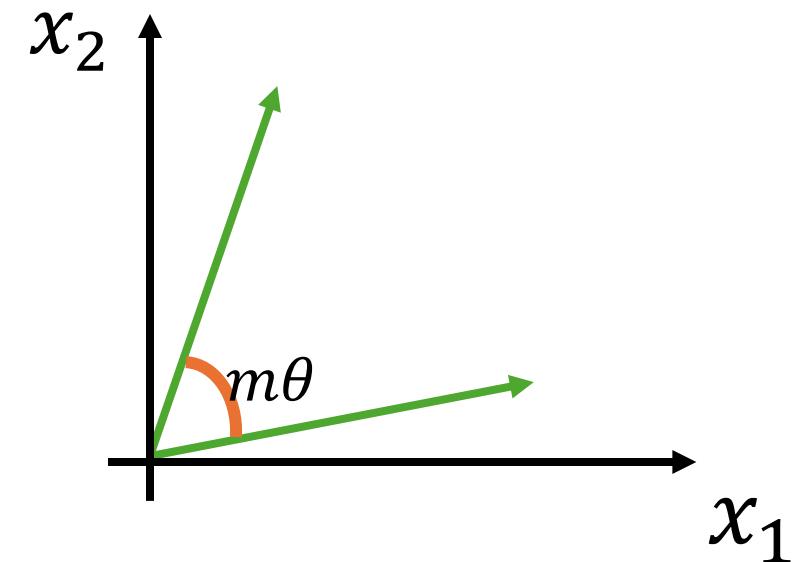
$$R_{\Theta}^m W_q x_m$$

query vector for the token at position  $m$

$$R_{\Theta}^n W_k x_n$$

key vector for the token at position  $n$

$$R_{\Theta}^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$



$$R_{\Theta}^m \mathbf{W}_q \mathbf{x}_m$$

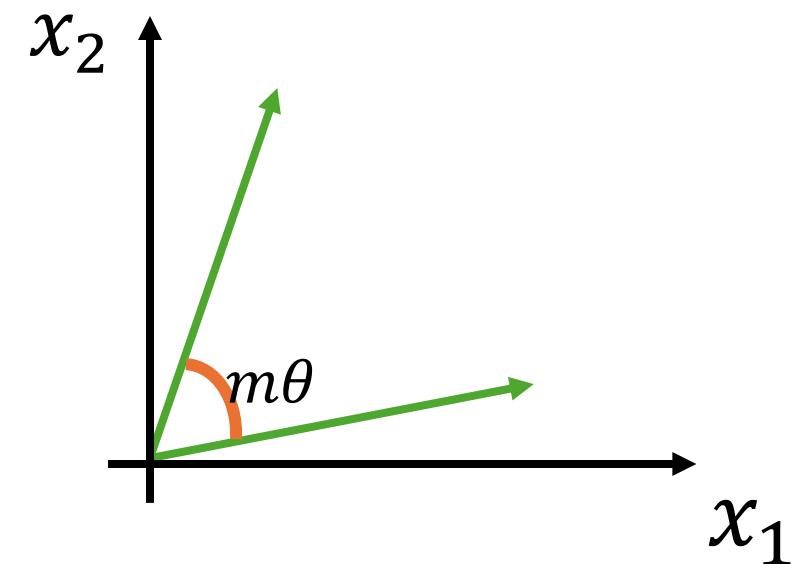
query vector for the token at position  $m$

$$R_{\Theta}^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$

$$R_{\Theta}^n \mathbf{W}_k \mathbf{x}_n$$

key vector for the token at position  $n$

$$\alpha_{m,n} = \mathbf{x}_n^T \mathbf{W}_k^T (R_{\Theta}^n)^T R_{\Theta}^m \mathbf{W}_q \mathbf{x}_m$$



$$R_{\Theta}^m \mathbf{W}_q \mathbf{x}_m$$

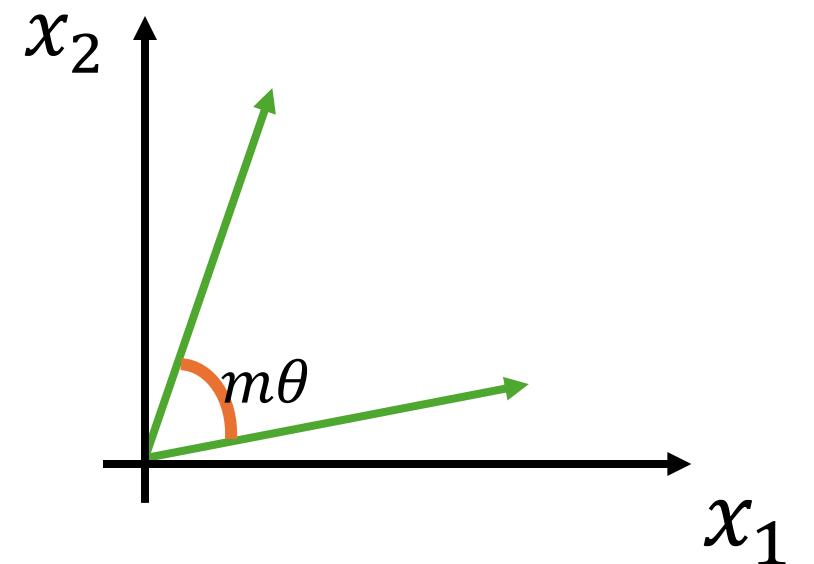
query vector for the token at position  $m$

$$R_{\Theta}^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$

$$R_{\Theta}^n \mathbf{W}_k \mathbf{x}_n$$

key vector for the token at position  $n$

$$\alpha_{m,n} = \mathbf{x}_n^T \mathbf{W}_k^T (R_{\Theta}^n)^T R_{\Theta}^m \mathbf{W}_q \mathbf{x}_m$$



$$R_{\Theta}^m \mathbf{W}_q \mathbf{x}_m$$

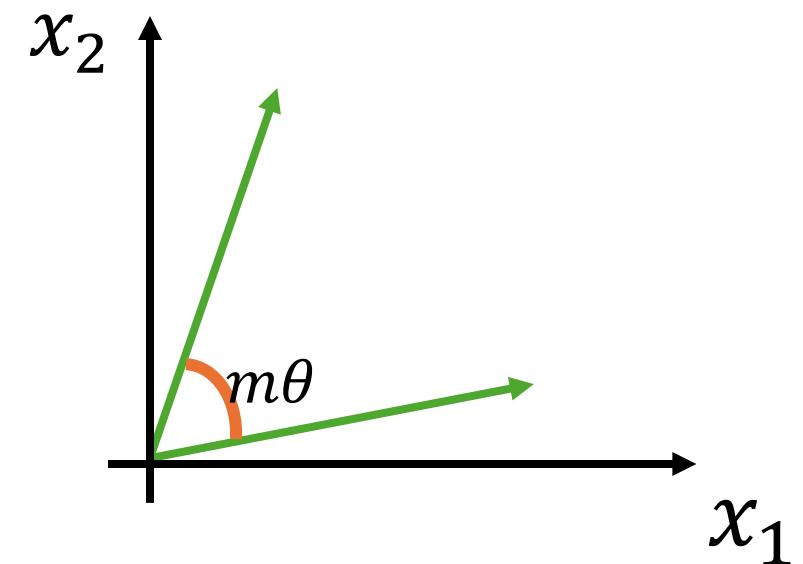
query vector for the token at position  $m$

$$R_{\Theta}^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$

$$R_{\Theta}^n \mathbf{W}_k \mathbf{x}_n$$

key vector for the token at position  $n$

$$\alpha_{m,n} = \mathbf{x}_n^T \mathbf{W}_k^T (R_{\Theta}^n)^T R_{\Theta}^m \mathbf{W}_q \mathbf{x}_m$$



$$R_{\Theta}^m W_q x_m$$

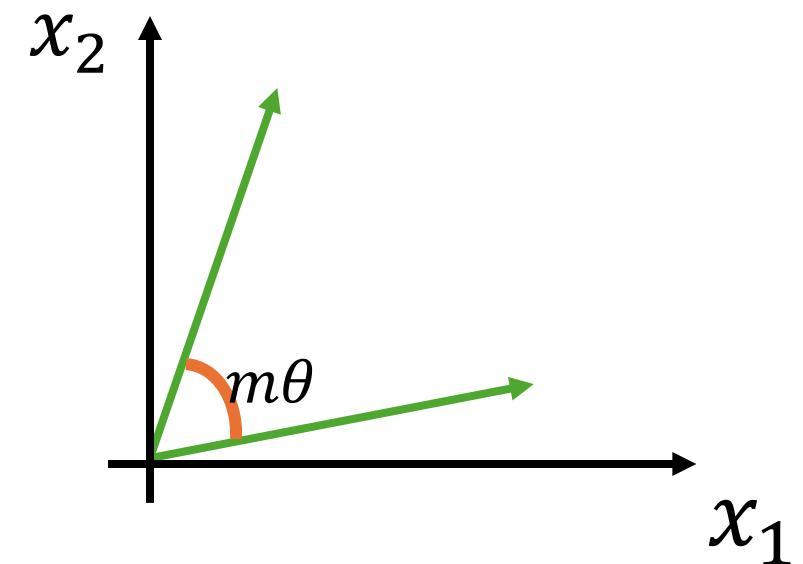
query vector for the token at position  $m$

$$R_{\Theta}^n W_k x_n$$

key vector for the token at position  $n$

$$\alpha_{m,n} = x_n^T W_k^T R_{\Theta}^{-n} R_{\Theta}^m W_q x_m$$

$$R_{\Theta}^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$



$$R_{\Theta}^m \mathbf{W}_q \mathbf{x}_m$$

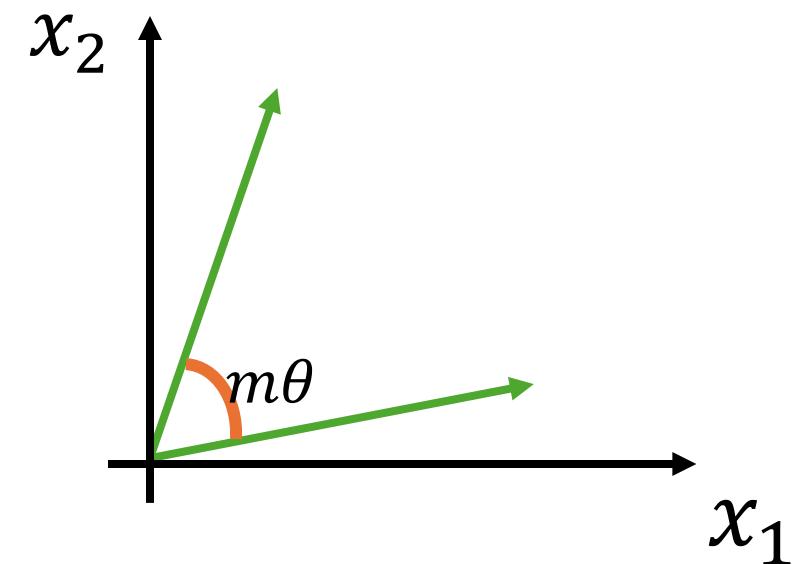
query vector for the token at position  $m$

$$R_{\Theta}^n \mathbf{W}_k \mathbf{x}_n$$

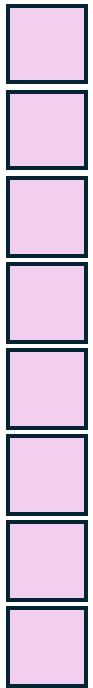
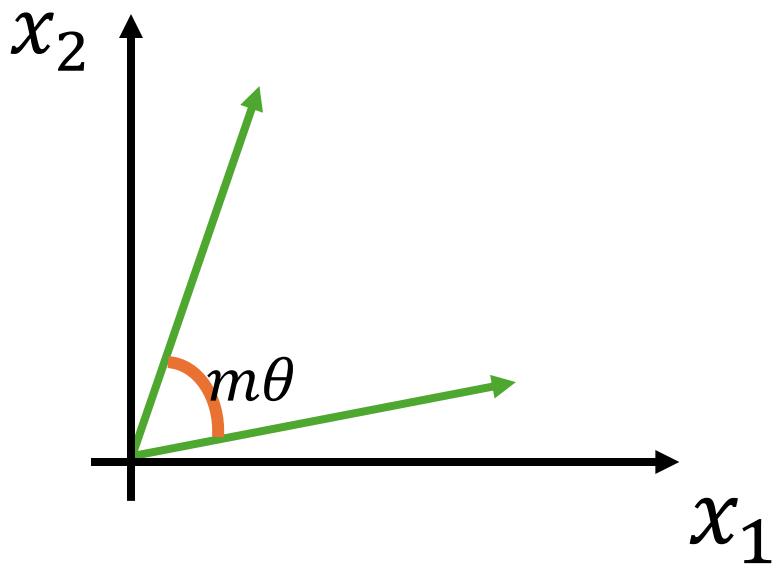
key vector for the token at position  $n$

$$\alpha_{m,n} = \mathbf{x}_n^T \mathbf{W}_k^T R_{\Theta}^{m-n} \mathbf{W}_q \mathbf{x}_m$$

$$R_{\Theta}^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$

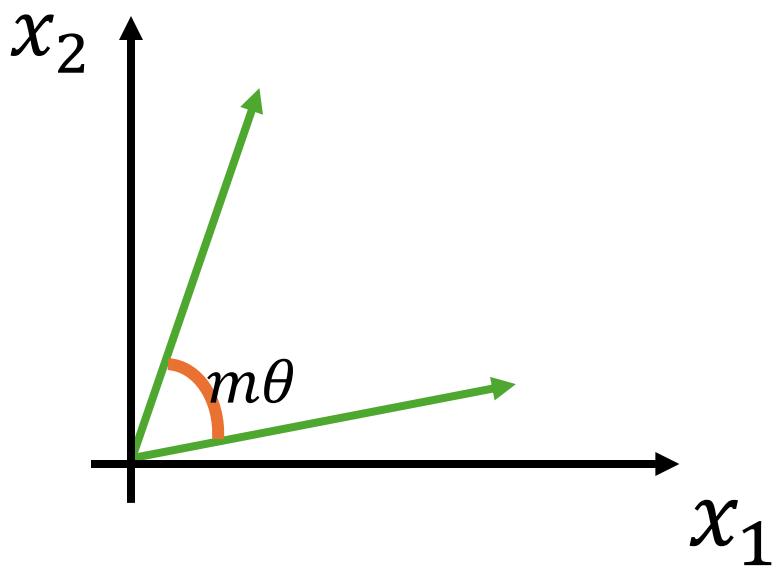


$$R_\Theta^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$



$W_q x_m$

$$R_\Theta^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$



$$\begin{array}{c|c} \text{[green square]} & = \begin{bmatrix} \cos(m\theta_1) & -\sin(m\theta_1) \\ \sin(m\theta_1) & \cos(m\theta_1) \end{bmatrix} \begin{array}{c|c} \text{[blue square]} \\ \text{[blue square]} \end{array} \\ \hline \end{array}$$

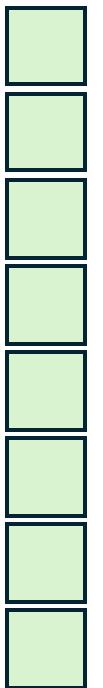
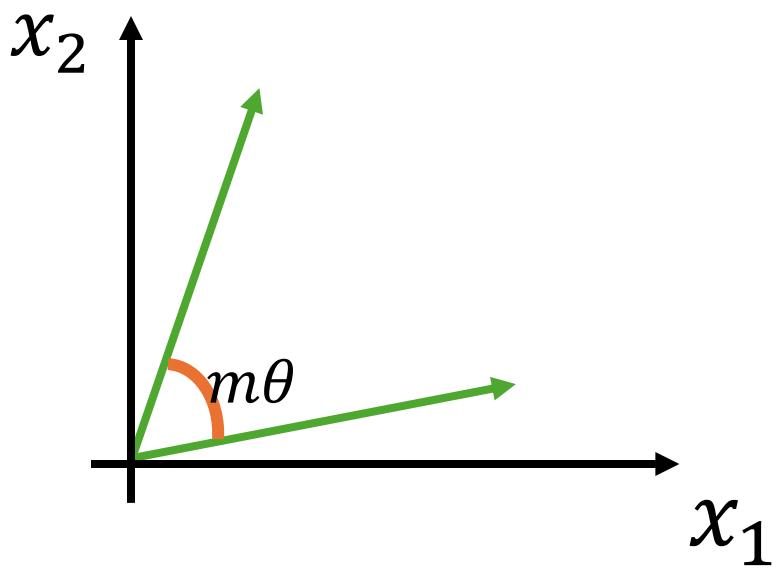
$$\begin{array}{c|c} \text{[green square]} & = \begin{bmatrix} \cos(m\theta_2) & -\sin(m\theta_2) \\ \sin(m\theta_2) & \cos(m\theta_2) \end{bmatrix} \begin{array}{c|c} \text{[blue square]} \\ \text{[blue square]} \end{array} \\ \hline \end{array}$$

$$\begin{array}{c|c} \text{[green square]} & = \begin{bmatrix} \cos(m\theta_3) & -\sin(m\theta_3) \\ \sin(m\theta_3) & \cos(m\theta_3) \end{bmatrix} \begin{array}{c|c} \text{[blue square]} \\ \text{[blue square]} \end{array} \\ \hline \end{array}$$

$$\begin{array}{c|c} \text{[green square]} & = \begin{bmatrix} \cos(m\theta_4) & -\sin(m\theta_4) \\ \sin(m\theta_4) & \cos(m\theta_4) \end{bmatrix} \begin{array}{c|c} \text{[blue square]} \\ \text{[blue square]} \end{array} \\ \hline \end{array}$$

$W_q x_m$

$$R_\Theta^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$

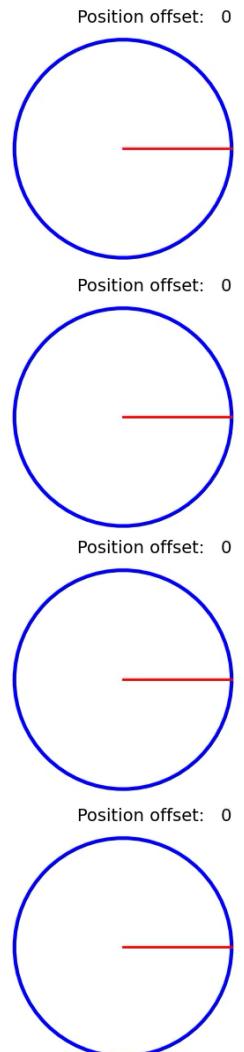
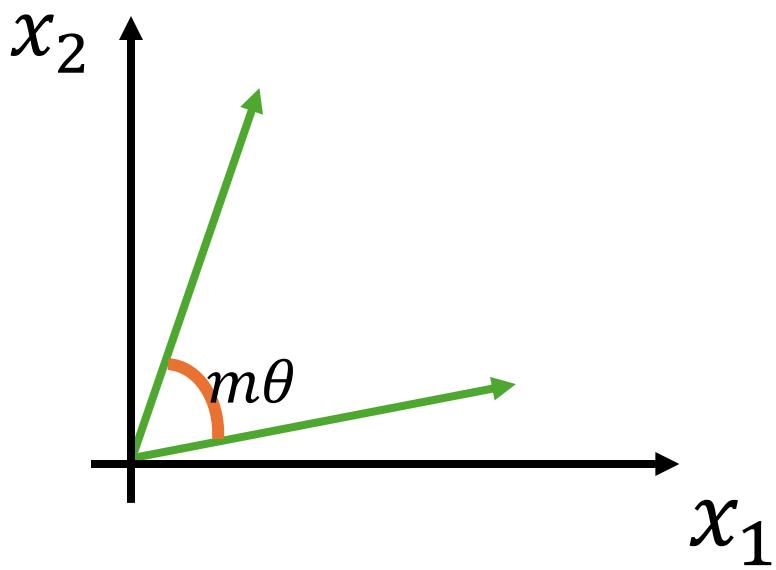


$$f_q(\mathbf{x}_m, m)$$

$$\begin{array}{c|c} \begin{array}{c} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{array} & = \begin{bmatrix} \cos(m\theta_1) & -\sin(m\theta_1) \\ \sin(m\theta_1) & \cos(m\theta_1) \end{bmatrix} \begin{array}{c} \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} \end{array} \\ \hline \begin{array}{c} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{array} & = \begin{bmatrix} \cos(m\theta_2) & -\sin(m\theta_2) \\ \sin(m\theta_2) & \cos(m\theta_2) \end{bmatrix} \begin{array}{c} \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} \end{array} \\ \hline \begin{array}{c} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{array} & = \begin{bmatrix} \cos(m\theta_3) & -\sin(m\theta_3) \\ \sin(m\theta_3) & \cos(m\theta_3) \end{bmatrix} \begin{array}{c} \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} \end{array} \\ \hline \begin{array}{c} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{array} & = \begin{bmatrix} \cos(m\theta_4) & -\sin(m\theta_4) \\ \sin(m\theta_4) & \cos(m\theta_4) \end{bmatrix} \begin{array}{c} \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} \end{array} \end{array}$$

$$\mathbf{W}_q \mathbf{x}_m$$

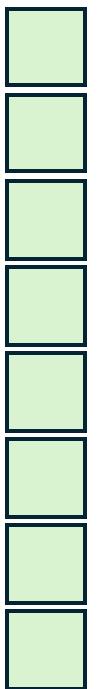
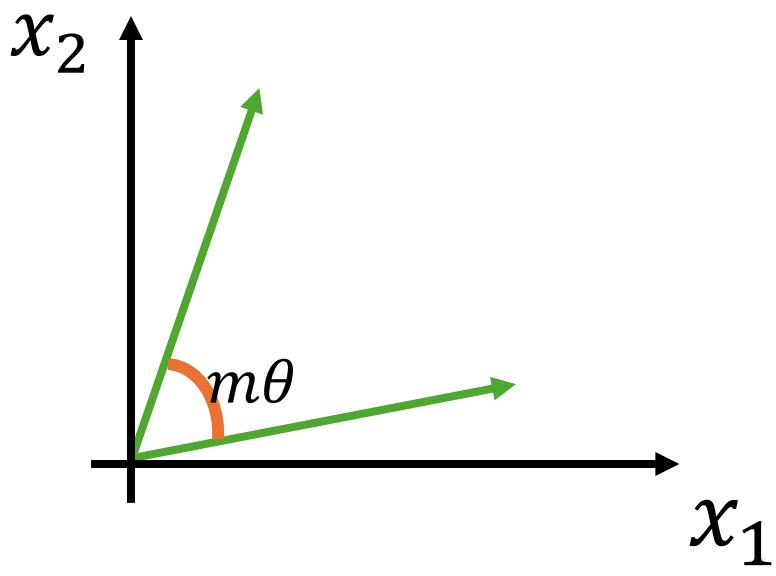
$$R_\Theta^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$



$$\begin{array}{lcl} \begin{array}{c} \text{---} \\ \text{---} \end{array} & = & \begin{bmatrix} \cos(m\theta_1) & -\sin(m\theta_1) \\ \sin(m\theta_1) & \cos(m\theta_1) \end{bmatrix} \begin{array}{c} \text{---} \\ \text{---} \end{array} \\ \begin{array}{c} \text{---} \\ \text{---} \end{array} & = & \begin{bmatrix} \cos(m\theta_2) & -\sin(m\theta_2) \\ \sin(m\theta_2) & \cos(m\theta_2) \end{bmatrix} \begin{array}{c} \text{---} \\ \text{---} \end{array} \\ \begin{array}{c} \text{---} \\ \text{---} \end{array} & = & \begin{bmatrix} \cos(m\theta_3) & -\sin(m\theta_3) \\ \sin(m\theta_3) & \cos(m\theta_3) \end{bmatrix} \begin{array}{c} \text{---} \\ \text{---} \end{array} \\ \begin{array}{c} \text{---} \\ \text{---} \end{array} & = & \begin{bmatrix} \cos(m\theta_4) & -\sin(m\theta_4) \\ \sin(m\theta_4) & \cos(m\theta_4) \end{bmatrix} \begin{array}{c} \text{---} \\ \text{---} \end{array} \end{array}$$

$W_q x_m$

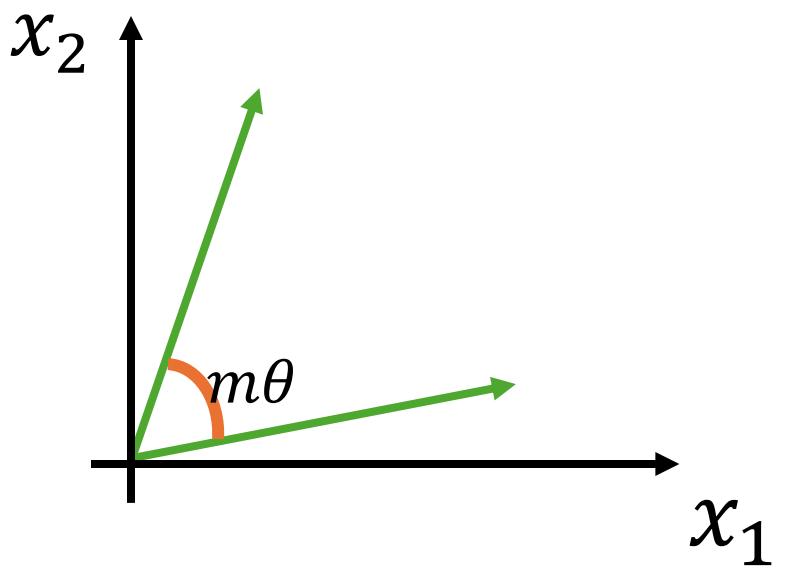
$$R_\Theta^m = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix}$$



$$\begin{array}{lcl} \begin{array}{c} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{array} & = & \begin{bmatrix} \cos(m\theta_1) & -\sin(m\theta_1) \\ \sin(m\theta_1) & \cos(m\theta_1) \end{bmatrix} \begin{array}{c} \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} \end{array} \\ \begin{array}{c} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{array} & = & \begin{bmatrix} \cos(m\theta_2) & -\sin(m\theta_2) \\ \sin(m\theta_2) & \cos(m\theta_2) \end{bmatrix} \begin{array}{c} \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} \end{array} \\ \begin{array}{c} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{array} & = & \begin{bmatrix} \cos(m\theta_3) & -\sin(m\theta_3) \\ \sin(m\theta_3) & \cos(m\theta_3) \end{bmatrix} \begin{array}{c} \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} \end{array} \\ \begin{array}{c} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{array} & = & \begin{bmatrix} \cos(m\theta_4) & -\sin(m\theta_4) \\ \sin(m\theta_4) & \cos(m\theta_4) \end{bmatrix} \begin{array}{c} \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} \end{array} \end{array}$$

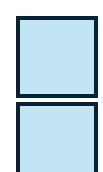
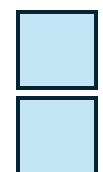
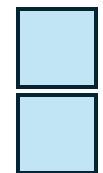
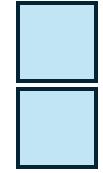
$$f_q(\mathbf{x}_m, m)$$

$$\mathbf{W}_q \mathbf{x}_m$$



$$f_q(\mathbf{x}_m, m)$$

$$\mathbf{W}_q \mathbf{x}_m$$

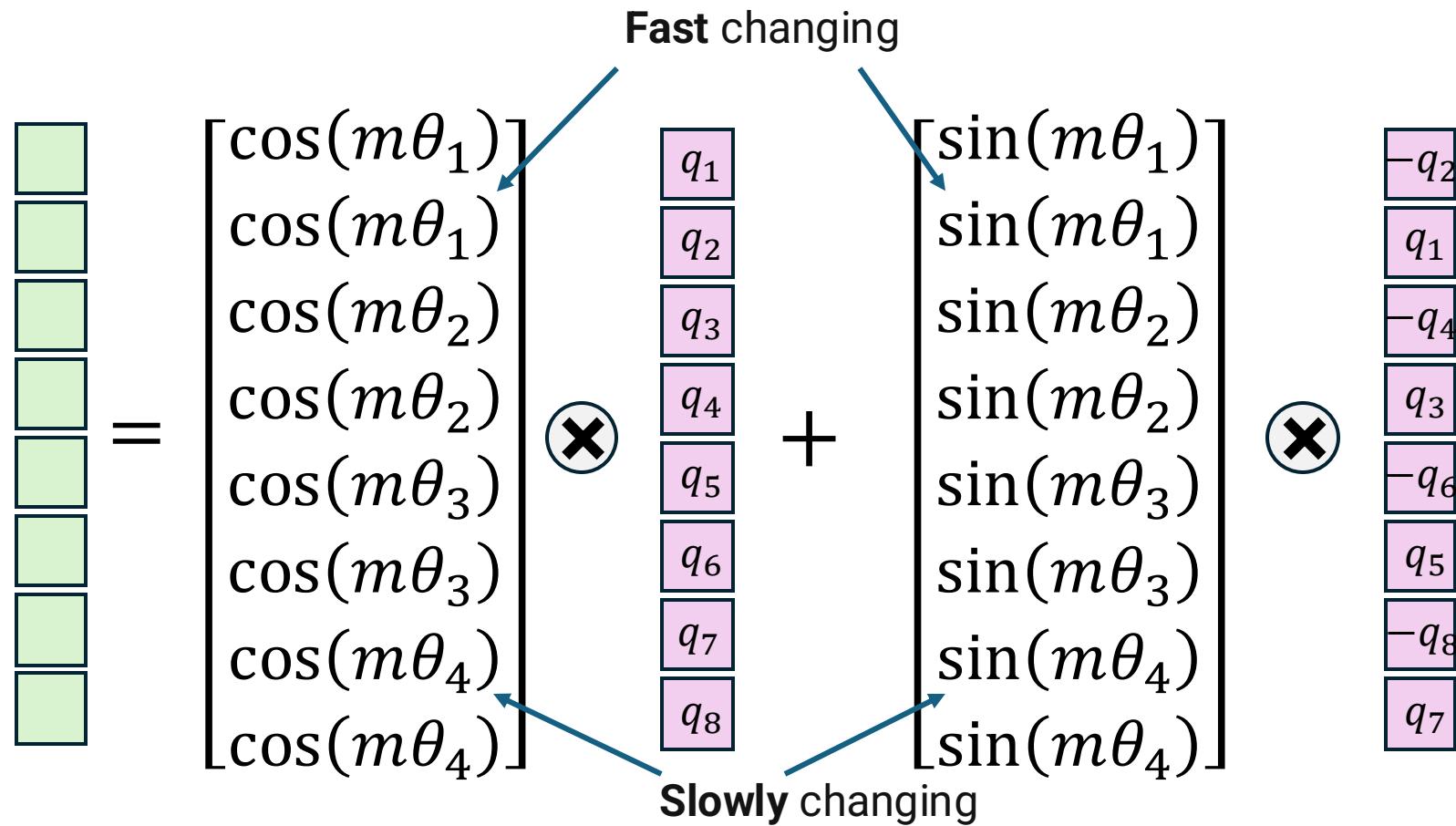


$$\begin{array}{c}
\begin{array}{|c|c|} \hline & \\ \hline \end{array} = \left( \begin{array}{ccccccc}
\cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\
\sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\
0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\
0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\
0 & 0 & 0 & 0 & \cdots & \cos m\theta_l & -\sin m\theta_l \\
0 & 0 & 0 & 0 & \cdots & \sin m\theta_l & \cos m\theta_l
\end{array} \right) \begin{array}{|c|c|} \hline q_1 \\ \hline q_2 \\ \hline q_3 \\ \hline q_4 \\ \hline q_5 \\ \hline q_6 \\ \hline \vdots \\ \hline q_8
\end{array}
\end{array}$$

$R_{\Theta}^m$

$f_q(x_m, m)$

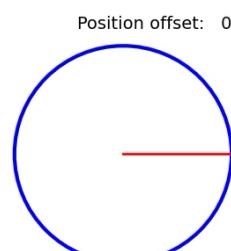
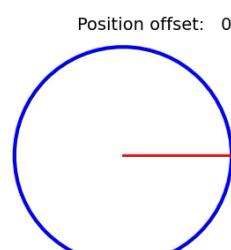
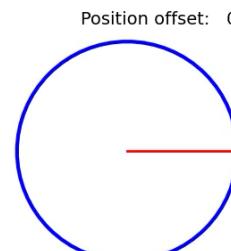
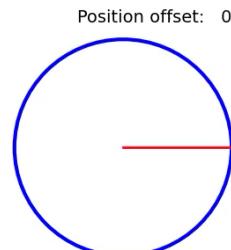
$W_q x_m$



$$f_q(x_m, m)$$

$$\theta_i = N^{-2i/d}$$

$$N = 10,000$$



$$\begin{bmatrix} \text{green square} \\ \text{green square} \end{bmatrix} = \begin{bmatrix} \cos(m\theta_1) & -\sin(m\theta_1) \\ \sin(m\theta_1) & \cos(m\theta_1) \end{bmatrix} \begin{bmatrix} \text{blue square} \\ \text{blue square} \end{bmatrix}$$

$$\begin{bmatrix} \text{green square} \\ \text{green square} \end{bmatrix} = \begin{bmatrix} \cos(m\theta_2) & -\sin(m\theta_2) \\ \sin(m\theta_2) & \cos(m\theta_2) \end{bmatrix} \begin{bmatrix} \text{blue square} \\ \text{blue square} \end{bmatrix}$$

$$\begin{bmatrix} \text{green square} \\ \text{green square} \end{bmatrix} = \begin{bmatrix} \cos(m\theta_3) & -\sin(m\theta_3) \\ \sin(m\theta_3) & \cos(m\theta_3) \end{bmatrix} \begin{bmatrix} \text{blue square} \\ \text{blue square} \end{bmatrix}$$

$$\begin{bmatrix} \text{green square} \\ \text{green square} \end{bmatrix} = \begin{bmatrix} \cos(m\theta_4) & -\sin(m\theta_4) \\ \sin(m\theta_4) & \cos(m\theta_4) \end{bmatrix} \begin{bmatrix} \text{blue square} \\ \text{blue square} \end{bmatrix}$$

$f_q(\mathbf{x}_m, m)$

$\mathbf{W}_q \mathbf{x}_m$

# 2D ROPE



Patch (m, n)

$$\begin{bmatrix} \cos(\theta_1 x) & -\sin(\theta_1 x) & 0 & 0 & \dots & 0 & 0 \\ \sin(\theta_1 x) & \cos(\theta_1 x) & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \cos(\theta_2 x) & -\sin(\theta_2 x) & \dots & 0 & 0 \\ 0 & 0 & \sin(\theta_2 x) & \cos(\theta_2 x) & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \cos(\theta_l x) & -\sin(\theta_l x) \\ 0 & 0 & 0 & 0 & \dots & \sin(\theta_l x) & \cos(\theta_l x) \end{bmatrix}$$

$R_x$

$$\begin{bmatrix} \cos(\theta_1 y) & -\sin(\theta_1 y) & 0 & 0 & \dots & 0 & 0 \\ \sin(\theta_1 y) & \cos(\theta_1 y) & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \cos(\theta_2 y) & -\sin(\theta_2 y) & \dots & 0 & 0 \\ 0 & 0 & \sin(\theta_2 y) & \cos(\theta_2 y) & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \cos(\theta_l y) & -\sin(\theta_l y) \\ 0 & 0 & 0 & 0 & \dots & \sin(\theta_l y) & \cos(\theta_l y) \end{bmatrix}$$

$R_y$

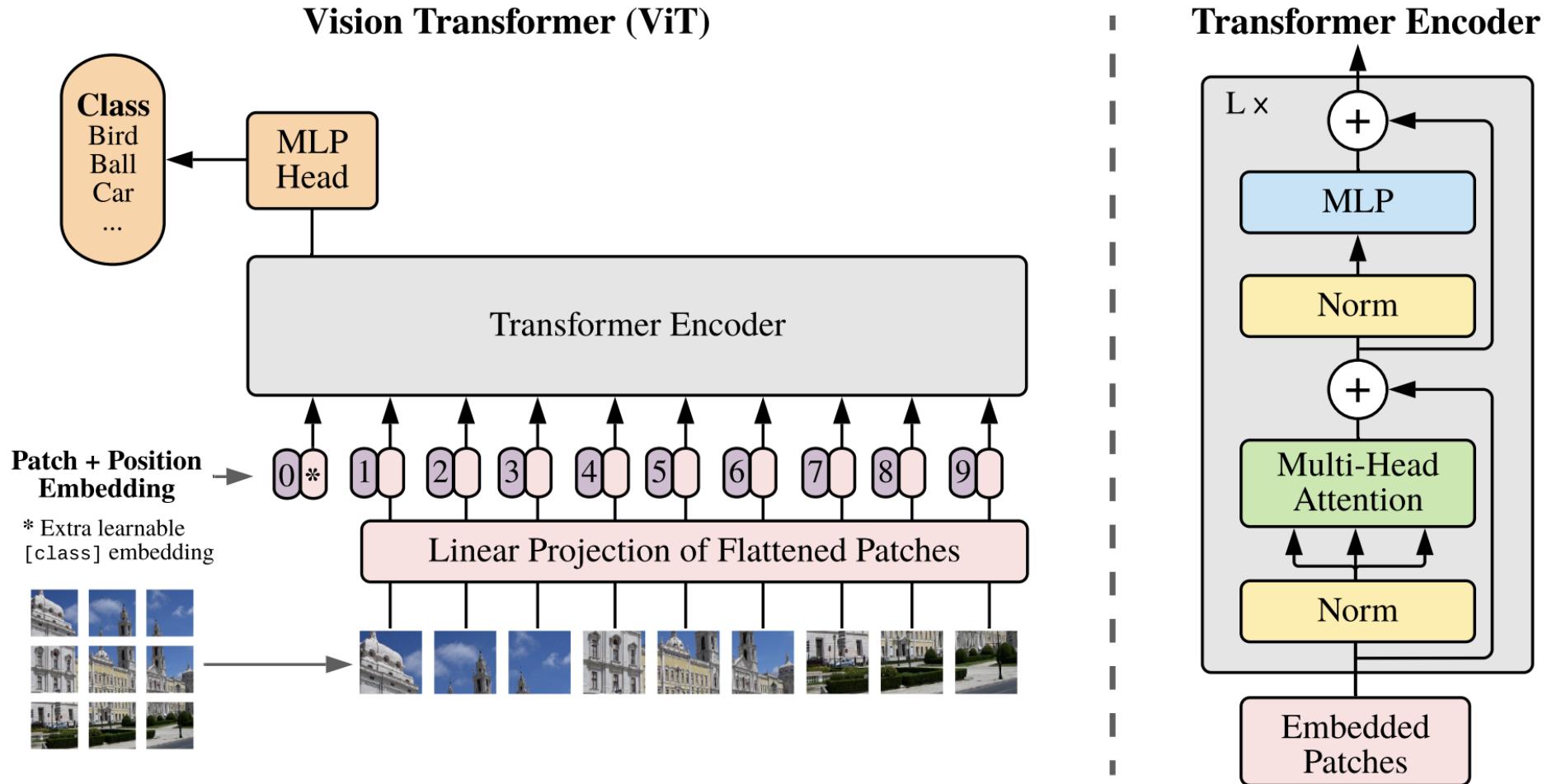
Rope on X  
values

Rope on Y  
values

# Switching Gears

Using ViT!

# Vision Transformer



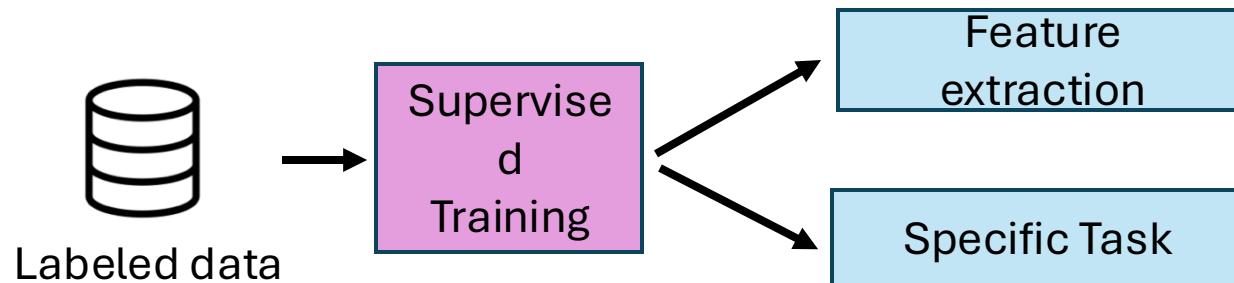
# Is semantic supervision necessary to learn good representations?

- Manual labeling doesn't scale, suffers from biases
- Plenty of unlabeled visual data already, and growing really fast
- And subject of the *Gelato Bet*:
- *Made on Sept 23 2014*
- *If, by the first day of autumn (Sept 23) of 2015, a method will exist that can match or beat the performance of R-CNN on Pascal VOC detection, without the use of any extra, human annotations (e.g. ImageNet) as pre-training, Mr. Malik promises to buy Mr. Efros one (1) gelato (2 scoops: one chocolate, one vanilla).*

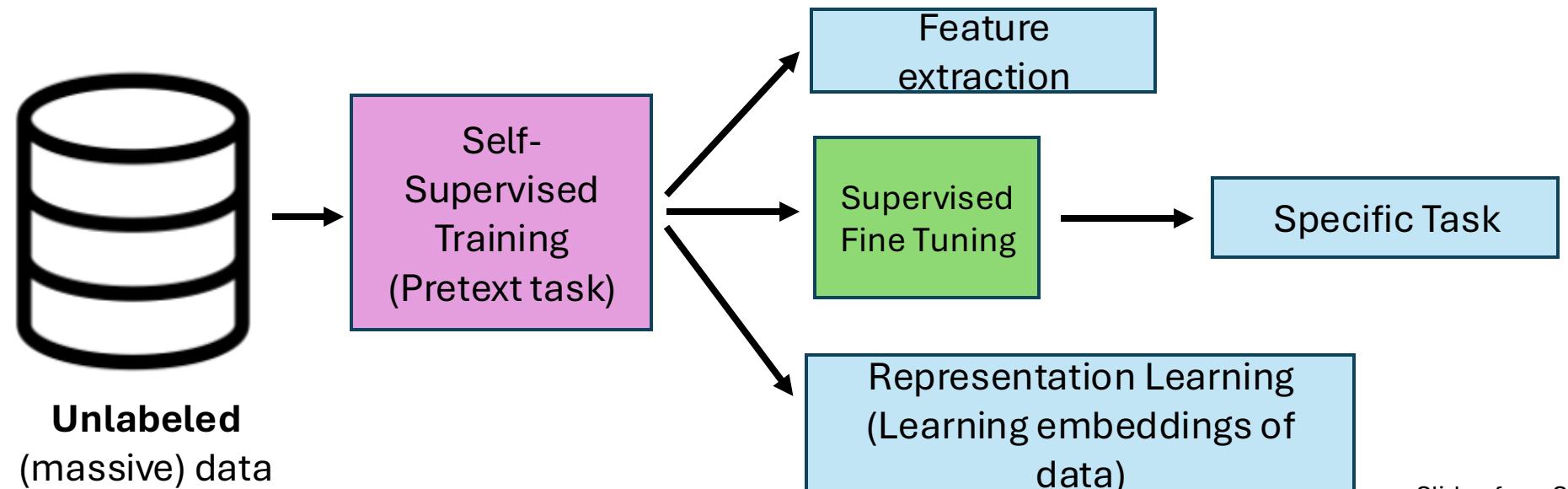


# Self supervised learning with ViTs

Supervised Learning



Self-Supervised Learning



# Pre-train representations on a pre-text task

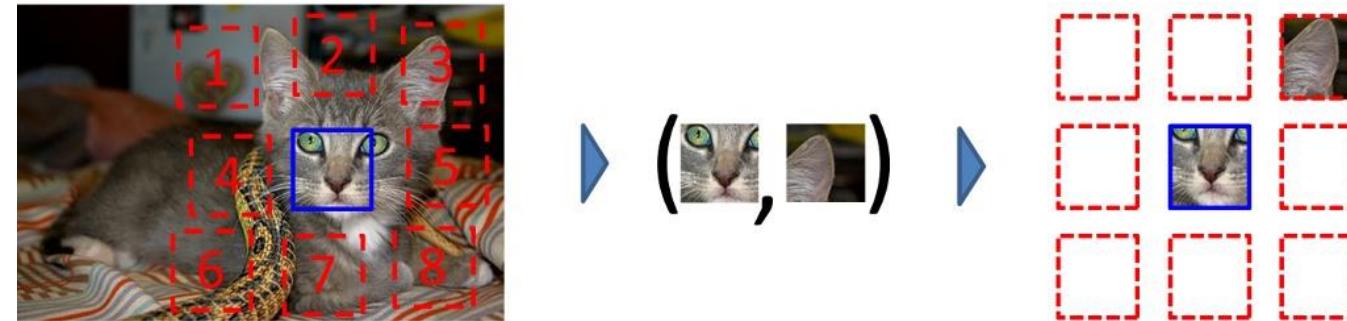
E.g. Colorization



After pre-training, use representation for down-stream tasks.

Many other possibilities,

- Spatial relationship between pair of patches



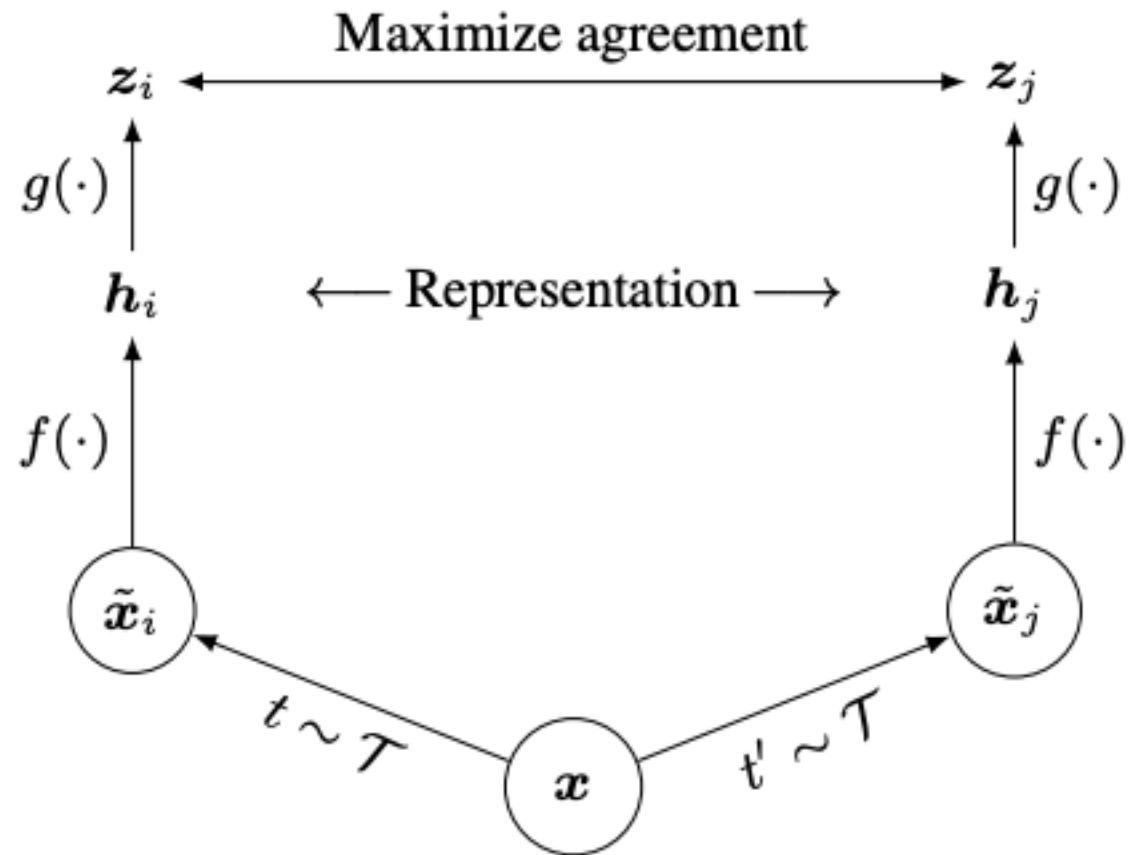
- Predict sound / frame ordering in a video
- Encourage two augmentations of same image to be closer to each other than to another image
- Predict hidden image patches from context

[Split-Brain Autoencoders: Unsupervised Learning by Cross-Channel Prediction](#), Zhang et al. CVPR 2017  
[Context as Supervisory Signal: Discovering Objects with Predictable Context](#), Doersch et al. ICCV 2015

# Contrastive Learning

- Encourage two augmentations of an image to be close.
- Using a contrastive loss:

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$



[A Simple Framework for Contrastive Learning of Visual Representations](#), Chen et al. ICML 2020

See also: [Momentum Contrast for Unsupervised Visual Representation Learning](#), He et al. CVPR 2020

# Augmentations



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)

(f) Rotate  $\{90^\circ, 180^\circ, 270^\circ\}$ 

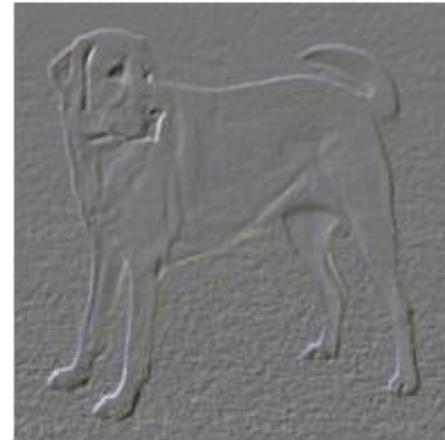
(g) Cutout



(h) Gaussian noise

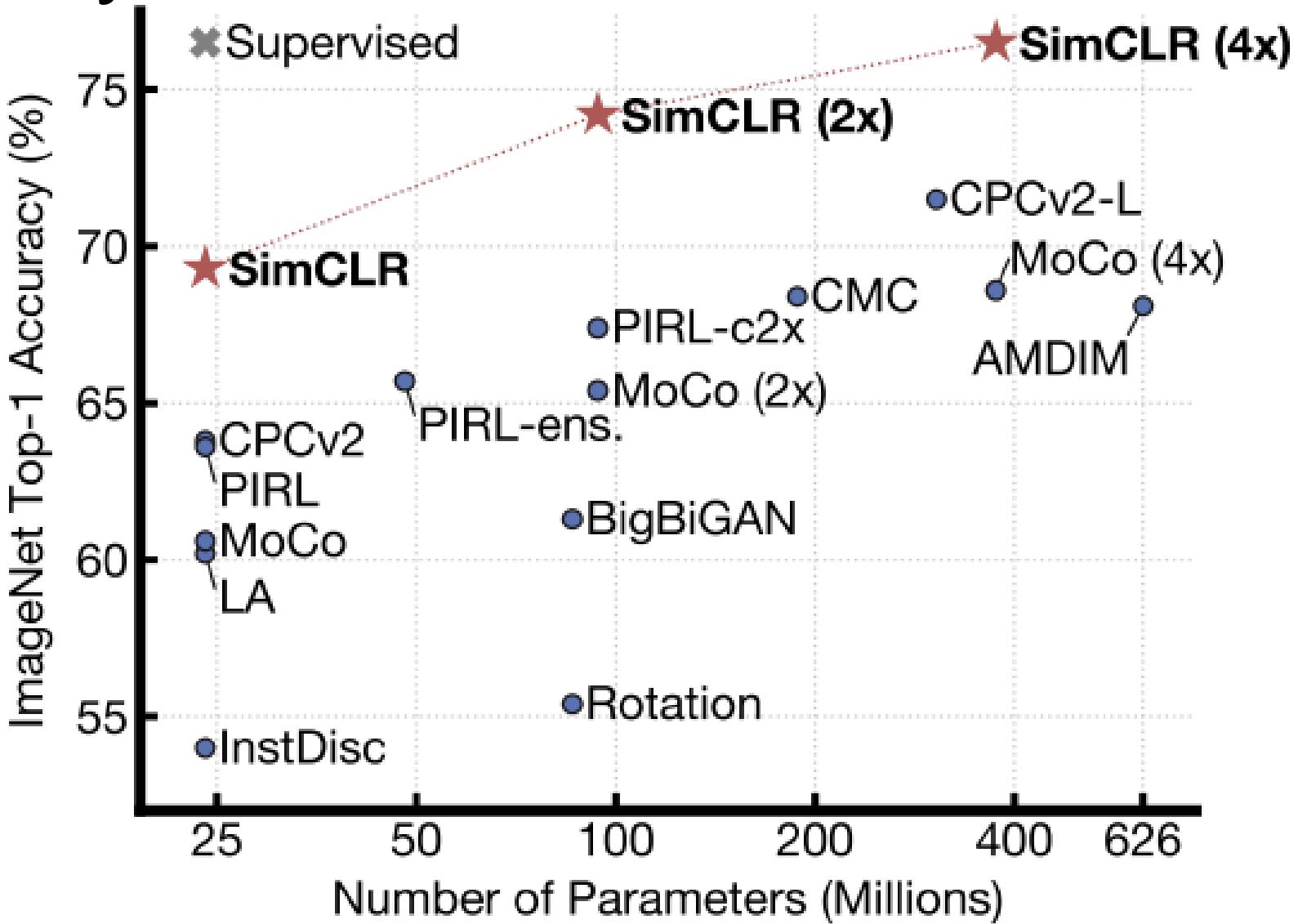


(i) Gaussian blur



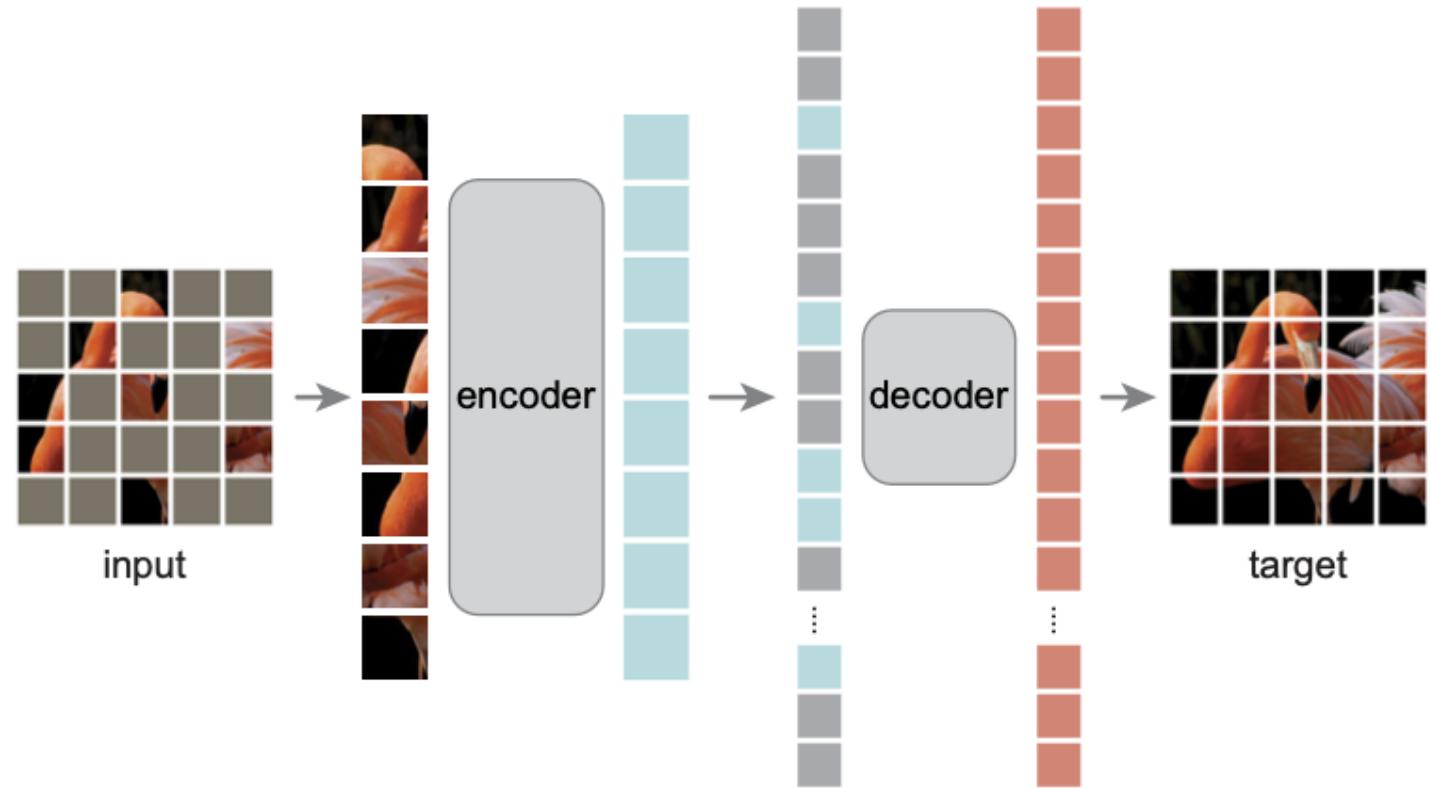
(j) Sobel filtering

# Was a very active area



# Masked Auto-Encoders

- Mask out image patches, predict masked patches from visible patches.
- Pre-train encoder & decoder.
- Use encoder as an image representation.
- Used as frozen backbone for various vision tasks like detection/segmentation (later in class)



# Better than semantic supervision on ImageNet 1K!

method	pre-train data	AP <sup>box</sup>		AP <sup>mask</sup>	
		ViT-B	ViT-L	ViT-B	ViT-L
supervised	IN1K w/ labels	47.9	49.3	42.9	43.9
MoCo v3	IN1K	47.9	49.3	42.7	44.0
BEiT	IN1K+DALLE	49.8	<b>53.3</b>	44.4	47.1
MAE	IN1K	<b>50.3</b>	<b>53.3</b>	<b>44.9</b>	<b>47.2</b>

Table 4. **COCO object detection and segmentation** using a ViT Mask R-CNN baseline. All entries are based on our implementation. Self-supervised entries use IN1K data *without* labels. Mask AP follows a similar trend as box AP.

Scaling behavior →

method	pre-train data	ViT-B	ViT-L
supervised	IN1K w/ labels	47.4	49.9
MoCo v3	IN1K	47.3	49.1
BEiT	IN1K+DALLE	47.1	53.3
MAE	IN1K	<b>48.1</b>	<b>53.6</b>

Table 5. **ADE20K semantic segmentation** (mIoU) using Uper-Net. BEiT results are reproduced using the official code. Other entries are based on our implementation. Self-supervised entries use IN1K data *without* labels.

dataset	ViT-B	ViT-L	ViT-H	ViT-H <sub>448</sub>	prev best
iNat 2017	70.5	75.7	79.3	<b>83.4</b>	75.4 [55]
iNat 2018	75.4	80.1	83.0	<b>86.8</b>	81.2 [54]
iNat 2019	80.5	83.4	85.7	<b>88.3</b>	84.1 [54]
Places205	63.9	65.8	65.9	<b>66.8</b>	66.0 [19] <sup>†</sup>
Places365	57.9	59.4	59.8	<b>60.3</b>	58.0 [40] <sup>‡</sup>

Table 6. **Transfer learning accuracy on classification datasets**, using MAE pre-trained on IN1K and then fine-tuned. We provide system-level comparisons with the previous best results.

<sup>†</sup>: pre-trained on 1 billion images. <sup>‡</sup>: pre-trained on 3.5 billion images.

# Improves performance on ImageNet itself

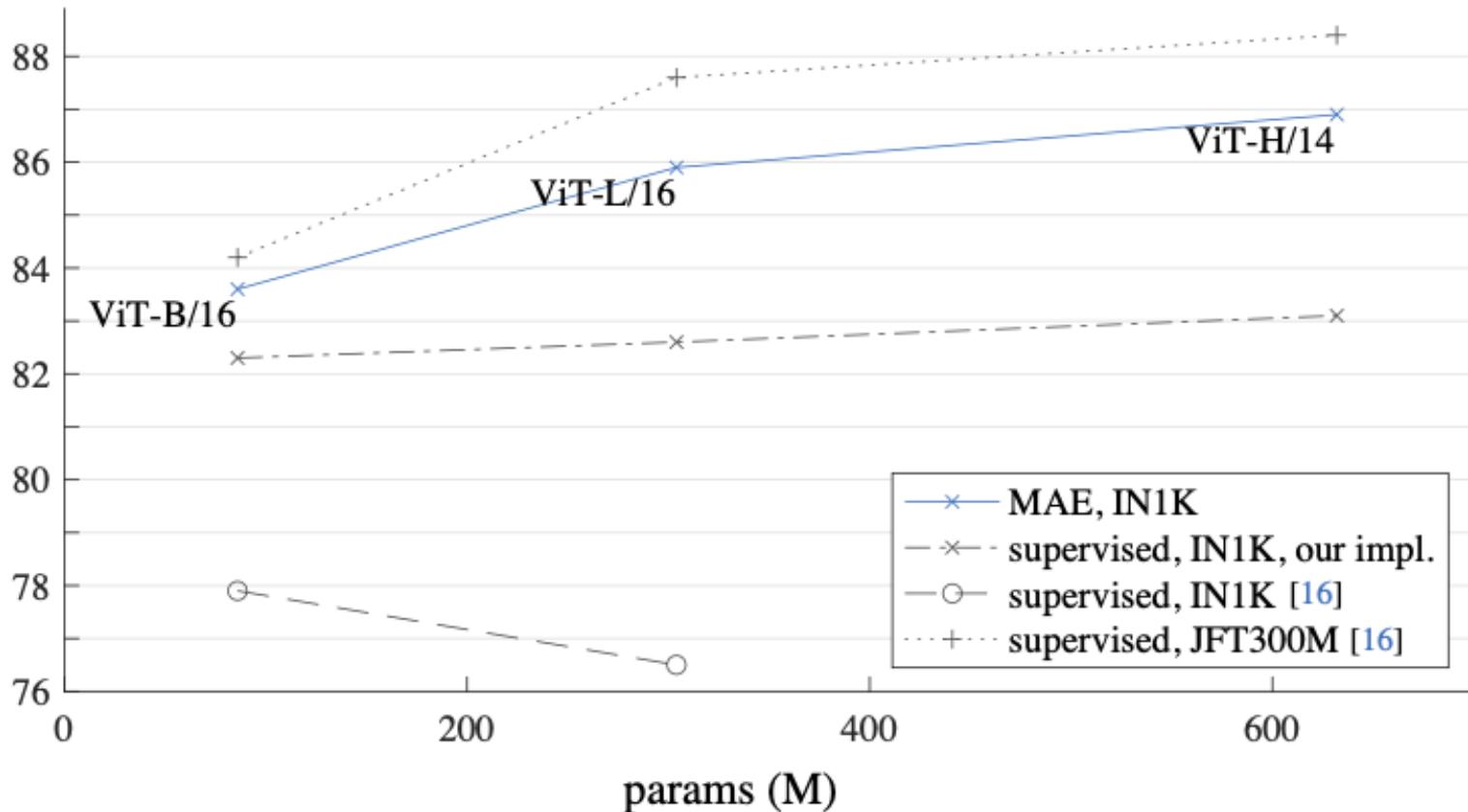
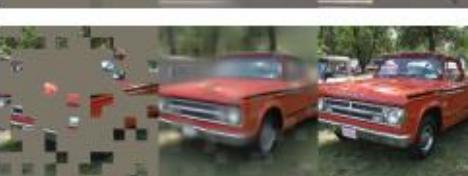
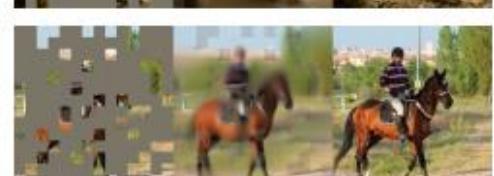
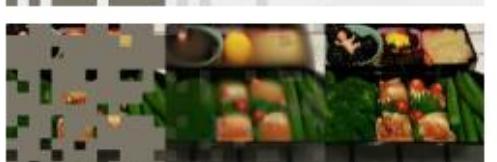


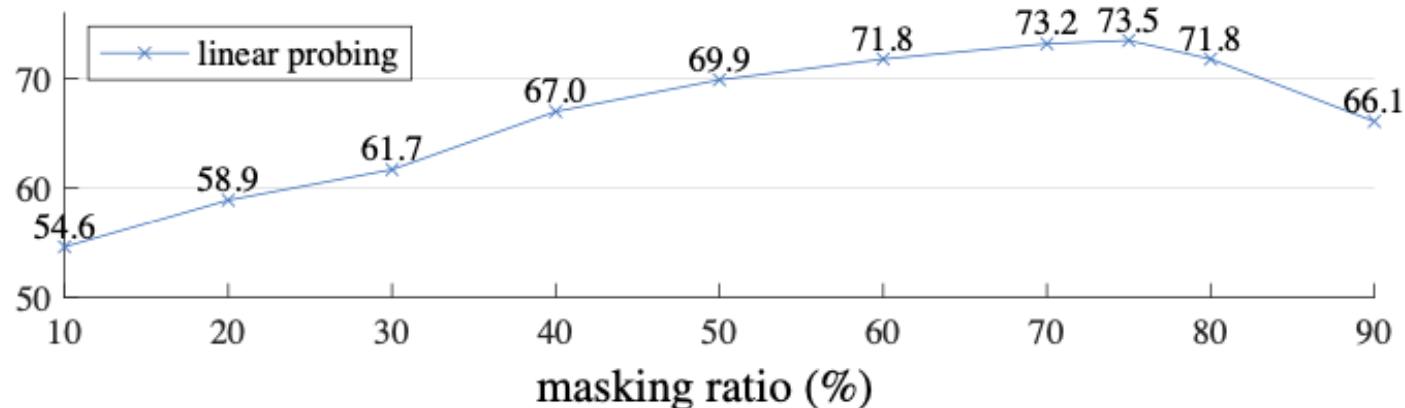
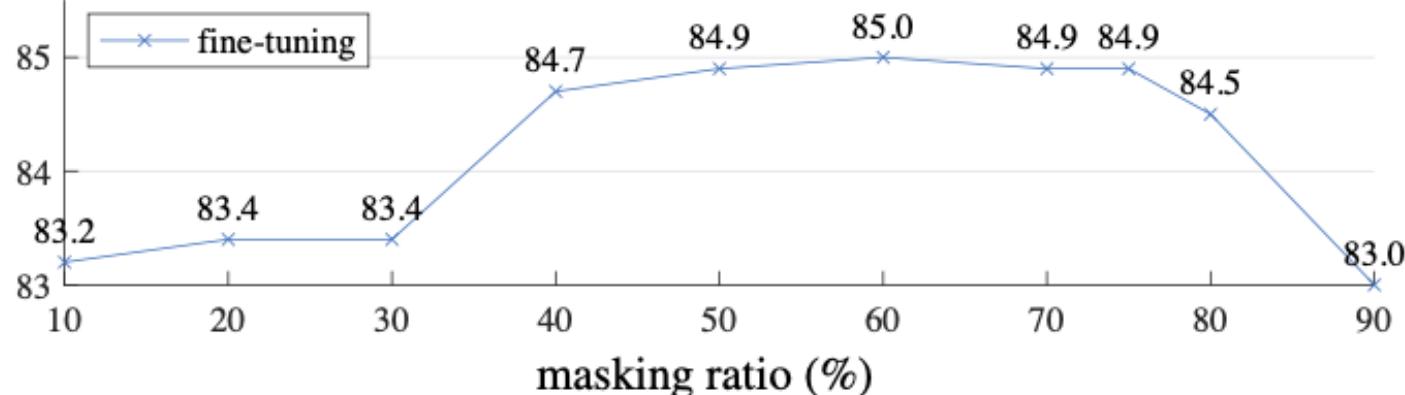
Figure 8. **MAE pre-training vs. supervised pre-training**, evaluated by fine-tuning in ImageNet-1K (224 size). We compare with the original ViT results [16] trained in IN1K or JFT300M.

# Better than past self-supervision approaches

method	pre-train data	ViT-B	ViT-L	ViT-H	ViT-H <sub>448</sub>
scratch, our impl.	-	82.3	82.6	83.1	-
DINO [5]	IN1K	82.8	-	-	-
MoCo v3 [9]	IN1K	83.2	84.1	-	-
BEiT [2]	IN1K+DALLE	83.2	85.2	-	-
MAE	IN1K	<u>83.6</u>	<u>85.9</u>	<u>86.9</u>	<b>87.8</b>



# Ablations



Need high masking ratio for good learning.  
NLP models use 15-20% masking ratio.

case	ft	lin	FLOPs
encoder w/ [M]	84.2	59.6	3.3×
encoder w/o [M]	<b>84.9</b>	<b>73.5</b>	<b>1×</b>

Faster and better to not input masked out patches to encoder

case	ft	lin
pixel (w/o norm)	84.9	73.5
pixel (w/ norm)	<b>85.4</b>	<b>73.9</b>
PCA	84.6	72.3
dVAE token	85.3	71.6

Normalized pixels are a better target than discrete tokens / PCA coefficients

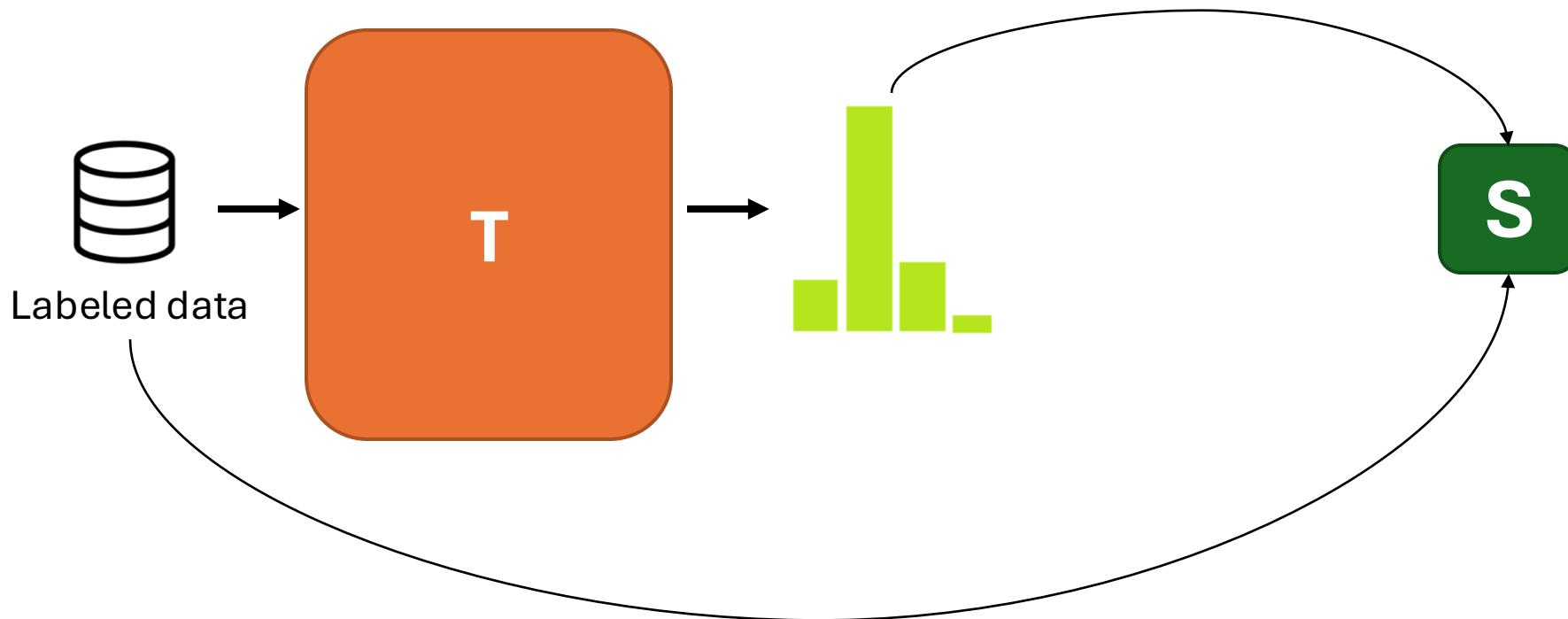
# Take away

- Able to pre-train data-hungry ViT model
- Very well structured experiments! Read and learn from it.



# DiNO - Approach

- Self supervised learning as a special case of **knowledge distillation**

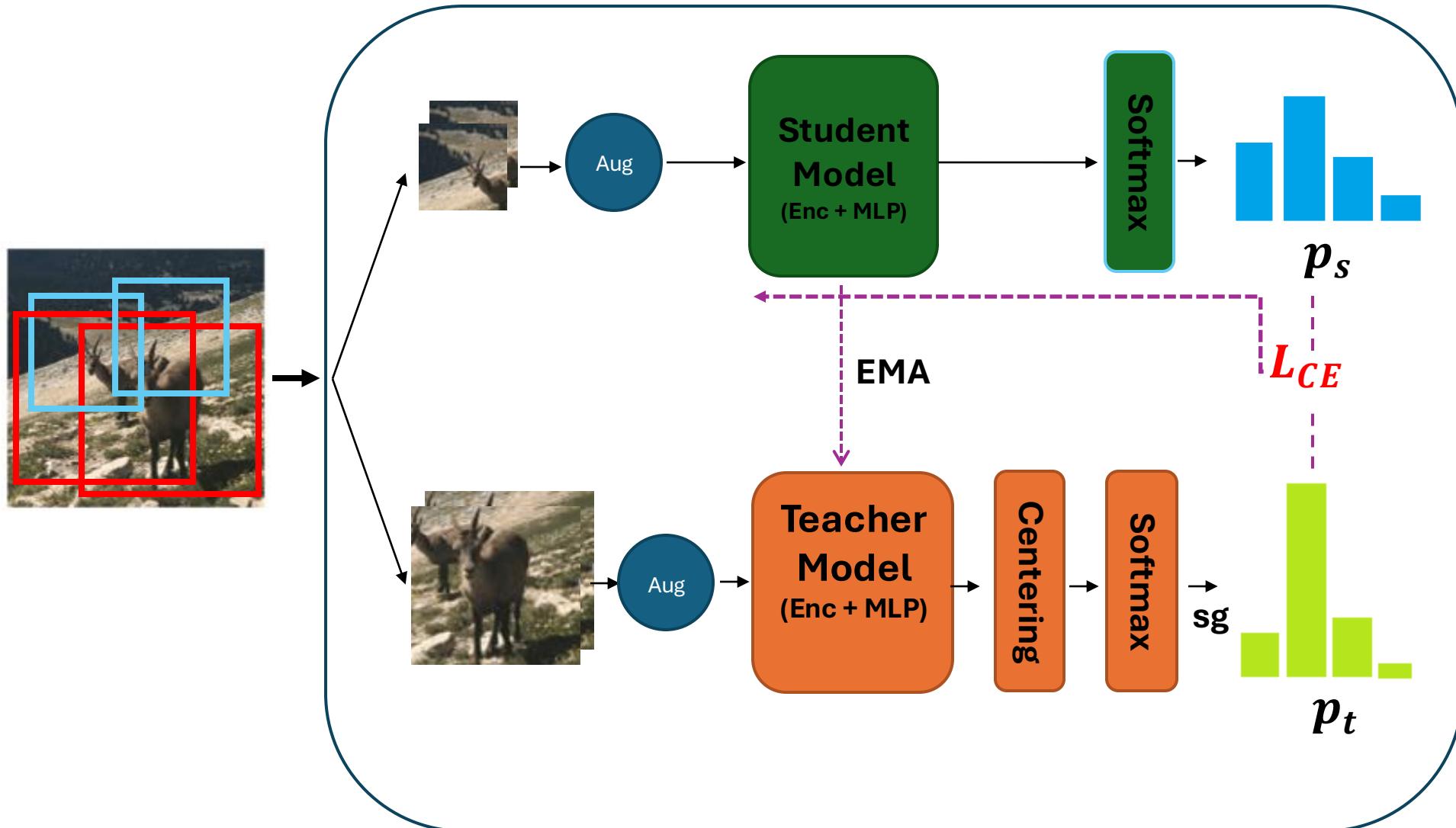


# DiNO - Training

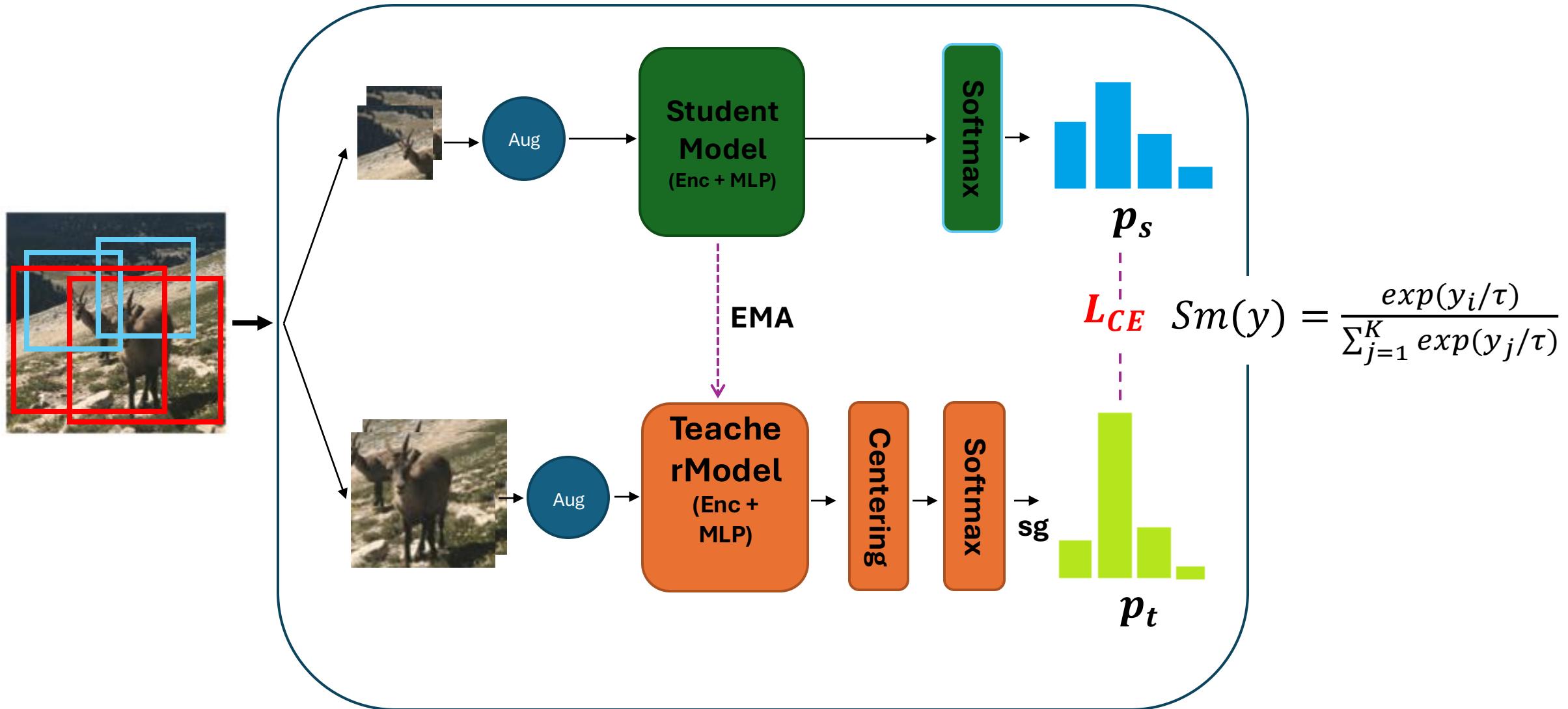


DiNO

# DiNO - Training



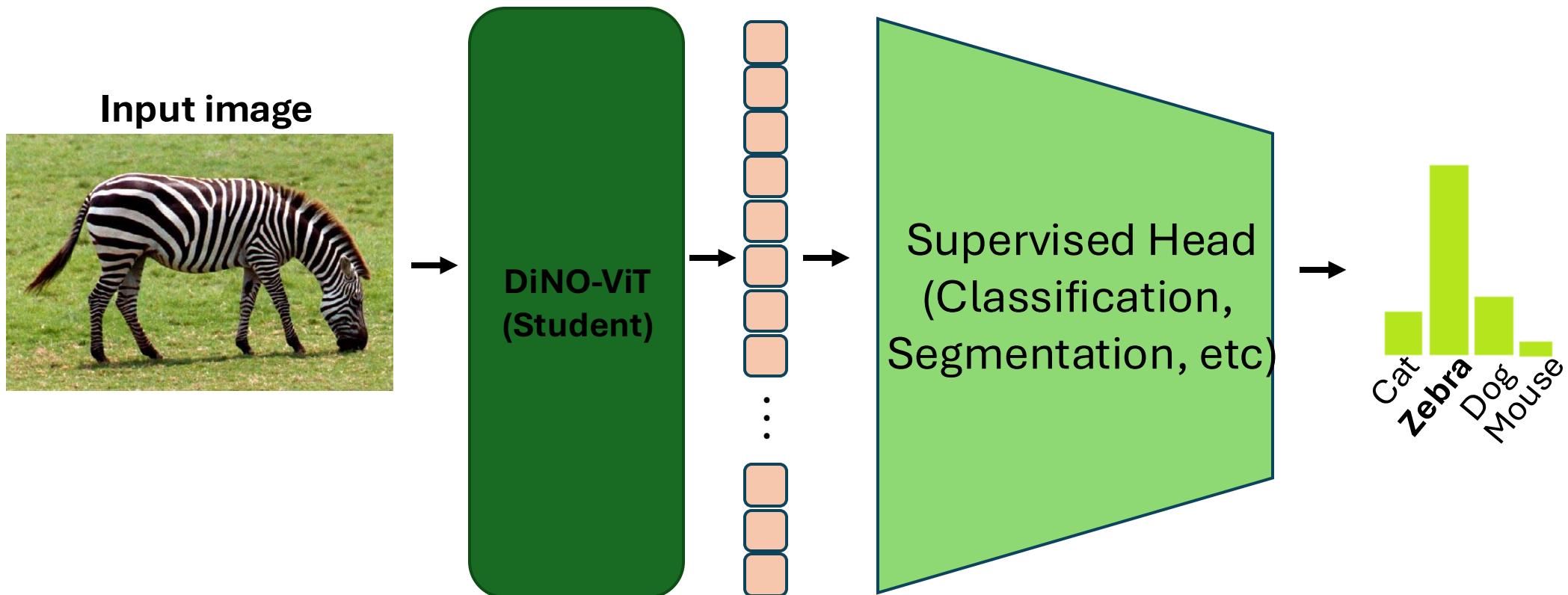
# DiNO - Training



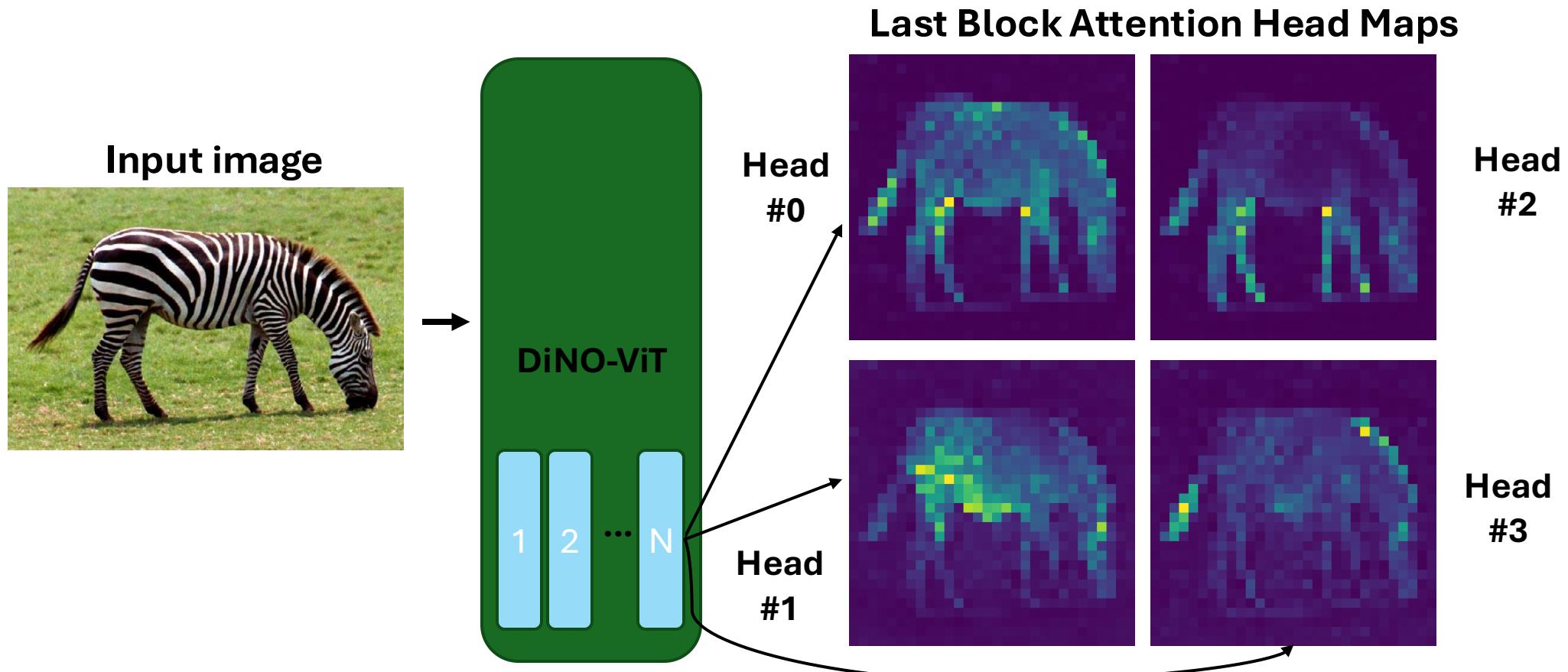
# DiNO - Explanation

- Augmentations
  - Tells the model what to ignore
    - Color jitter, Gaussian Blur, Solarize
    - Acts as a data prior
- “Global – local” cropping
- Teacher out-distribution sharpening via centering & Low-temperature in softmax
- The student encoder learns “abstract representations”
  - No awareness of “class labels” or meaning behind logits

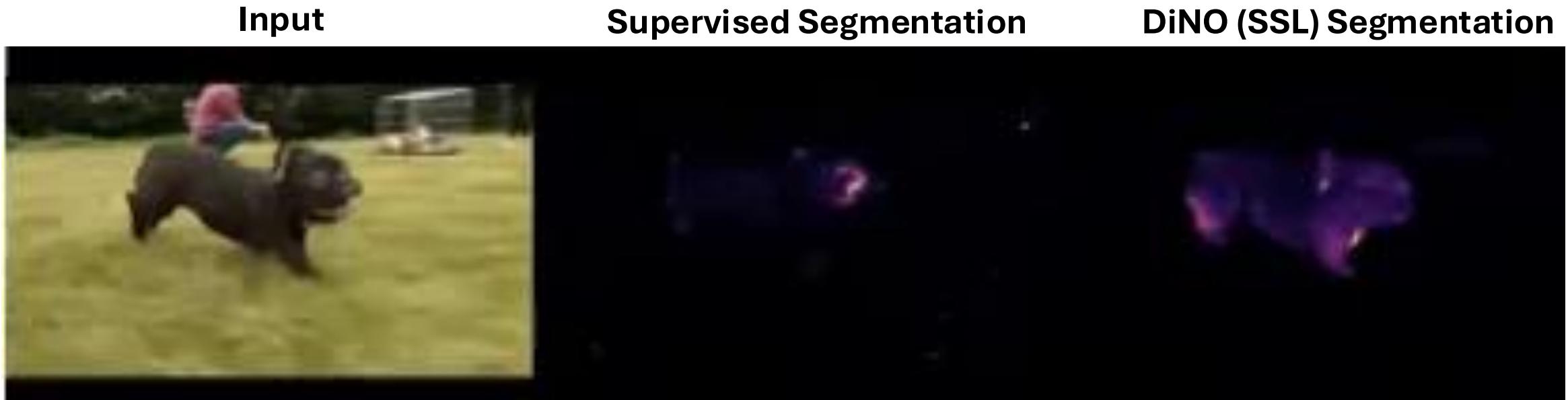
# DiNO - Inference #1



# DiNO - Inference #2



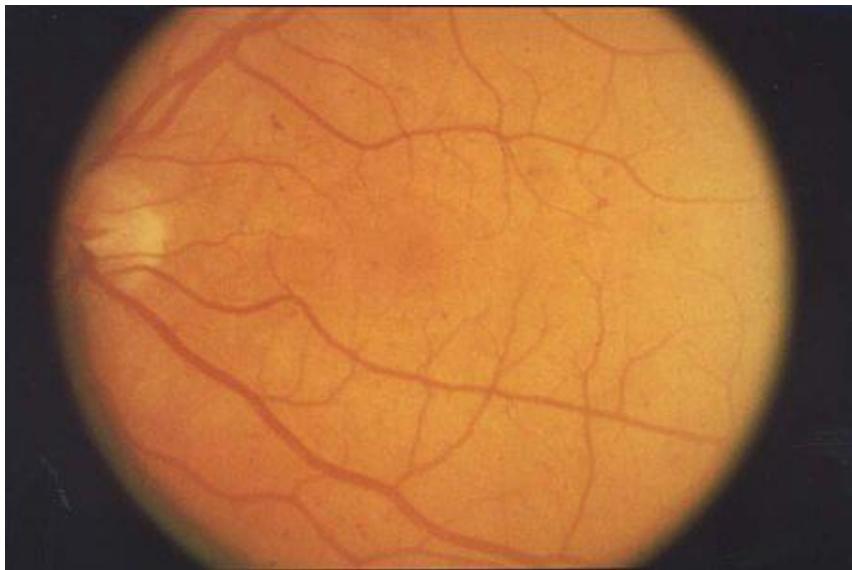
# DiNO - Inference #2



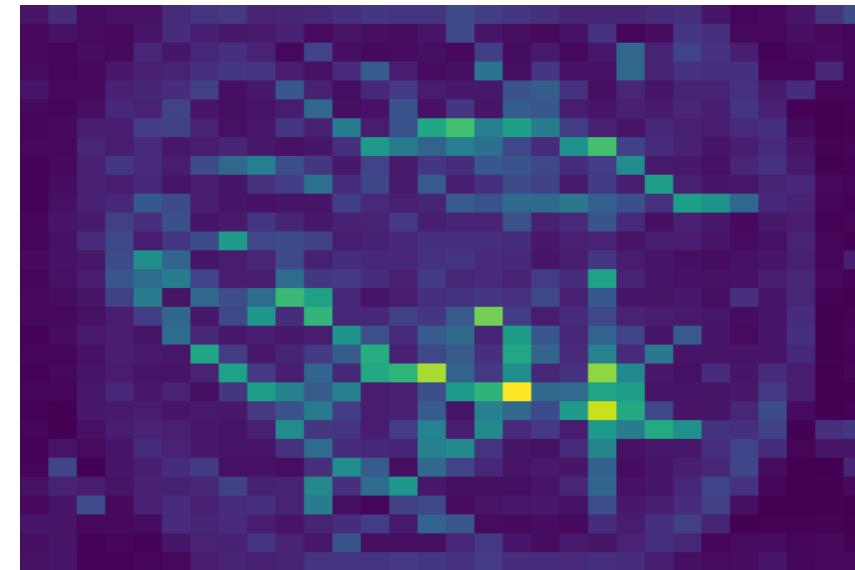
# DiNO - Inference #2

- Self supervised learning also makes learned representations applicable to out-of-distribution data

**Input image**



**Last block attention map**

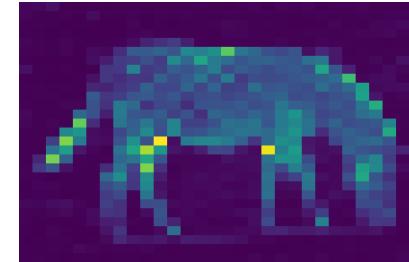


# Usage

```
vit_model = torch.hub.load('facebookresearch/dino:main',
                           'dino_vits16', pretrained=True)
img = imread('zebra.png')

x = vit_model.prepare_tokens(img)
for blk in vit_model.blocks[:-1]:
    x = blk(x)
attn_maps = vit_model.blocks[-1](x, return_attention=True)

# Choose head, Get attention map of class token
attn_map = attn_maps[0, HEAD, 0, 1:].reshape((1, 1, H_PATCHES, W_PATCHES))
attn_map = F.interpolate(attn_map, scale_factor=16, mode="nearest")
```



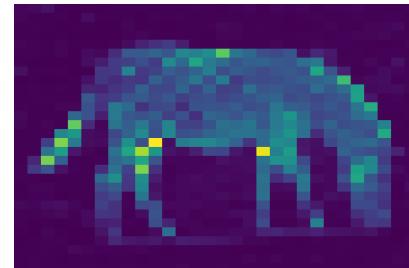
# Usage

```
vit_model = torch.hub.load('facebookresearch/dino:main',
                           'dino_vits16', pretrained=True)

img = imread('zebra.png')

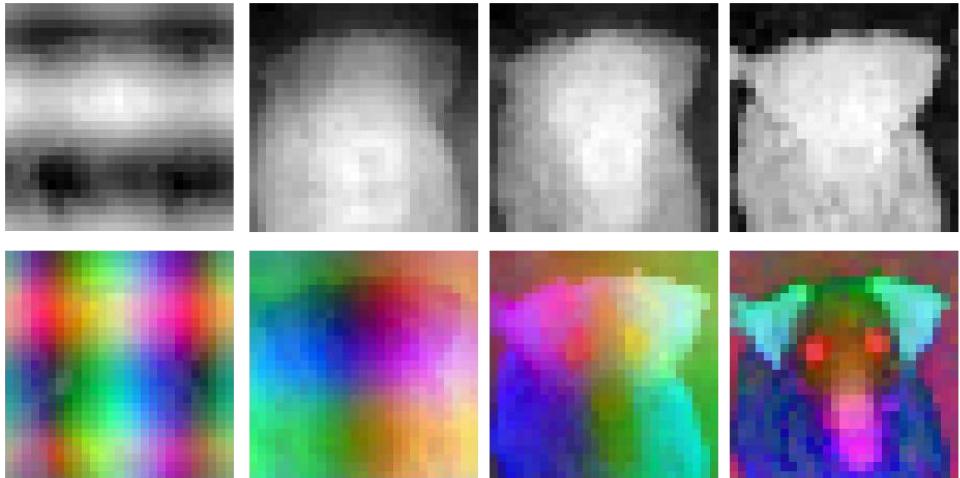
attn_maps = vit_model.get_last_selfattention(img)

# Choose head, Get attention map of class token
attn_map = attn_maps[0, HEAD, 0, 1:].reshape((1, 1, H_PATCHES, W_PATCHES))
attn_map = F.interpolate(attn_map, scale_factor=16, mode="nearest")
```



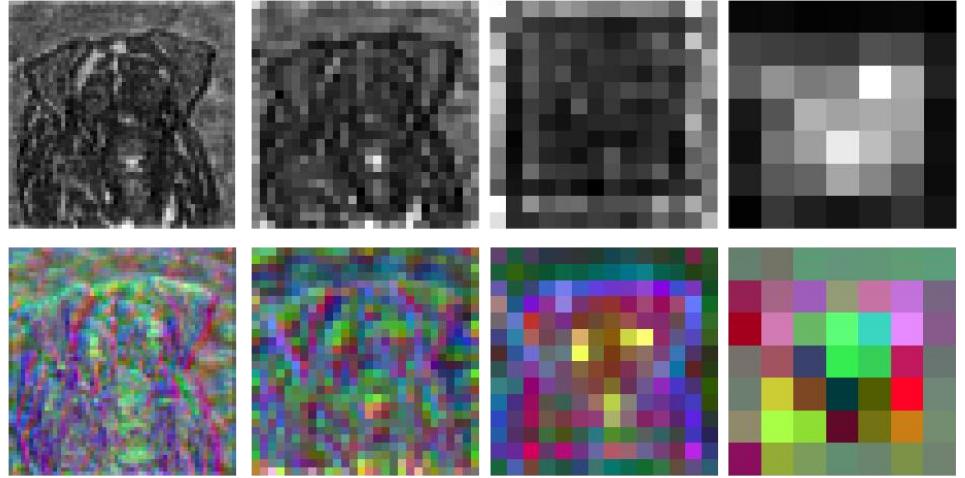
# PCA (Keys) across layers

Self-supervised ViT (DINO-ViT)



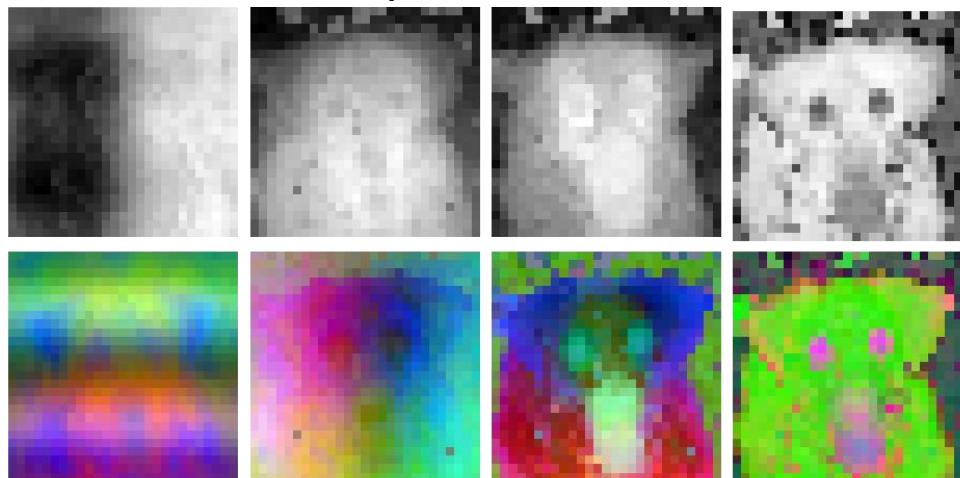
shallow deep

Self-supervised ResNet (DINO-ResNet)

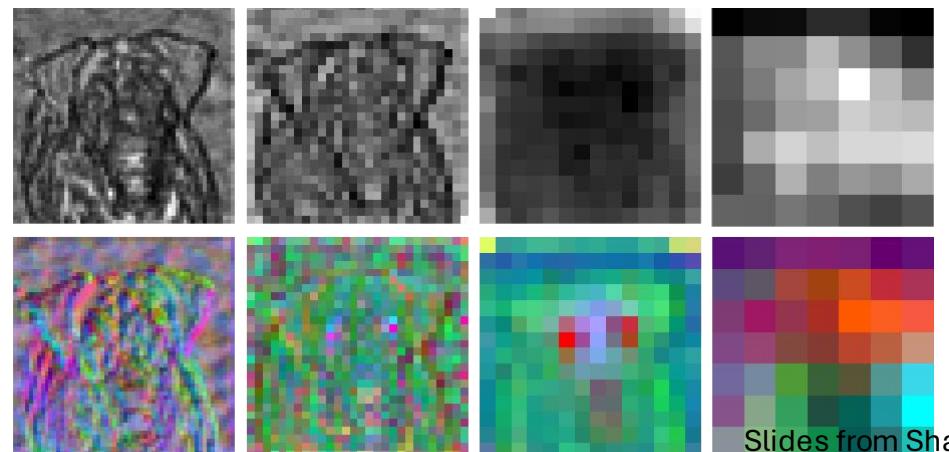


shallow deep

Supervised ViT

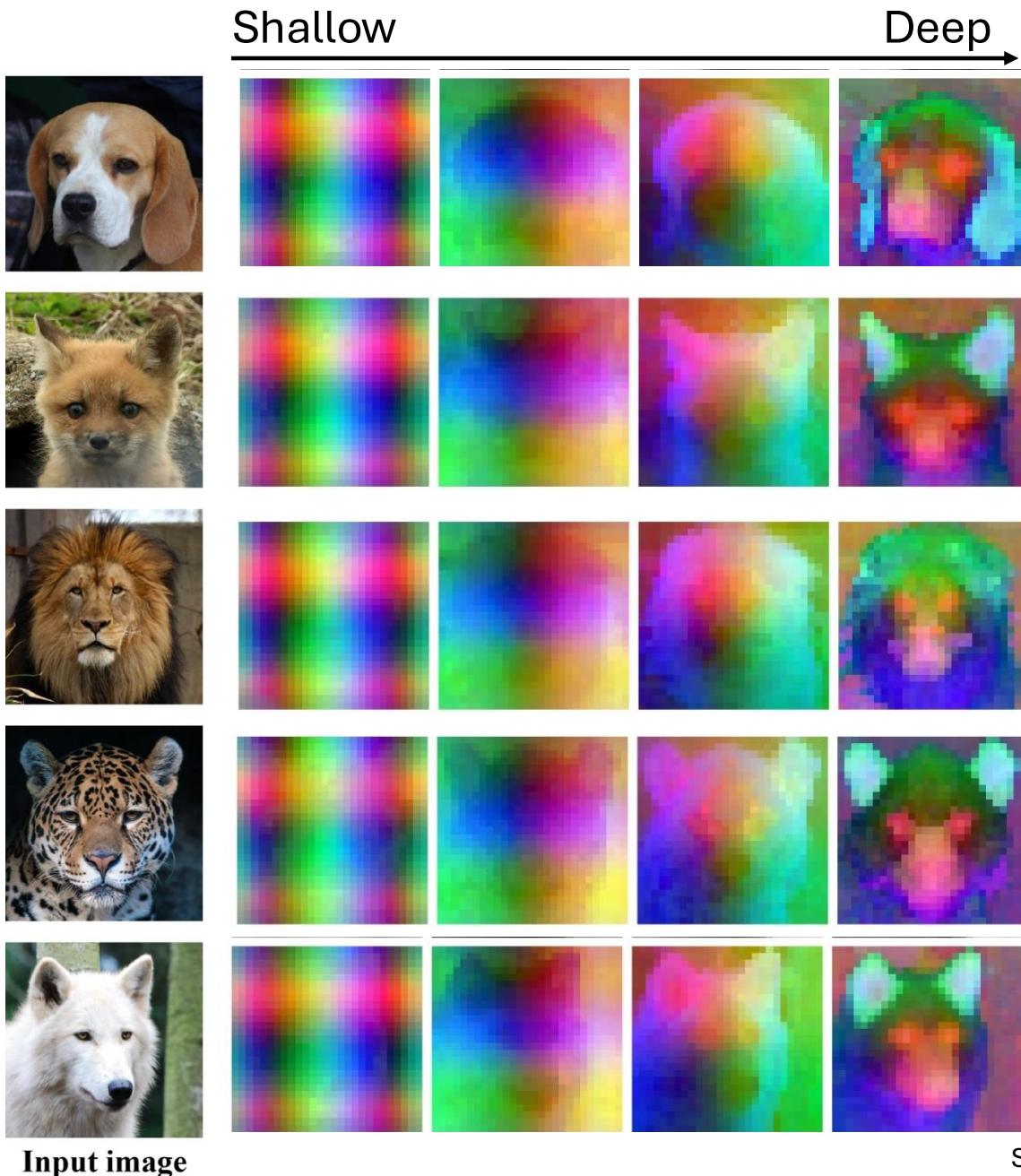


Supervised ResNet



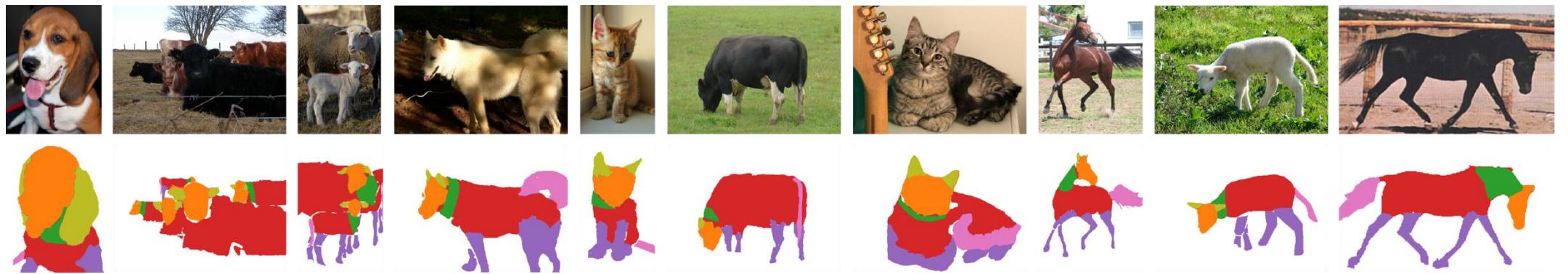
Slides from Shai Bagon

# PCA – DiNO ViT



Amir, S., Gandelsman, Y., Bagon, S., Dekel, T.,  
“Deep ViT Features as dense visual descriptors”

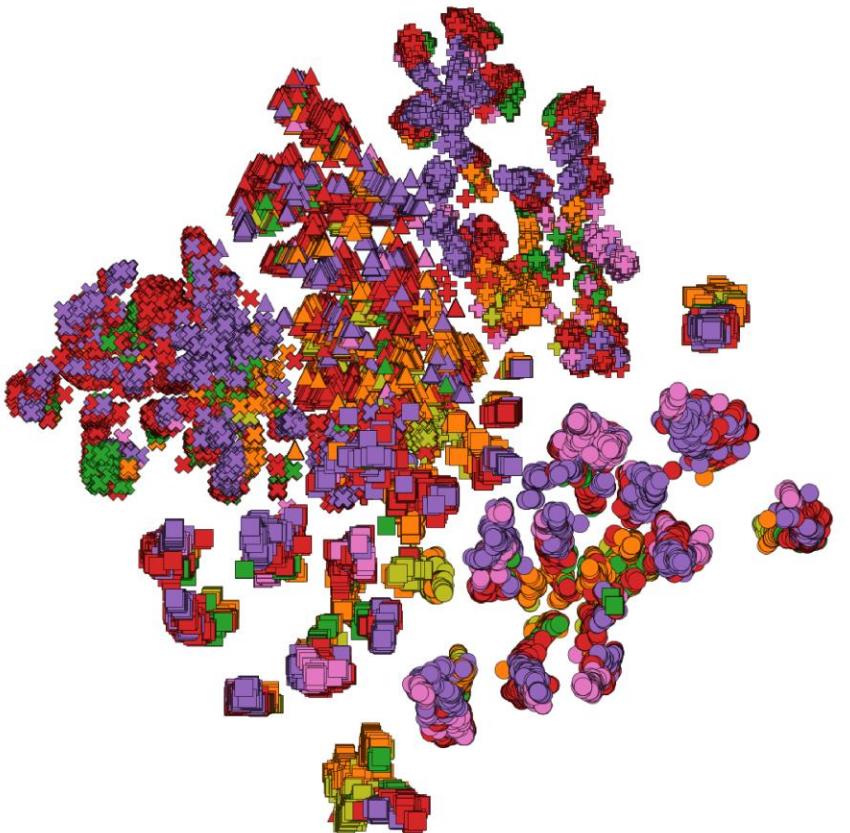
Slides from Shai Bagon



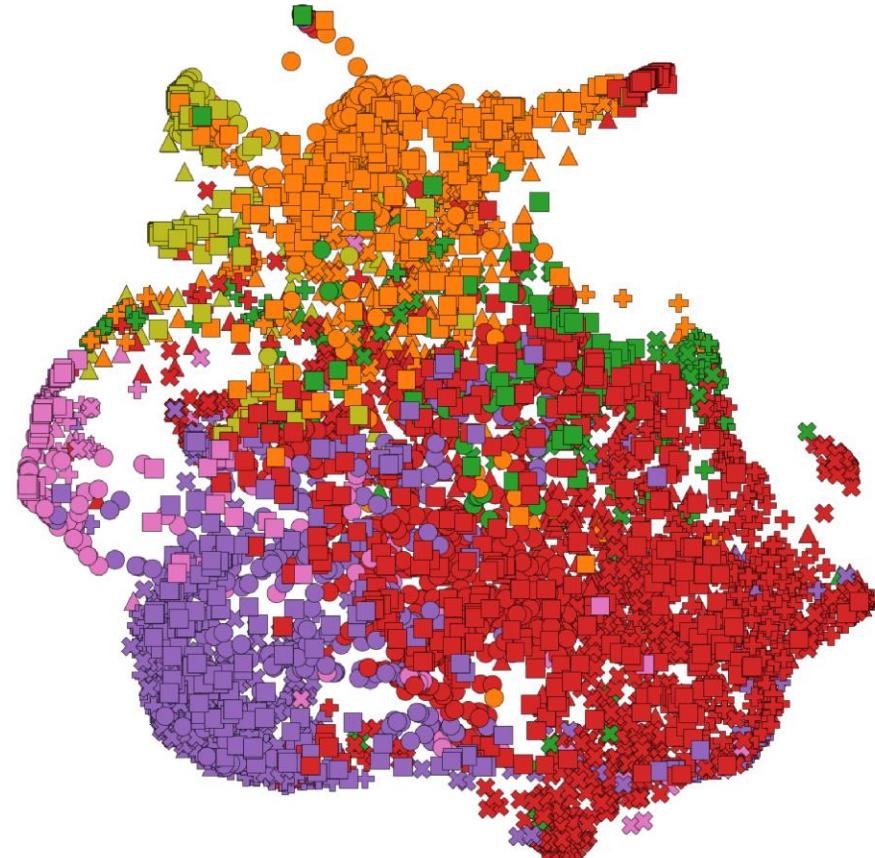
(a) Sample images and ground truth parts

**Torso**  
**Neck**  
**Head**  
**Ears**  
**Tail**  
**Limbs**

- **Cat**
- **Dog**
- + **Horse**
- ◊ **Sheep**
- △ **Cow**



(c) Supervised ViT



(b) Self-Supervised ViT (DINO-ViT)

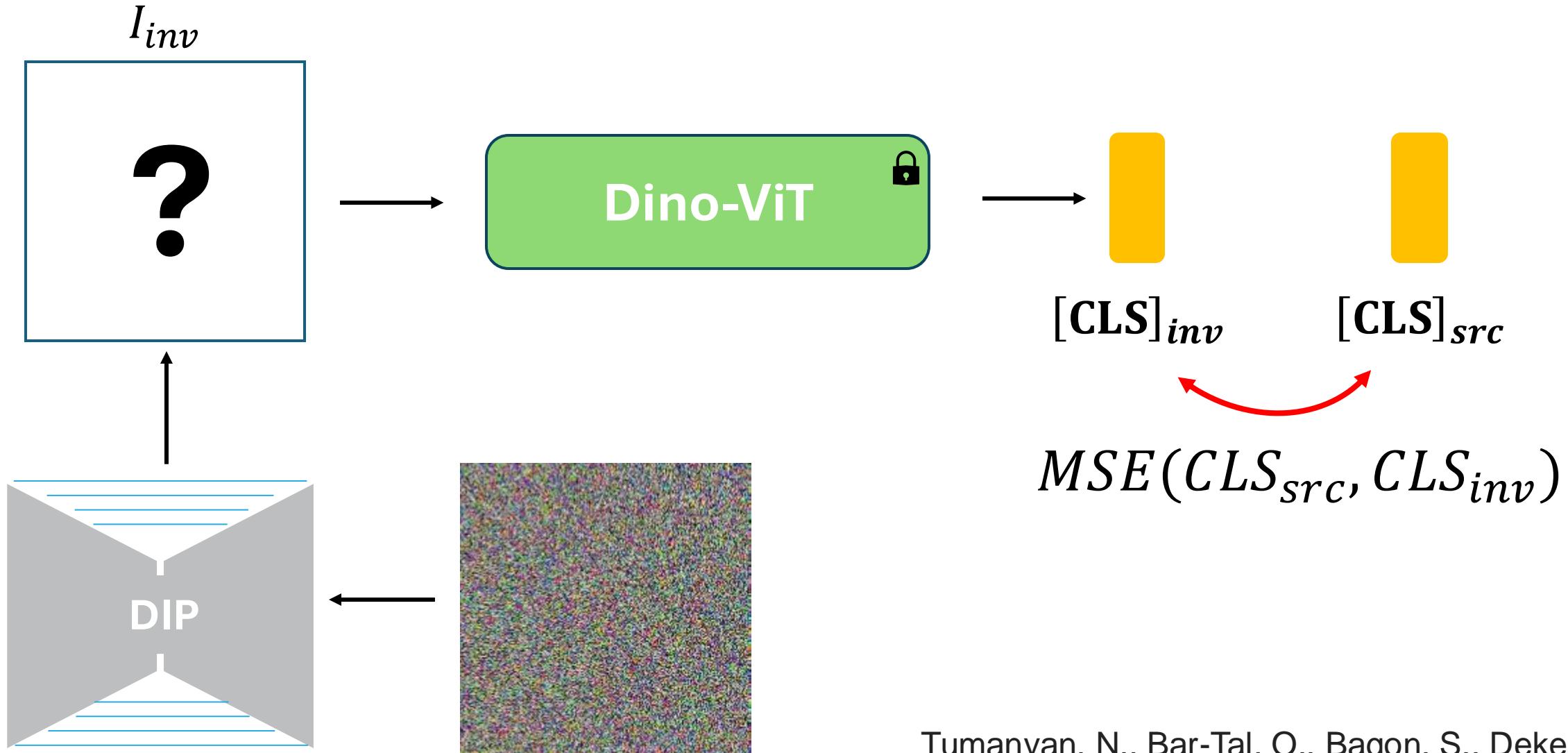
Slides from Shai Bagon

T-SNE of keys  
visualization

# [CLS] as a global appearance representation



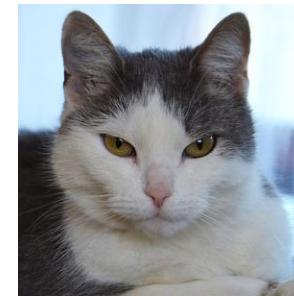
# [CLS] as a global appearance representation



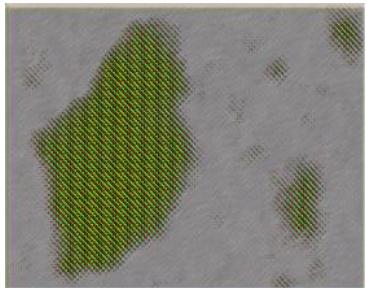
Tumanyan, N., Bar-Tal, O., Bagon, S., Dekel, T.  
Splicing vit features for semantic appearance transfer (CVPR 2022)  
Slides from Shai Bagon

# [CLS] as a global appearance representation

Input



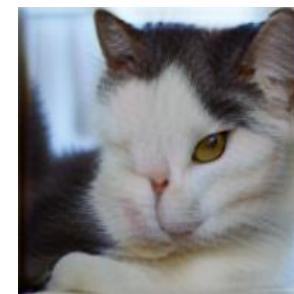
layer 0



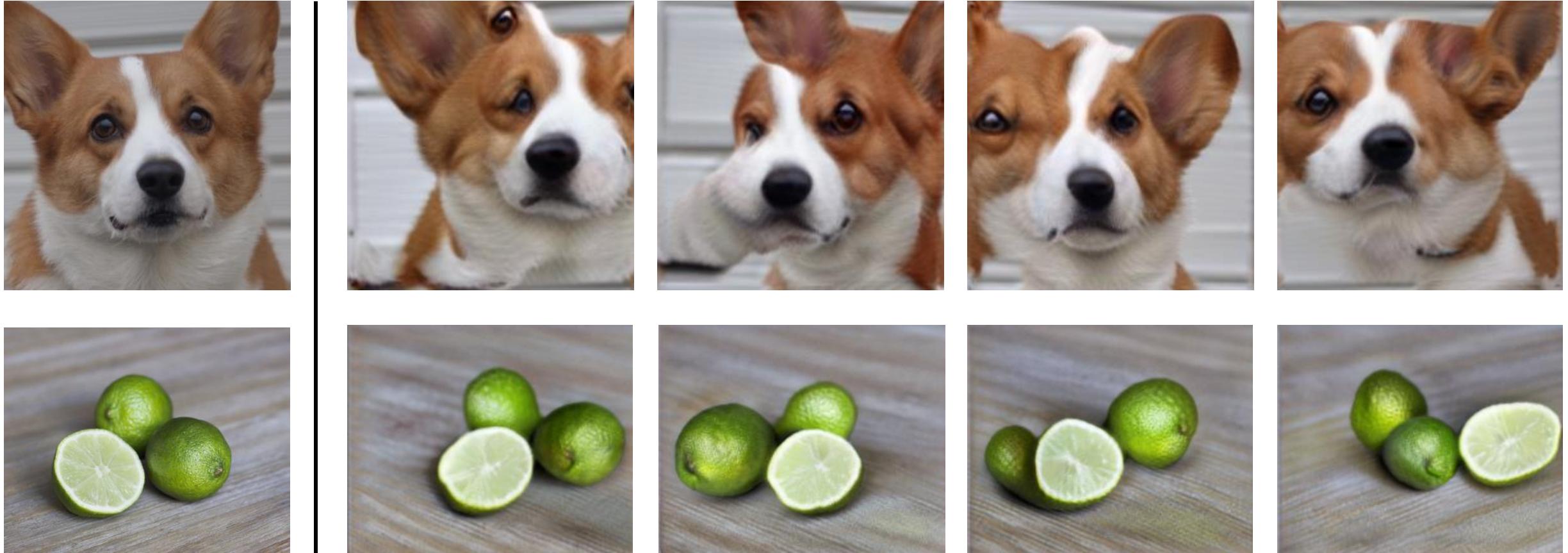
layer 3



layer 11



# [CLS] as a global appearance representation



Inversion run 1

Inversion run 2

Inversion run 3

Inversion run 4