# Lecture 23: Recent architectures

Transformer slides from S. Lazebnik

# Announcements

- Mon. Dec. 11: 10:30am - 1:30pm in FXB1109
- Mon. Dec. 11:   2:30pm - 4:30pm over Zoom
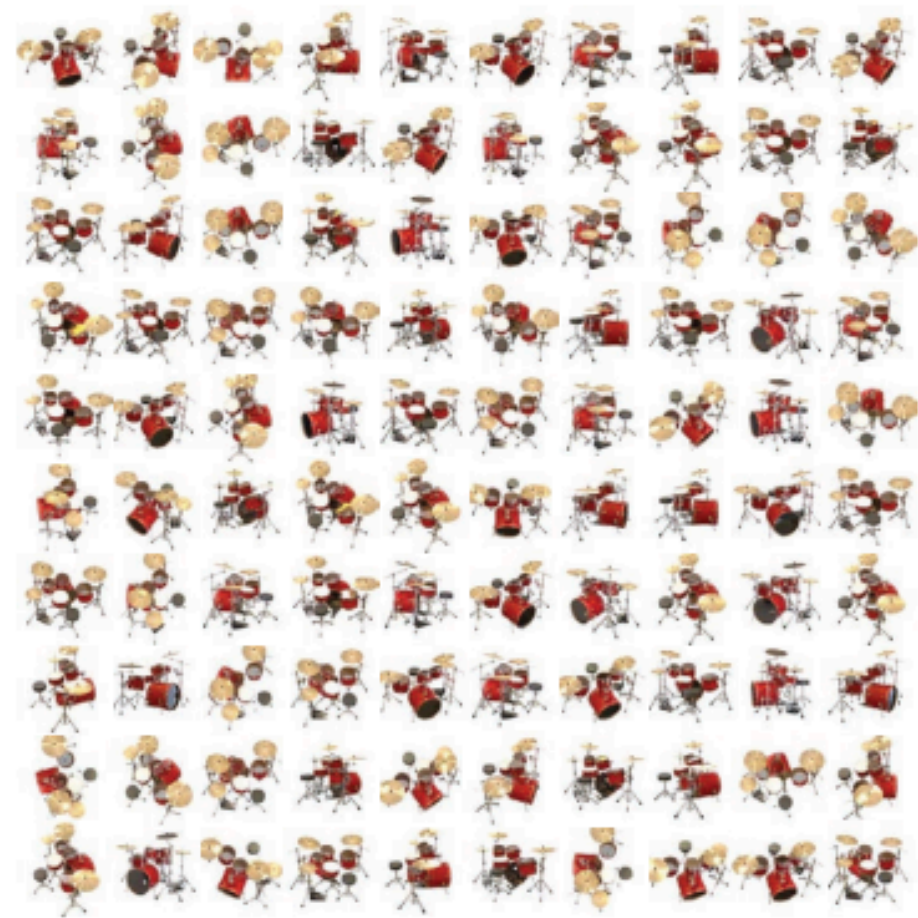- Tues. Dec. 12: 10:30am - 1:30pm in FXB1109

- In-person split in two halves (e.g., 10:30-12pm, 12pm - 1:30pm), so you can leave for lunch if you want.

- Video submission option for those who can't make it (due on Dec. 11 at **noon**), with explanation for why you can't come.

# Today

- Neural fields
- Transformers for vision

# 3D view synthesis



Input views → Create model → Render new views

## What representation should we use?

[Source: Mildenhall et al., "NeRF", 2020]

# Idea #1: Image-based rendering



→ View from a different angle

Point cloud
(reconstructed with SfM + multi-view stereo)
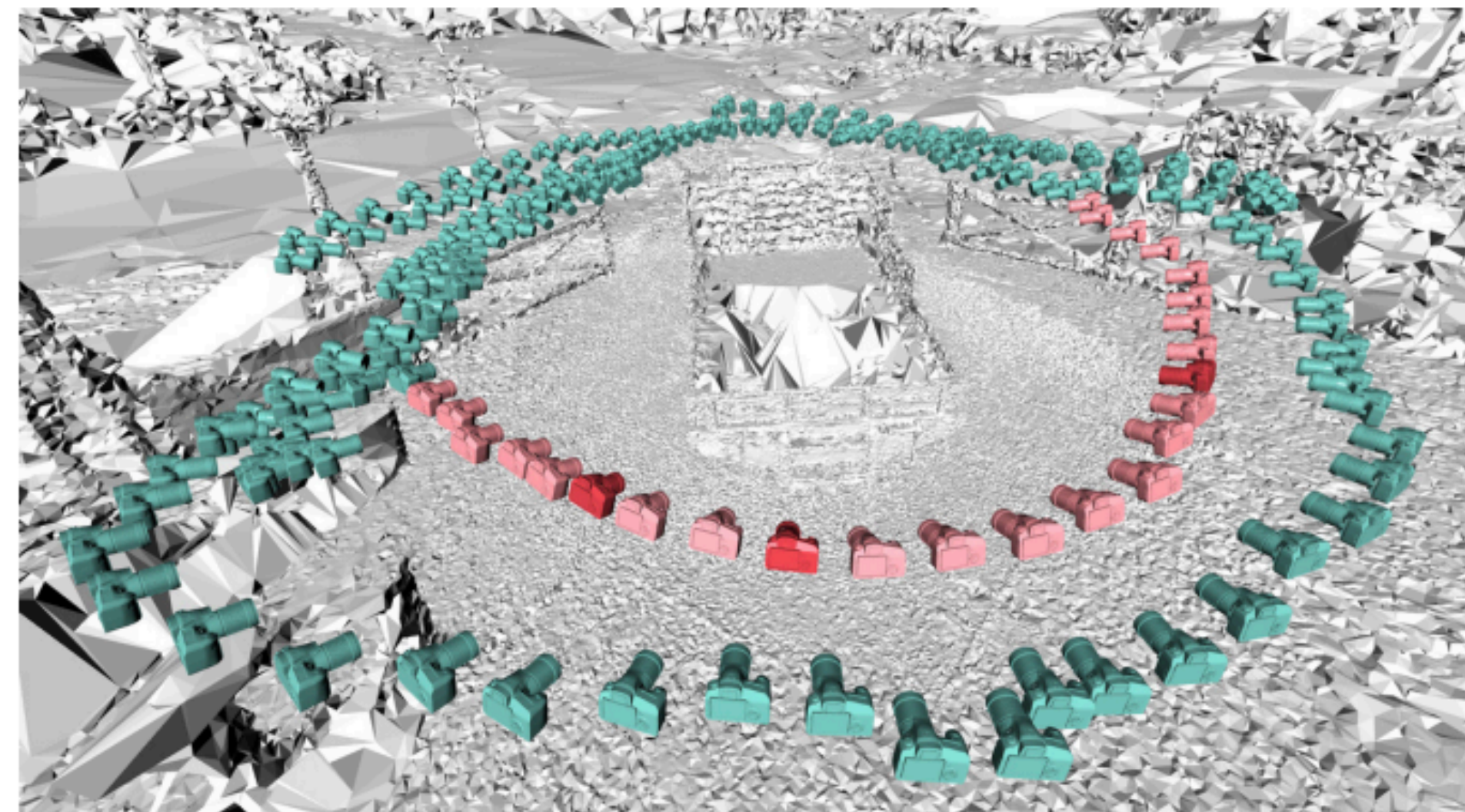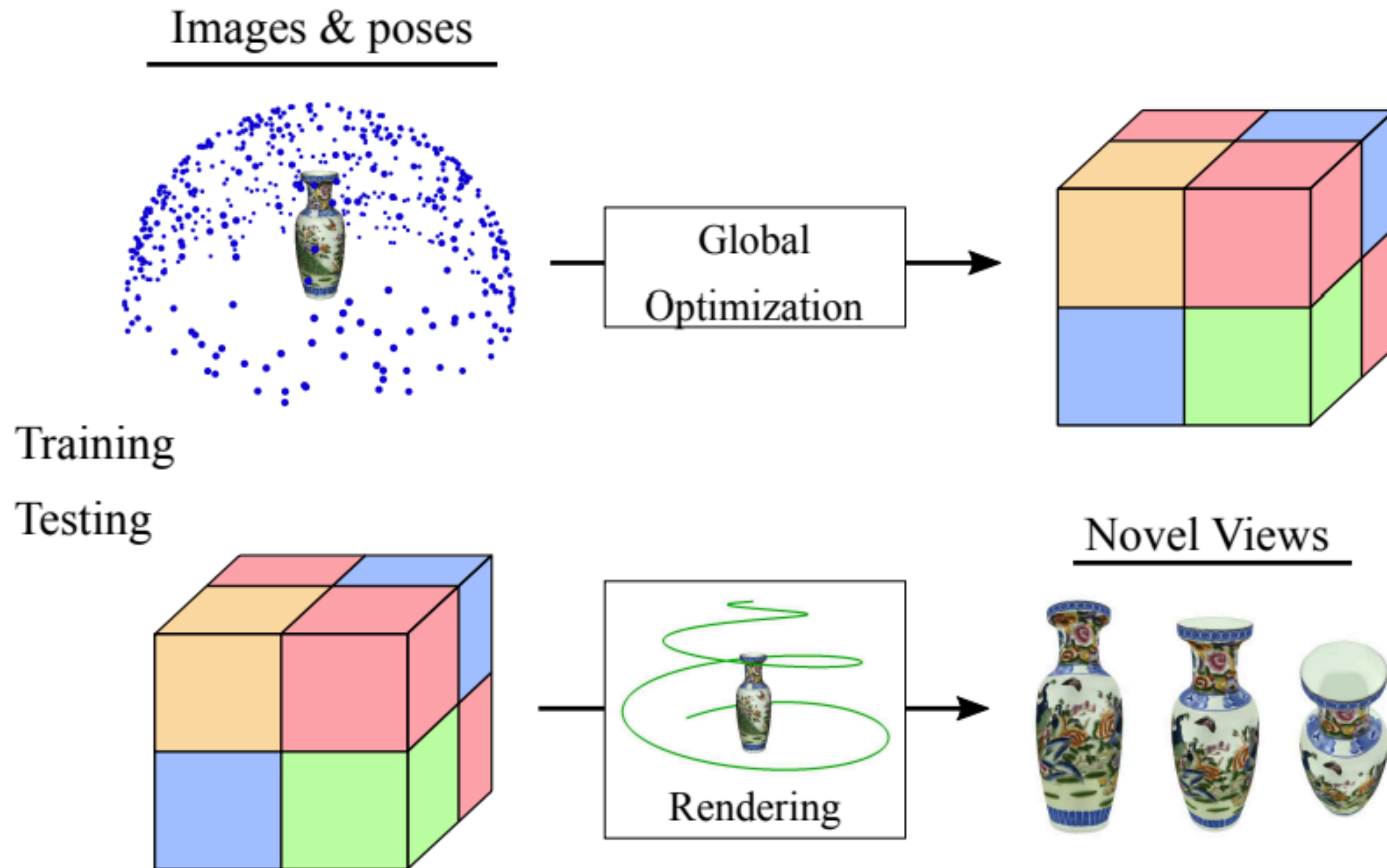
[Riegler and Koltun, 2020]

# Idea #1: Image-based rendering



Point cloud

Proxy geometry (a mesh)

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

6

[Riegler and Koltun, 2020]

# Idea #1: Image-based rendering

[Riegler and Koltun, 2020]

# Idea #2: voxel representation



Images & poses

Training

Testing

Global Optimization

Rendering

Novel Views

[Source: Sitzmann et al., "DeepVoxels", 2019]

# Idea #2: voxel representation
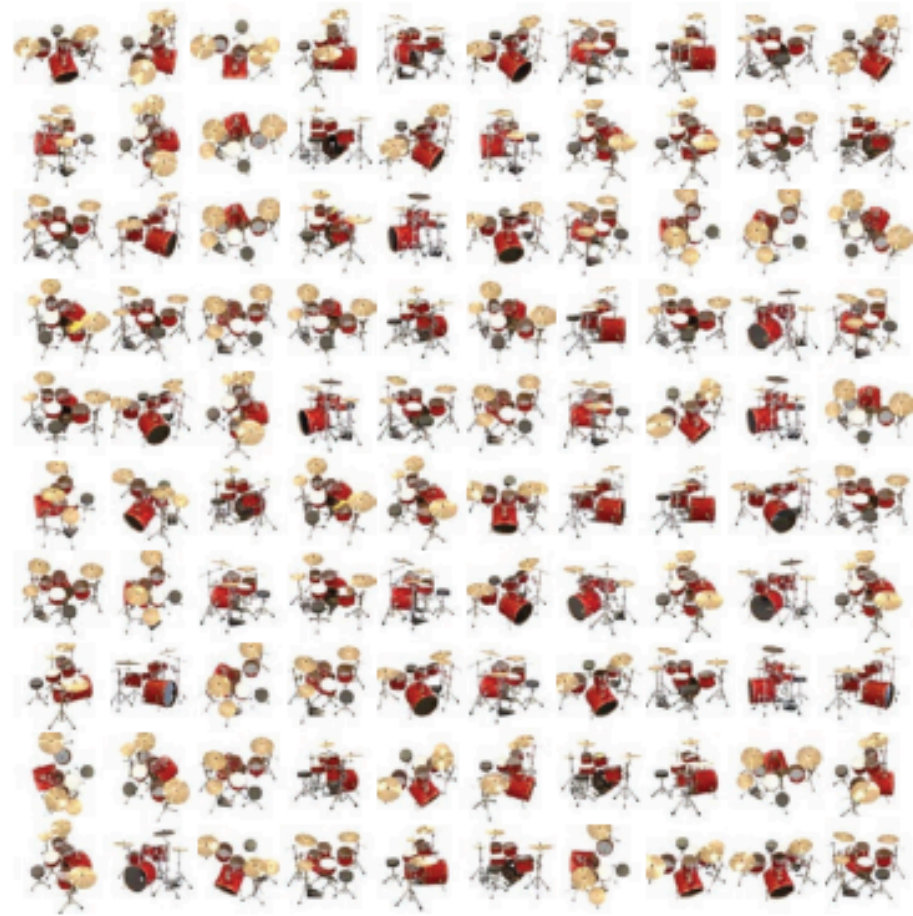


Input views

Position  Viewing direction  Color  Density

$$V[x, y, z, \theta, \phi] = (R, G, B, \sigma)$$
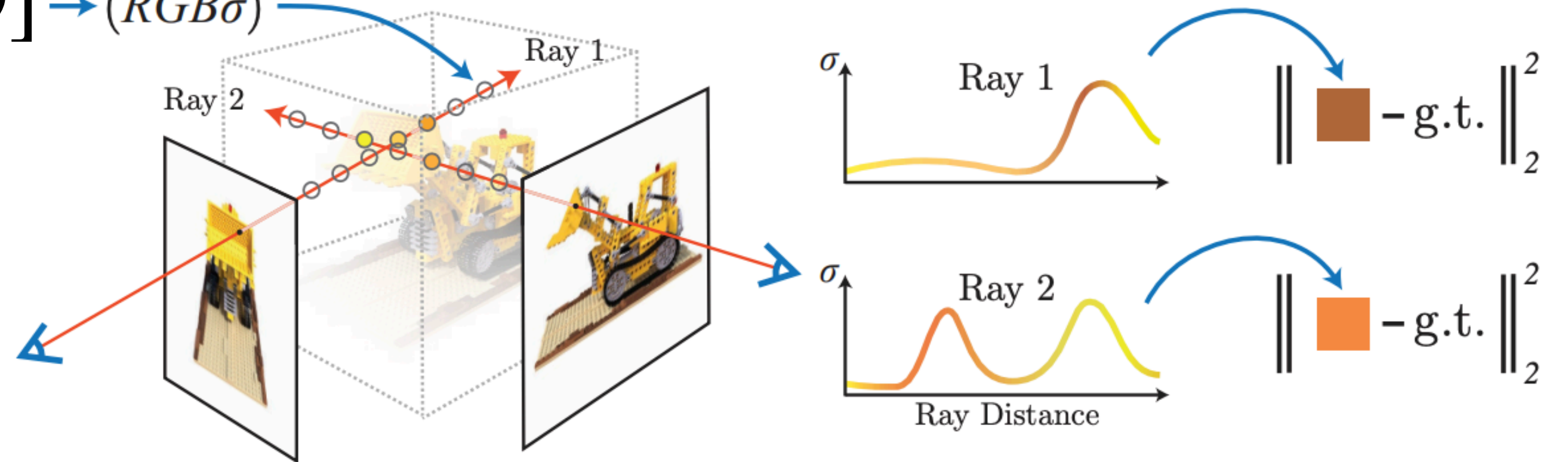
# Idea #2: voxel representation

Position  Viewing direction     Color      Density

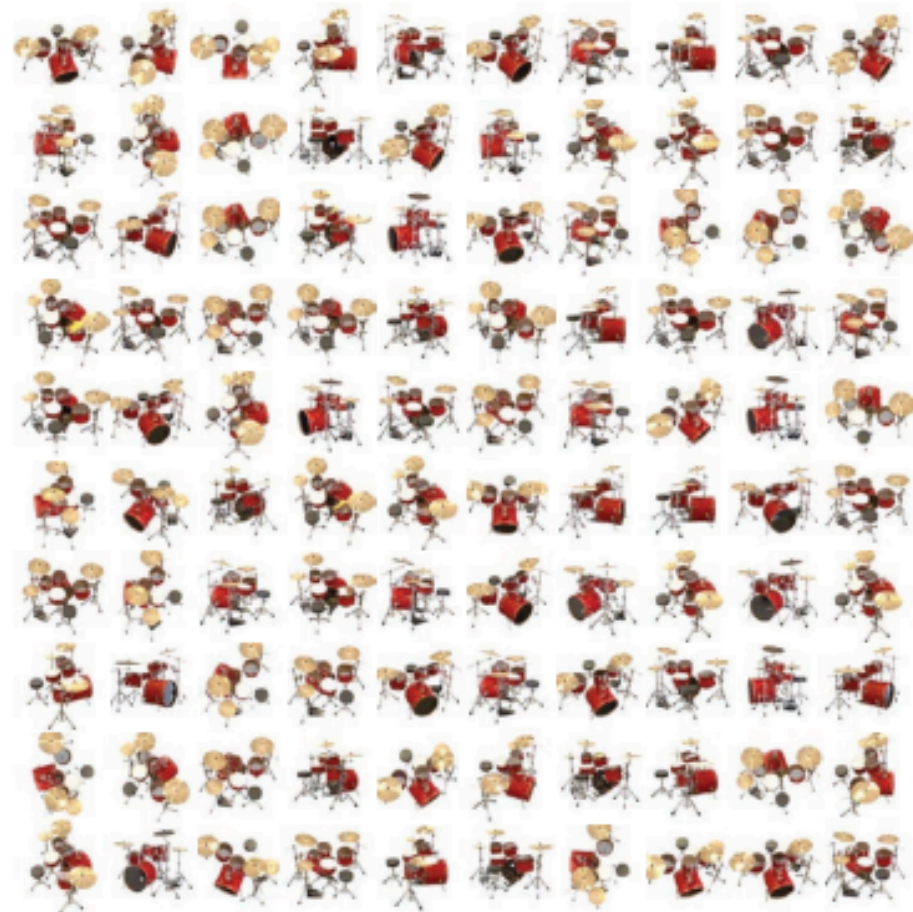$$V[x, y, z, \theta, \phi] = (R, G, B, \sigma)$$

Input views

**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)



Input views

$\longrightarrow$

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

- 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
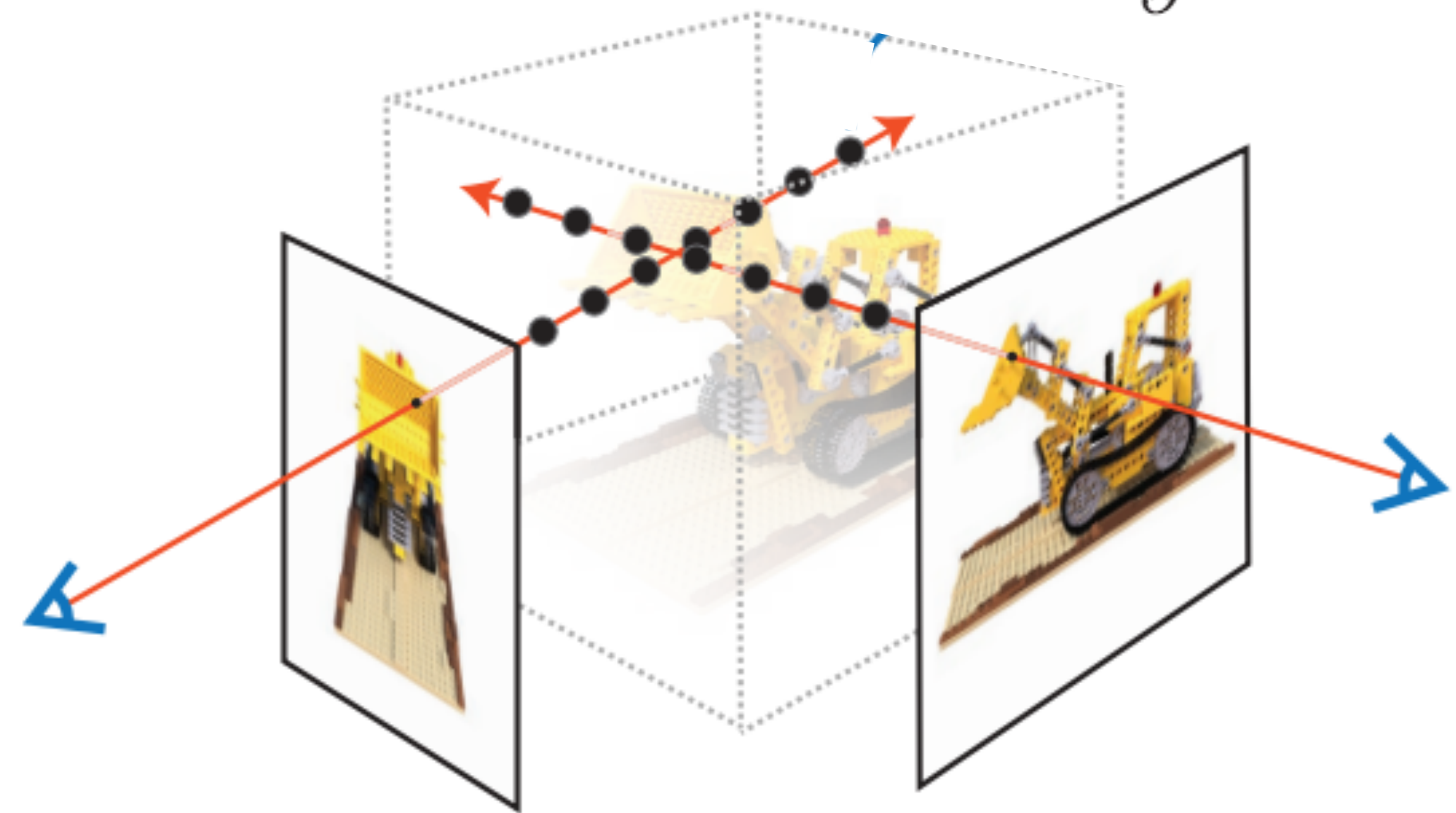
# Idea #3: neural radiance field (NeRF)

Learn volume:
color + occupancy

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

3D scene

Viewpoints

[Mildenhall*, Srinivasan*, Tanick*, et al., Neural radiance fields, 2020]

# Learning a NeRF



5D Input
Position + Direction

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

Output
Color + Density

Ray 1
Ray 2

Volume
Rendering

$\sigma$
Ray 1

$\sigma$
Ray 2
Ray Distance

Rendering
Loss

$\left\| \, \blacksquare - \text{g.t.} \, \right\|_2^2$

$\left\| \, \blacksquare - \text{g.t.} \, \right\|_2^2$

13

# Neural rendering



5D Input
Position + Direction

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

Output
Color + Density

Volume
Rendering

Ray 1

Ray 2

Ray Distance

Rendering
Loss

$\left\| \quad - \text{g.t.} \right\|_2^2$

$\left\| \quad - \text{g.t.} \right\|_2^2$

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

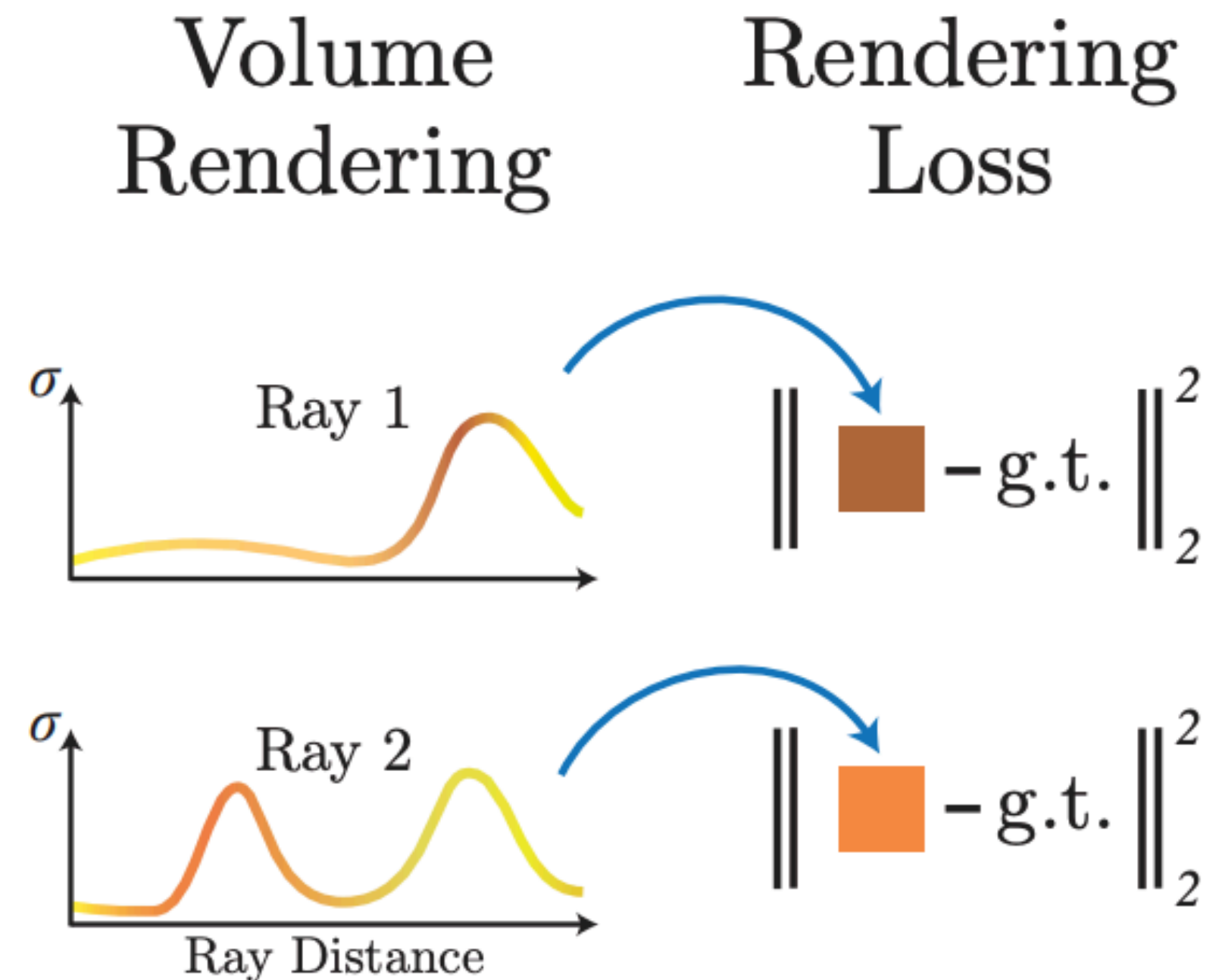Ray: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  For color $c$ and density $\sigma$.

14

# Neural rendering

Color for ray $\mathbf{r}$

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

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

A point distance $t$ along $\mathbf{r}$, centered at

# Neural rendering

Color for ray $\mathbf{r}$

Weight

Color at 3D point $\mathbf{r}(t)$ and direction $\mathbf{d}$

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

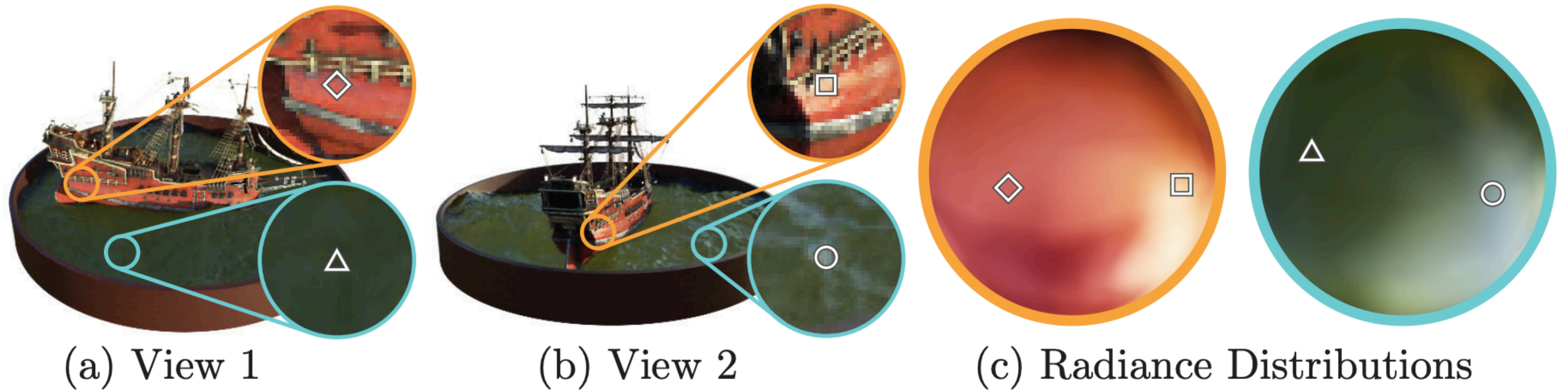Ray: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$     For color $c$ and density $\sigma$.

16

# Neural rendering

Density at
point $\mathbf{r}(t)$

Probability that ray hasn't
been absorbed

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

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

17

# Loss function



$$\mathcal{L} = \sum_{\mathbf{r} \in \mathscr{R}} \|C(\mathbf{r}) - C_{gt}(\mathbf{r})\|_2^2$$

Minimize difference between predicted and observed colors.

In practice: coarse-to-fine and other tricks.

18

# Implementation details

# Why is it good to be view-dependent?



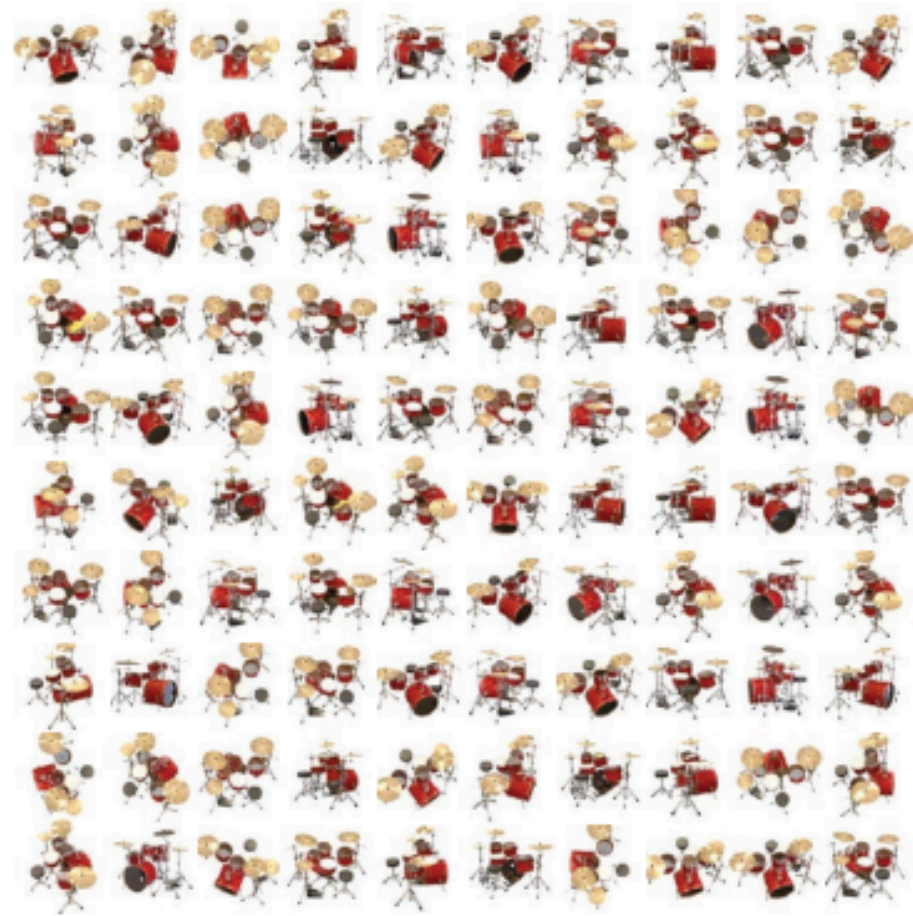(a) View 1    (b) View 2    (c) Radiance Distributions

# Representing the inputs



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

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 [Tancik et al., 2020]

# Fourier features



Input views

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

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

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

- Plug in the coordinate to sinusoids at different frequencies (e.g. L = 10).

22

# MLP architecture



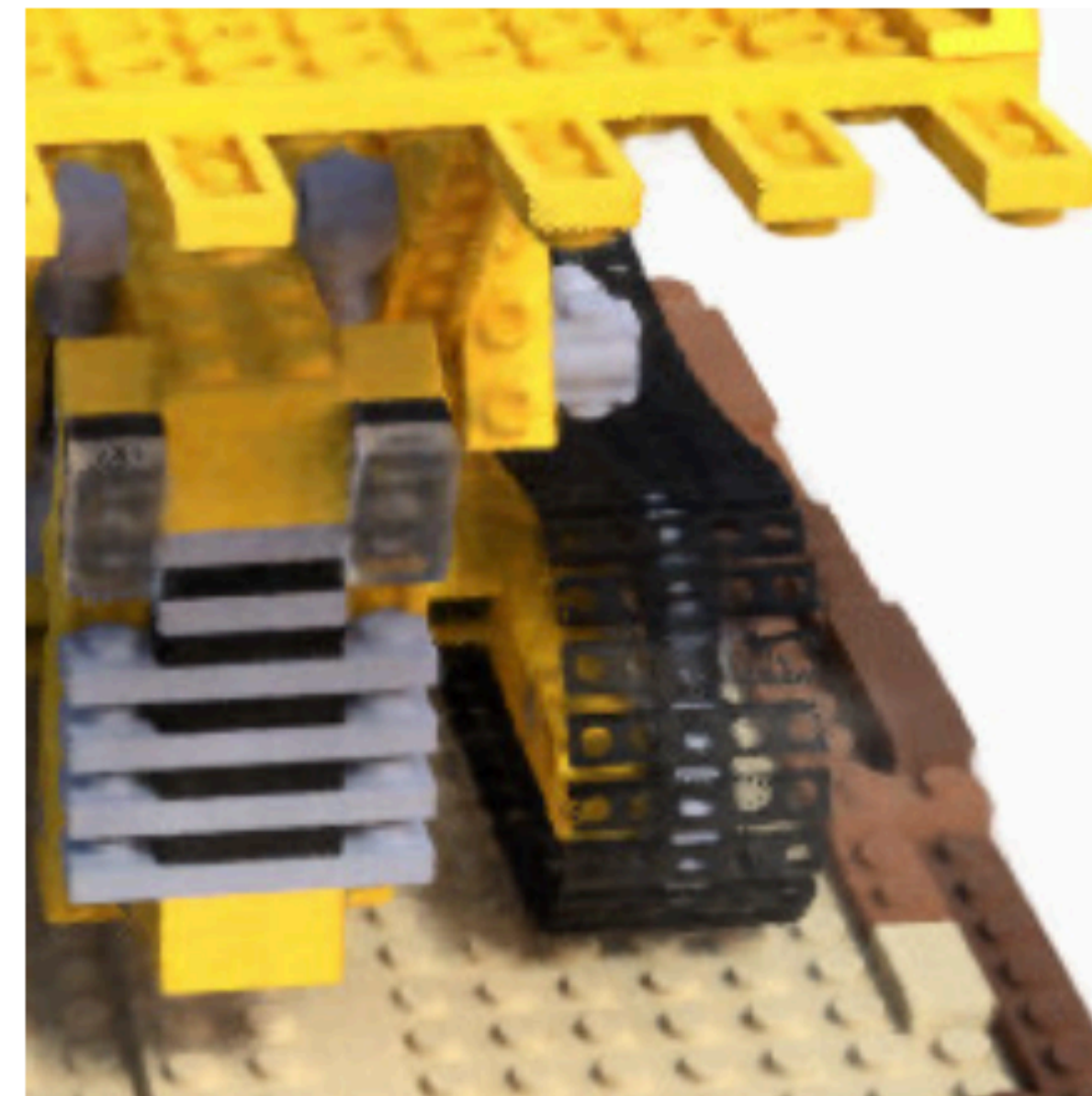$$\gamma(p) = \left( \sin(2^0 \pi p), \cos(2^0 \pi p), \cdots, \sin(2^{L-1} \pi p), \cos(2^{L-1} \pi p) \right)$$

[Source: Mildenhall et al., "NeRF", 2020]

# Results for a novel viewpoint



Ground Truth     Complete Model     No View Dependence     No Positional Encoding

24

[Source: Mildenhall et al., "NeRF", 2020]

# Fourier features



(a) Coordinate-based MLP

$\mathbf{v}$ $x$ $y$

$\gamma(\mathbf{v})$

$\mathbf{y}$ $R$ $G$ $B$

No Fourier features $\gamma(\mathbf{v}) = \mathbf{v}$

$(x, y)$

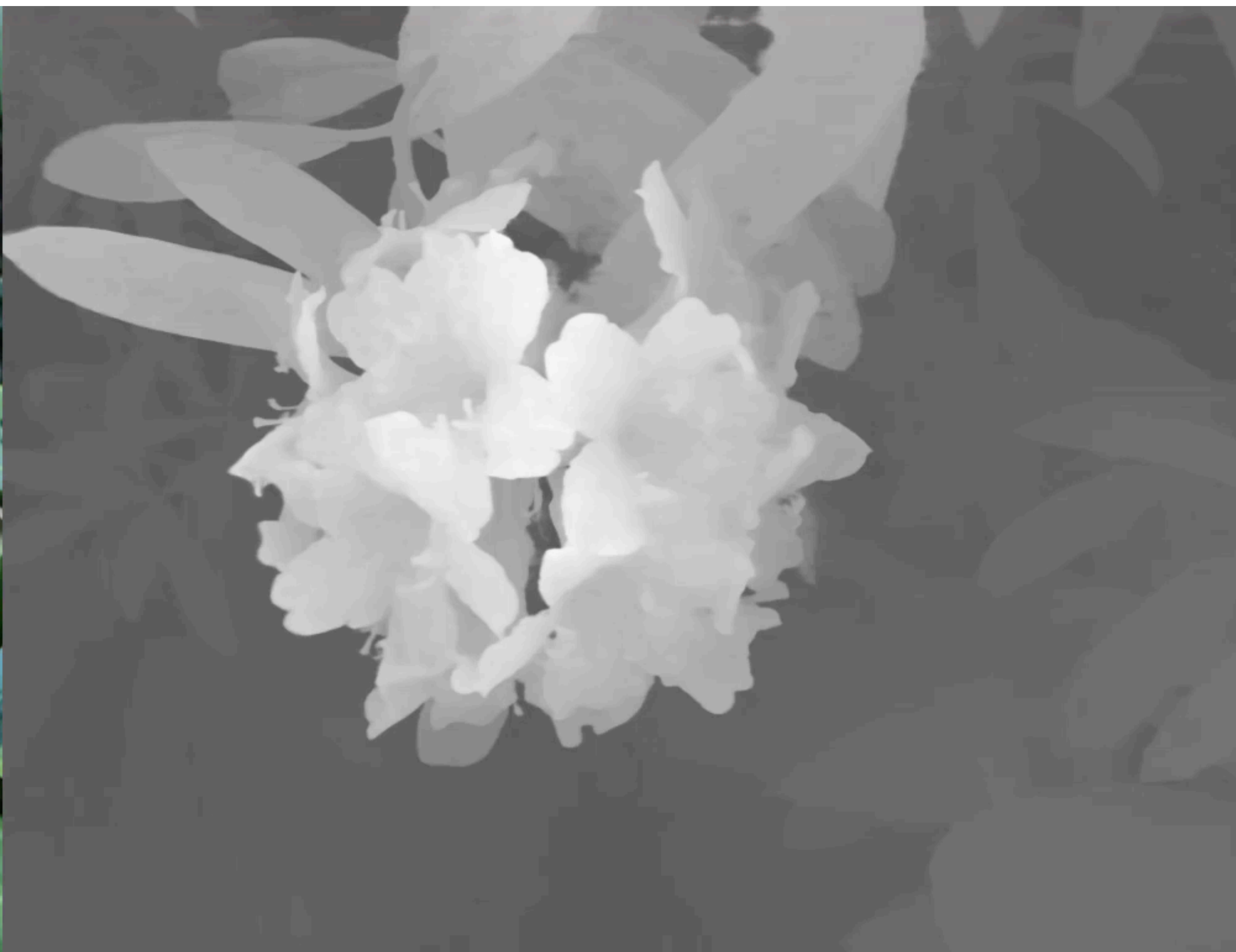With Fourier features $\gamma(\mathbf{v}) = \mathrm{FF}(\mathbf{v})$

(b) Image regression
$(x, y) \rightarrow \mathrm{RGB}$

- Neural nets have trouble learning high frequency functions

- This mapping explicitly represents different frequencies (forces net to pay more attention high frequencies)



freq weighting↑

No mapping
Target signal
Training points

See [Tancik et al., "Fourier Features Let Networks Learn High Frequency Functions in Low Dimensional Domains", 2020]

# Results



[Mildenhall*, Srinivasan*, Tanick*, et al. 2020]

# Results

[Mildenhall*, Srinivasan*, Tanick*, et al. 2020]

# Results

[Mildenhall*, Srinivasan*, Tanick*, et al. 2020]

# NeRF state-of-the-art

[Barron et al., "Zip-NeRF: Anti-Aliased Grid-Based Neural Radiance Fields. 2023]

# Extension: internet photo collections

[Martin-Brualla, Radwan et al. "NeRF in the Wild", 2020]

# Extension: internet photo collections

[Martin-Brualla, Radwan et al. "NeRF in the Wild", 2020]

# Lots of other applications



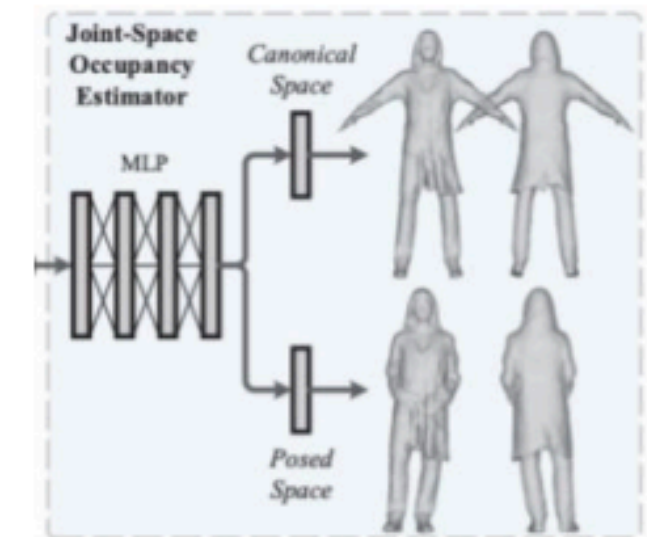$f_\theta : \mathbb{R}^m \rightarrow \mathbb{R}^n$

Neural field

2D and 3D Reconstruction

Generative Models

Digital Humans

903.63 KB

Compression

Robotics

…and Beyond!

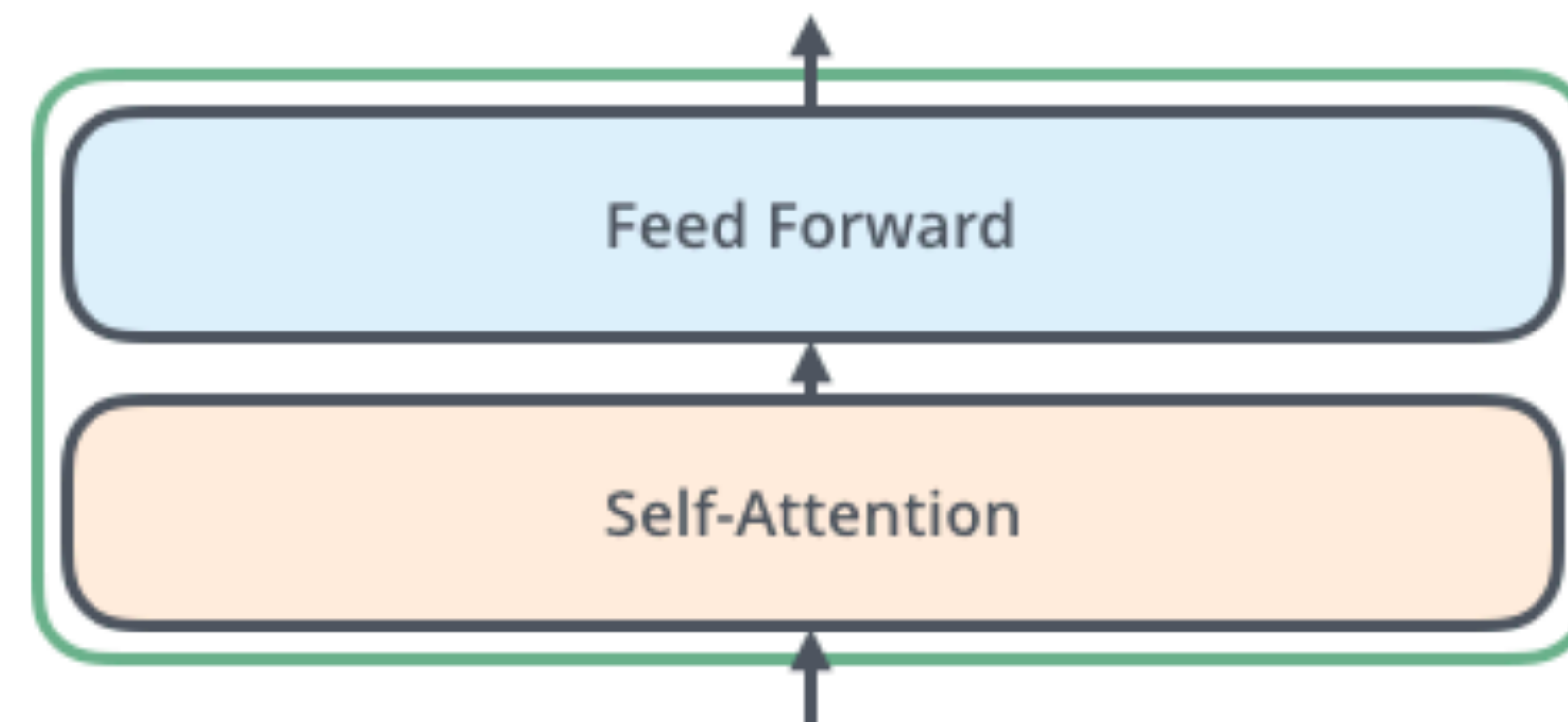[Xie et al., "Neural Fields in Visual Computing and Beyond", 2021]

# 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)



A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, NeurIPS 2017

Source: S. Lazebnik

Image source
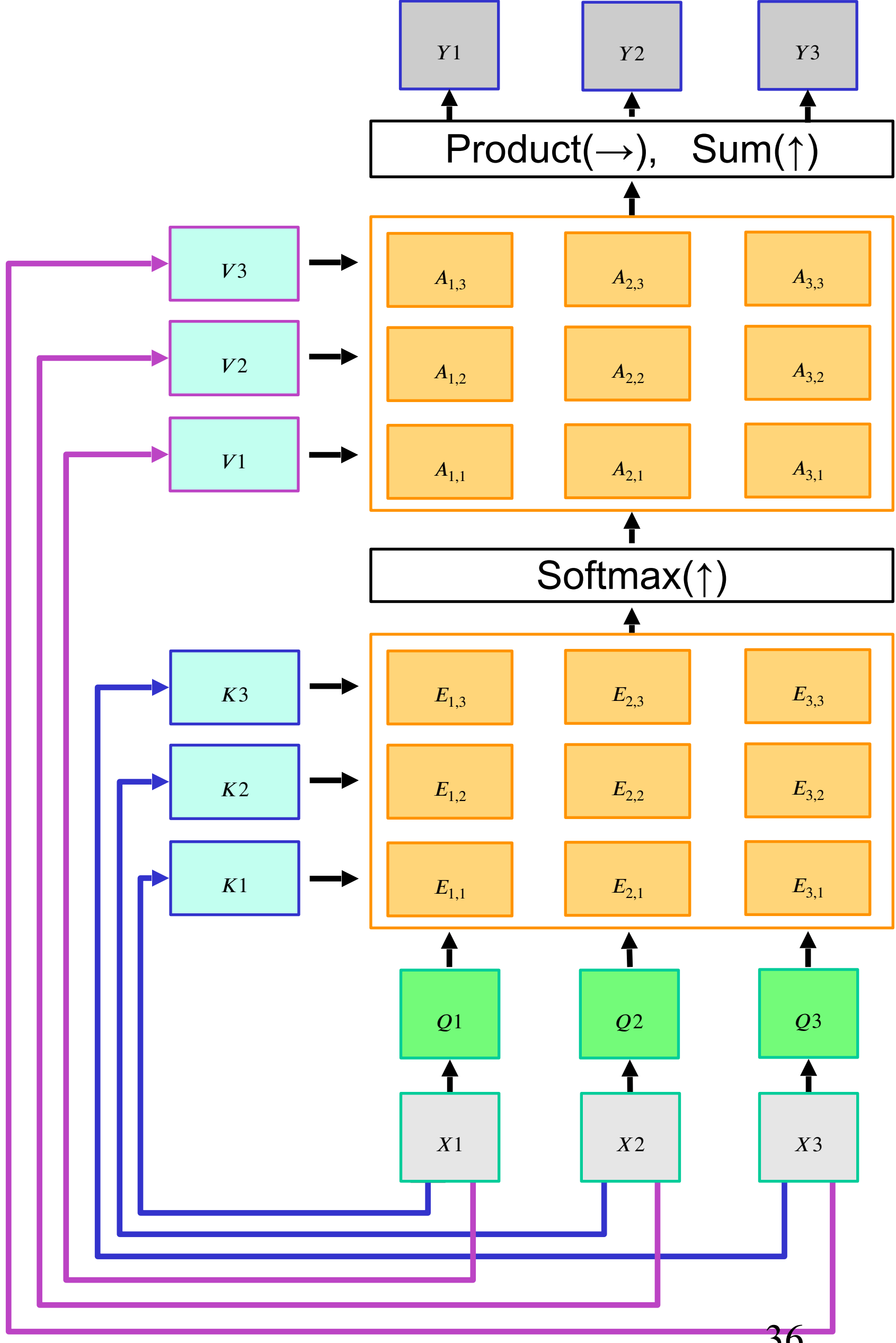
# Self-attention layer

- 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:

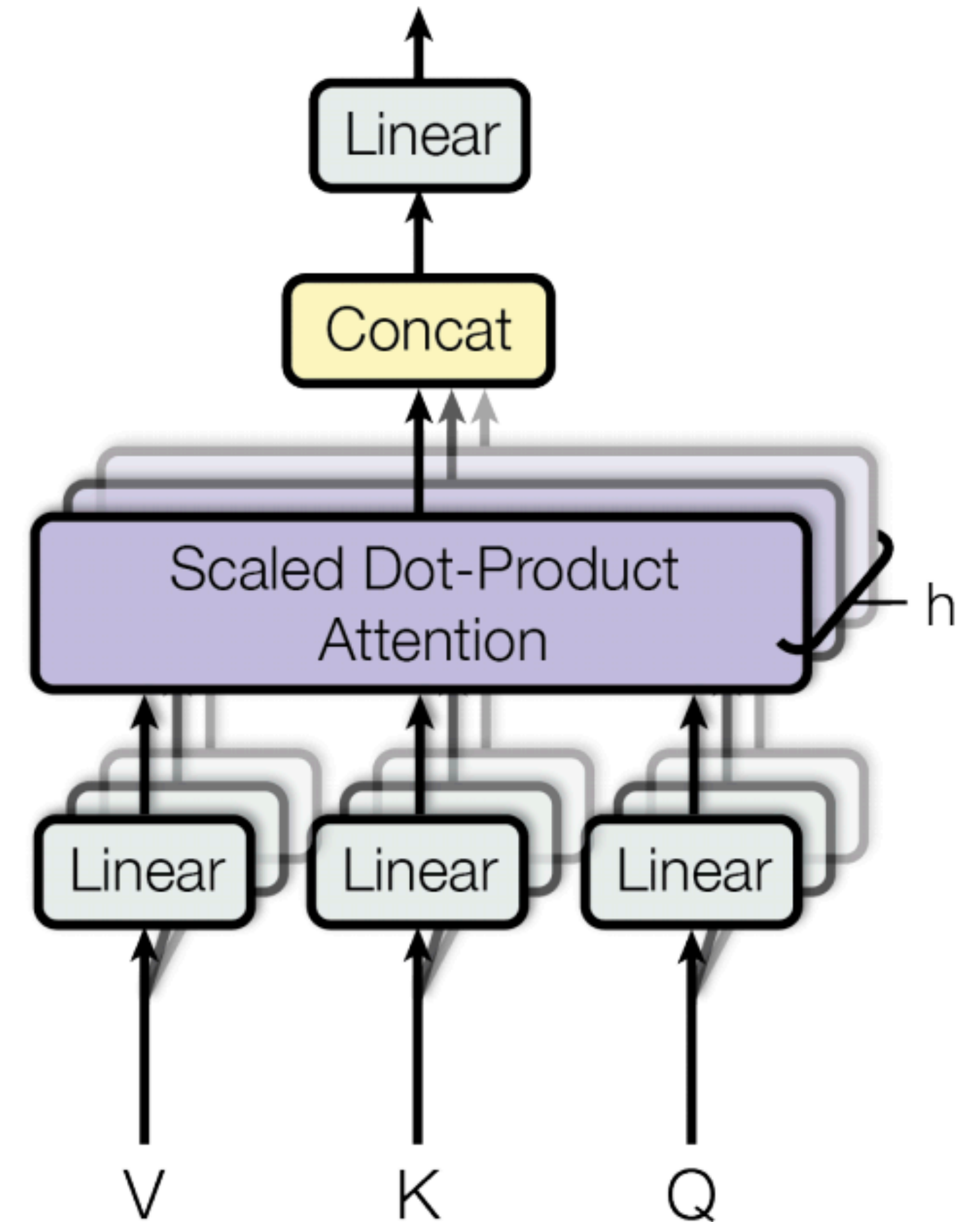$$Y_i = \sum_j A_{i,j} V_j \quad \text{or} \quad Y = AV$$



**One query per input vector**

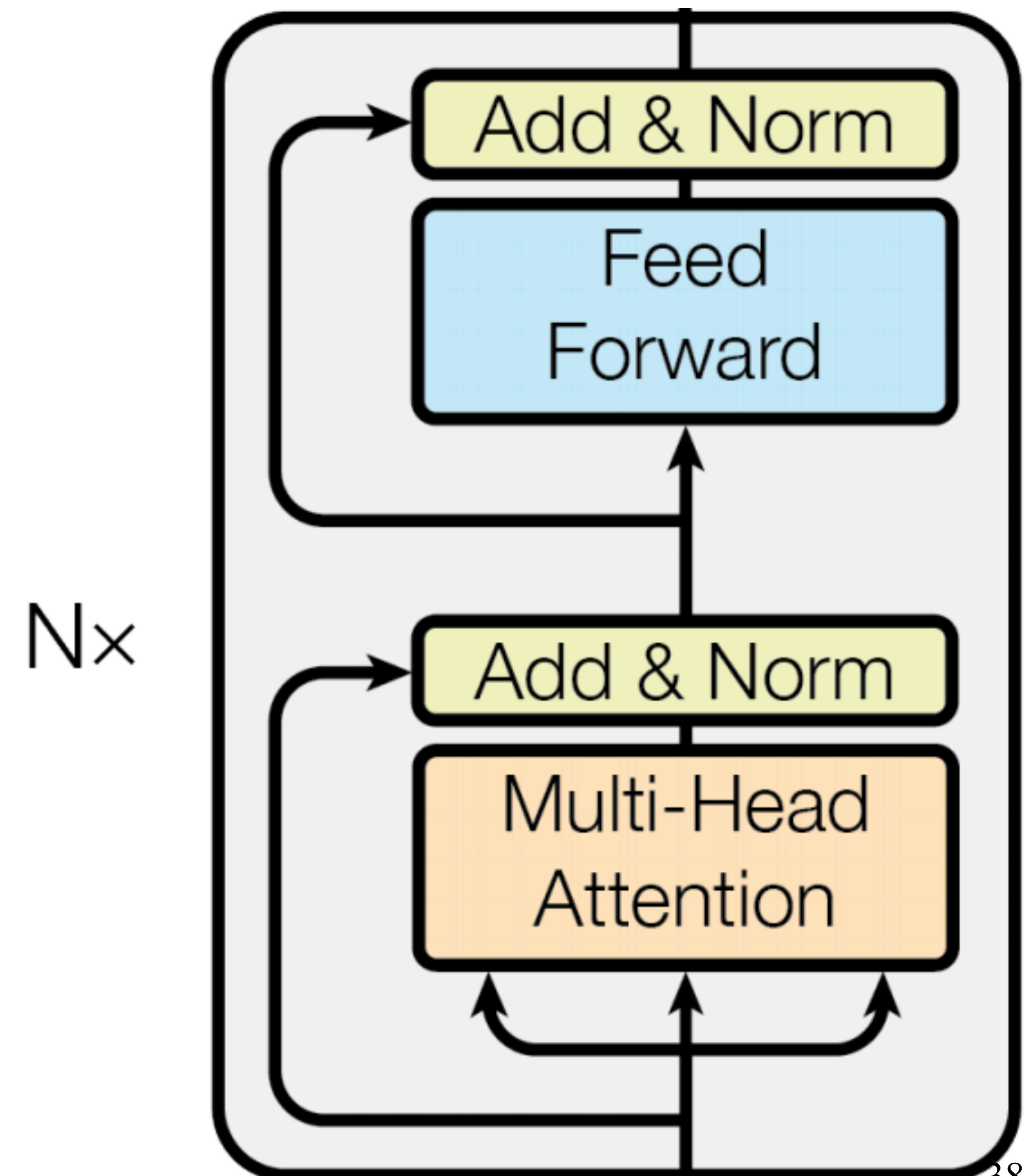Adapted from J. Johnson and S. Lazebnik.

# 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

Source: S. Lazebnik

# 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!



Nx

Source: S. Lazebnik

# 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
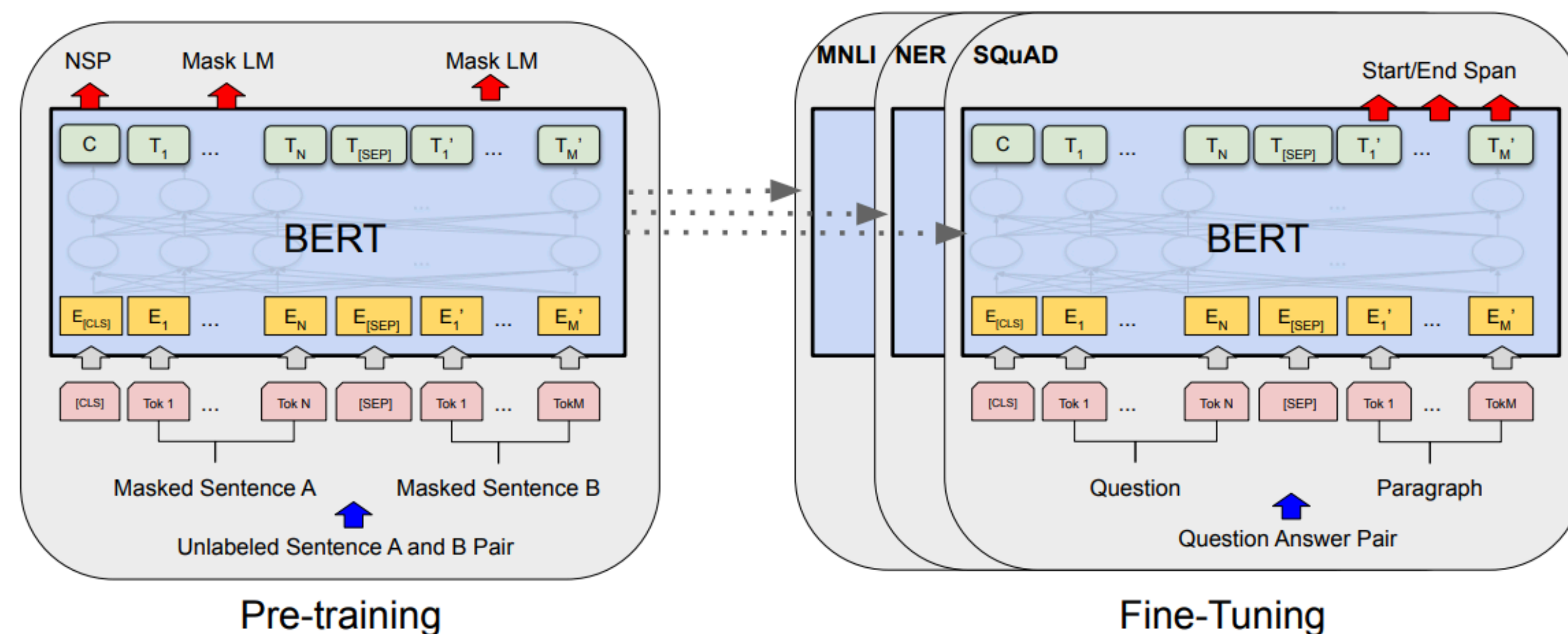3. Fine-tune the transformer on desired NLP task

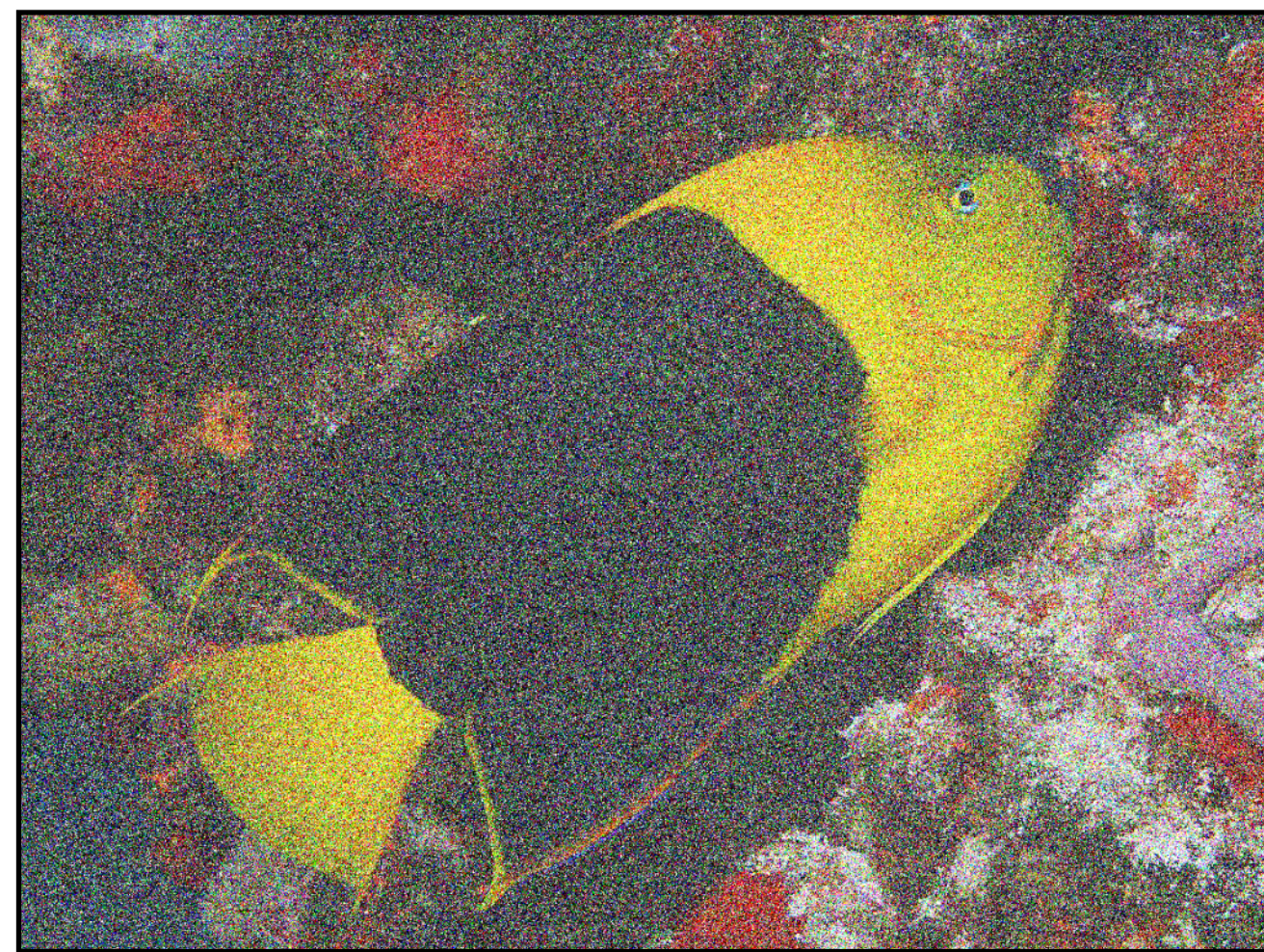| Model Alias | Org. | Article Reference |
|---|---|---|
| ULMfit | fast.ai | *Universal Language Model Fine-tuning for Text Classification*<br>Howard and Ruder |
| ELMo | AllenNLP | *Deep contextualized word representations*<br>Peters et al. |
| OpenAI GPT | OpenAI | *Improving Language Understanding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

Source: S. Lazebnik

Image source

# 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

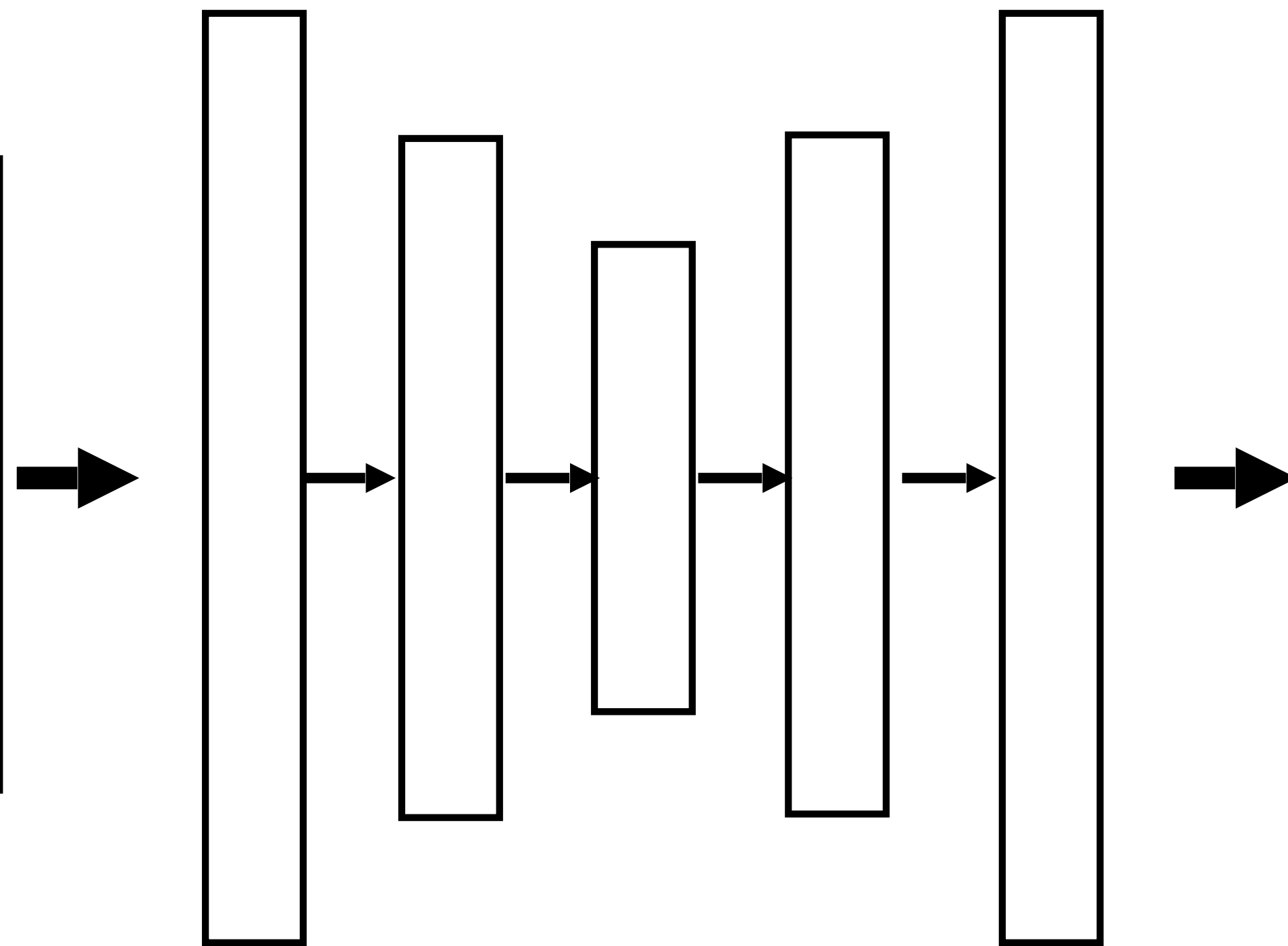3.  Fine-tune the transformer on desired NLP task

Bidirectional Encoder Representations from Transformers (BERT)

J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, EMNLP 2018

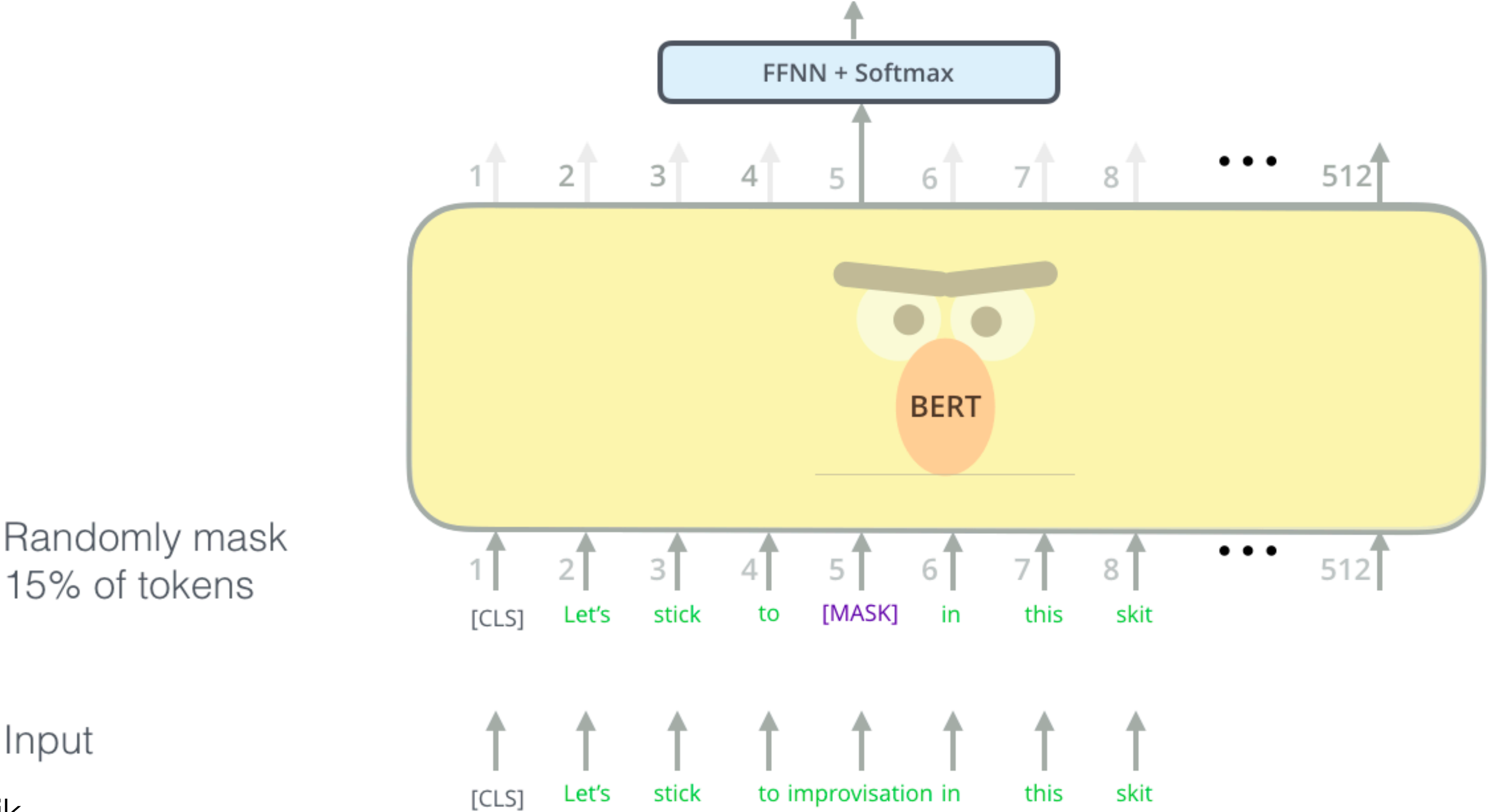# Recall: denoising autoencoder
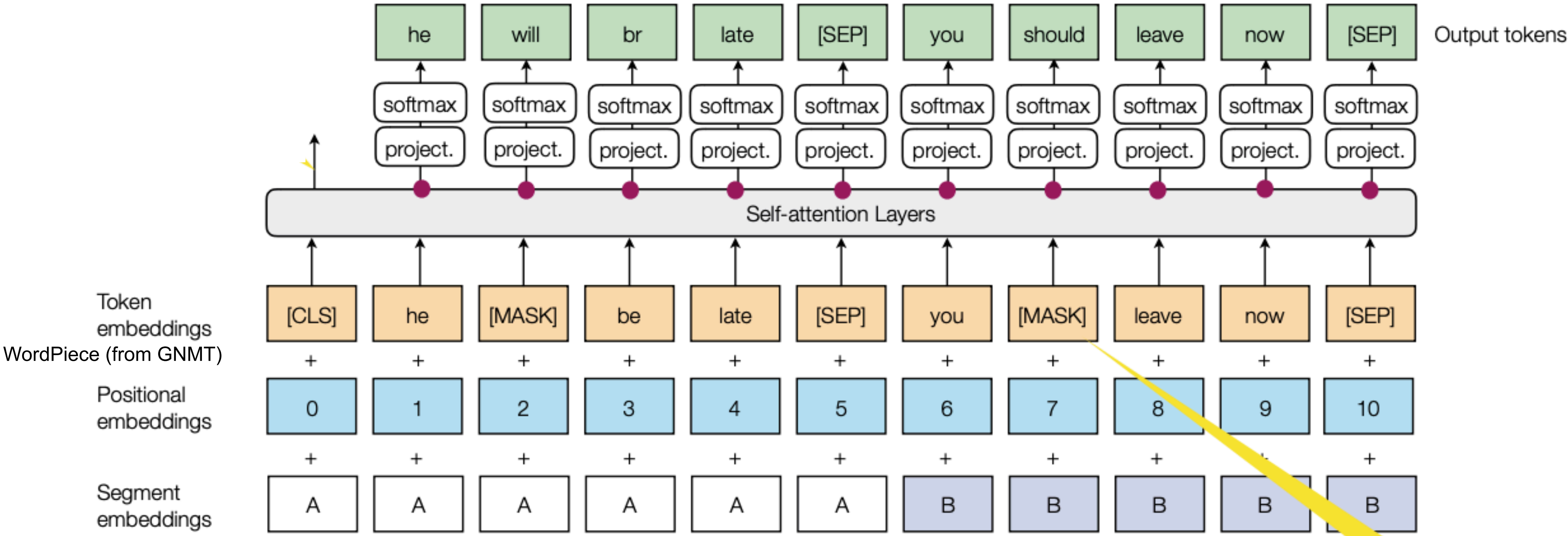


Noisy image

Reconstructed image

[Vincent et al., 2008]

# BERT: Pretext tasks

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

42

# BERT: More detailed view



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

WordPiece (from GNMT)

15% of tokens get masked

Source: S. Lazebnik

Image source

# BERT: Evaluation

- General Language Understanding Evaluation (GLUE) benchmark (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$_{BASE}$ | 84.6/83.4 | 71.2 | 90.5 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT$_{LARGE}$ | **86.7/85.9** | **72.1** | **92.7** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **82.1** |

Source: S. Lazebnik

# Image GPT

- Image resolution up to 64x64, color values quantized to 512 levels (9 bits), dense attention

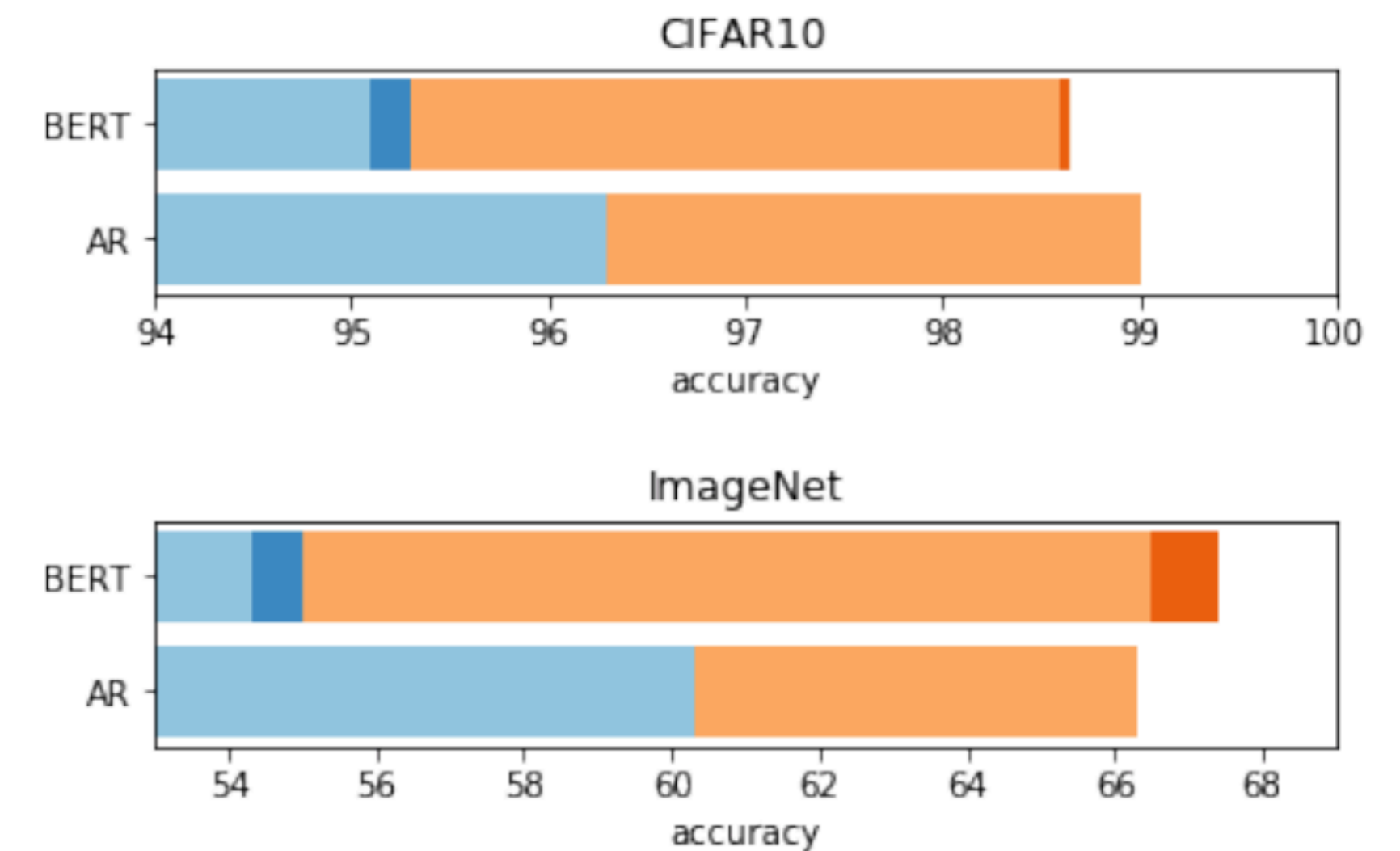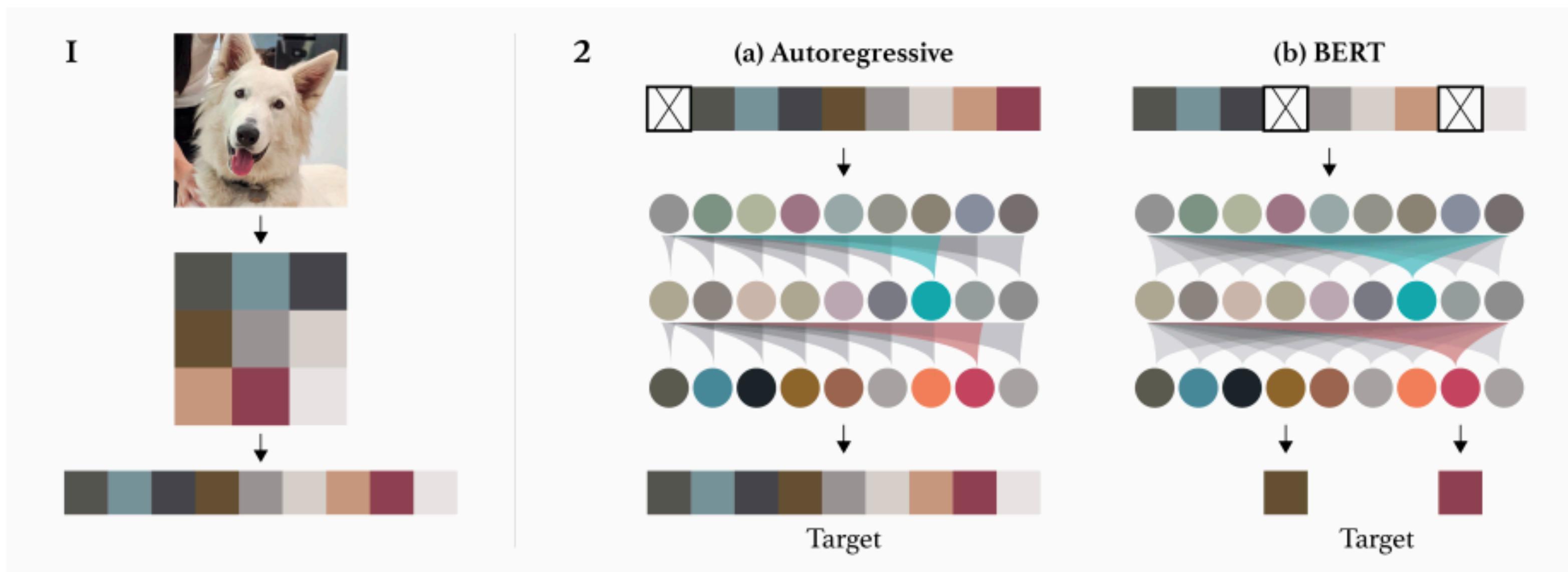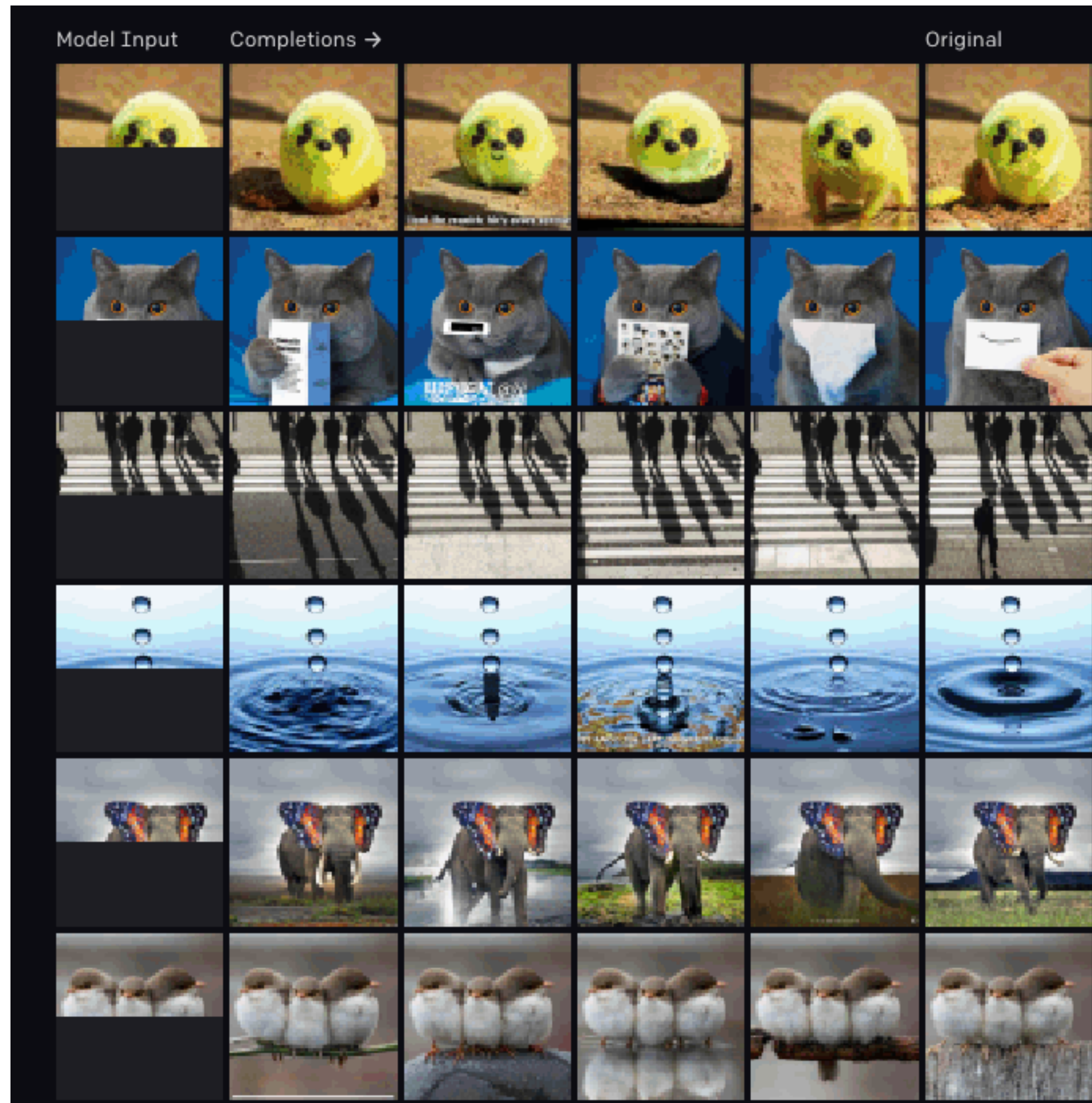- For transfer learning, average-pool encoded features across all positions



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.

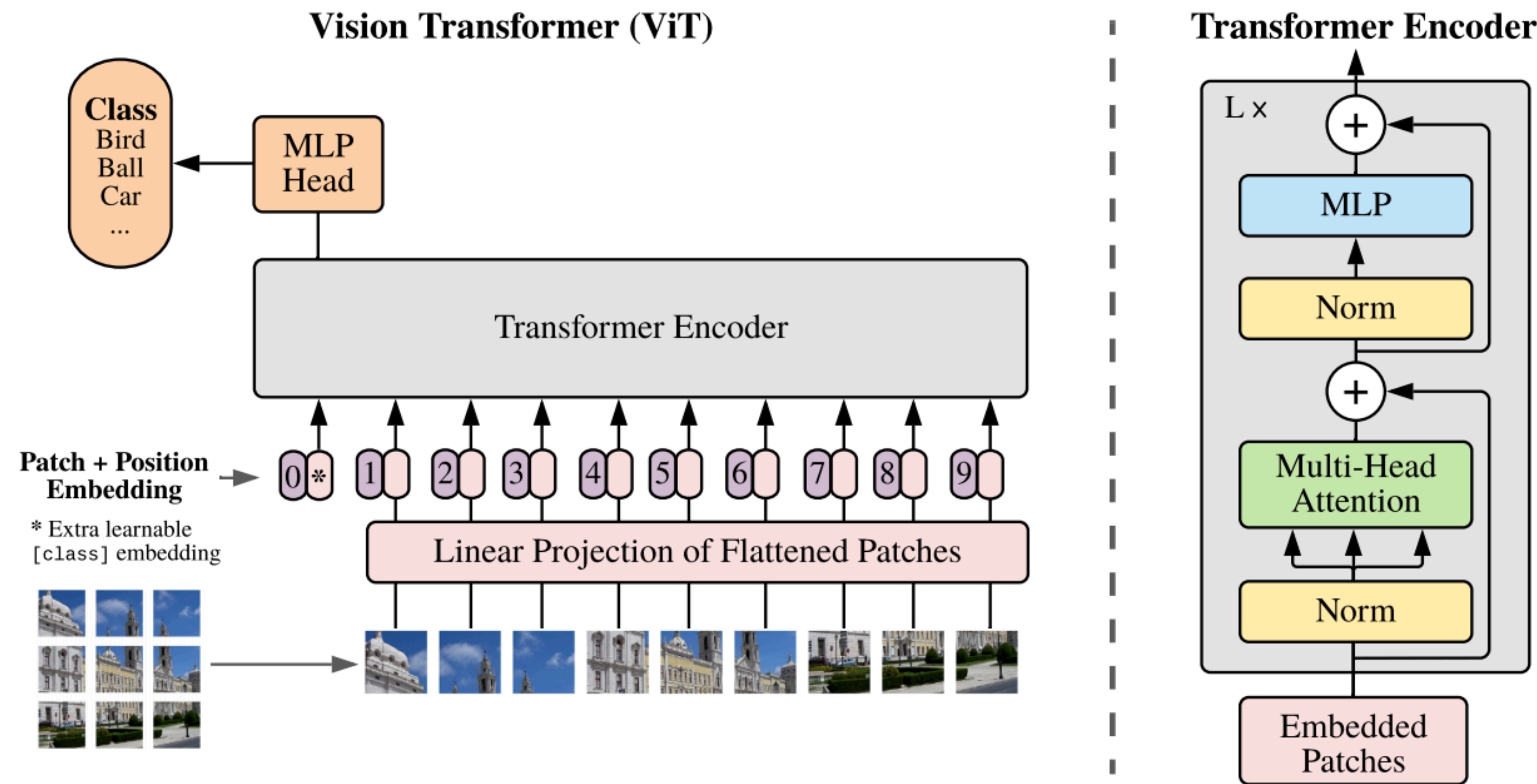M. Chen et al., Generative pretraining from pixels, ICML 2020

# Image GPT – OpenAI

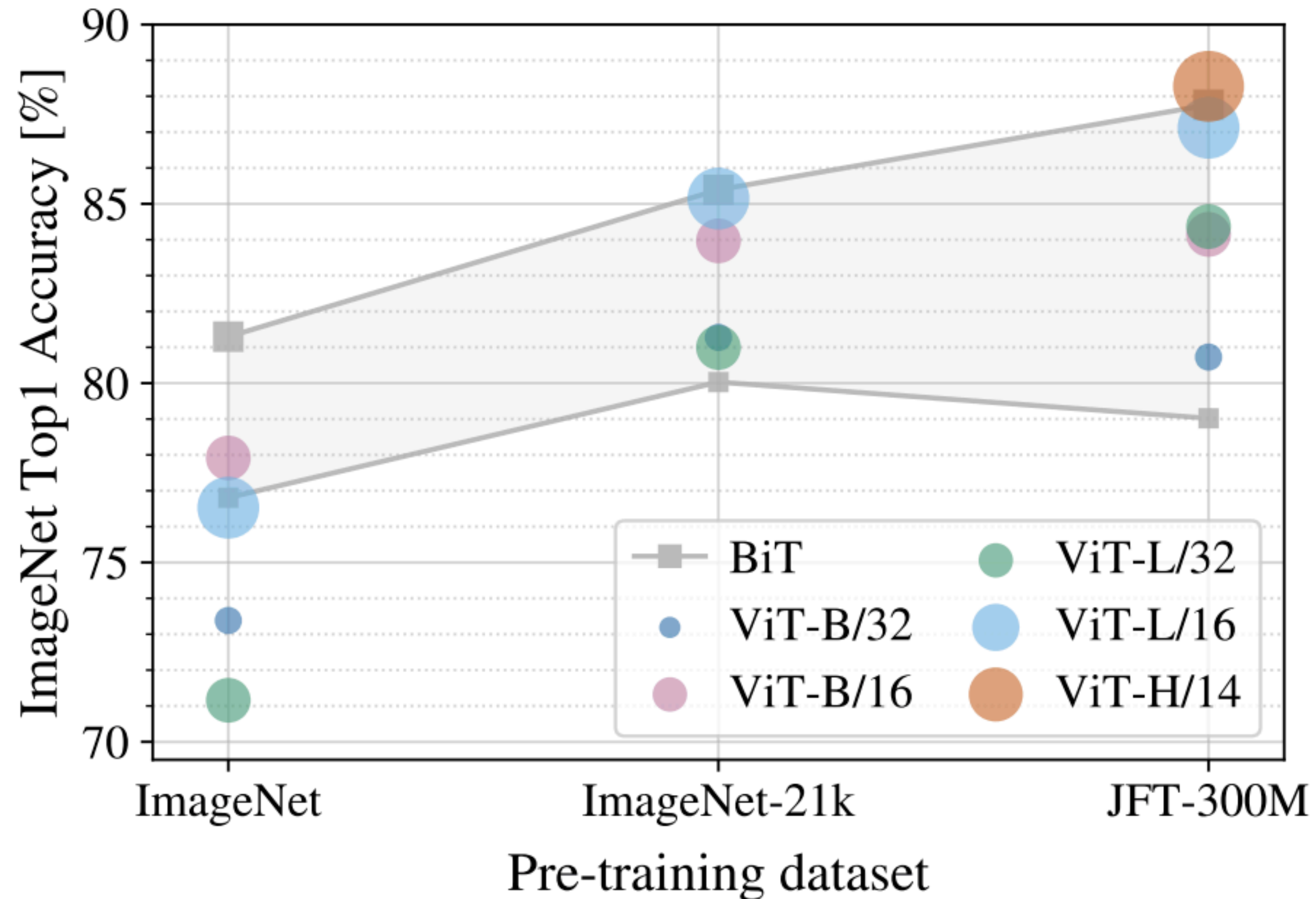M. Chen et al., Generative pretraining from pixels, ICML 2020

# Vision transformer (ViT)

- Split an image into patches, feed linearly projected patches into standard transformer encoder
  - With patches of 14x14 pixels, you need 16x16=256 patches to represent 224x224 images



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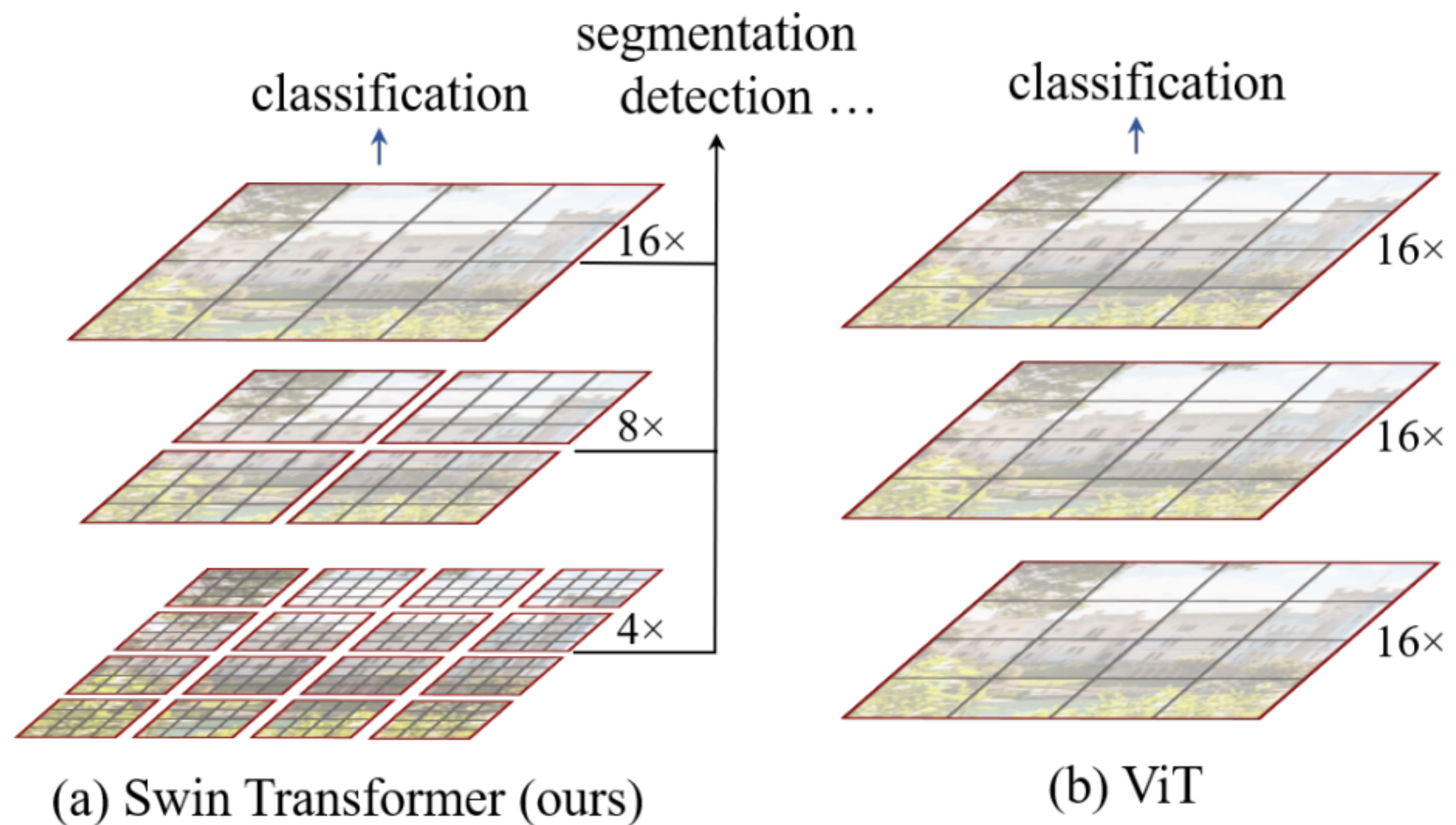
# Vision transformer (ViT)



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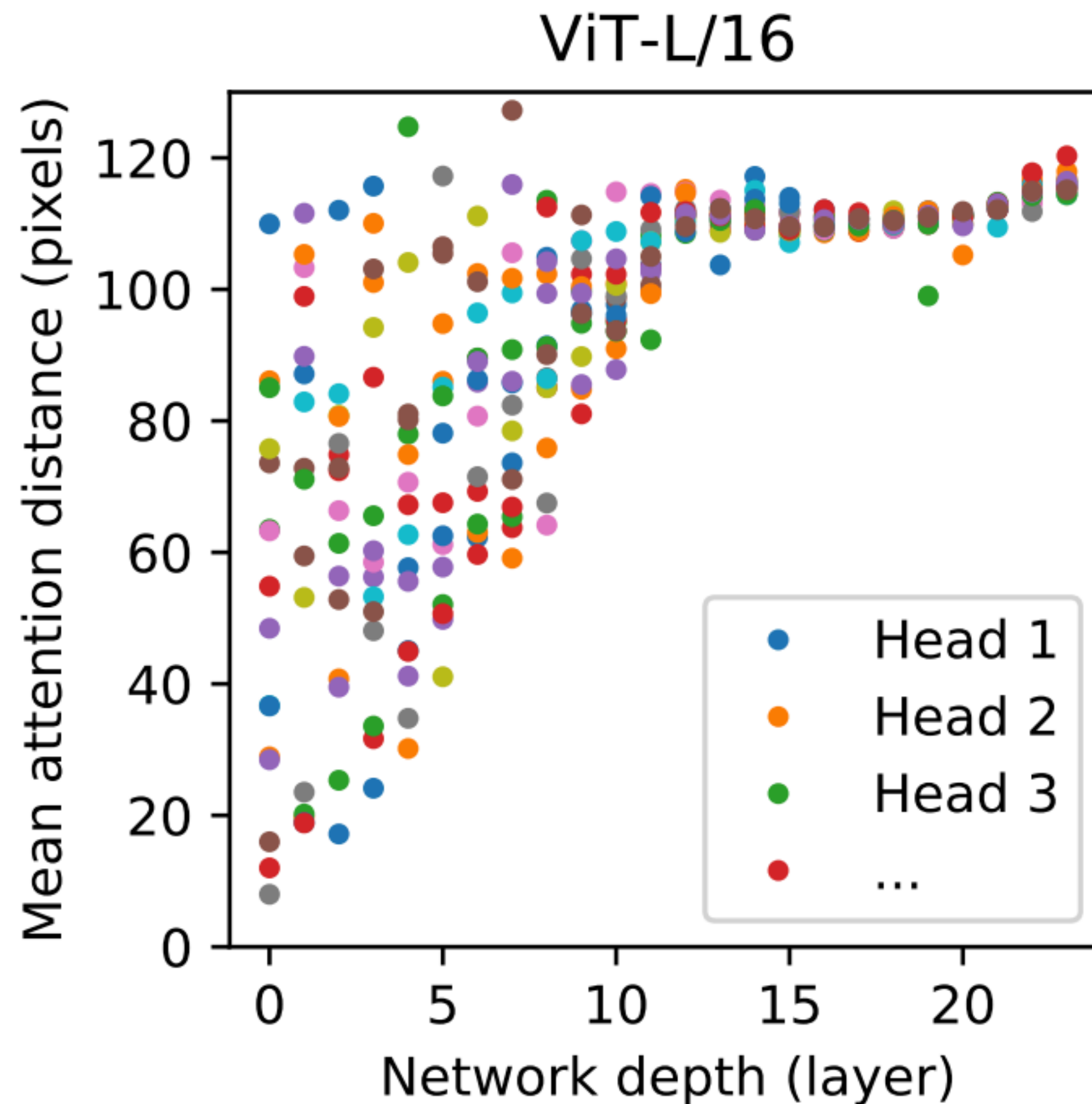ViT: Vision Transformer (Base/Large/Huge, patch size of 14x14, 16x16, or 32x32)

Internal Google dataset (not public)

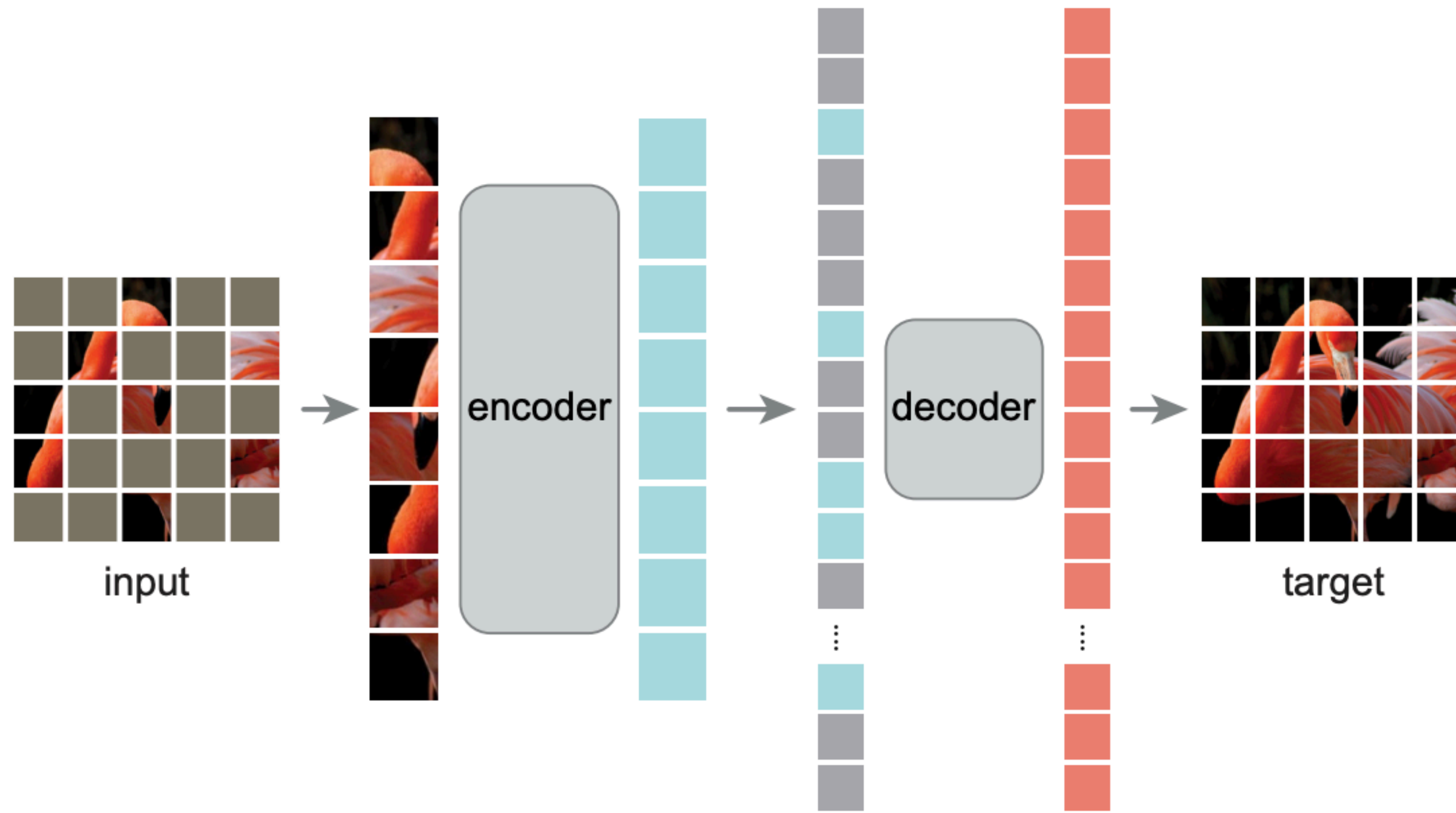A. Dosovitskiy et al. An image is worth 16x16 words: Transformers for image recognition at scale. ICLR 2021

Source: S. Lazebnik

# Swin Transformer: windowed attention



(a) Swin Transformer (ours)

(b) ViT

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

# Vision transformer (ViT)



ViT-L/16

A. Dosovitskiy et al. An image is worth 16x16 words: Transformers for image recognition at scale. ICLR 2021
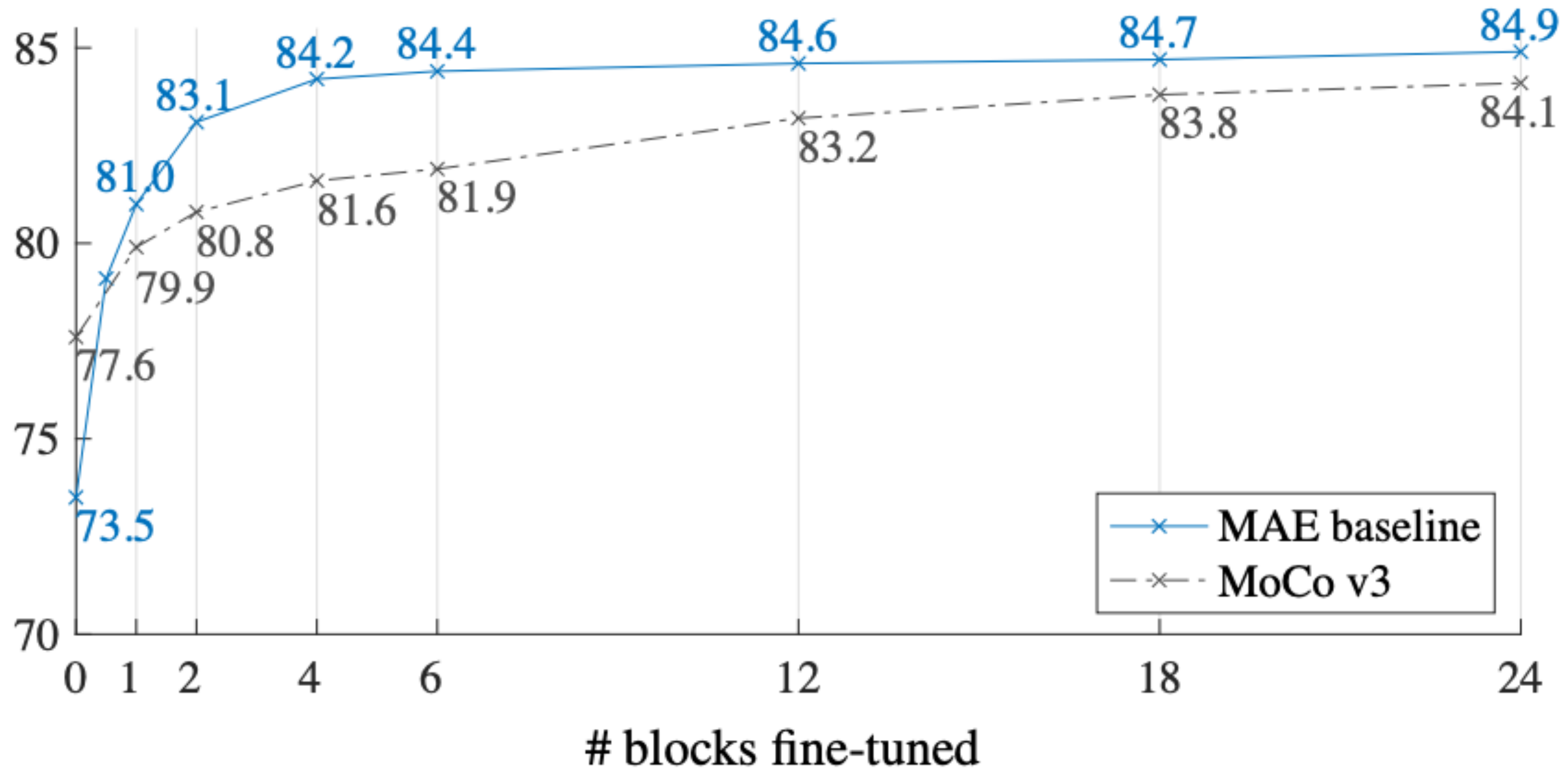
# Application: self-supervised learning with masked autoencoders



input

encoder

decoder

target

Source: S. Lazebnik

[Kaiming He et al., "Masked Autoencoders Are Scalable Vision Learners", 2021]

# Application: self-supervised learning

[Kaiming He et al., "Masked Autoencoders Are Scalable Vision Learners", 2021]

[Kaiming He et al., "Masked Autoencoders Are Scalable Vision Learners", 2021]

Space Attention (S)        Joint Space-Time        Divided Space-Time
                           Attention (ST)          Attention (T+S)

Source: [Gedas Bertasius et al, "Is Space-Time Attention All You Need for Video Understanding?", 2021]

# Application: video self-supervised learning

Source: [Feichtenhofer et al., "Masked Autoencoders As Spatiotemporal Learners", 2022]

[Arsha Nagrani et al, "Attention Bottlenecks for Multimodal Fusion", 2021]

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



N. Carion et al., End-to-end object detection with transformers, ECCV 2020

Source: S. Lazebnik

# CNNs vs. Transformer performance: pretty similar

| Improvements | Top-1 | Δ |
|---|---|---|
| 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** |

Speed-Accuracy Pareto Curve

2.7x Speedup

+ Minor Architectural Changes

+ Improved Training Strategies

- - - ResNet-RS
- - - EfficientNet
● ResNet

Top-1 ImageNet Accuracy

Time Per Training Step (Sec)

[Irwan Bello et al. "Revisiting ResNets: Improved Training and Scaling Strategies", 2021]

# Improving CNNs



[Liu et al. "A ConvNet for the 2020s", 2022]

Also: [Smith et al., ConvNets Match Vision Transformers at Scale, 2023]

# Next class: Bias and ethics