

# Lecture 25: Recent architectures

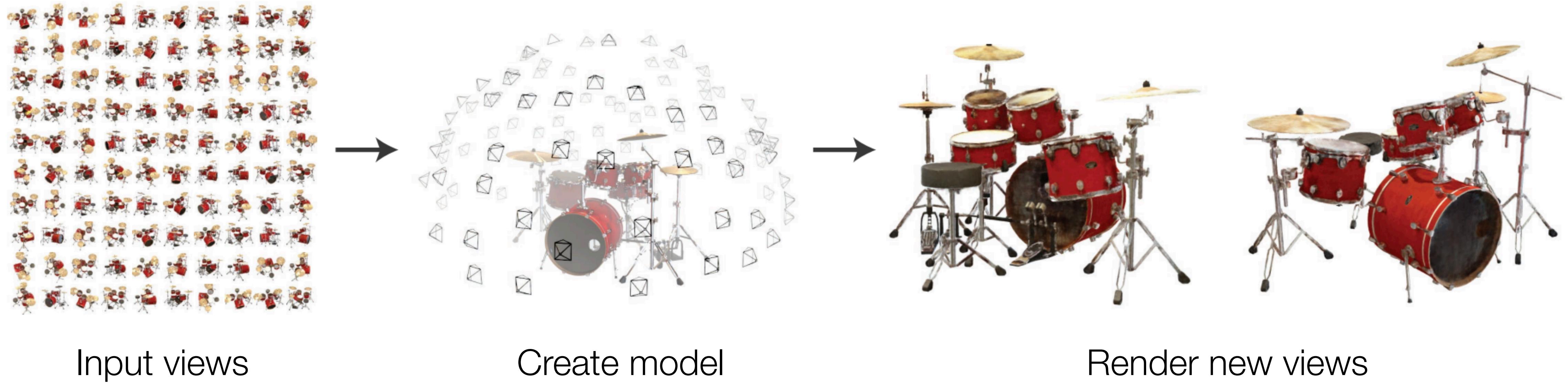
# Reminders

- Sign up for final presentation (see Piazza)

# Today

- Neural fields
- Transformers for vision

# 3D view synthesis



What representation should we use?

# Idea #1: Image-based rendering



Point cloud  
(reconstructed with SfM + multi-view stereo)



View from a different angle

# Idea #1: Image-based rendering



Point cloud



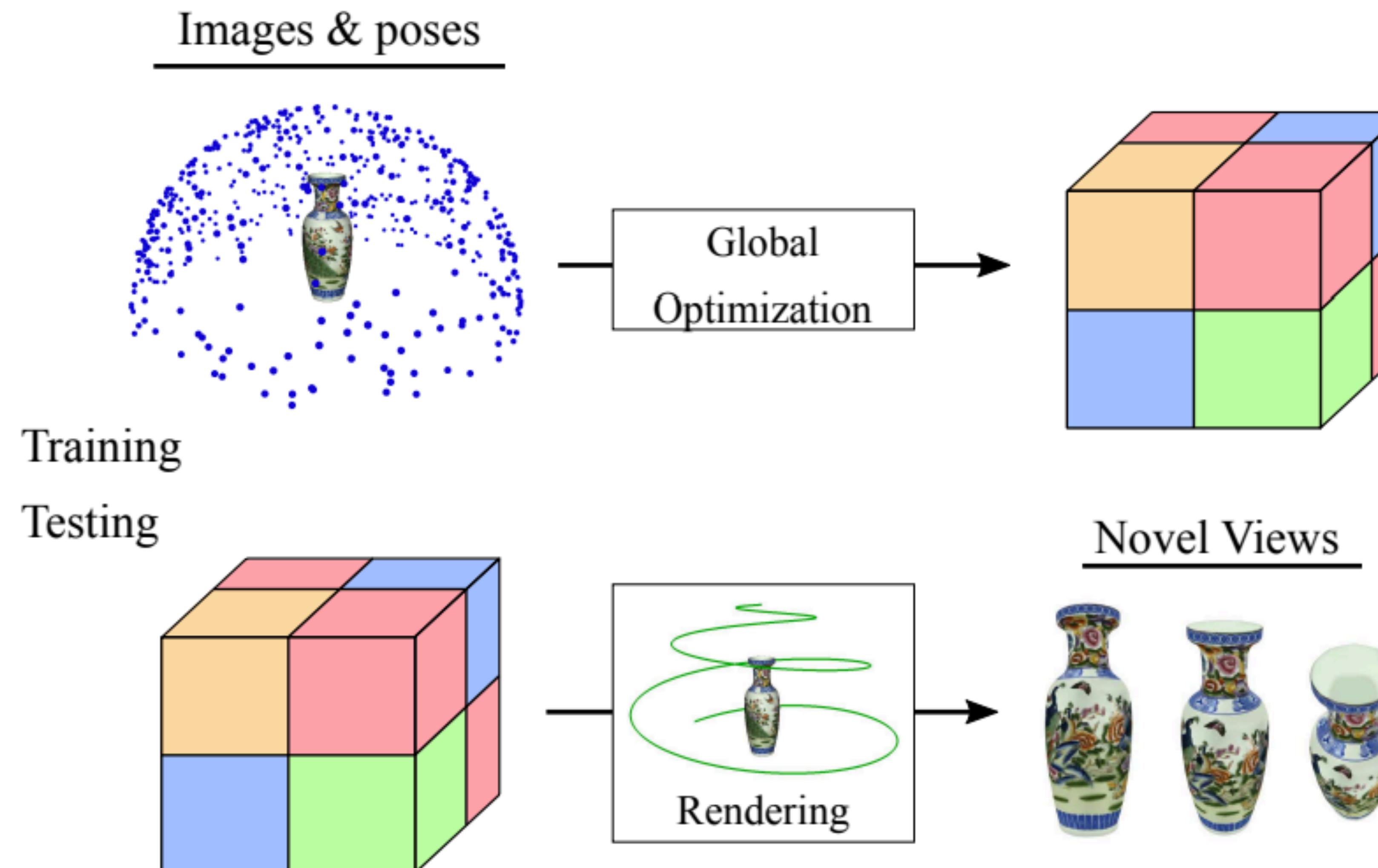
Proxy geometry (a mesh)

To synthesize a new view, select colors from existing views using proxy geometry.

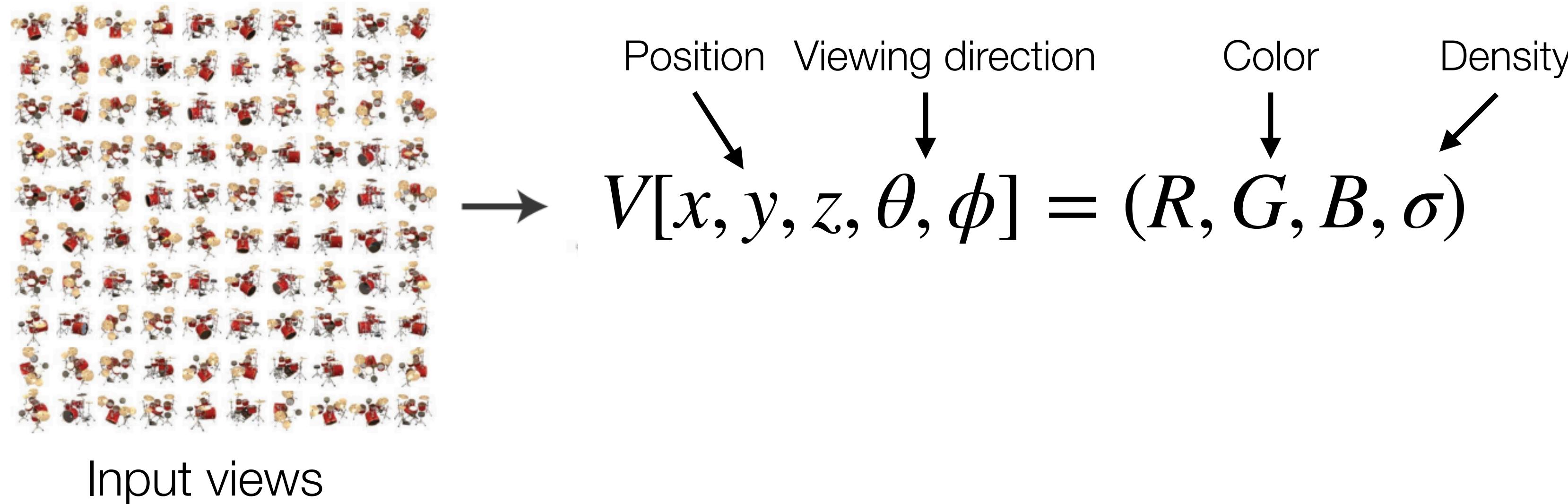
# Idea #1: Image-based rendering



# Idea #2: voxel representation



# Idea #2: voxel representation



# Idea #2: voxel representation



Input views

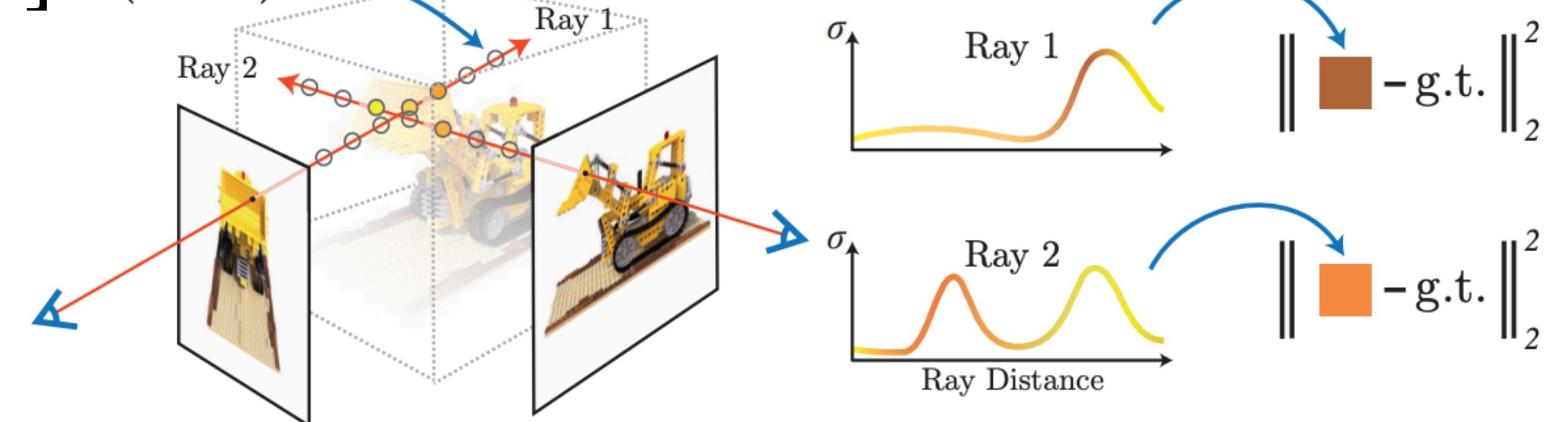
Position   Viewing direction      Color      Density

↓                  ↓                  ↓                  ↓

$$\rightarrow V[x, y, z, \theta, \phi] = (R, G, B, \sigma)$$

## Training:

$$V[x, y, z, \theta, \phi] \rightarrow (RGB\sigma)$$



**Problem:** A huge table!  $\mathcal{O}(D^3 A^2)$

# Idea #3: neural radiance field (NeRF)



$$\rightarrow F_{\Theta}(x, y, z, \theta, \phi) = (R, G, B, \sigma)$$

- Represent using a **neural radiance field**.
- Function that maps a  $(x, y, z, \theta, \phi)$  to a color and density.
- Typically parameterized as a multi-layer perceptron (MLP)
- Goal: find parameters  $\Theta$  for MLP that explain the images

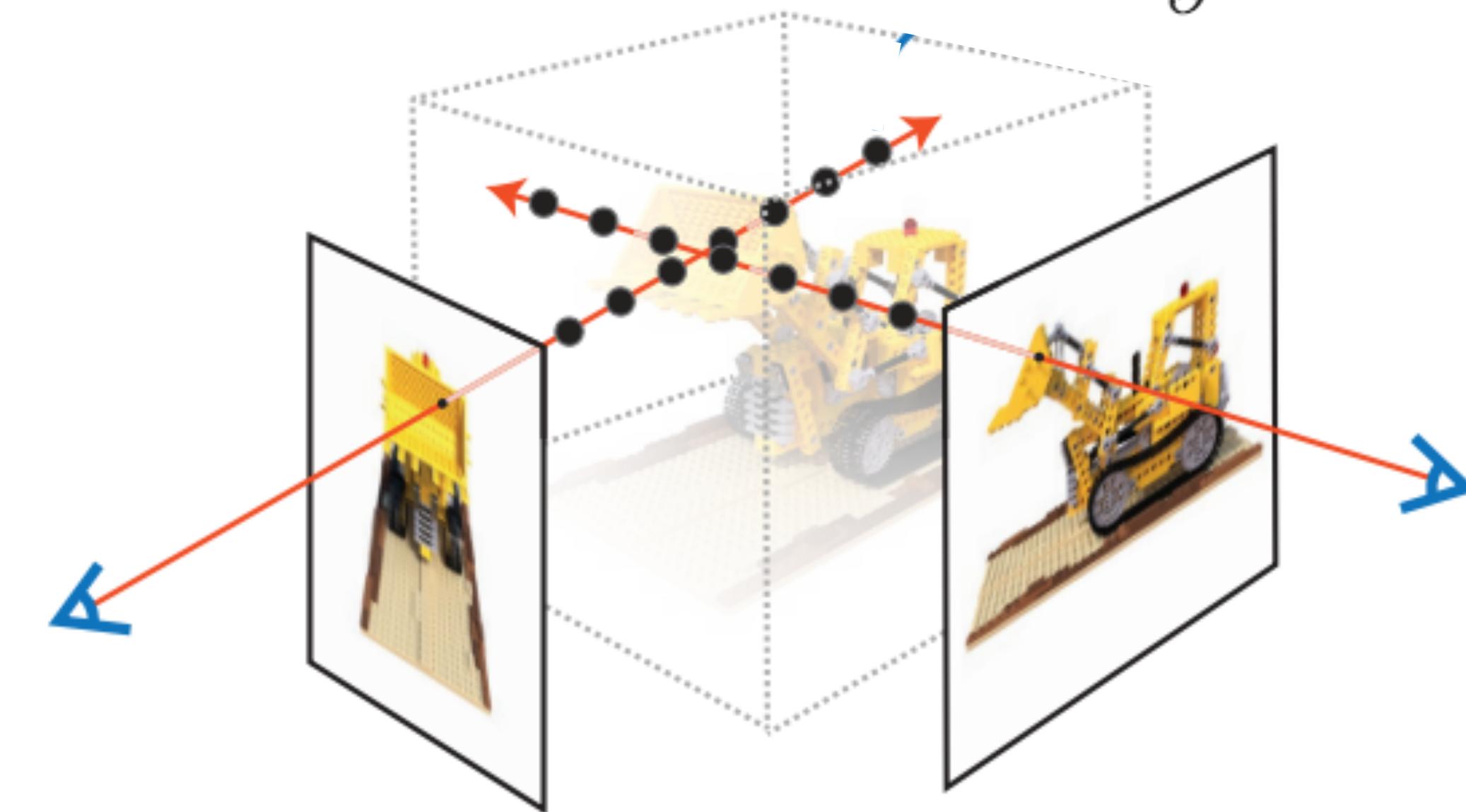
# Idea #3: neural radiance field (NeRF)



3D scene

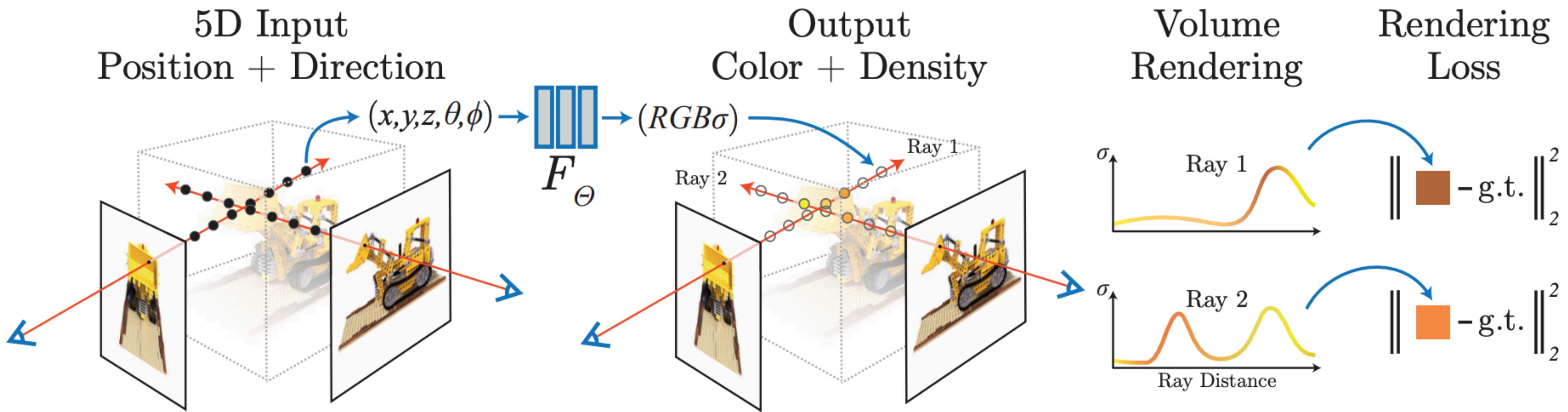
Learn volume:  
color + occupancy

$$(x,y,z,\theta,\phi) \rightarrow F_{\Theta} \rightarrow (RGB\sigma)$$

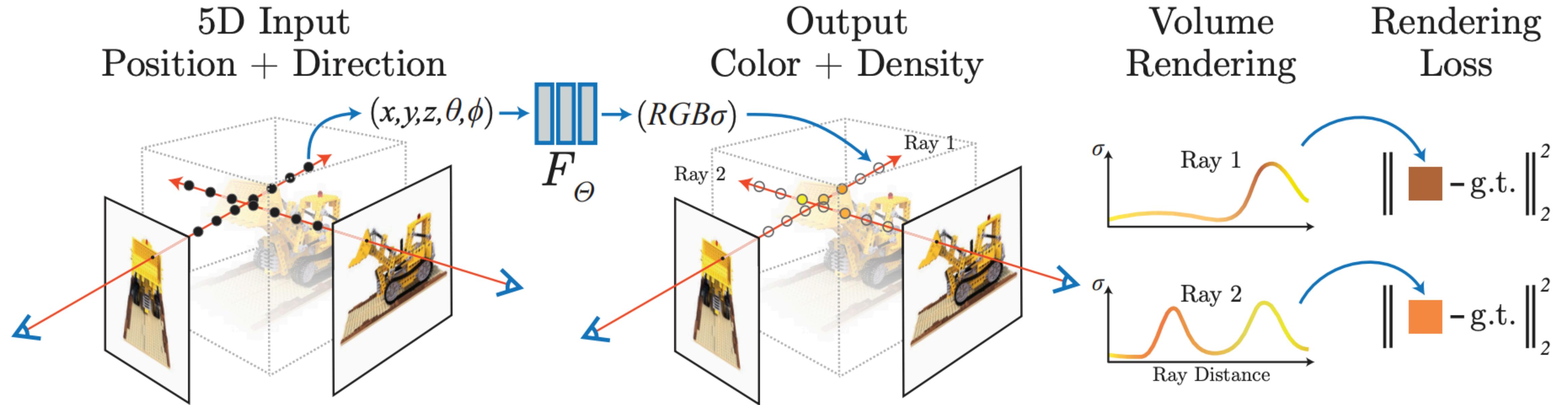


Viewpoints

# Learning a NeRF



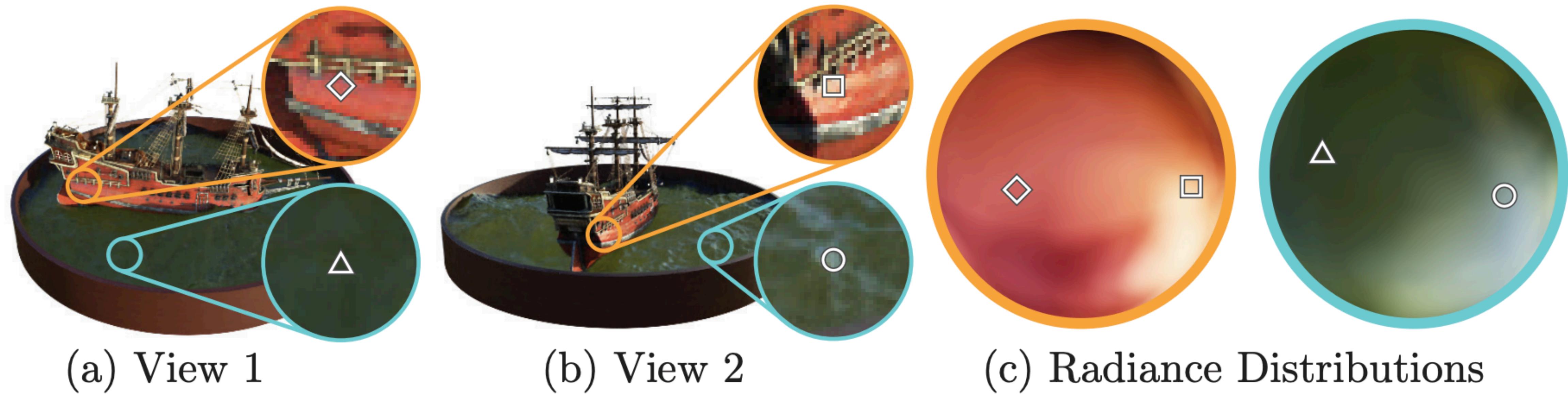
# Neural rendering



$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \text{ where } T(t) = \exp\left(- \int_{t_n}^t \sigma(\mathbf{r}(s)) ds\right)$$

Ray:  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  For color  $c$  and density  $\sigma$ .

# Why is it good to be view-dependent?



# Representing the inputs



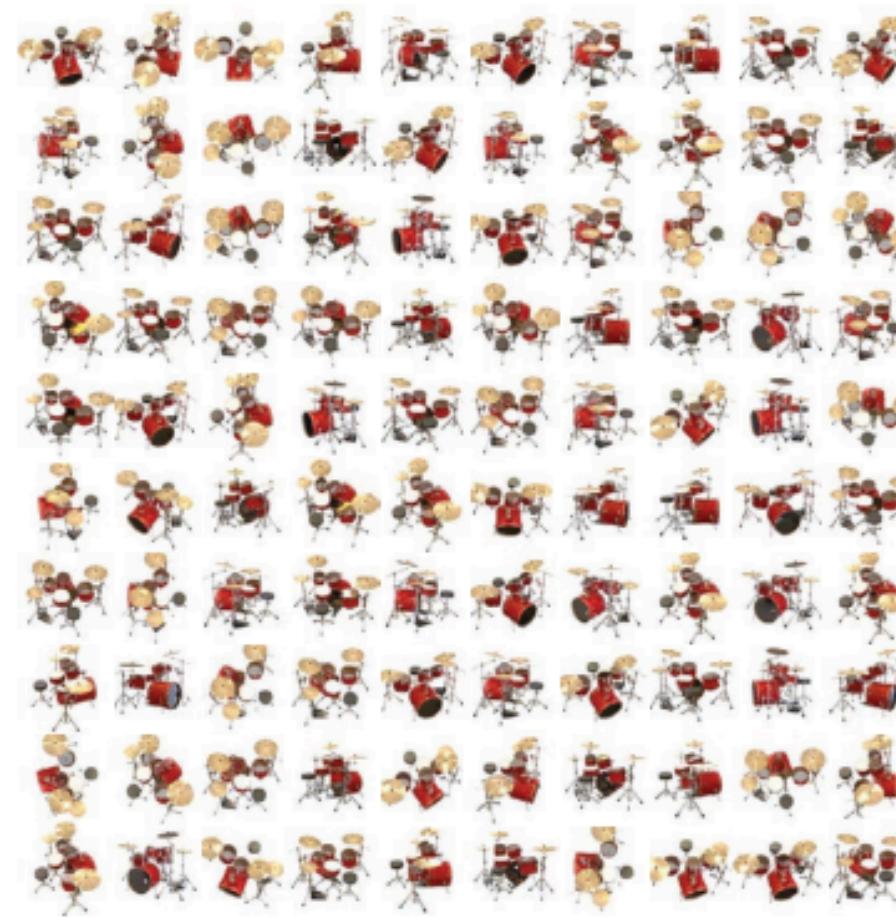
$$\rightarrow F_{\Theta}(x, y, z, \theta, \phi) = (R, G, B, \sigma)$$

Input views

- In theory, could just plug in 4 inputs  $x, y, z, \theta, \phi$
- However, this leads to blurry results.
- Neural nets show a bias toward low frequency functions.



# Representing the inputs



$$\rightarrow F_{\Theta}(x, y, z, \theta, \phi) = (R, G, B, \sigma)$$

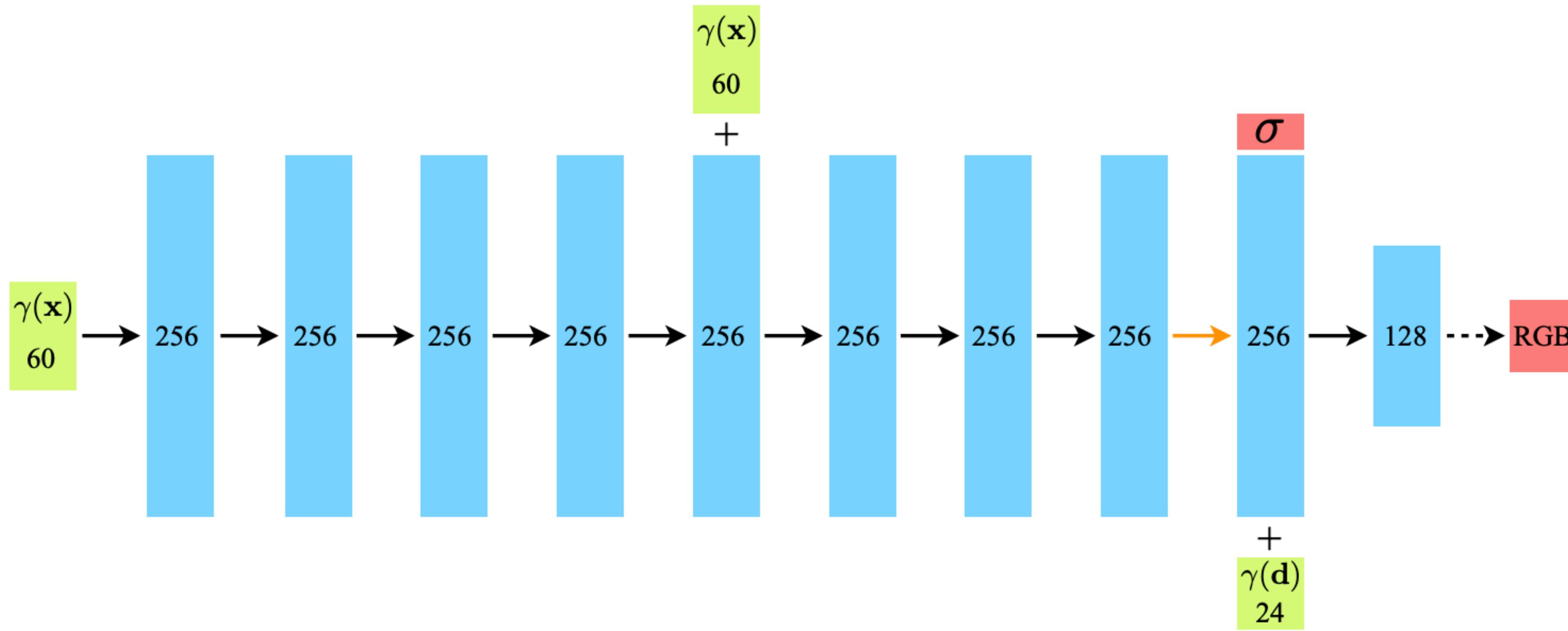
Input views

- Use a **positional encoding**. Given a scalar  $p$ , compute:

$$\gamma(p) = (\sin(2^0\pi p), \cos(2^0\pi p), \dots, \sin(2^{L-1}\pi p), \cos(2^{L-1}\pi p))$$

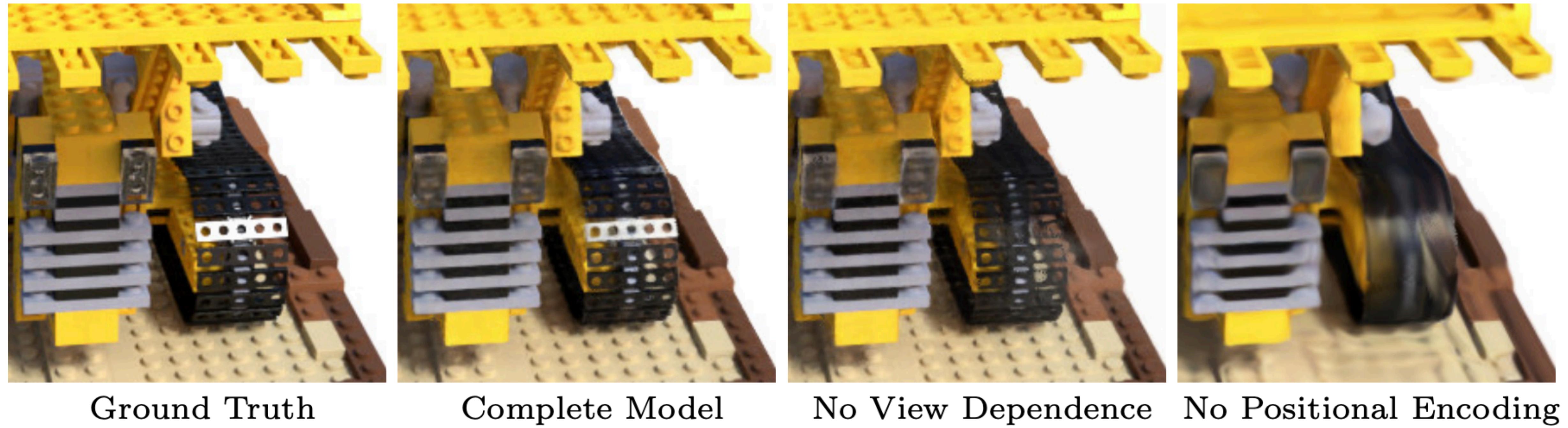
- Plug in the coordinate to sinusoids at different frequencies (e.g.  $L = 10$ ). Creates a high-frequency input.

# MLP architecture



$$\gamma(p) = (\sin(2^0 \pi p), \cos(2^0 \pi p), \dots, \sin(2^{L-1} \pi p), \cos(2^{L-1} \pi p))$$

# Results for a novel viewpoint



# Results



[Mildenhall\*, Srinivasan\*,<sup>20</sup> Tanick\*, et al. 2020]

# Results



[Mildenhall\*, Srinivasan\*,<sup>21</sup> Tanick\*, et al. 2020]

# Results



[Mildenhall\*, Srinivasan\*,<sup>22</sup> Tanick\*, et al. 2020]

# Extension: internet photo collections



# Extension: internet photo collections



# BLOCK-NERF

## RESULTS

Matthew Tancik<sup>1\*</sup>

Sabeek Pradhan<sup>2</sup>

Jonathan T. Barron<sup>3</sup> Henrik Kretzschmar<sup>2</sup>

Vincent Casser<sup>2</sup>

Ben Mildenhall<sup>3</sup>

Xinchen Yan<sup>2</sup>

Pratul Srinivasan<sup>3</sup>

<sup>1</sup>UC Berkeley

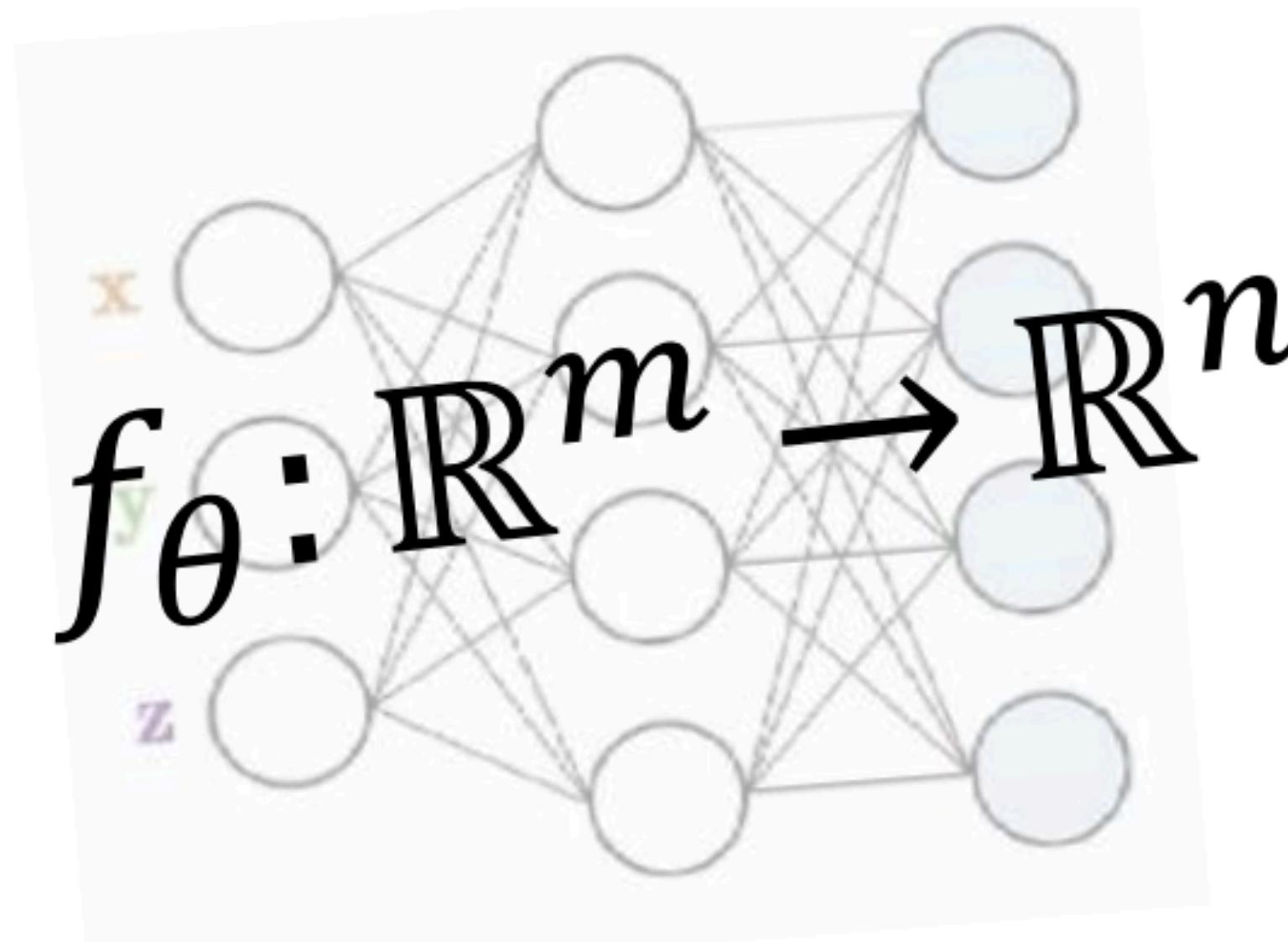
<sup>2</sup>Waymo

<sup>3</sup>Google Research

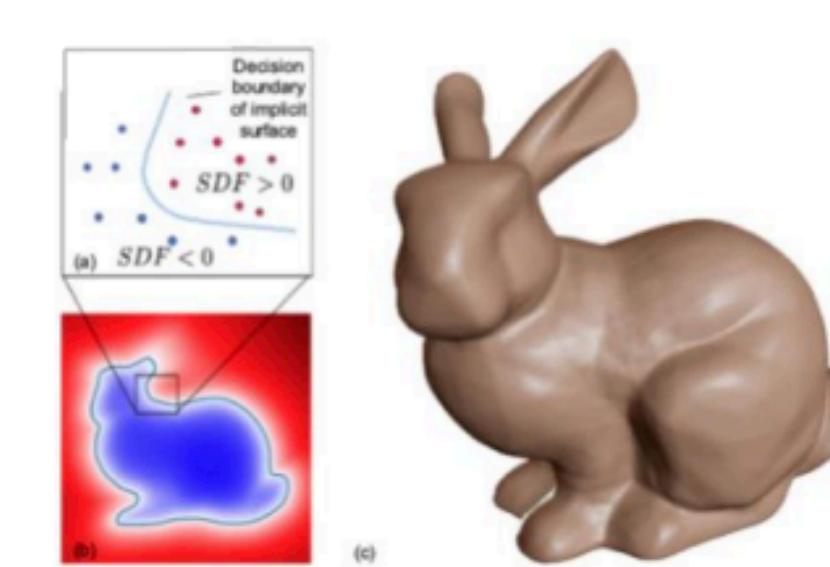


\*Work done as an intern at Waymo

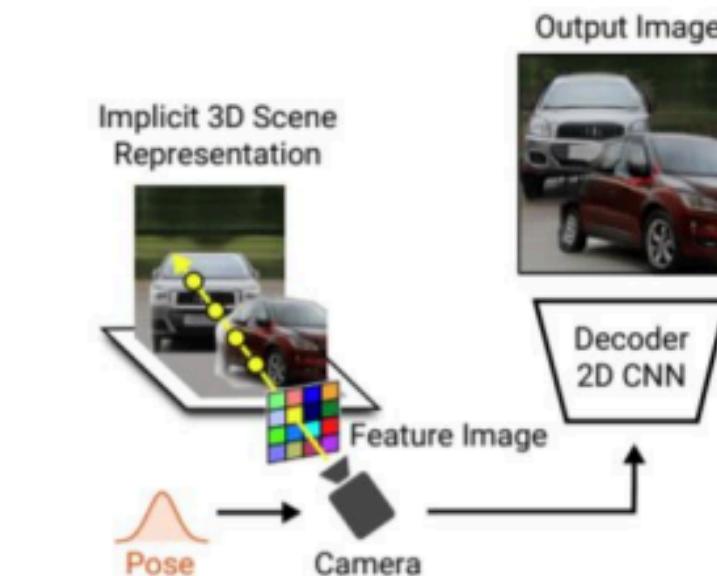
# Lots of other applications



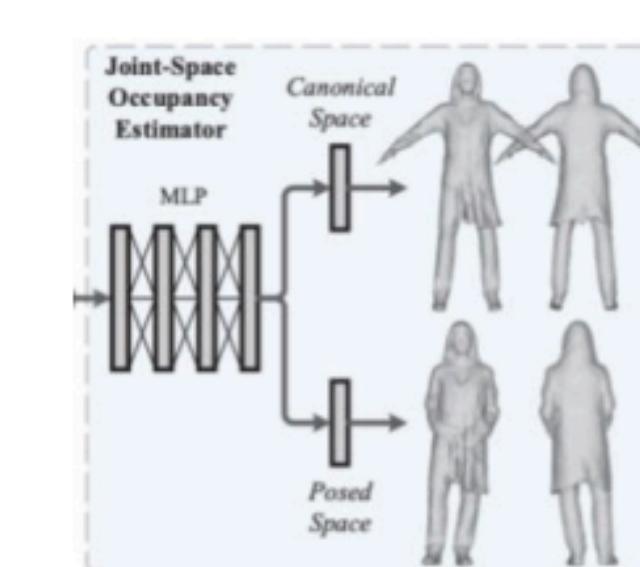
Neural field



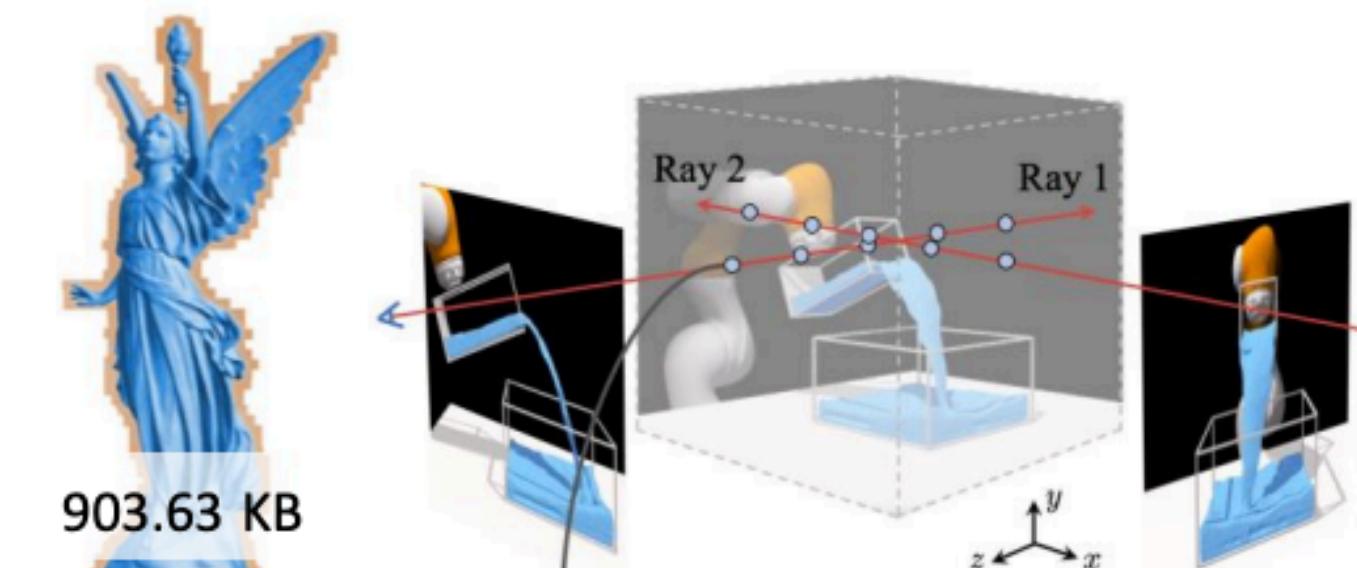
2D and 3D Reconstruction



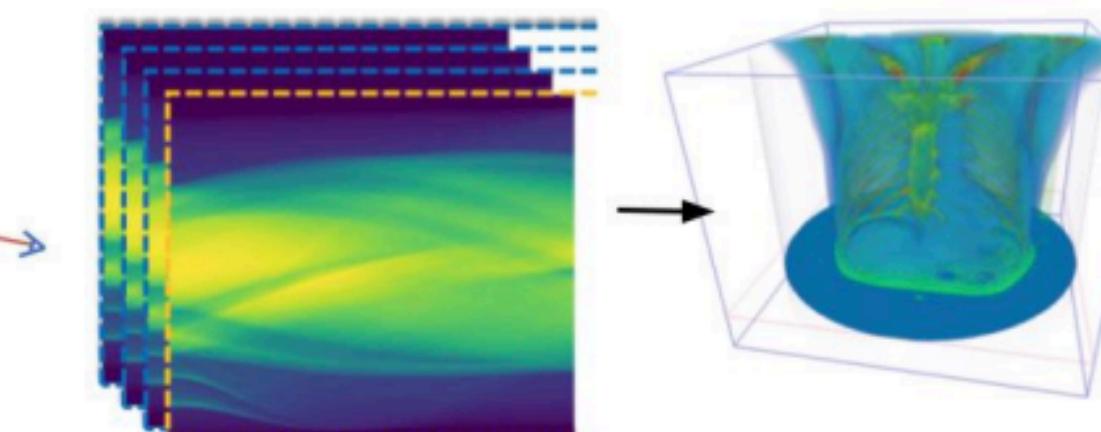
Generative Models



Digital Humans



Compression



Robotics

...and Beyond!

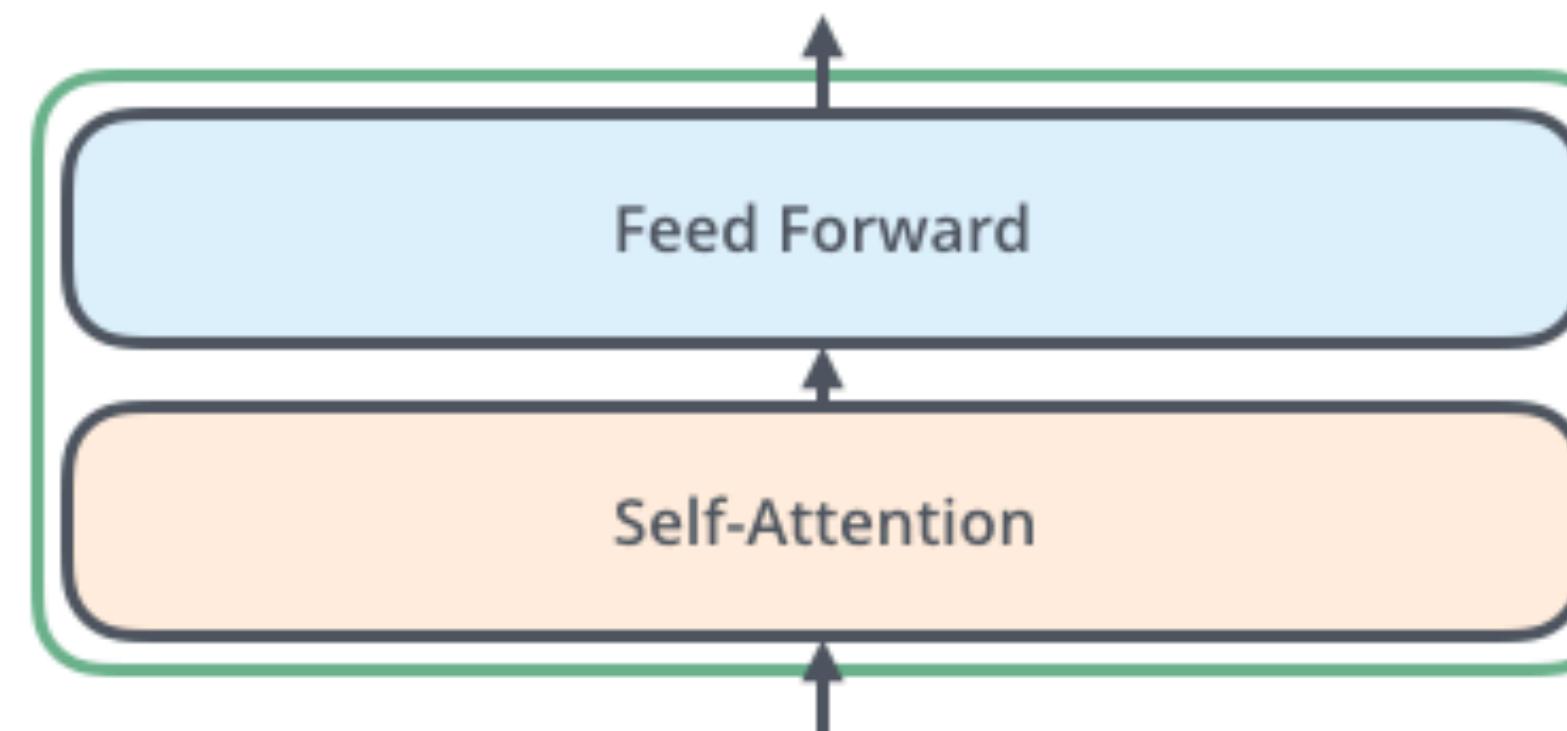
# Today

- Neural fields
- **Transformers for vision**

# Recall: Transformers

---

- Build whole model out of self-attention
- Uses only point-wise processing and attention (no recurrent units or convolutions)



# Self-attention layer

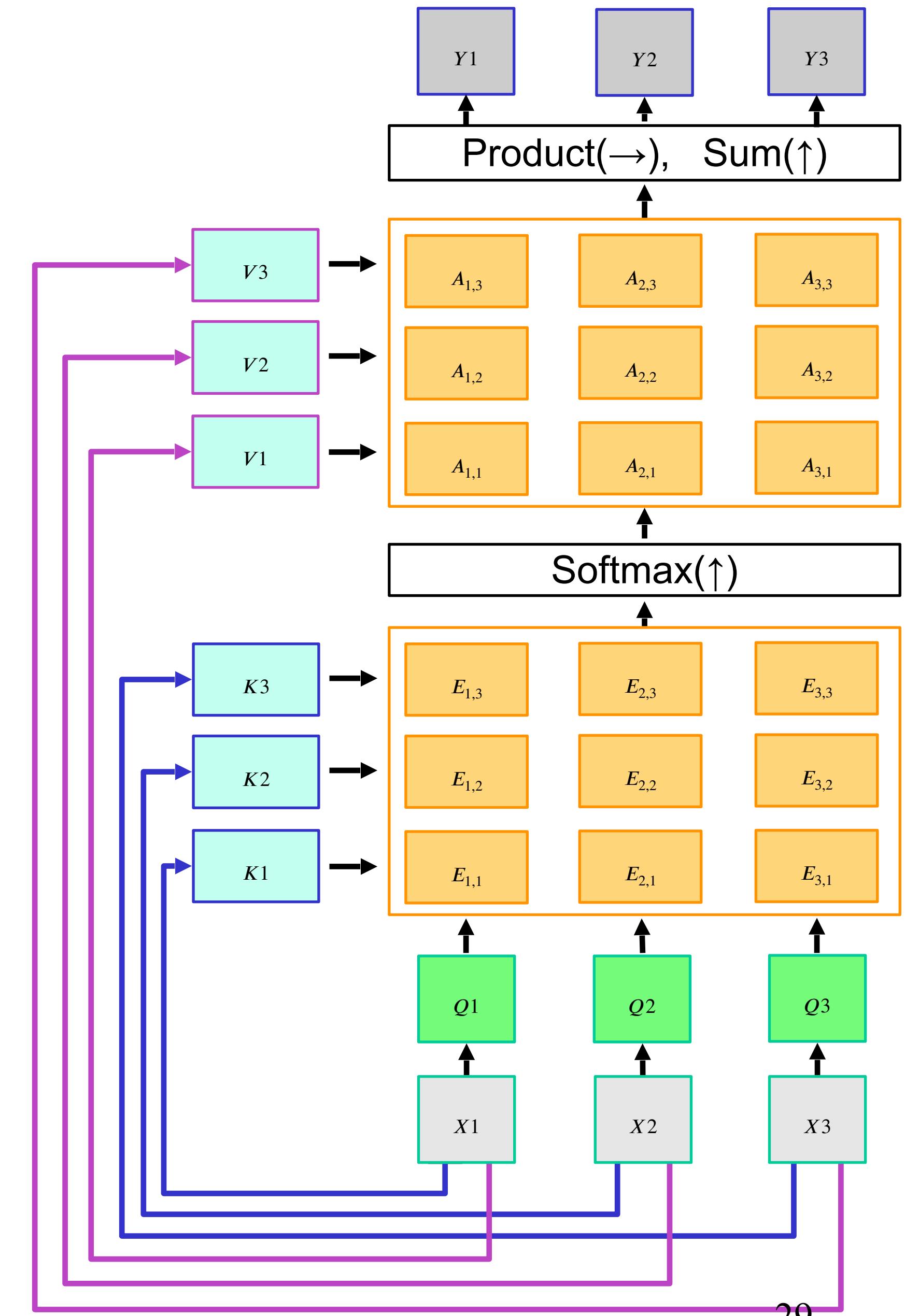
- Query vectors:  $Q = XW_Q$
- Key vectors:  $K = XW_K$
- Value vectors:  $V = XW_V$
- Similarities: *scaled dot-product attention*

$$E_{i,j} = \frac{(Q_i \cdot K_j)}{\sqrt{D}} \quad \text{or} \quad E = QK^T / \sqrt{D}$$

( $D$  is the dimensionality of the keys)

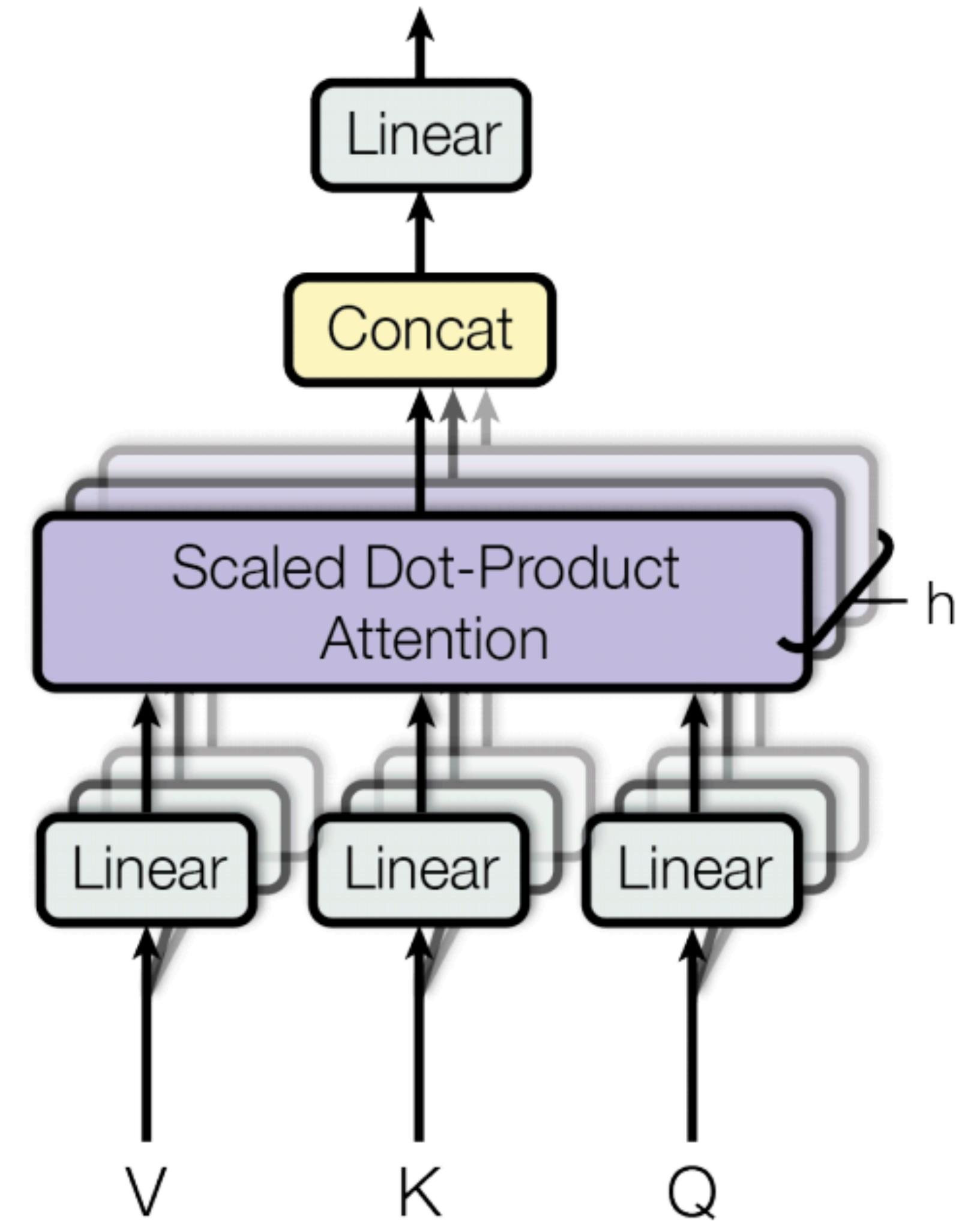
- Attn. weights:  $A = \text{softmax}(E, \text{ dim } = 1)$
- Output vectors:

$$Y_i = \sum_j A_{i,j} V_j \quad \text{or} \quad Y = AV$$



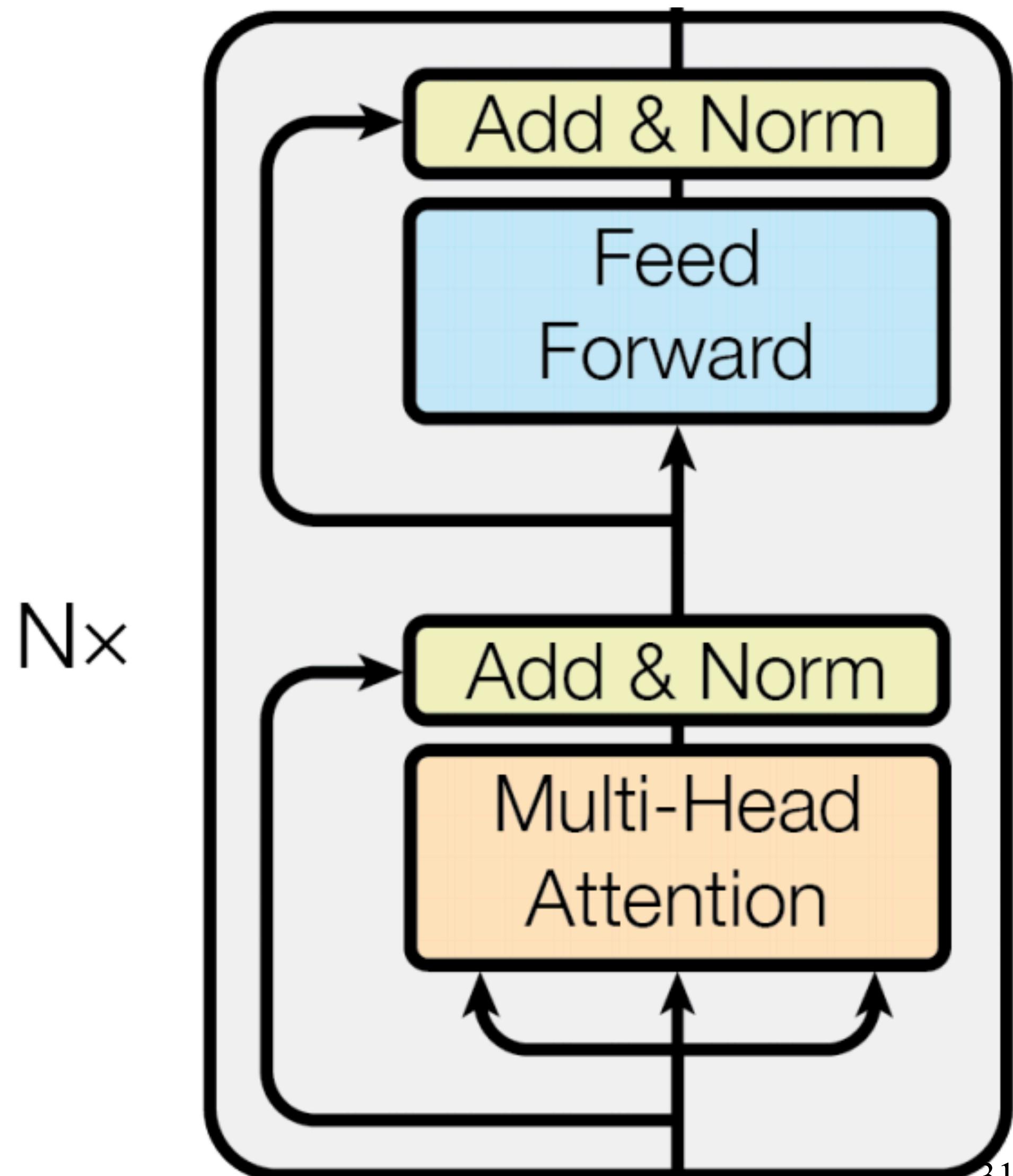
# Multi-head attention

- Run  $h$  attention models in parallel on top of different linearly projected versions of  $Q$ ,  $K$ ,  $V$ ; concatenate and linearly project the results
- Intuition: enables model to attend to different kinds of information at different positions



# Transformer blocks

- A **Transformer** is a sequence of transformer blocks
  - Vaswani et al.: N=12 blocks, embedding dimension = 512, 6 attention heads
  - **Add & Norm:** residual connection followed by layer normalization
  - **Feedforward:** two linear layers with ReLUs in between, applied independently to each vector
  - Attention is the only interaction between inputs!



# Self-supervised learning in Natural Language Processing

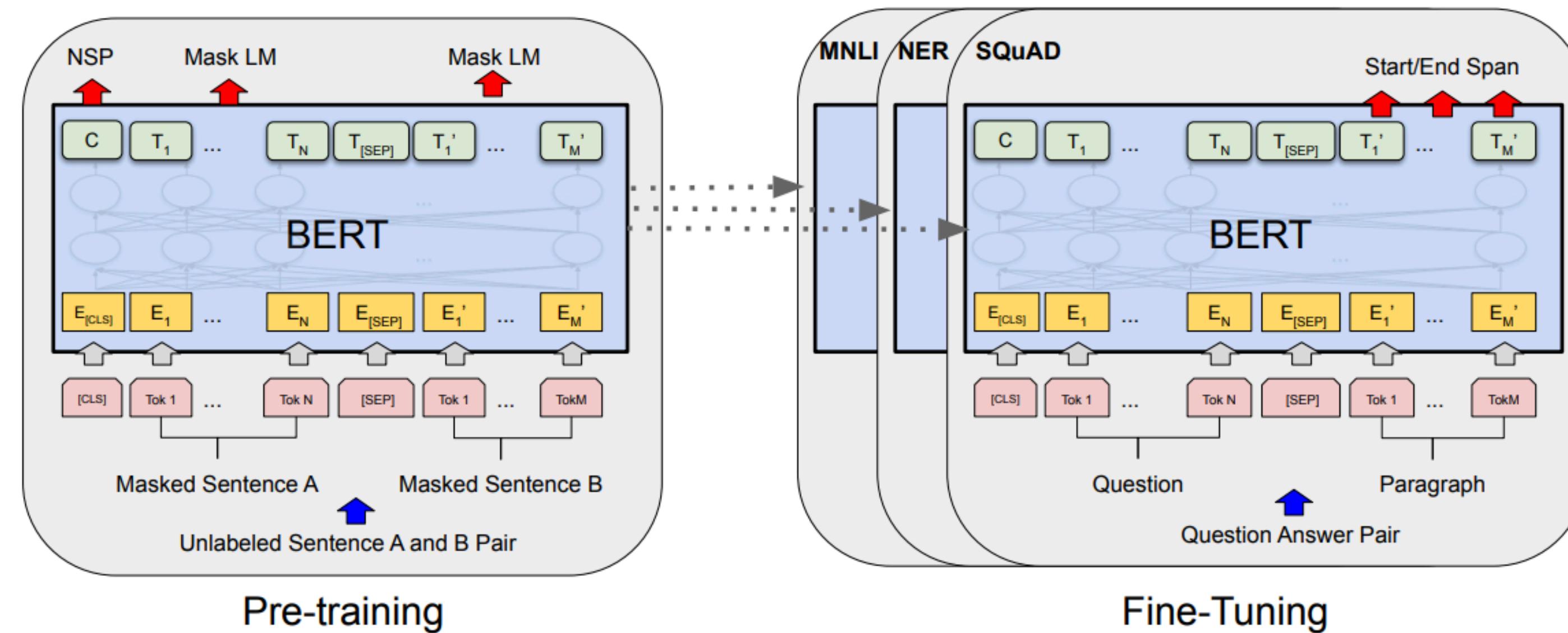
1. Download A LOT of text from the internet
2. Train a giant transformer using a suitable pretext task
3. Fine-tune the transformer on desired NLP task

Model Alias	Org.	Article Reference
ULMfit	fast.ai	<i>Universal Language Model Fine-tuning for Text Classification</i> Howard and Ruder
 ELMo	AllenNLP	<i>Deep contextualized word representations</i> Peters et al.
OpenAI GPT	OpenAI	<i>Improving Language Understanding by Generative Pre-Training</i> Radford et al.
 BERT	Google	<i>BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding</i> Devlin et al.
XLM	Facebook	<i>Cross-lingual Language Model Pretraining</i> Lample and Conneau

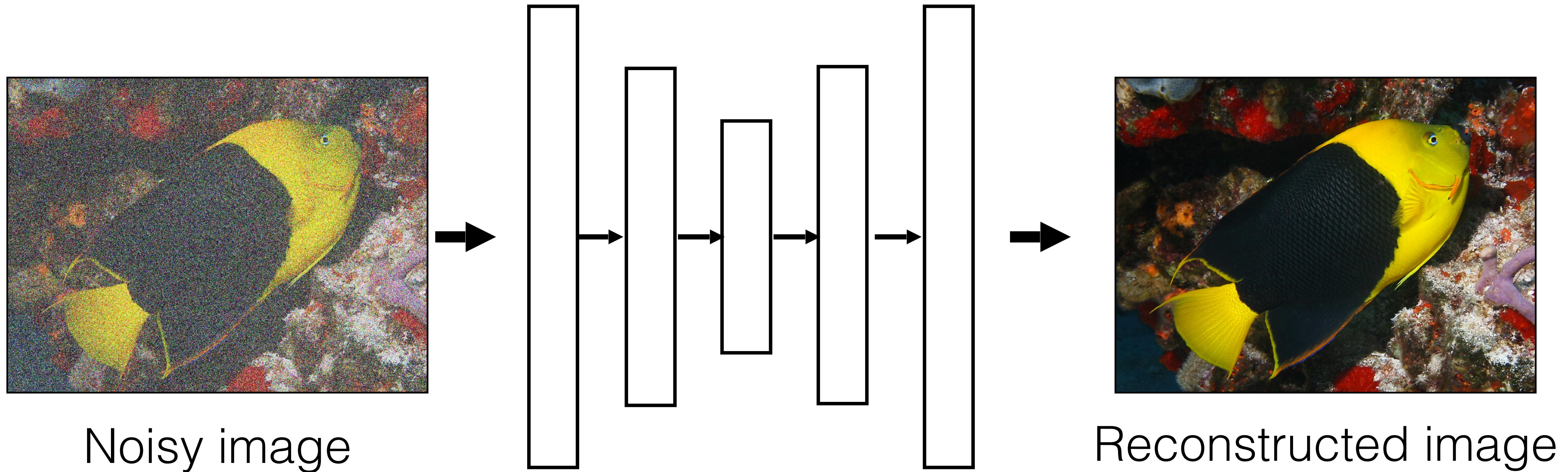
# Self-supervised language modeling with transformers

1. Download A LOT of text from the internet
2. Train a giant transformer using a suitable pretext task
3. Fine-tune the transformer on desired NLP task

## Bidirectional Encoder Representations from Transformers (BERT)

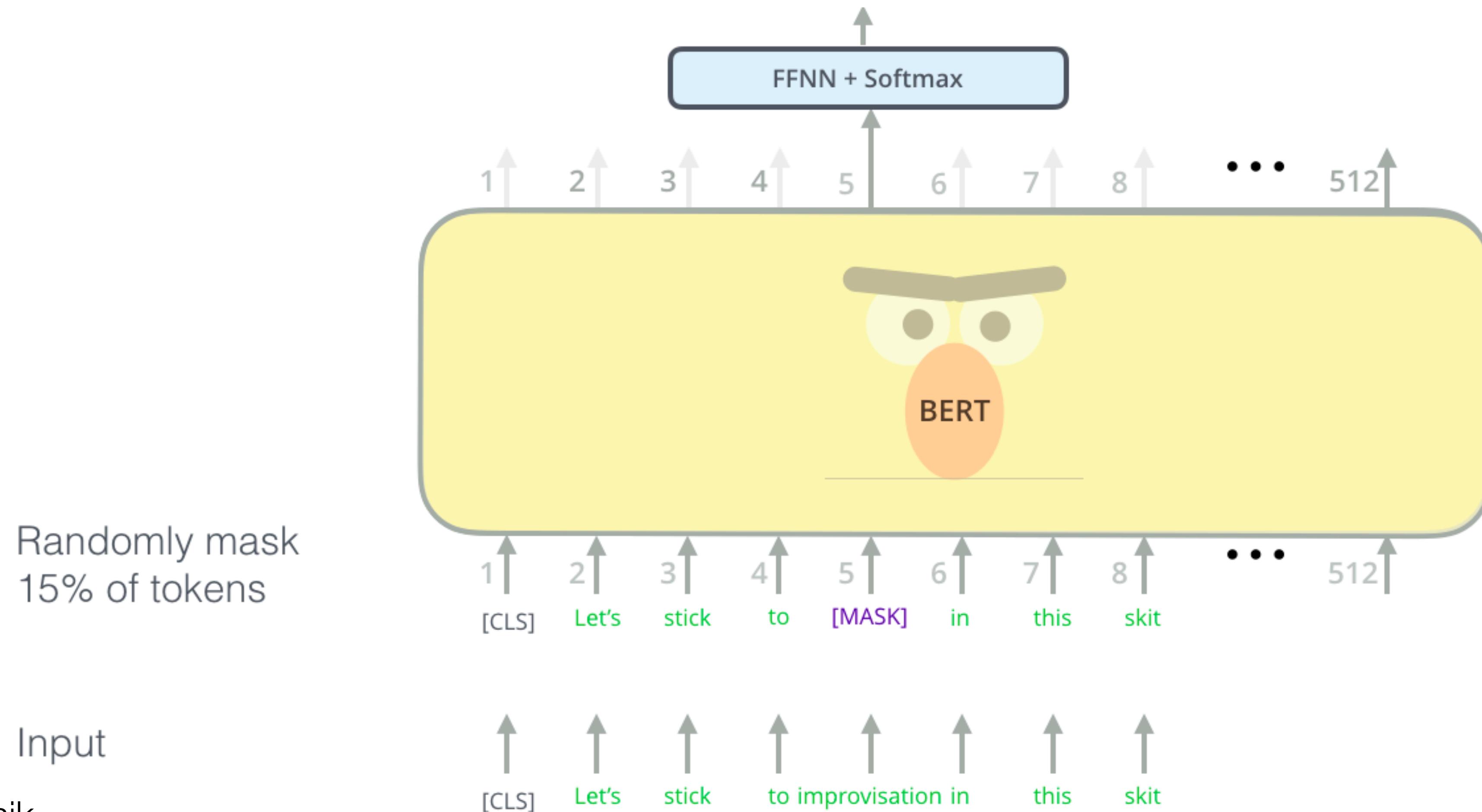


# Recall: denoising autoencoder

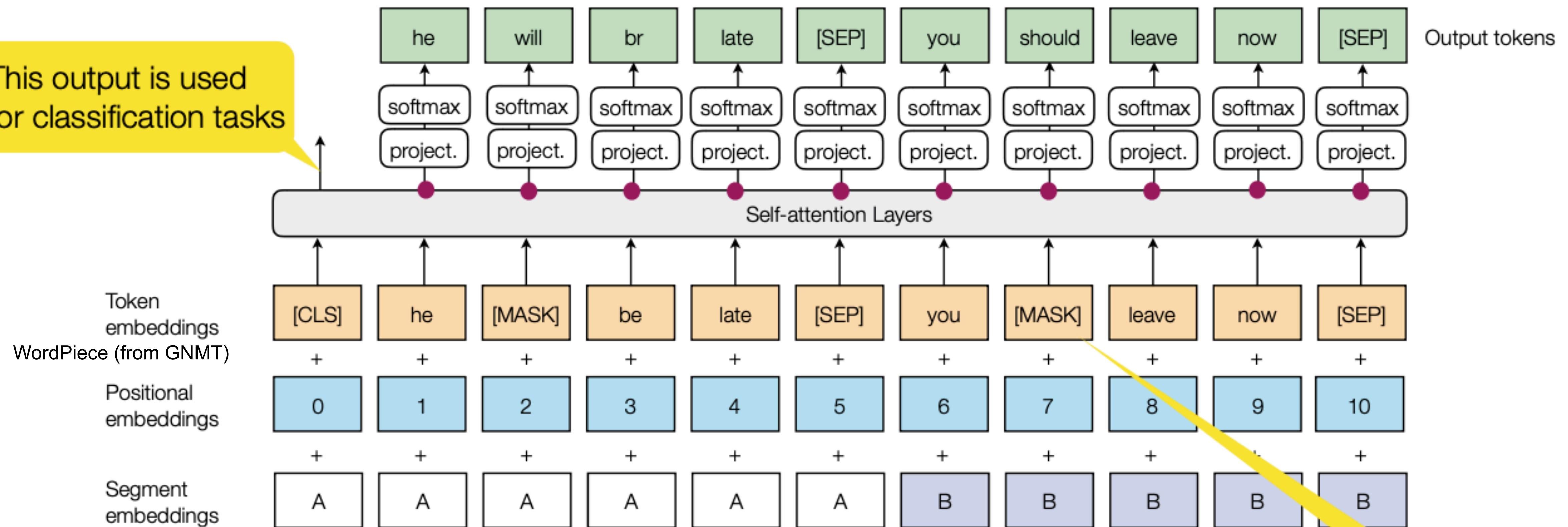


# BERT: Pretext tasks

- Masked language model (MLM)
  - Randomly mask 15% of tokens in input sentences, goal is to reconstruct them using bidirectional context



# BERT: More detailed view



Trained on Wikipedia (2.5B words) + BookCorpus (800M words)

15% of tokens  
get masked

# BERT: Evaluation

---

- General Language Understanding Evaluation (GLUE) benchmark ([gluebenchmark.com](http://gluebenchmark.com))

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

# Image GPT

- Image resolution up to 64x64, color values quantized to 512 levels (9 bits), dense attention
- For transfer learning, average-pool encoded features across all positions

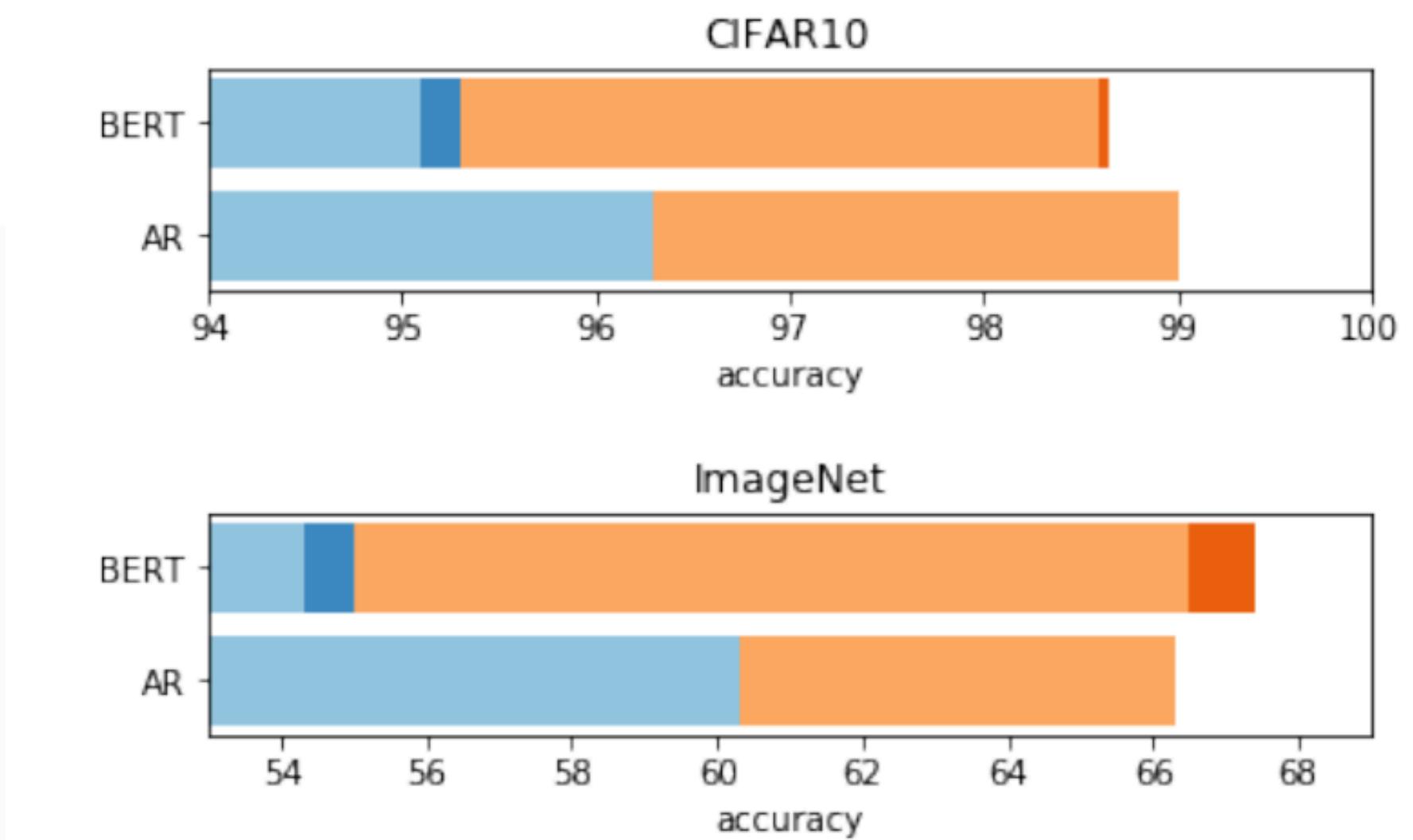
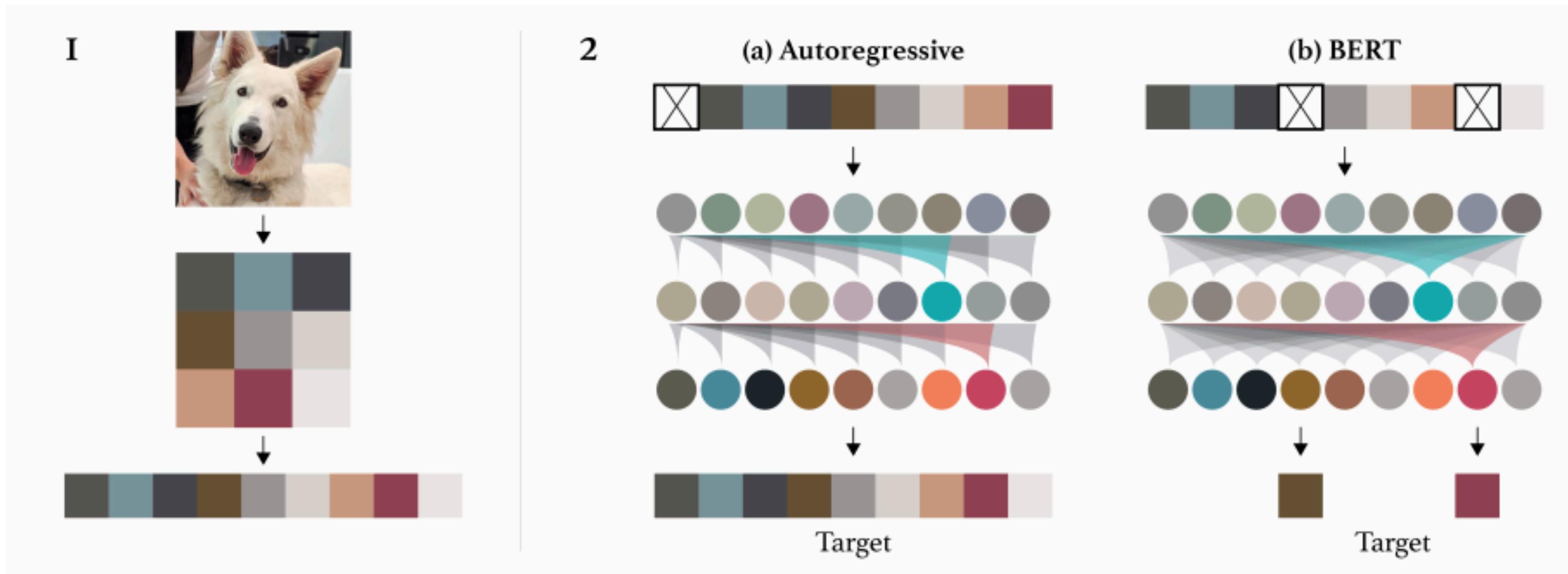
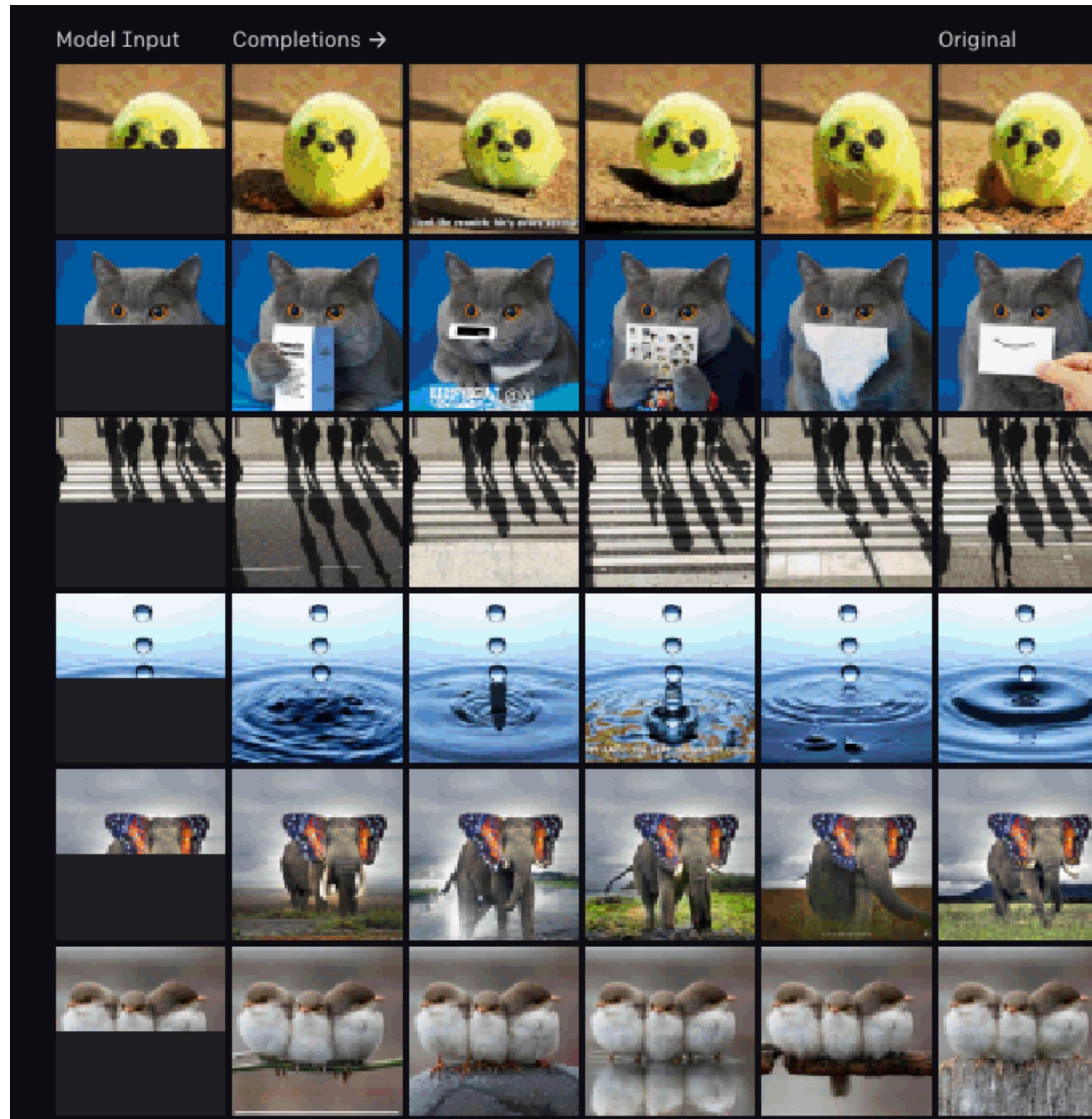


Figure 4. Comparison of auto-regressive pre-training with BERT pre-training using iGPT-L at an input resolution of  $32^2 \times 3$ . Blue bars display linear probe accuracy and orange bars display fine-tune accuracy. Bold colors show the performance boost from ensembling BERT masks. We see that auto-regressive models produce much better features than BERT models after pre-training, but BERT models catch up after fine-tuning.

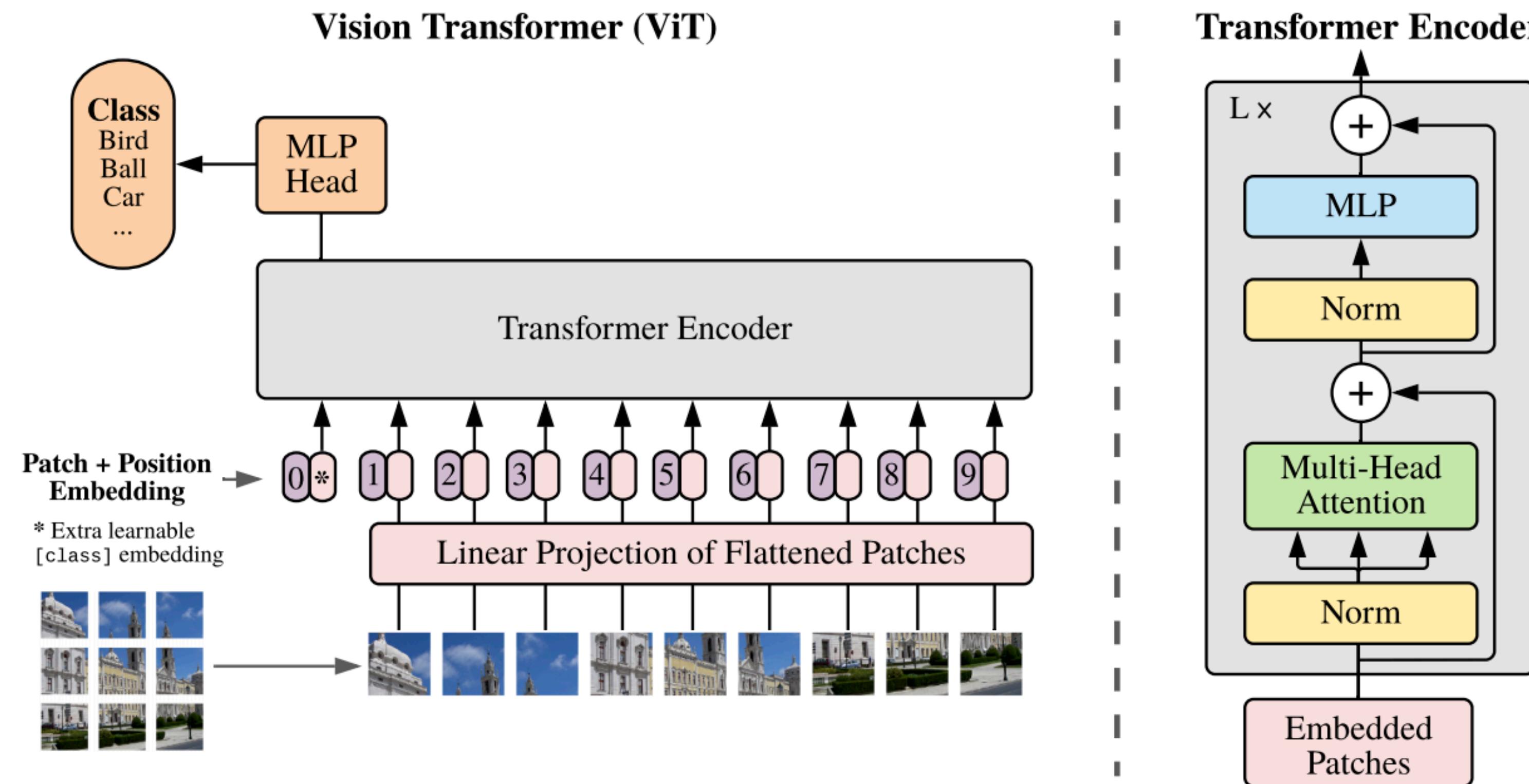
# Image GPT – OpenAI



<https://openai.com/blog/image-gpt/>

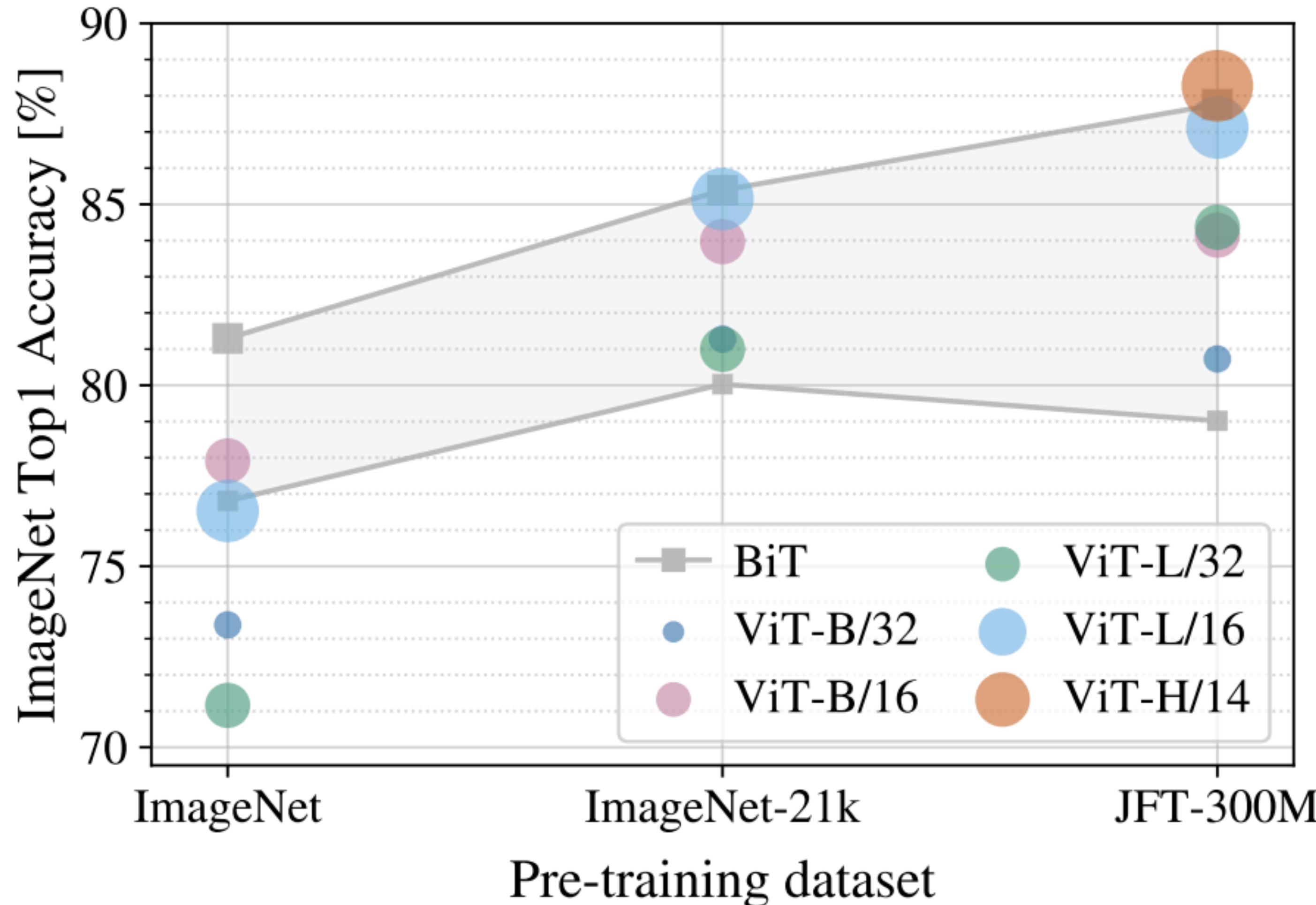
# Vision transformer (ViT)

- Split an image into patches, feed linearly projected patches into standard transformer encoder
  - With patches of 14x14 pixels, you need  $16 \times 16 = 256$  patches to represent 224x224 images



A. Dosovitskiy et al. [An image is worth 16x16 words: Transformers for image recognition at scale](#). ICLR 2021

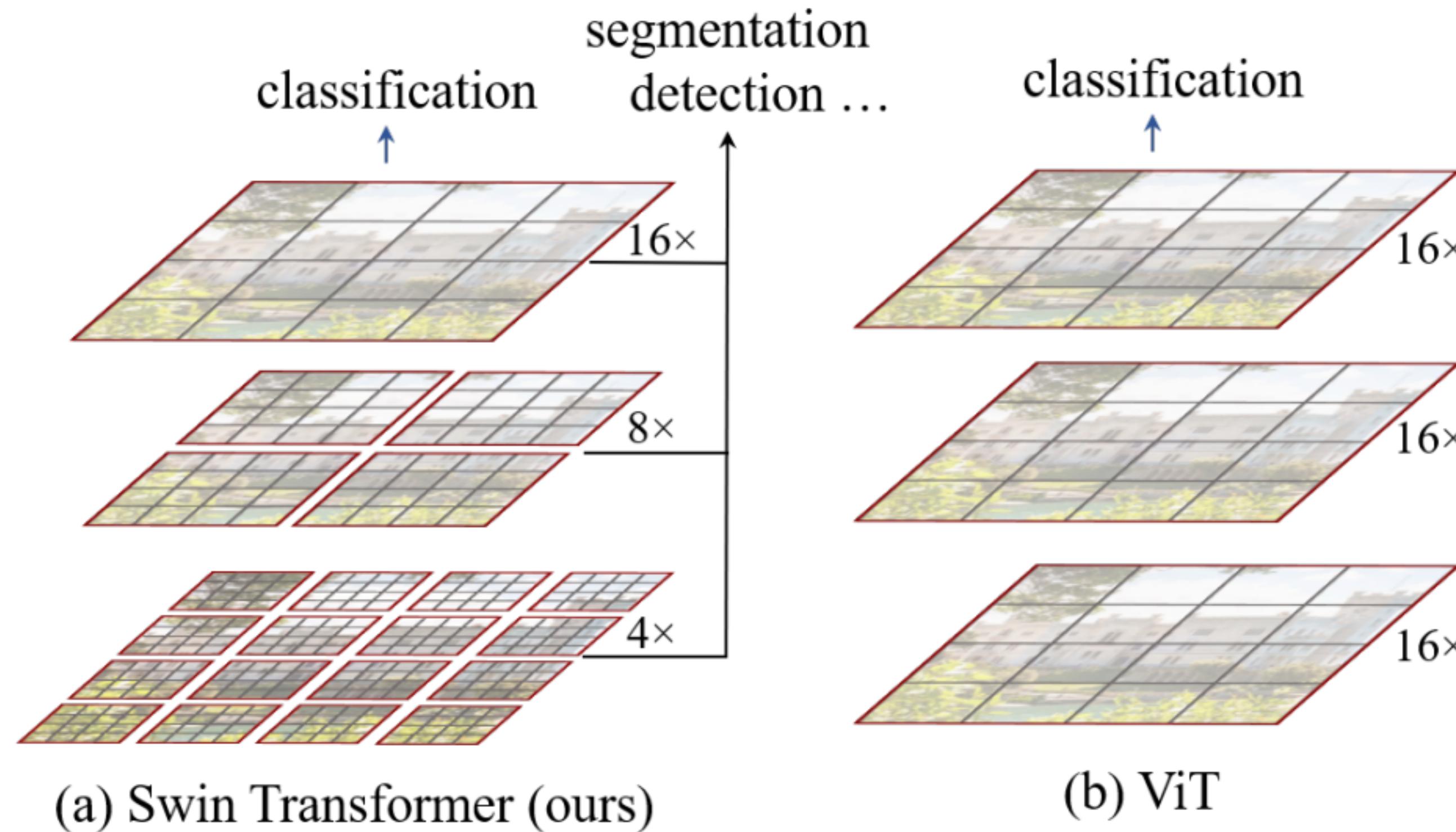
# Vision transformer (ViT)



BiT: [Big Transfer](#) (ResNet)  
ViT: Vision Transformer (Base/Large/Huge,  
patch size of 14x14, 16x16, or 32x32)

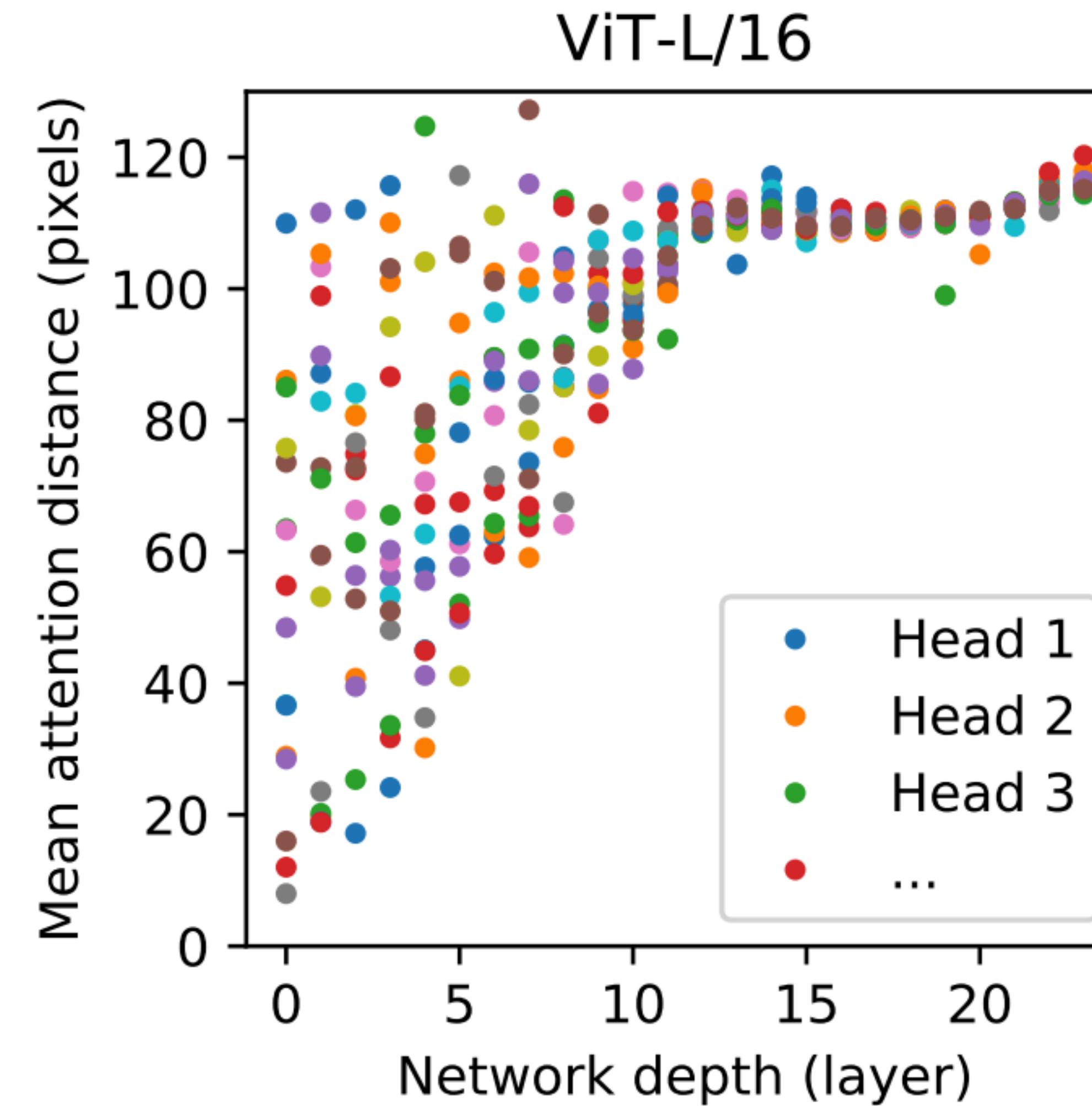
[Internal Google dataset](#) (not public)

# Swin Transformer: windowed attention

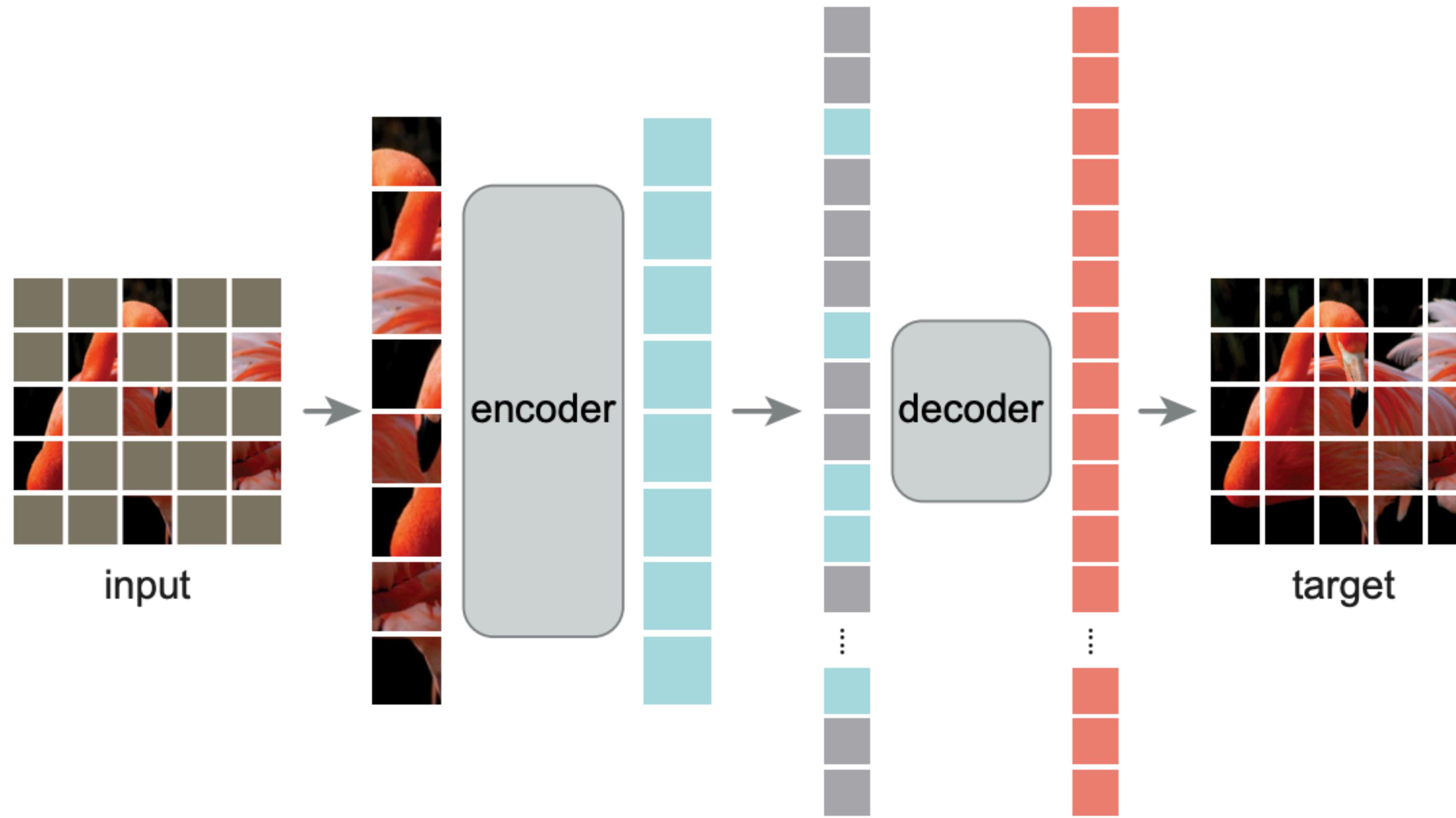


[Liu et al., “Swin Transformer: Hierarchical Vision Transformer using Shifted Windows”]

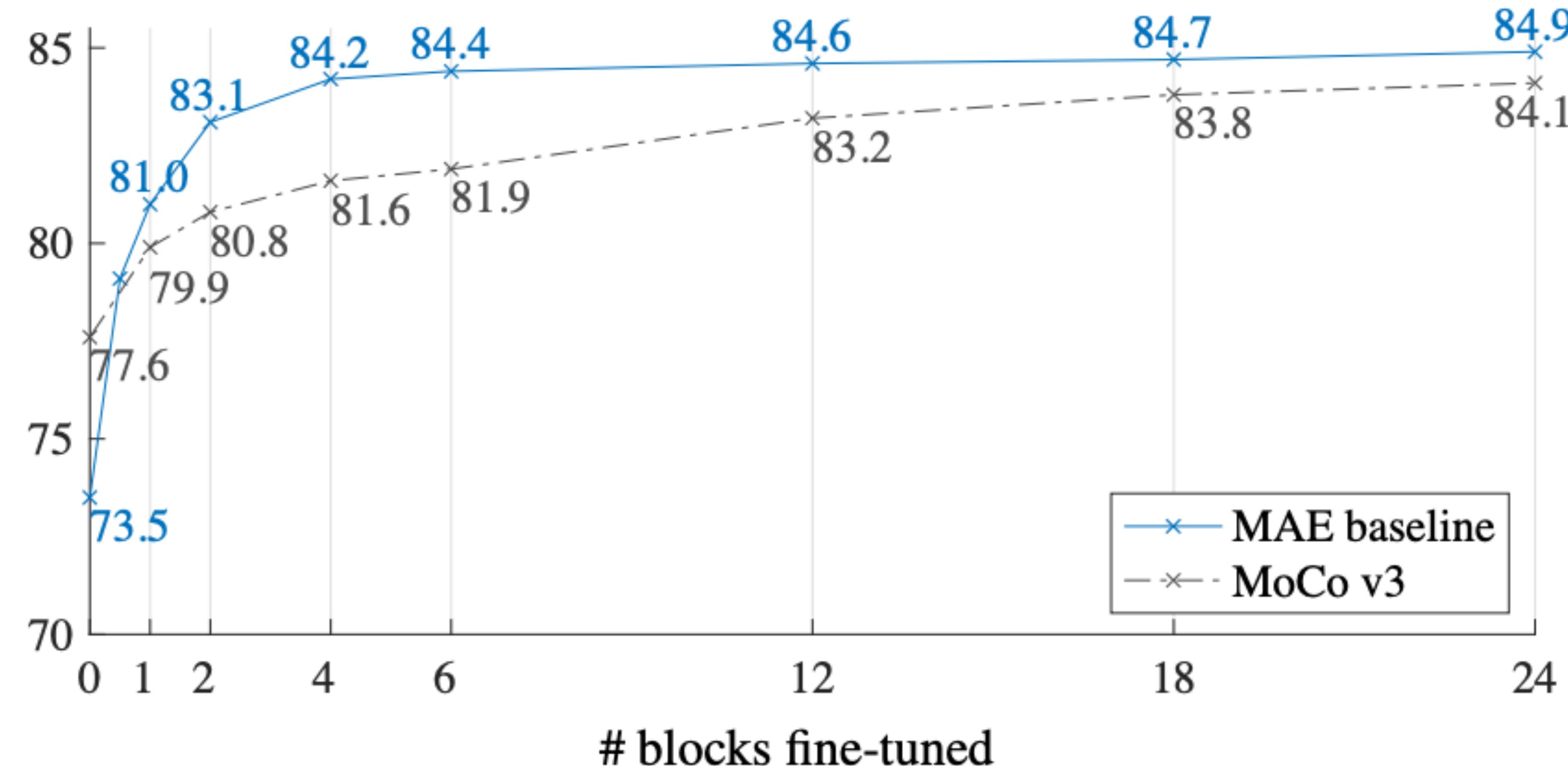
# Vision transformer (ViT)



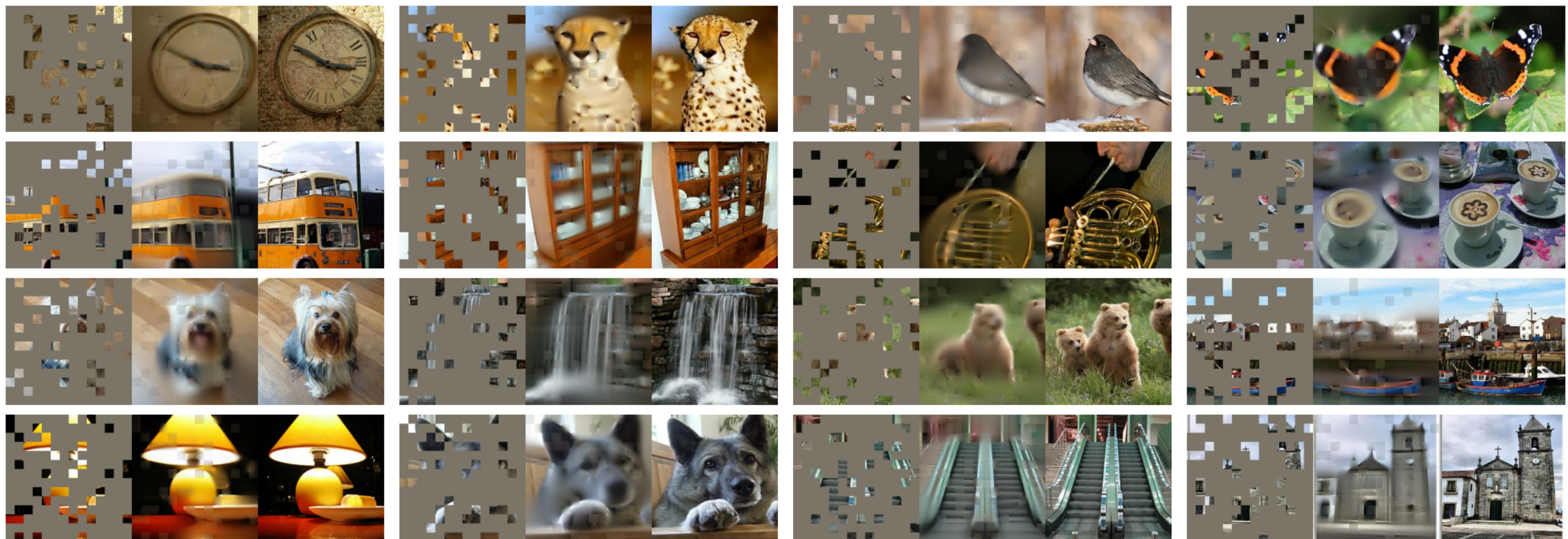
# Application: self-supervised learning



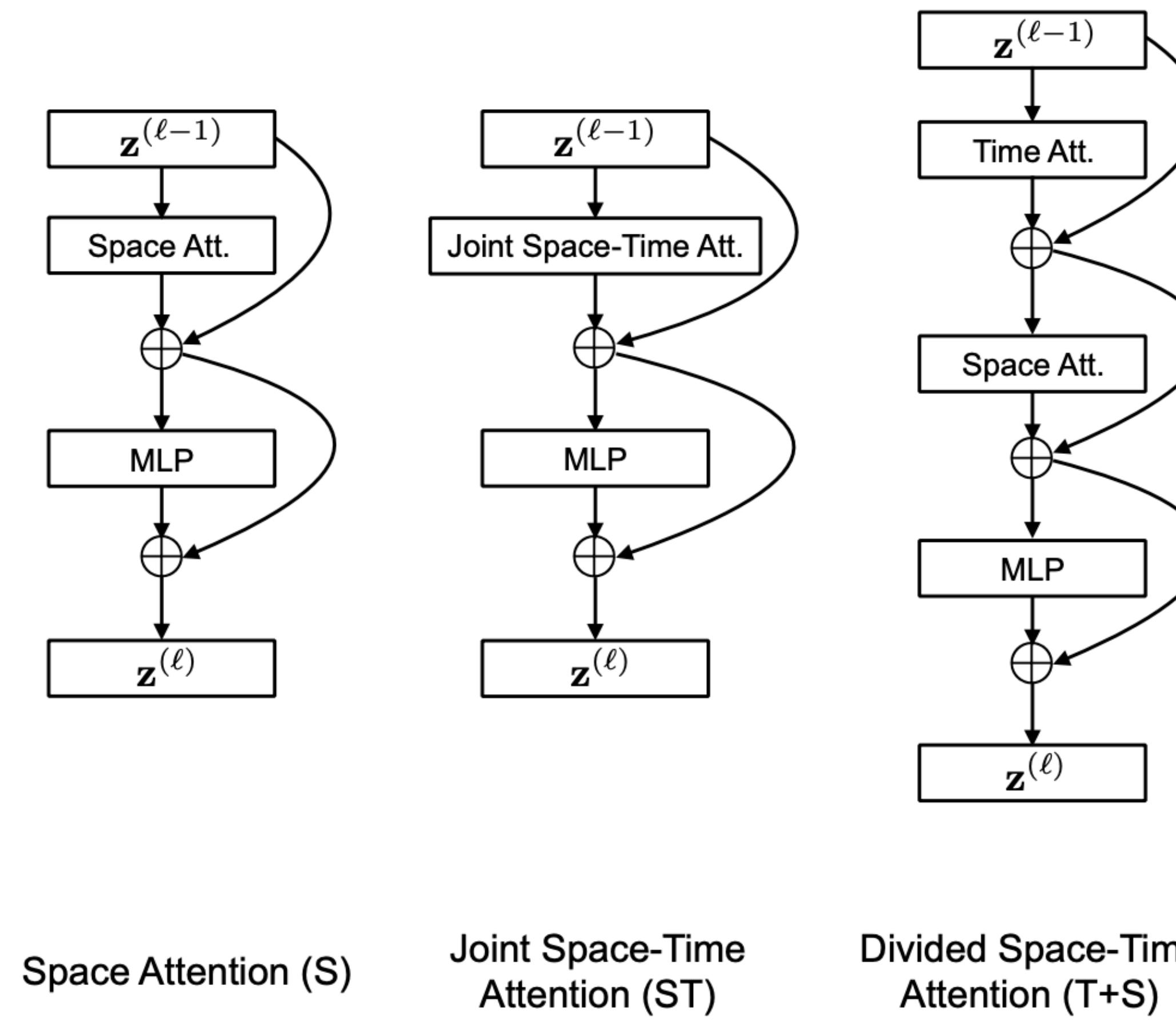
# Application: self-supervised learning



# Application: self-supervised learning



# Application: video

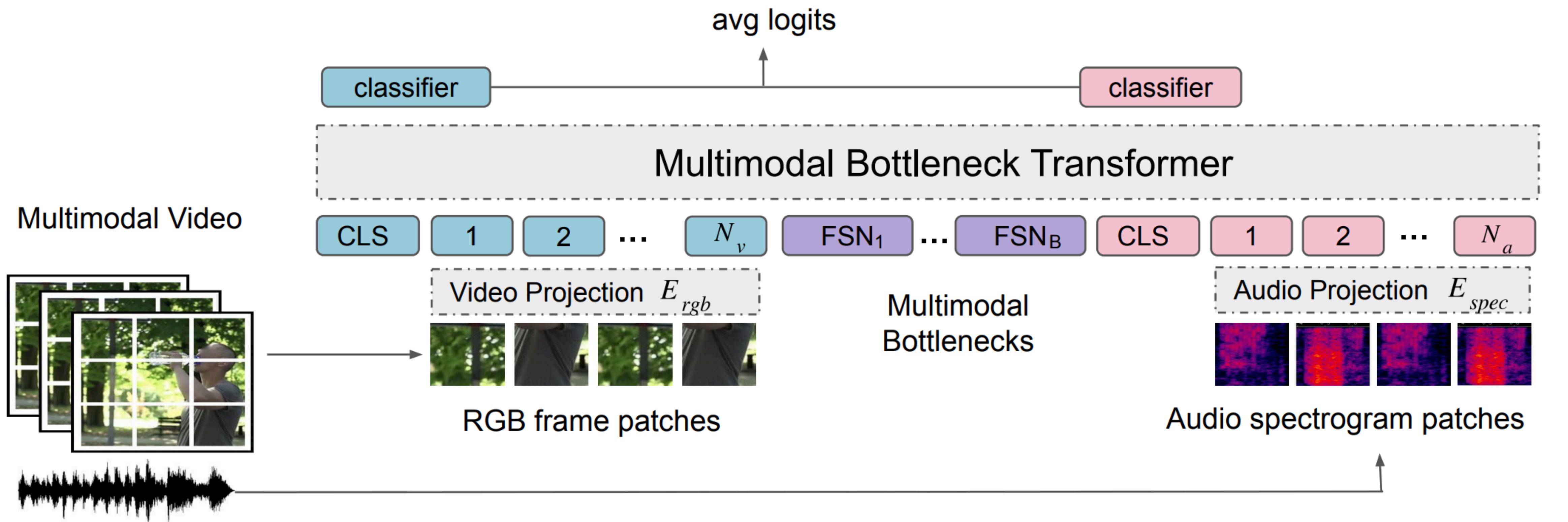


Space Attention (S)

Joint Space-Time  
Attention (ST)

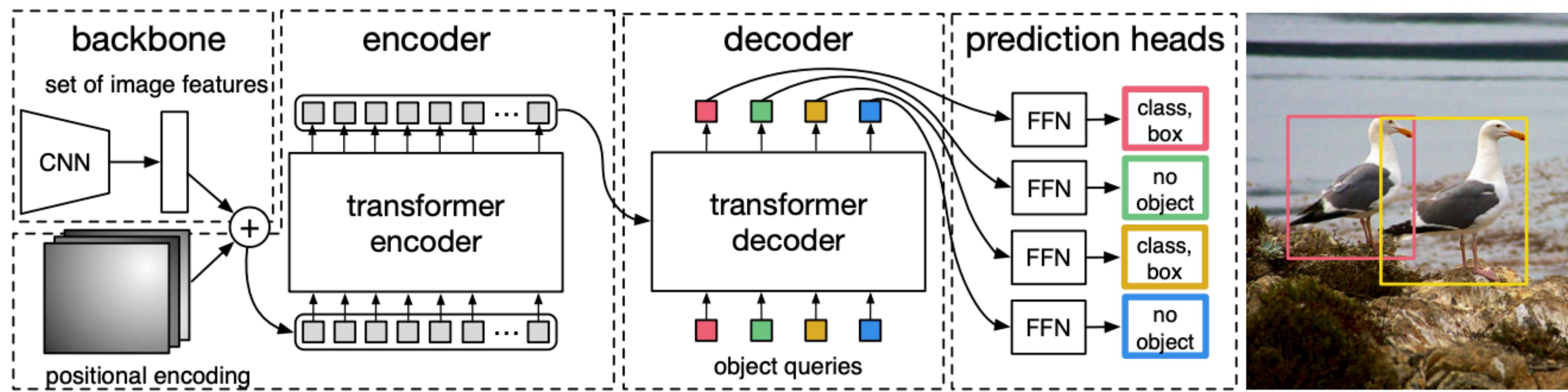
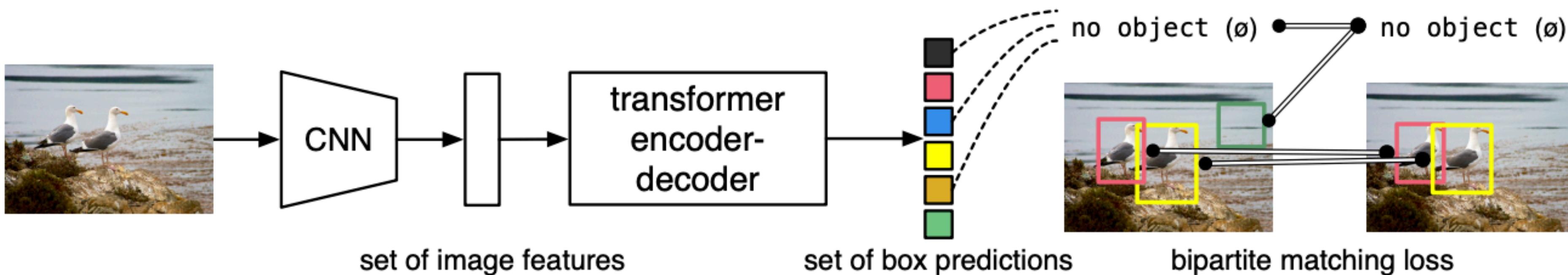
Divided Space-Time  
Attention (T+S)

# Application: multimodal models



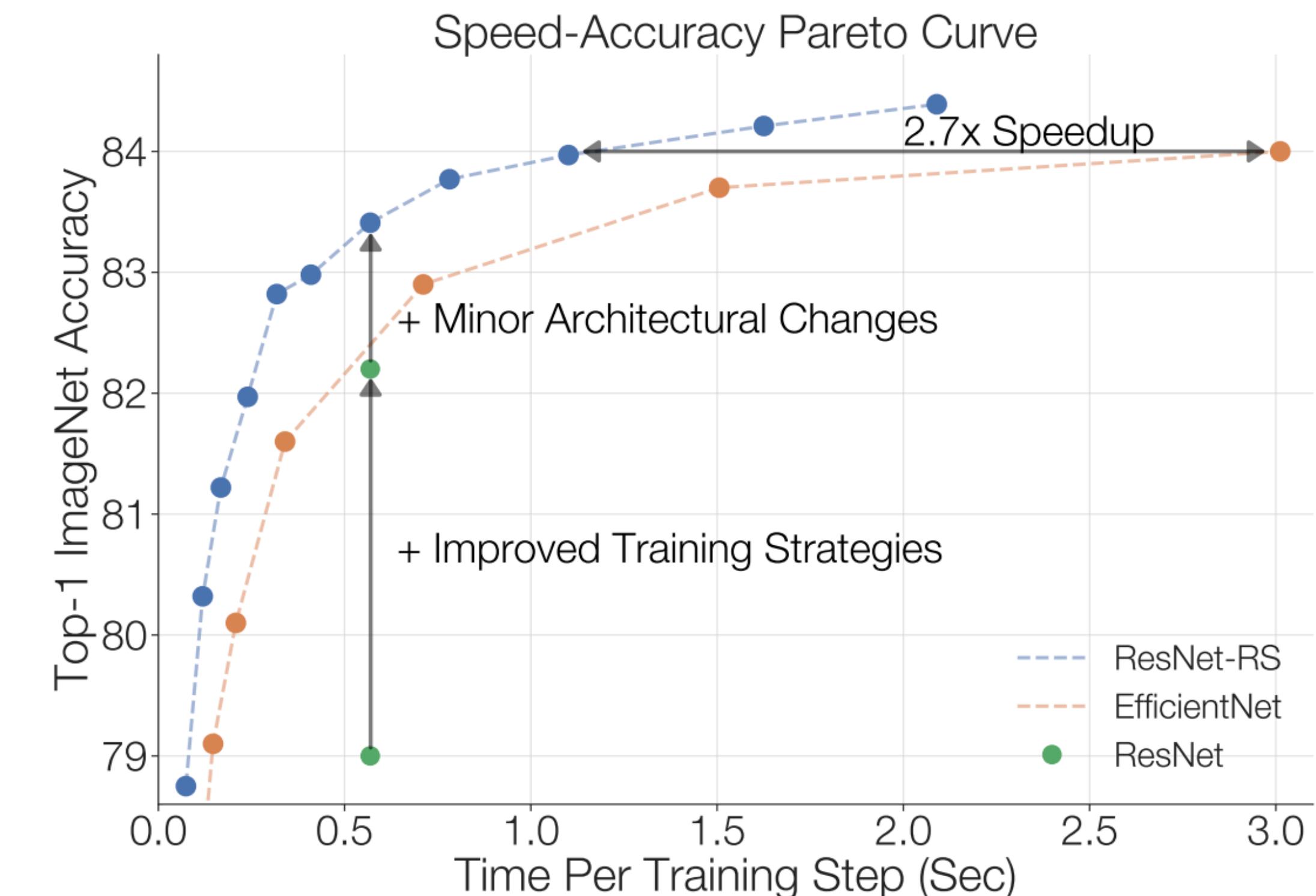
# Detection Transformer (DETR) – Facebook AI

- Hybrid of CNN and transformer, aimed at standard recognition task

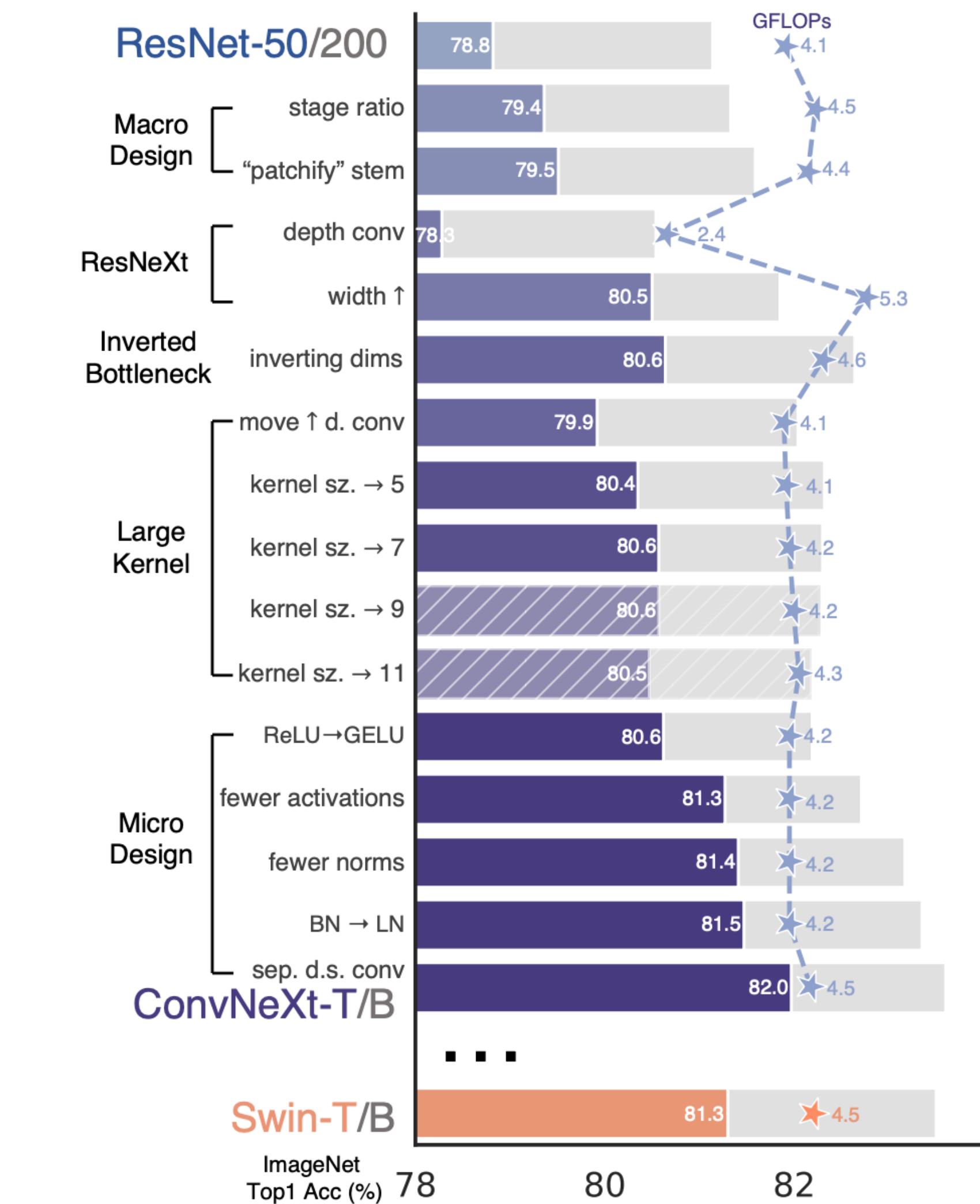
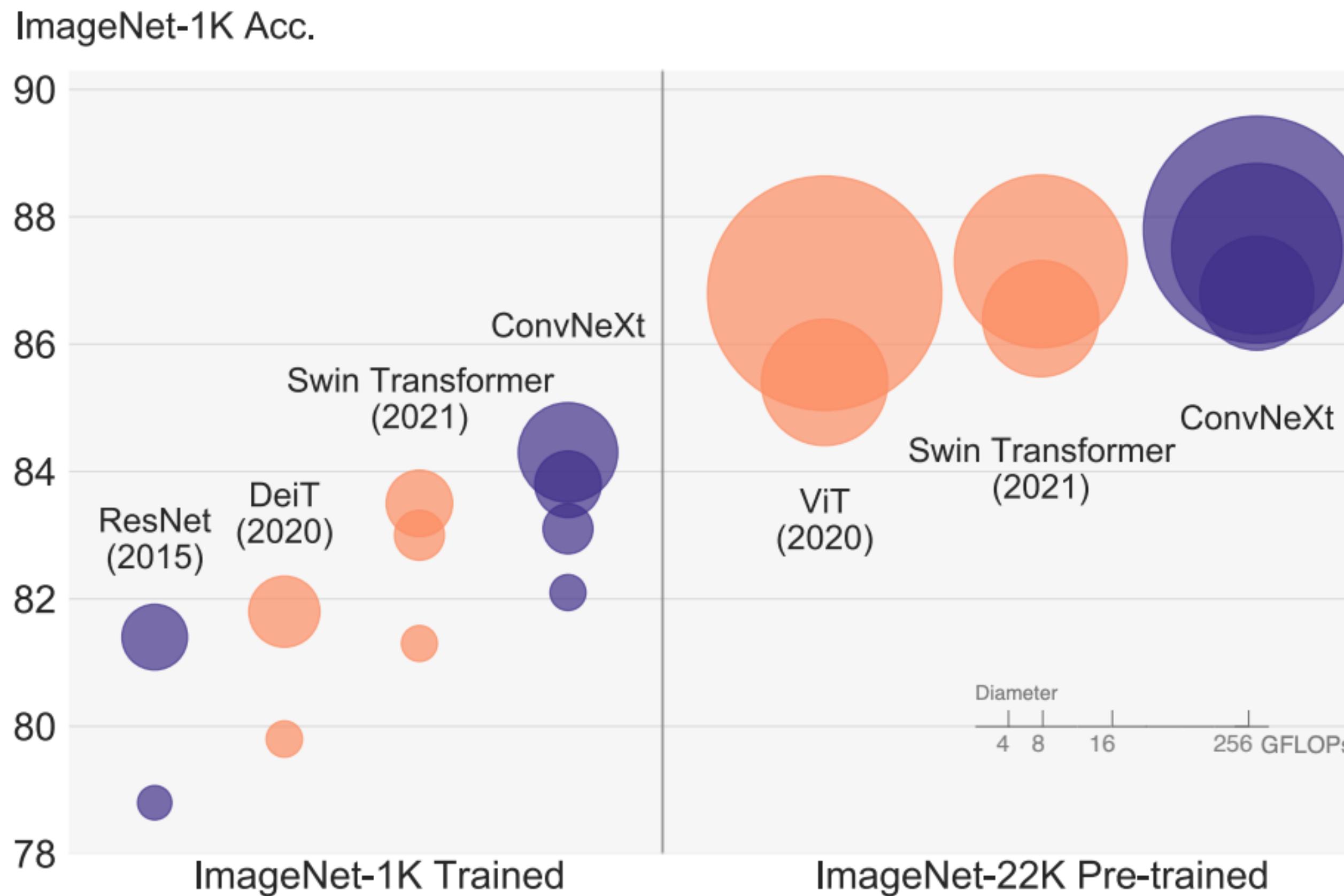


# CNNs can be improved too, though.

Improvements	Top-1	$\Delta$
ResNet-200	79.0	—
+ Cosine LR Decay	79.3	+0.3
+ Increase training epochs	78.8 †	-0.5
+ EMA of weights	79.1	+0.3
+ Label Smoothing	80.4	+1.3
+ Stochastic Depth	80.6	+0.2
+ RandAugment	81.0	+0.4
+ Dropout on FC	80.7 ‡	-0.3
+ Decrease weight decay	82.2	+1.5
+ Squeeze-and-Excitation	82.9	+0.7
+ ResNet-D	83.4	+0.5



# CNNs can be improved too, though.



# Next class: Bias and ethics