

# How to use Attention / Transformers for Vision?

## Last lecture

- Transformers – encoder and decoder blocks
- SA, Multihead SA, Masked SA

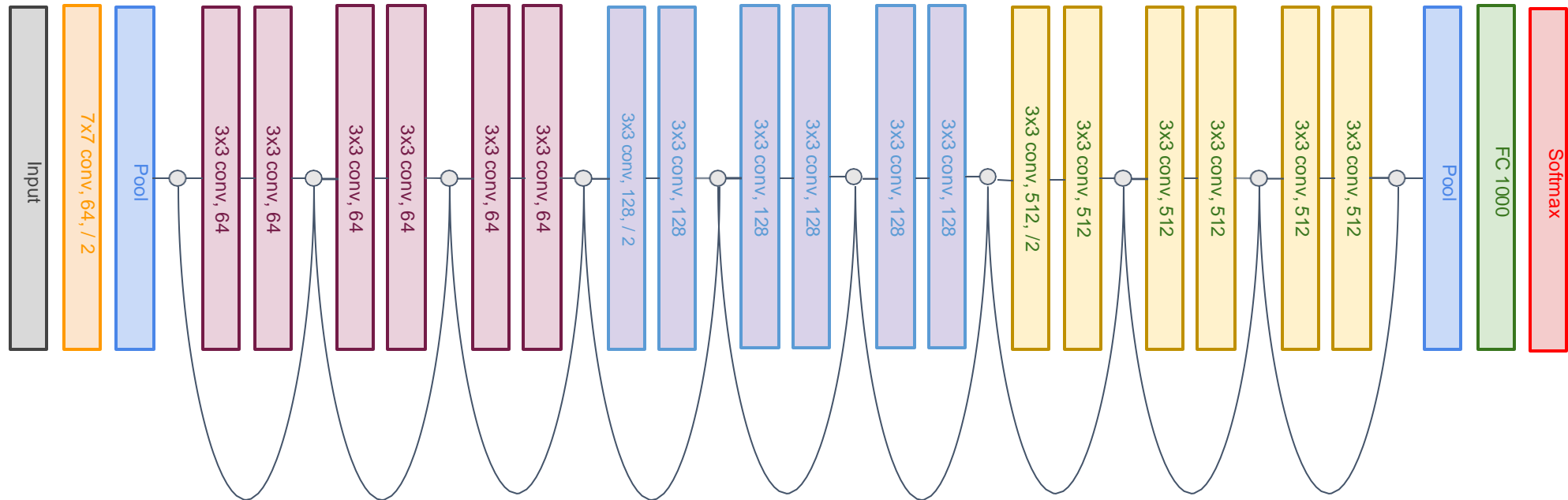
## This lecture

- ViT, SwinT
- Improving ViT training using distillation and augmentation/regularization

# Idea #1: Add attention to existing CNNs

Start from standard CNN architecture (e.g. ResNet)

How can you incorporate SA?



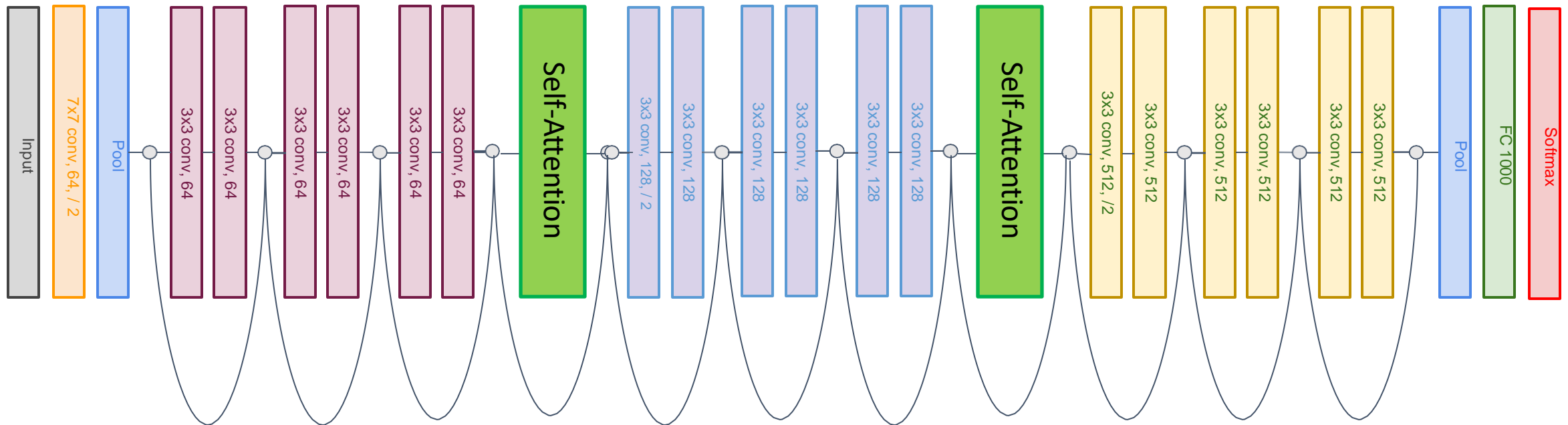
Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

Wang et al, "Non-local Neural Networks", CVPR 2018

# Idea #1: Add attention to existing CNNs

Start from standard CNN architecture (e.g. ResNet)

Add Self-Attention blocks between existing ResNet blocks



Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

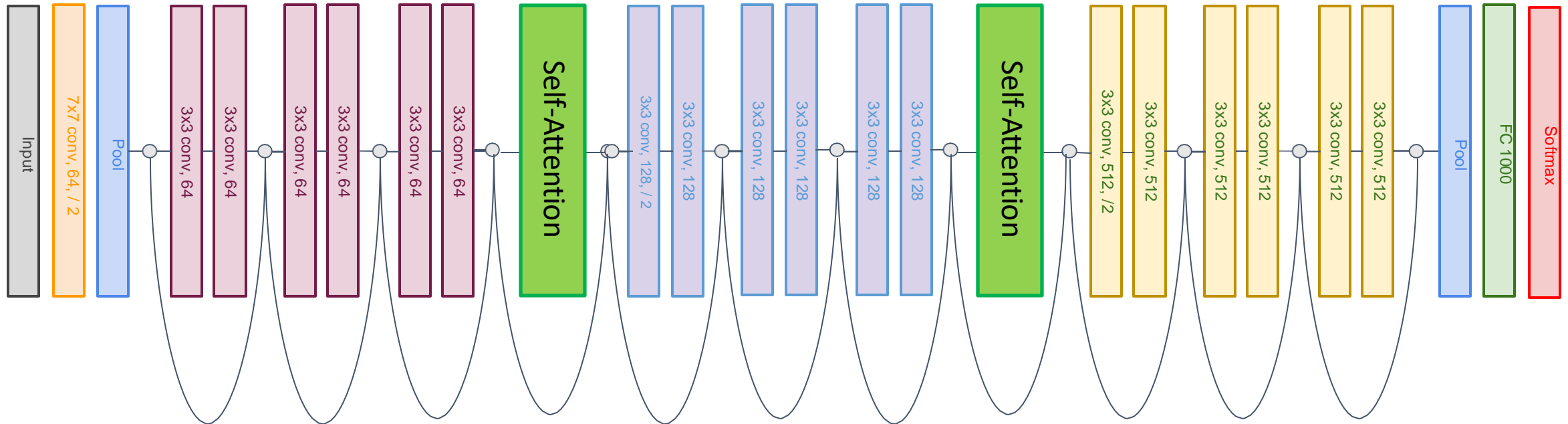
Wang et al, "Non-local Neural Networks", CVPR 2018

# Idea #1: Add attention to existing CNNs

## CNNs

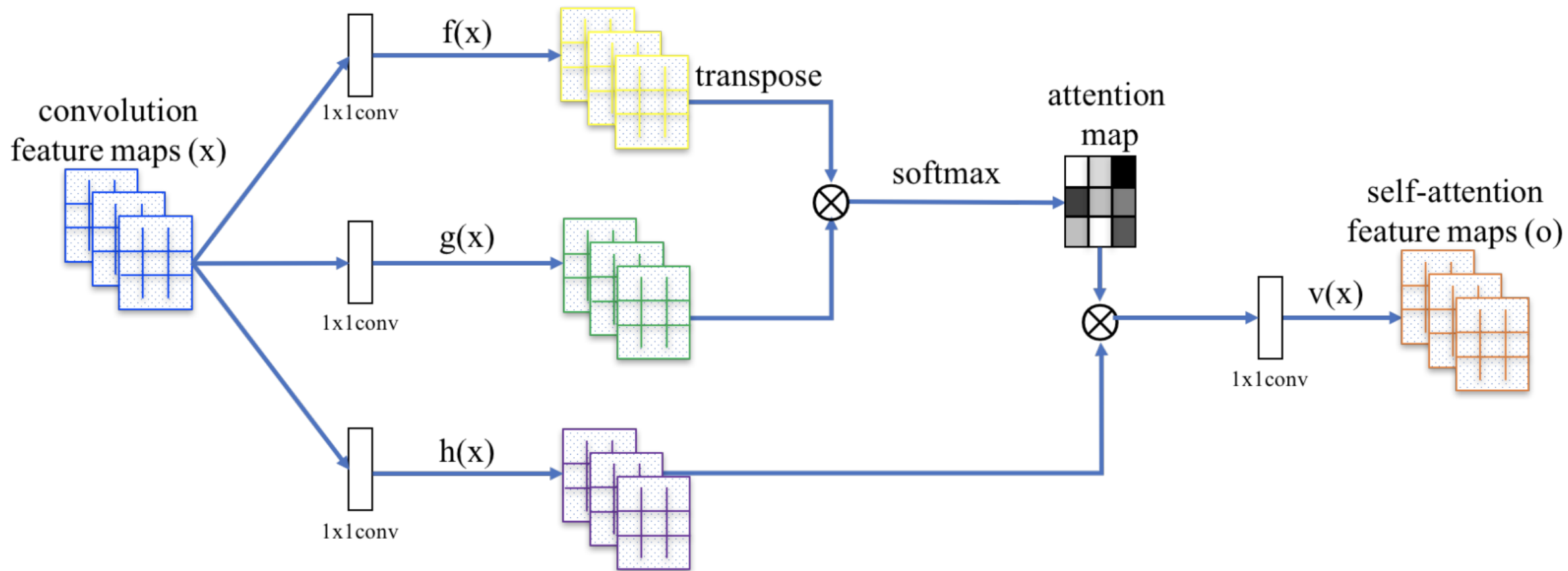
Model is still a CNN! Start from standard CNN architecture (e.g. ResNet)

Can we replace  
convolution entirely? Add Self-Attention blocks between existing ResNet blocks



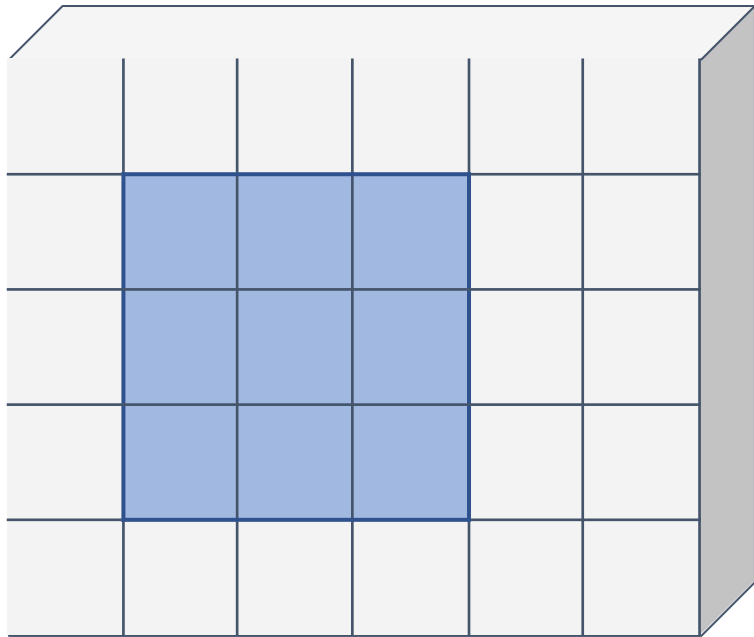
Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

Wang et al, "Non-local Neural Networks", CVPR 2018

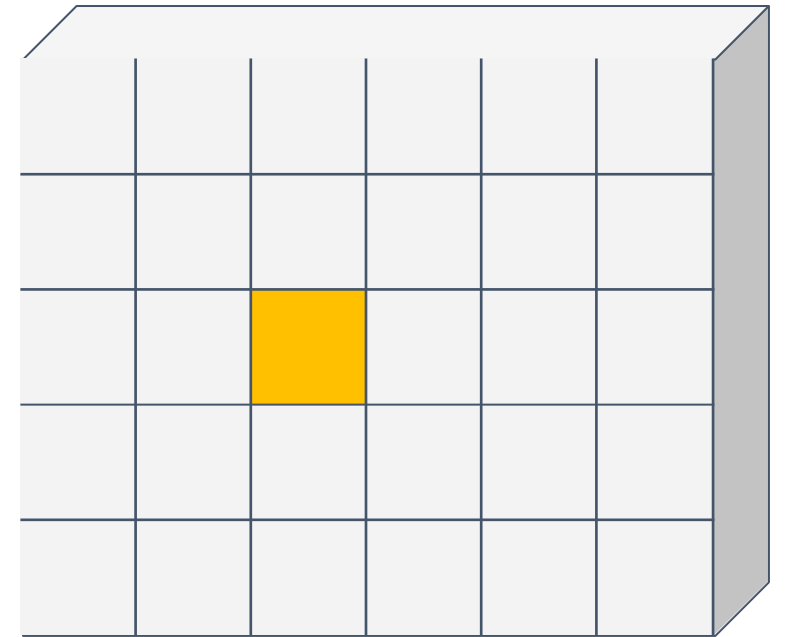


# Idea #2: Replace Convolution with “Local Attention”

Convolution: Output at each position is inner product of conv kernel with receptive field in input



Input:  $C \times H \times W$

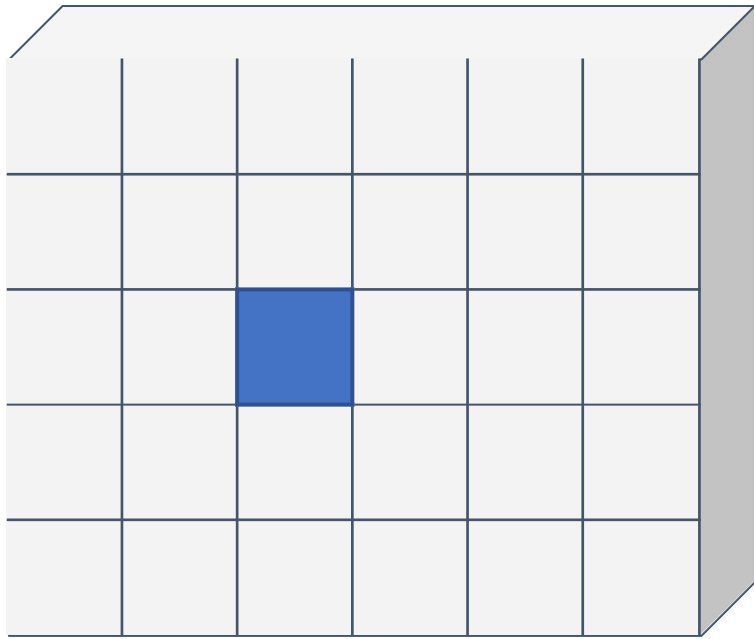


Output:  $C' \times H \times W$

Hu et al, “Local Relation Networks for Image Recognition”, ICCV 2019; Ramachandran et al, “Stand-Alone Self-Attention in Vision Models”, NeurIPS 2019

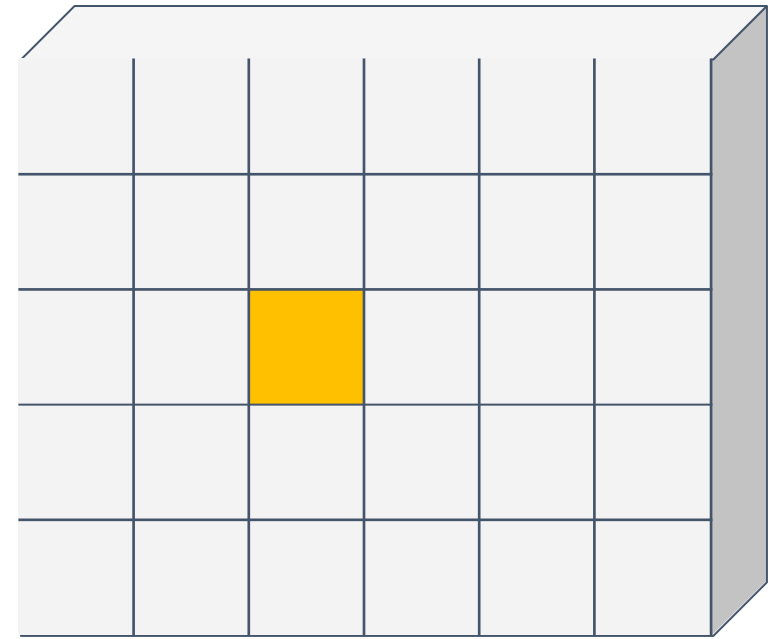
# Idea #2: Replace Convolution with “Local Attention”

Map center of receptive field to **query**



Input:  $C \times H \times W$

Query:  $D_Q$



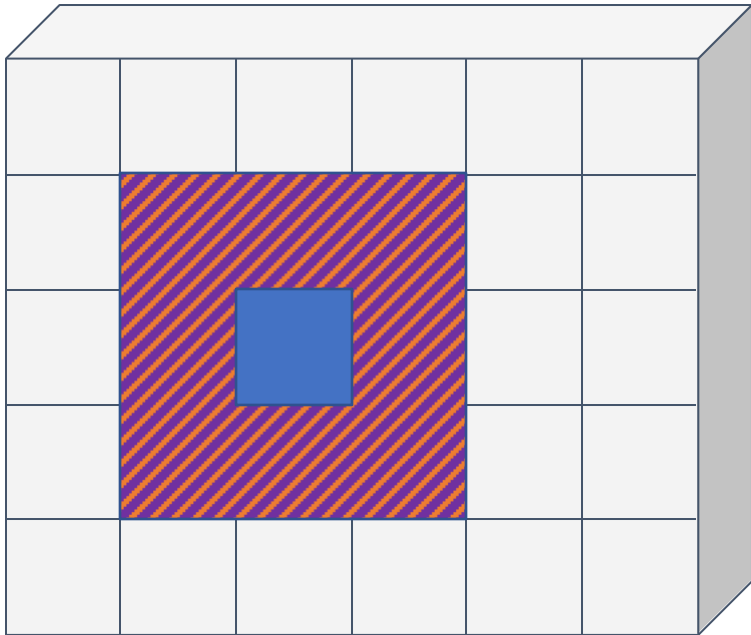
Output:  $C' \times H \times W$

Hu et al, “Local Relation Networks for Image Recognition”, ICCV 2019; Ramachandran et al, “Stand-Alone Self-Attention in Vision Models”, NeurIPS 2019

# Idea #2: Replace Convolution with “Local Attention”

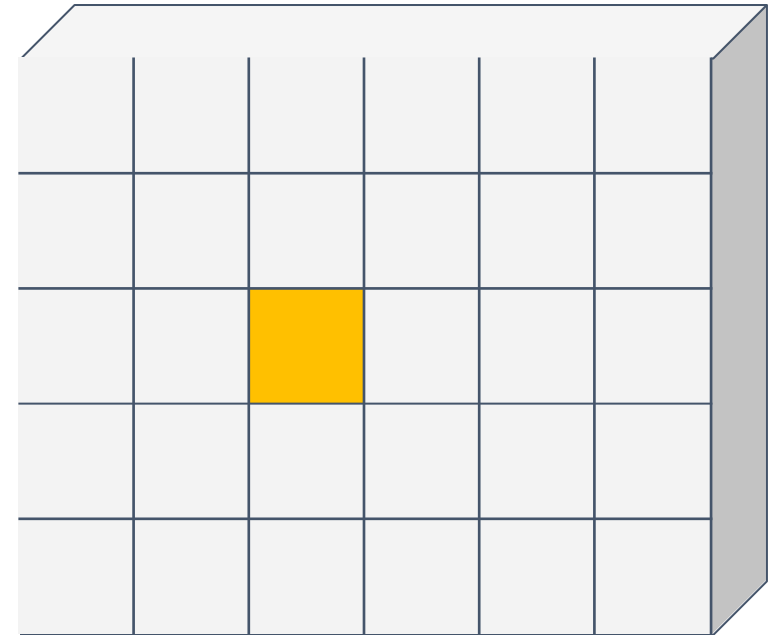
Map center of receptive field to **query**

Map each element in receptive field to **key** and **value**



Input:  $C \times H \times W$

Query:  $D_Q$   
Keys:  $R \times R \times D_Q$   
Values:  $R \times R \times C'$



Output:  $C' \times H \times W$

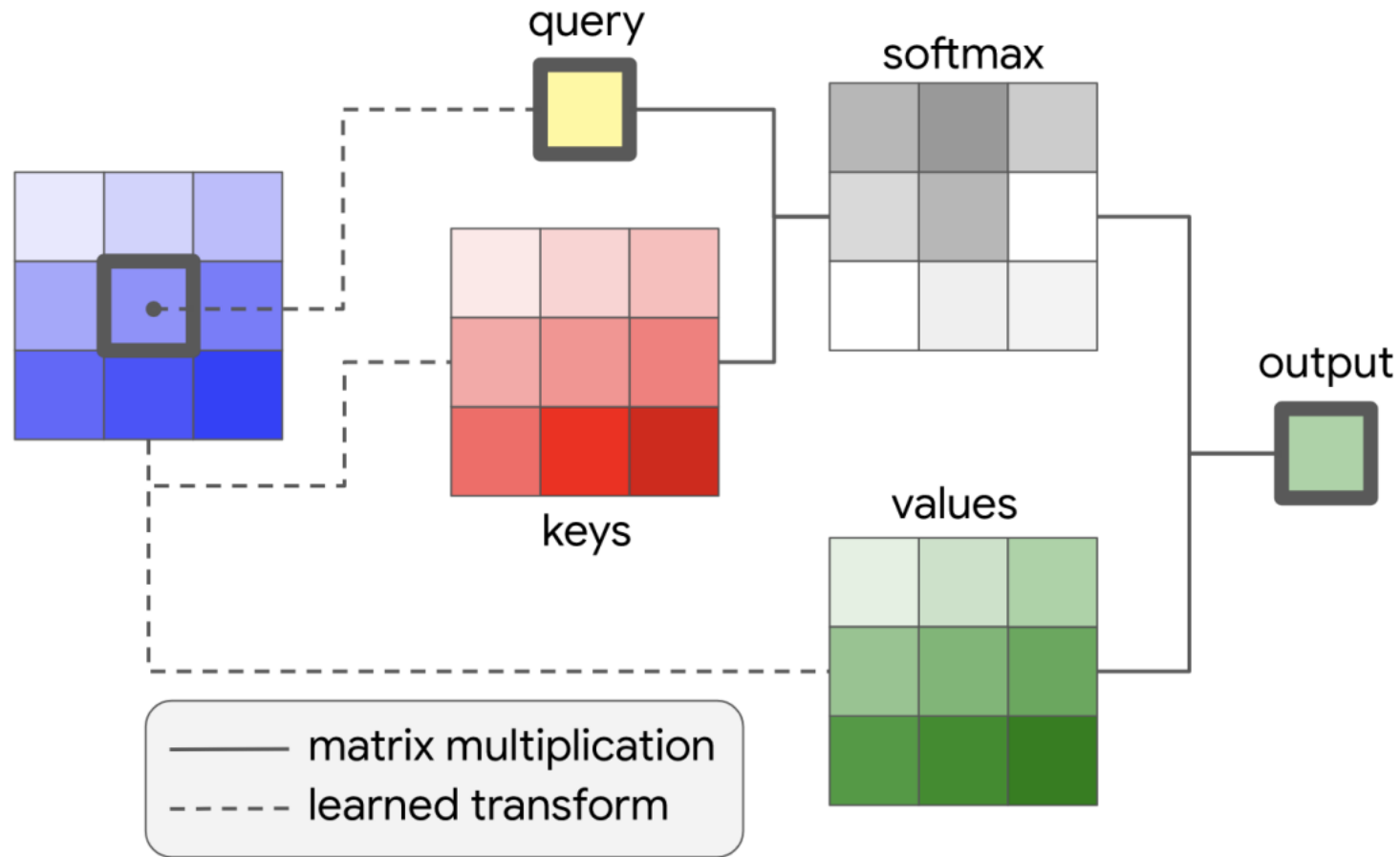
Hu et al, “Local Relation Networks for Image Recognition”, ICCV 2019; Ramachandran et al, “Stand-Alone Self-Attention in Vision Models”, NeurIPS 2019



# Idea #2: Replace Convolution with “Local Attention”

Map center of receptive field to **query**

Map each element in receptive field to **key** and **value**



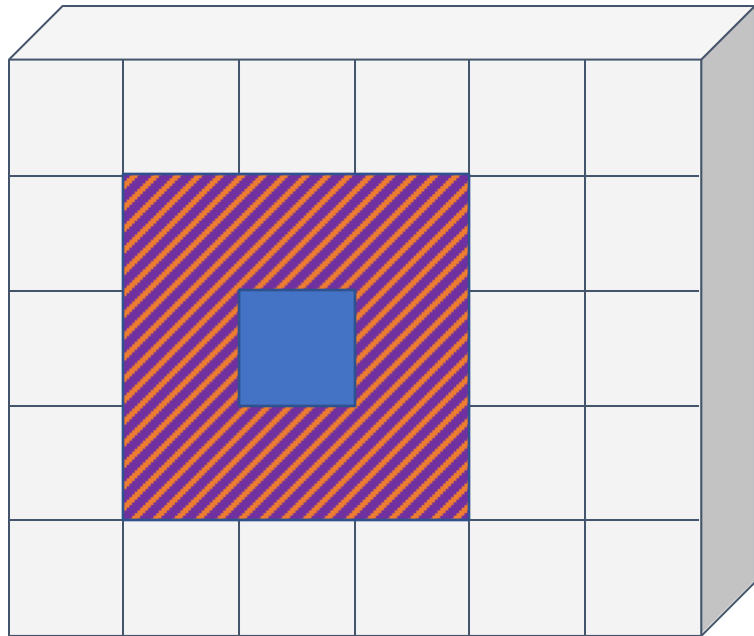
Ramachandran et al, “Stand-Alone Self-Attention in Vision Models”, NeurIPS 2019

# Idea #2: Replace Convolution with “Local Attention”

Map center of receptive field to **query**

Map each element in receptive field to **key** and **value**

Compute **output** using attention



Input:  $C \times H \times W$

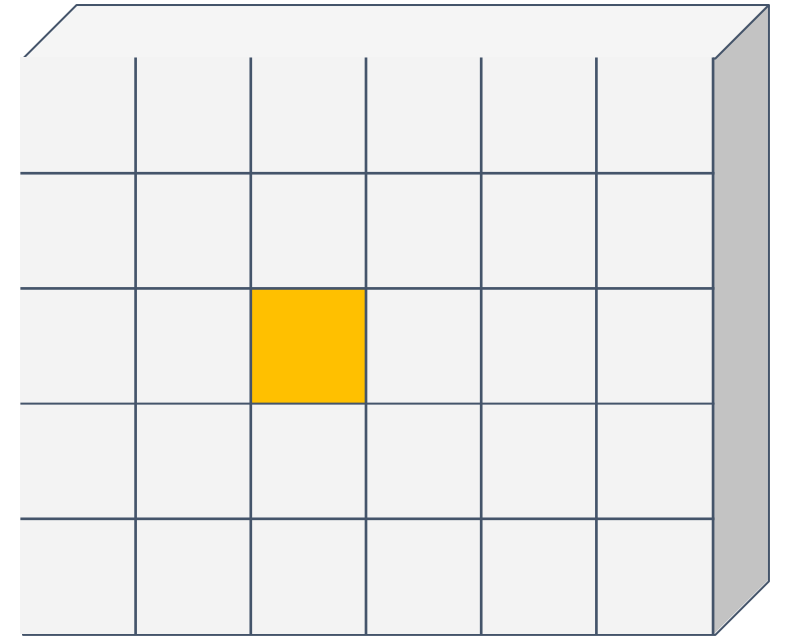
Query:  $D_Q$

Keys:  $R \times R \times D_Q$

Values:  $R \times R \times C'$

Output:  $C$

↓  
Attention  
↑



Output:  $C' \times H \times W$

Hu et al, “Local Relation Networks for Image Recognition”, ICCV 2019; Ramachandran et al, “Stand-Alone Self-Attention in Vision Models”, NeurIPS 2019

# Idea #2: Replace Convolution with “Local Attention”

Map center of receptive field to **query**

Map each element in receptive field to **key** and **value**

Compute **output** using attention

Replace all conv in ResNet with local attention

LR = “Local Relation”

stage	output	ResNet-50	LR-Net-50 ( $7\times 7, m=8$ )
res1	$112\times 112$	$7\times 7$ conv, 64, stride 2	$1\times 1, 64$ $7\times 7$ LR, 64, stride 2
res2	$56\times 56$	$3\times 3$ max pool, stride 2	$3\times 3$ max pool, stride 2
		$\begin{bmatrix} 1\times 1, 64 \\ 3\times 3 \text{ conv}, 64 \\ 1\times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 100 \\ 7\times 7 \text{ LR}, 100 \\ 1\times 1, 256 \end{bmatrix} \times 3$
res3	$28\times 28$	$\begin{bmatrix} 1\times 1, 128 \\ 3\times 3 \text{ conv}, 128 \\ 1\times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1\times 1, 200 \\ 7\times 7 \text{ LR}, 200 \\ 1\times 1, 512 \end{bmatrix} \times 4$
res4	$14\times 14$	$\begin{bmatrix} 1\times 1, 256 \\ 3\times 3 \text{ conv}, 256 \\ 1\times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1\times 1, 400 \\ 7\times 7 \text{ LR}, 400 \\ 1\times 1, 1024 \end{bmatrix} \times 6$
res5	$7\times 7$	$\begin{bmatrix} 1\times 1, 512 \\ 3\times 3 \text{ conv}, 512 \\ 1\times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 800 \\ 7\times 7 \text{ LR}, 800 \\ 1\times 1, 2048 \end{bmatrix} \times 3$
	$1\times 1$	global average pool 1000-d fc, softmax	global average pool 1000-d fc, softmax
# params		$25.5\times 10^6$	$23.3\times 10^6$
FLOPs		$4.3\times 10^9$	$4.3\times 10^9$

Hu et al, “Local Relation Networks for Image Recognition”, ICCV 2019;

Ramachandran et al, “Stand-Alone Self-Attention in Vision Models”, NeurIPS 2019

# Idea #2: Replace Convolution with “Local Attention”

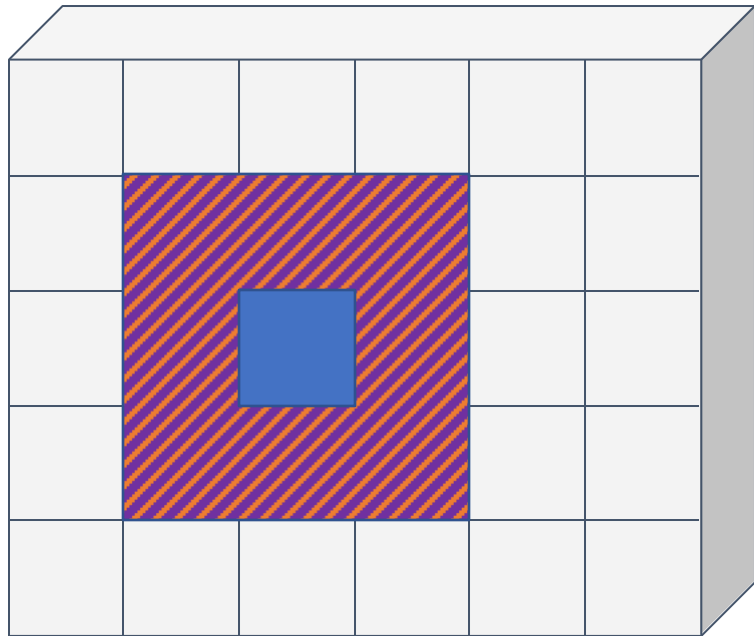
Map center of receptive field to **query**

Map each element in receptive field to **key** and **value**

Compute **output** using attention

Replace all conv in ResNet with local attention

Lots of tricky details,  
hard to implement,  
only marginally better  
than ResNets

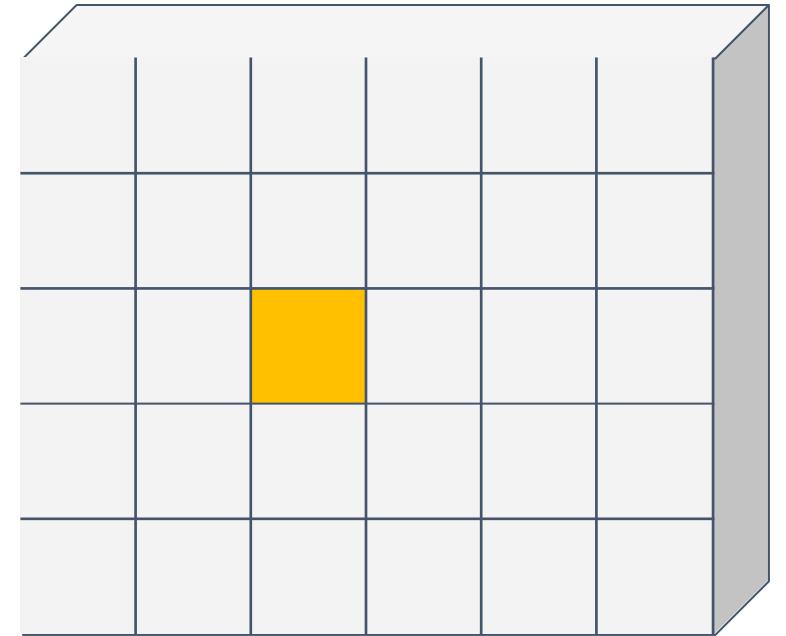


Input:  $C \times H \times W$

Query:  $D_Q$   
Keys:  $R \times R \times D_Q$   
Values:  $R \times R \times C'$

Output:  $C$

↓  
Attention  
↑

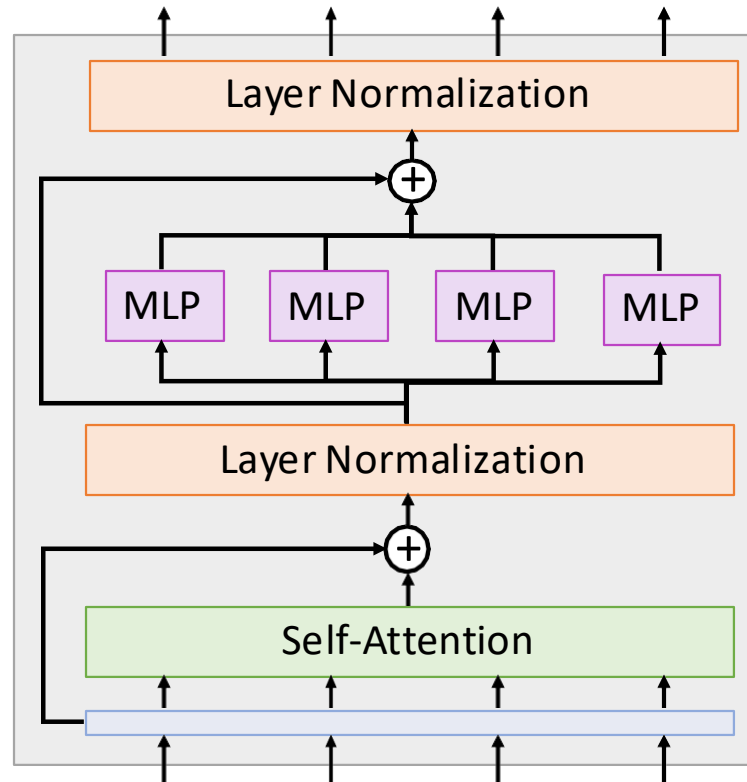
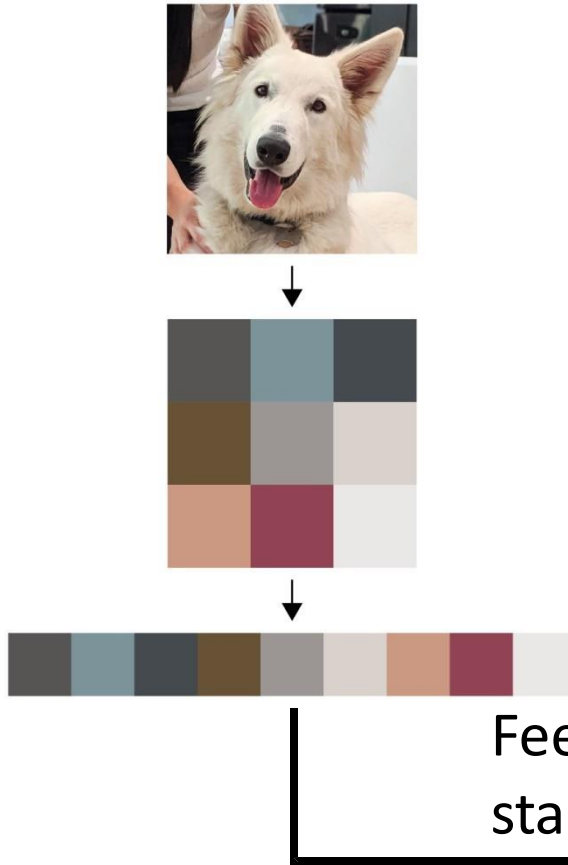


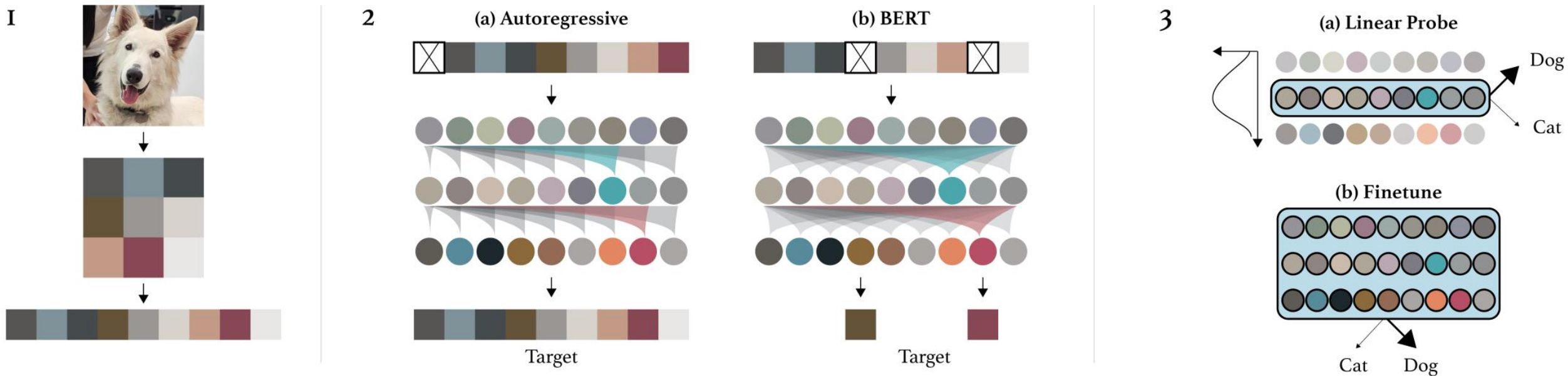
Output:  $C' \times H \times W$

Hu et al, “Local Relation Networks for Image Recognition”, ICCV 2019; Ramachandran et al, “Stand-Alone Self-Attention in Vision Models”, NeurIPS 2019

# Idea #3: Standard Transformer on Pixels

Treat an image as a set of pixel values

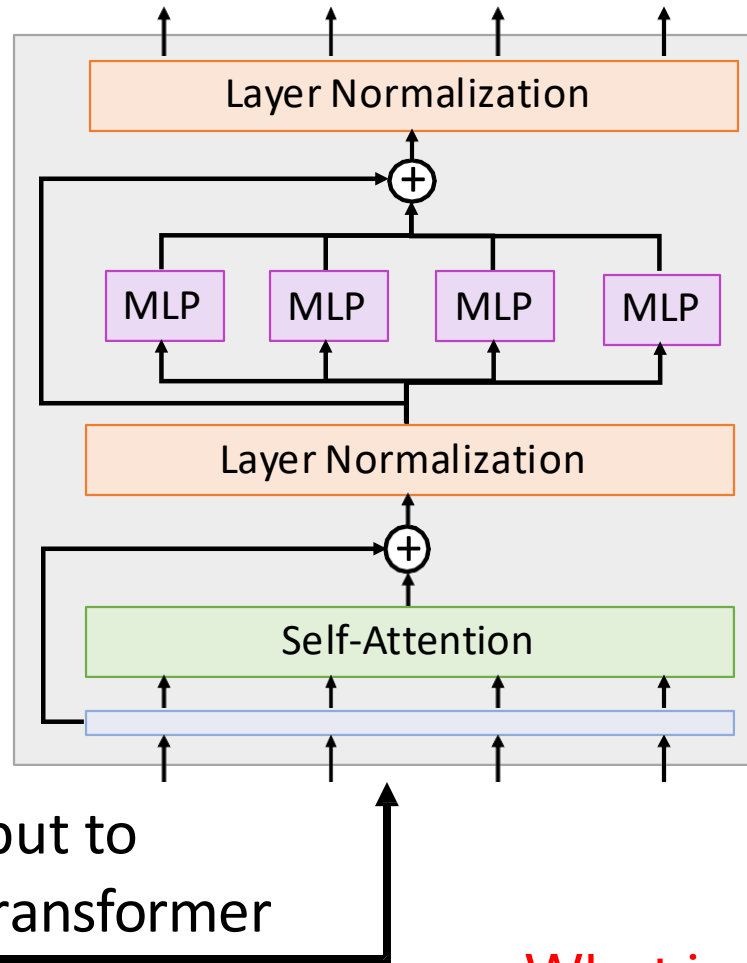
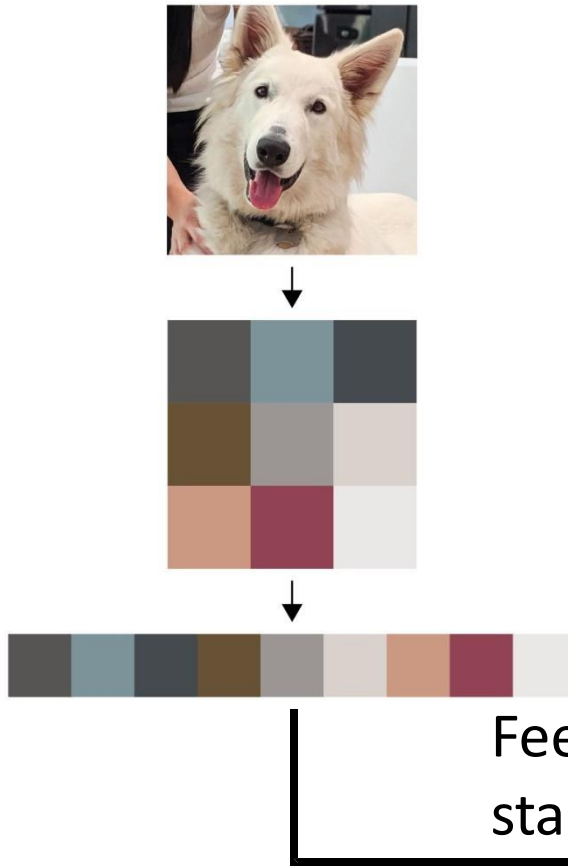




- First, pre-process raw images by resizing to a low resolution and reshaping into a 1D sequence.
- pre-training objectives, auto-regressive next pixel prediction or masked pixel prediction.
- Evaluate the representations learned by these objectives with linear probes or fine-tuning.

# Idea #3: Standard Transformer on Pixels

Treat an image as a set of pixel values



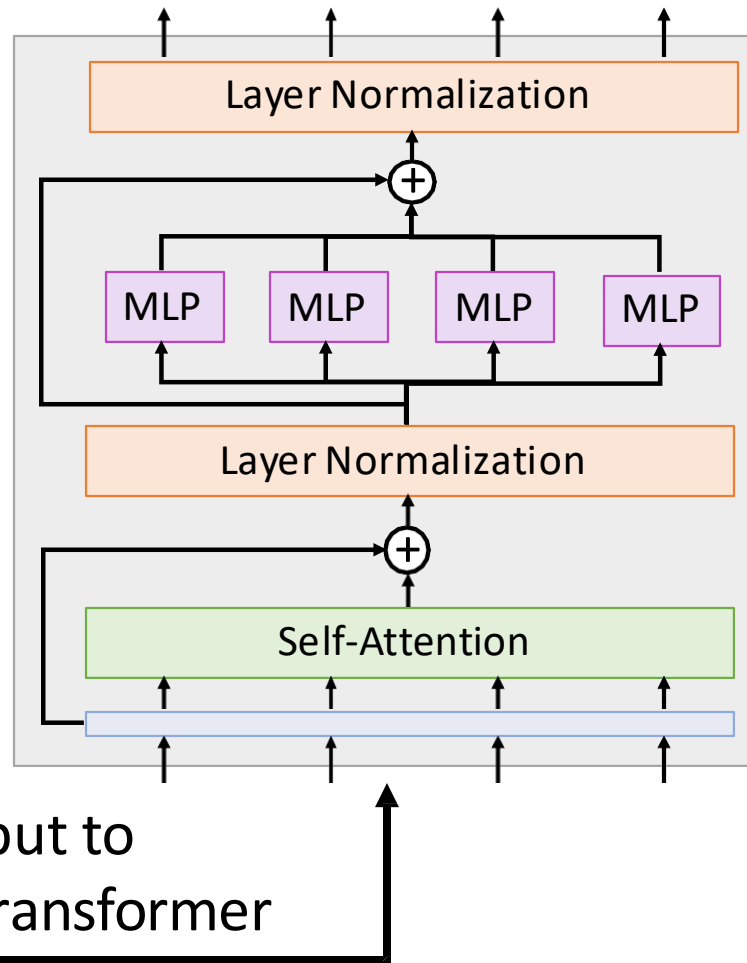
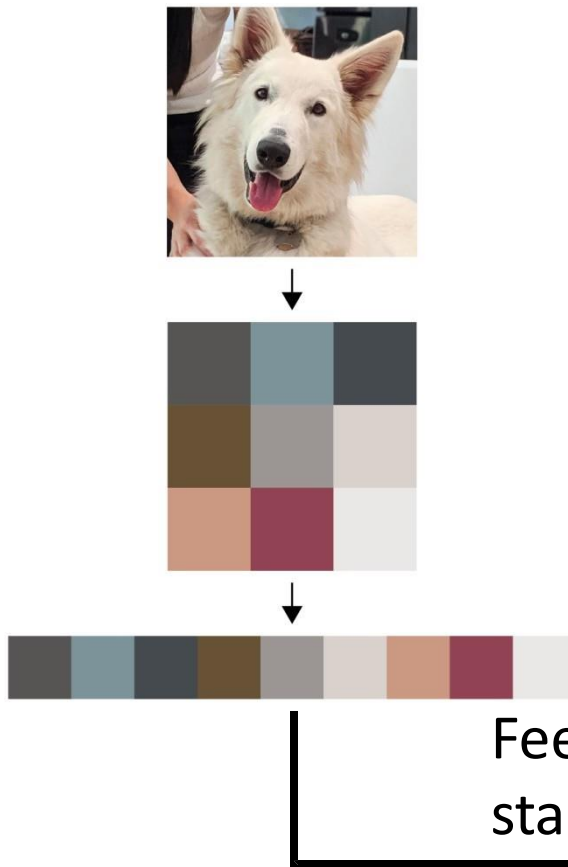
Problem: Memory use!

$R \times R$  image needs  $R^4$  elements per attention matrix

What is the problem?

# Idea #3: Standard Transformer on Pixels

Treat an image as a set of pixel values



Problem: Memory use!

$R \times R$  image needs  $R^4$  elements per attention matrix

$R=128$ , 48 layers, 16 heads per layer takes 768GB of memory for attention matrices for a single example...



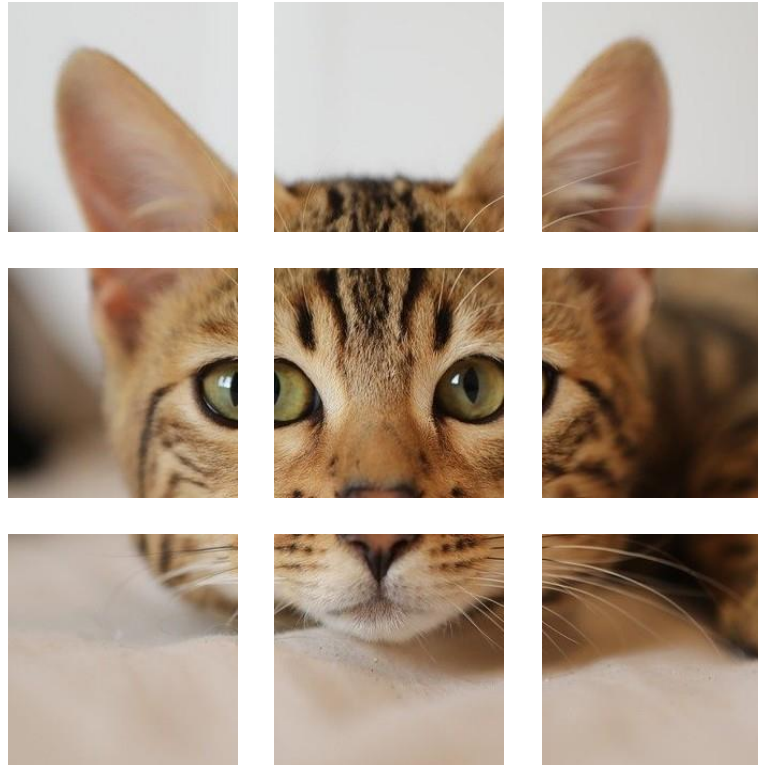
# Idea #4: Standard Transformer on Patches



Dosovitskiy et al, “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”, ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Idea #4: Standard Transformer on Patches

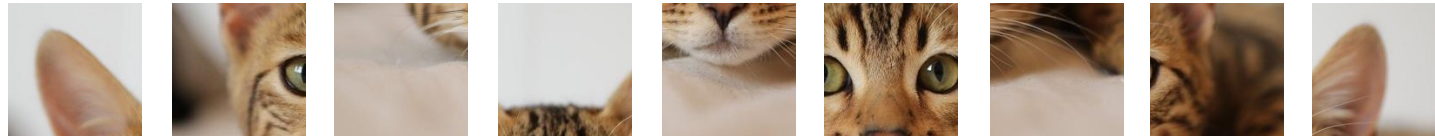


Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Idea #4: Standard Transformer on Patches

N input patches, each  
of shape 3x16x16



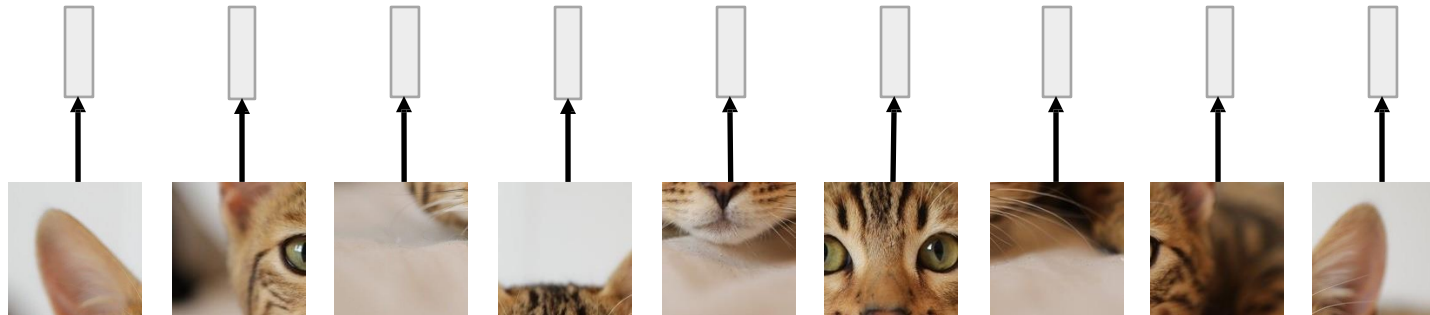
Dosovitskiy et al, “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”, ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

# Idea #4: Standard Transformer on Patches

Linear projection to  
D-dimensional vector

N input patches, each  
of shape 3x16x16



Dosovitskiy et al, “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”, ICLR 2021

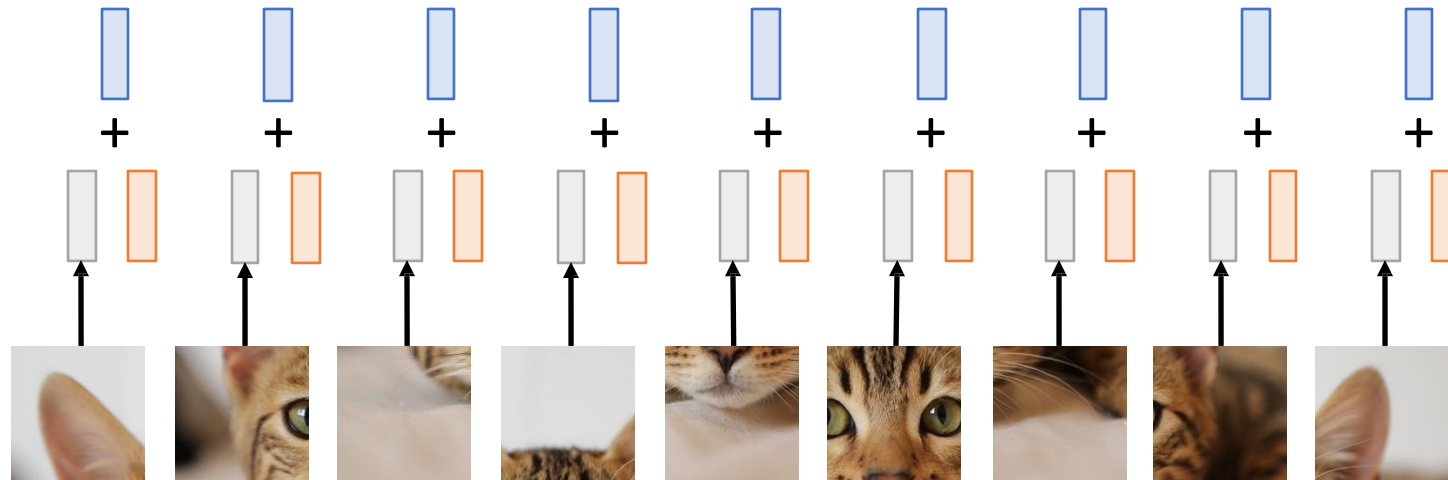
[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

# Idea #4: Standard Transformer on Patches

Add positional embedding: learned D-dim vector per position

Linear projection to D-dimensional vector

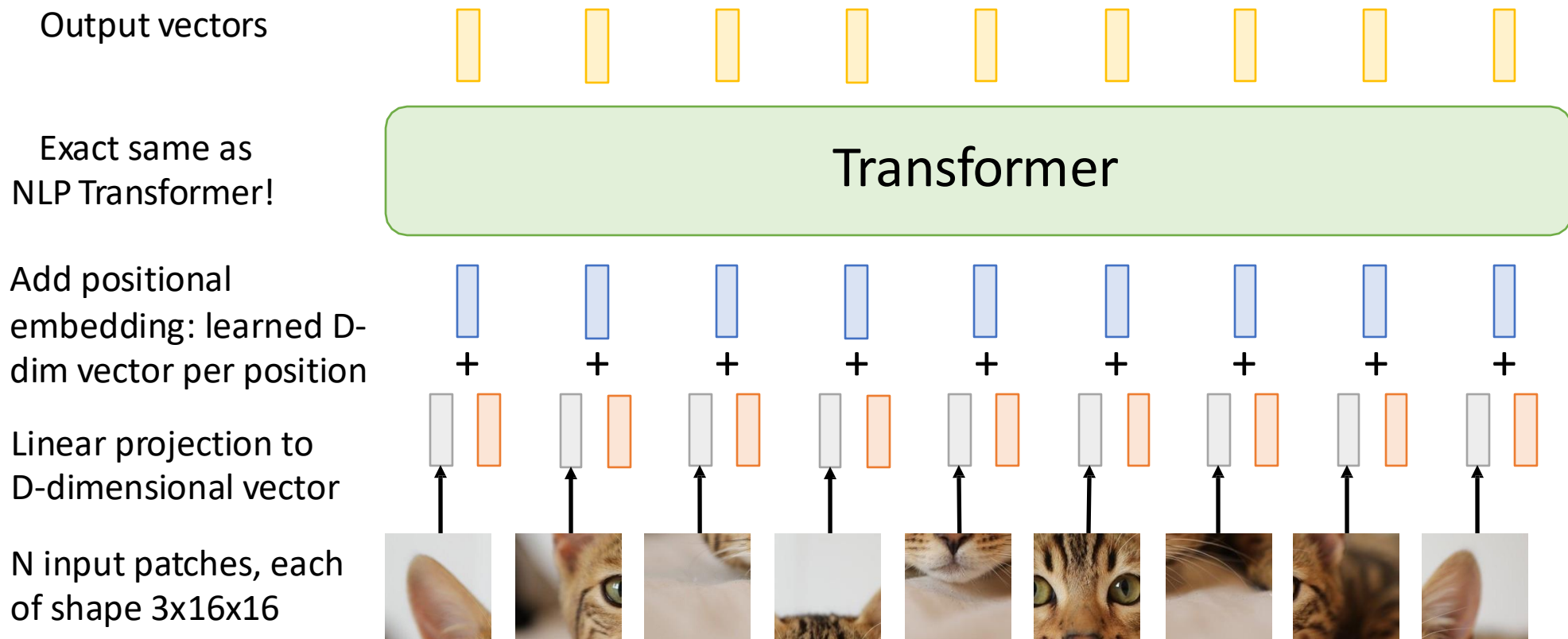
N input patches, each of shape 3x16x16



Dosovitskiy et al, “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”, ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

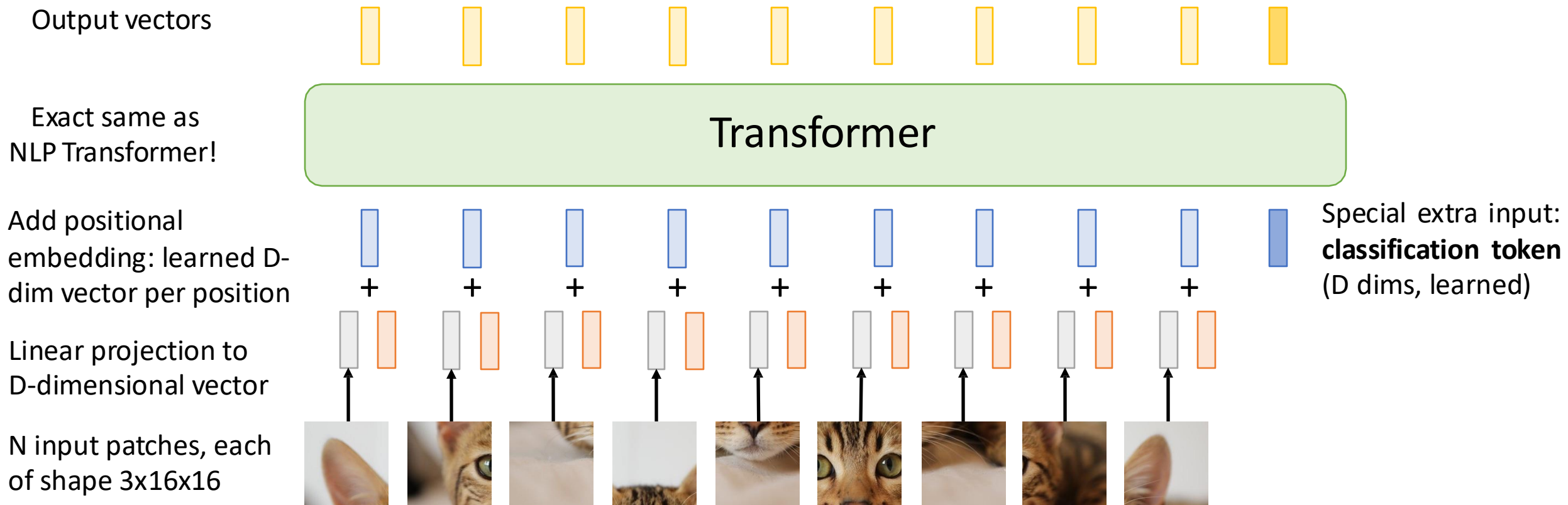
# Idea #4: Standard Transformer on Patches



Dosovitskiy et al, “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”, ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

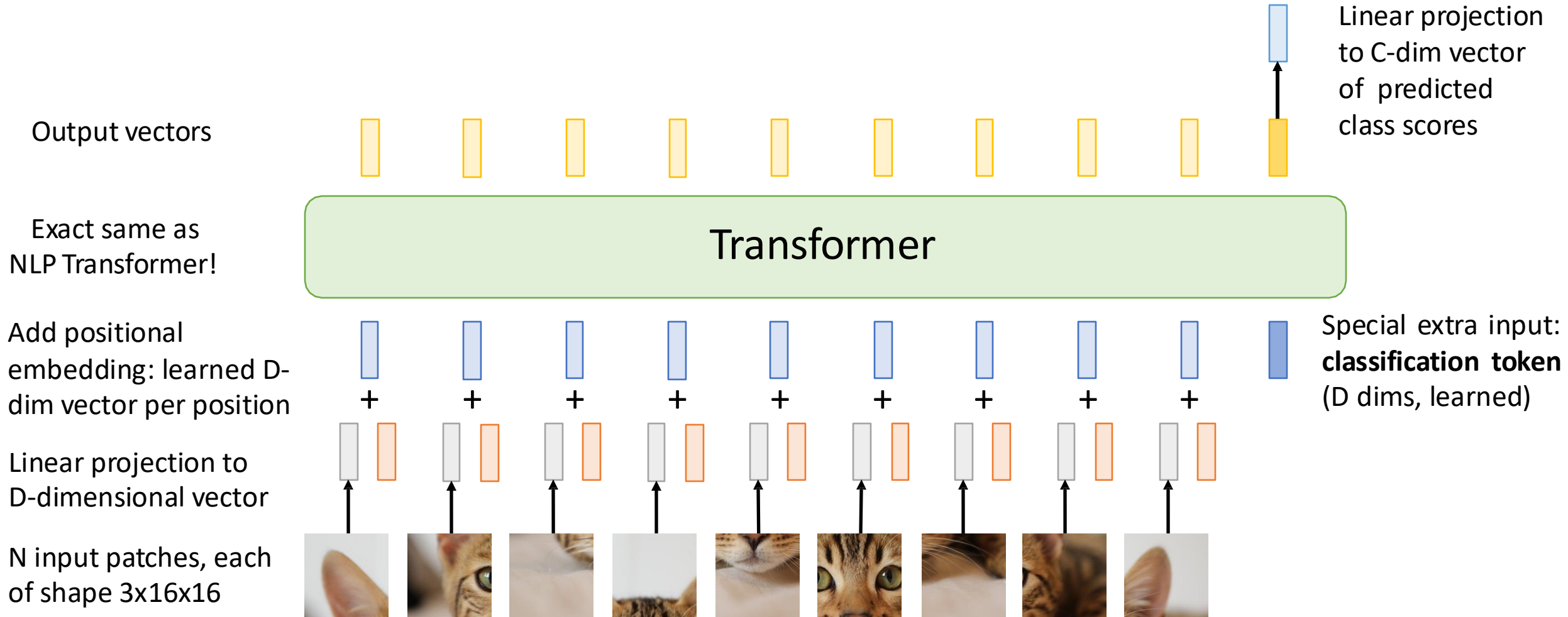
# Idea #4: Standard Transformer on Patches



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

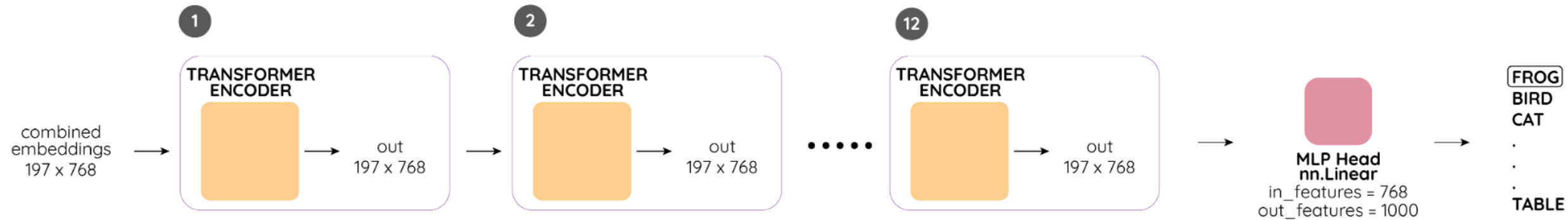
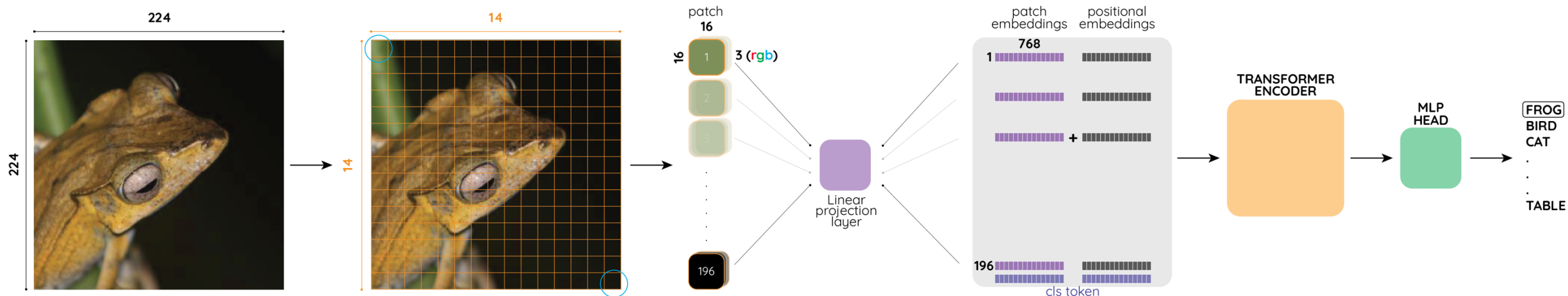
# Idea #4: Standard Transformer on Patches



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)





What will be the order of attention matrix?

# Vision Transformer

(ViT)

Computer vision model  
with no convolutions!

Output vectors

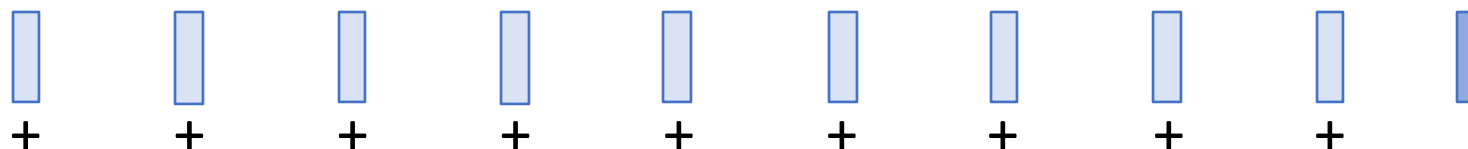


Linear projection  
to C-dim vector  
of predicted  
class scores

Exact same as  
NLP Transformer!



Add positional  
embedding: learned D-  
dim vector per position

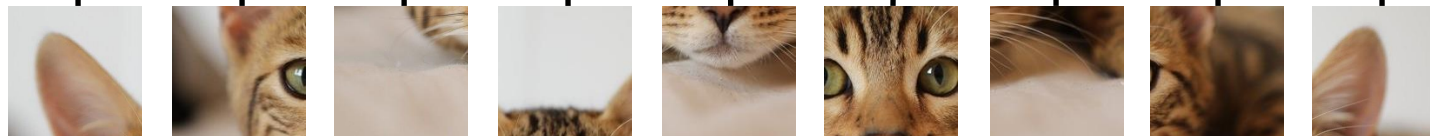


Special extra input:  
**classification token**  
(D dims, learned)

Linear projection to  
D-dimensional vector



N input patches, each  
of shape 3x16x16



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

# Vision Transformer (ViT)

In practice: take 224x224 input image, divide into 14x14 grid of 16x16 pixel patches (or 16x16 grid of 14x14 patches)

Each attention matrix has  $14^4 = 38,416$  entries, takes 150 KB (or 65,536 entries, takes 256 KB)

Output vectors

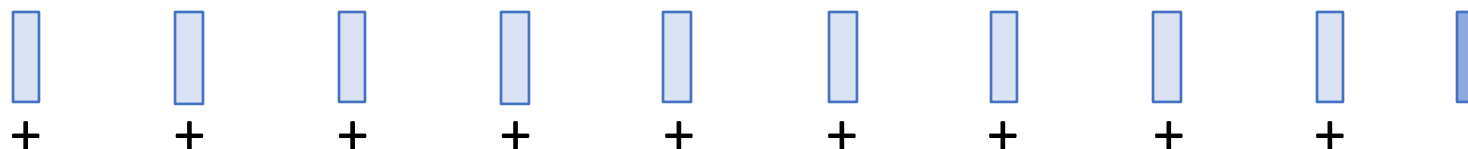


Linear projection to C-dim vector of predicted class scores

Exact same as NLP Transformer!

Transformer

Add positional embedding: learned D-dim vector per position

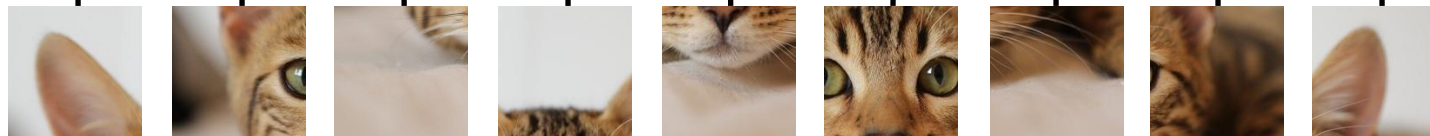


Special extra input: **classification token** (D dims, learned)

Linear projection to D-dimensional vector



N input patches, each of shape 3x16x16



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformer (ViT)

In practice: take 224x224 input image, divide into 14x14 grid of 16x16 pixel patches (or 16x16 grid of 14x14 patches)

With 48 layers, 16 heads per layer, all attention matrices take 112 MB (or 192MB)

Output vectors



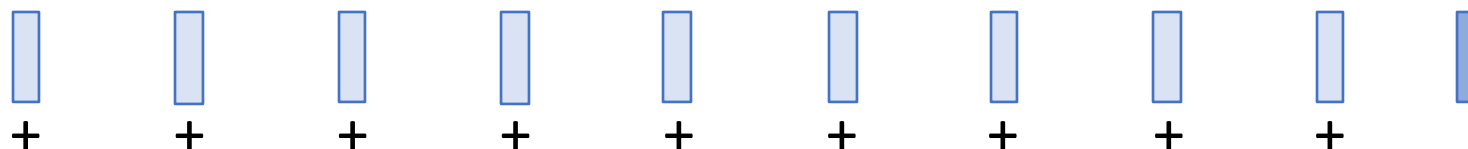
Linear projection to C-dim vector of predicted class scores



Exact same as NLP Transformer!



Add positional embedding: learned D-dim vector per position



Linear projection to D-dimensional vector



N input patches, each of shape 3x16x16



Special extra input: **classification token** (D dims, learned)



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformer (ViT)

In practice: take 224x224 input image, divide into 14x14 grid of 16x16 pixel patches (or 16x16 grid of 14x14 patches)

With 48 layers, 16 heads per layer, all attention matrices take 112 MB (or 192MB)

Output vectors



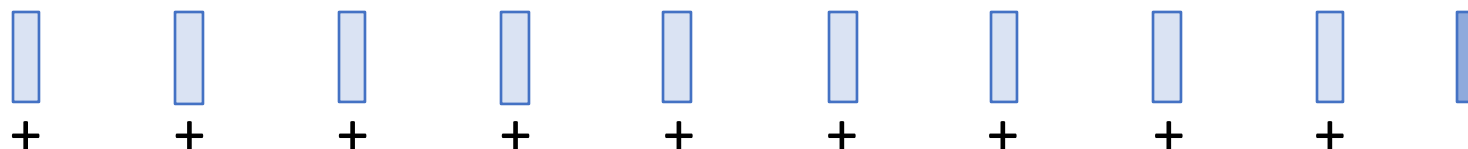
Linear projection to C-dim vector of predicted class scores



Exact same as NLP Transformer!



Add positional embedding: learned D-dim vector per position



Linear projection to D-dimensional vector



N input patches, each of shape 3x16x16



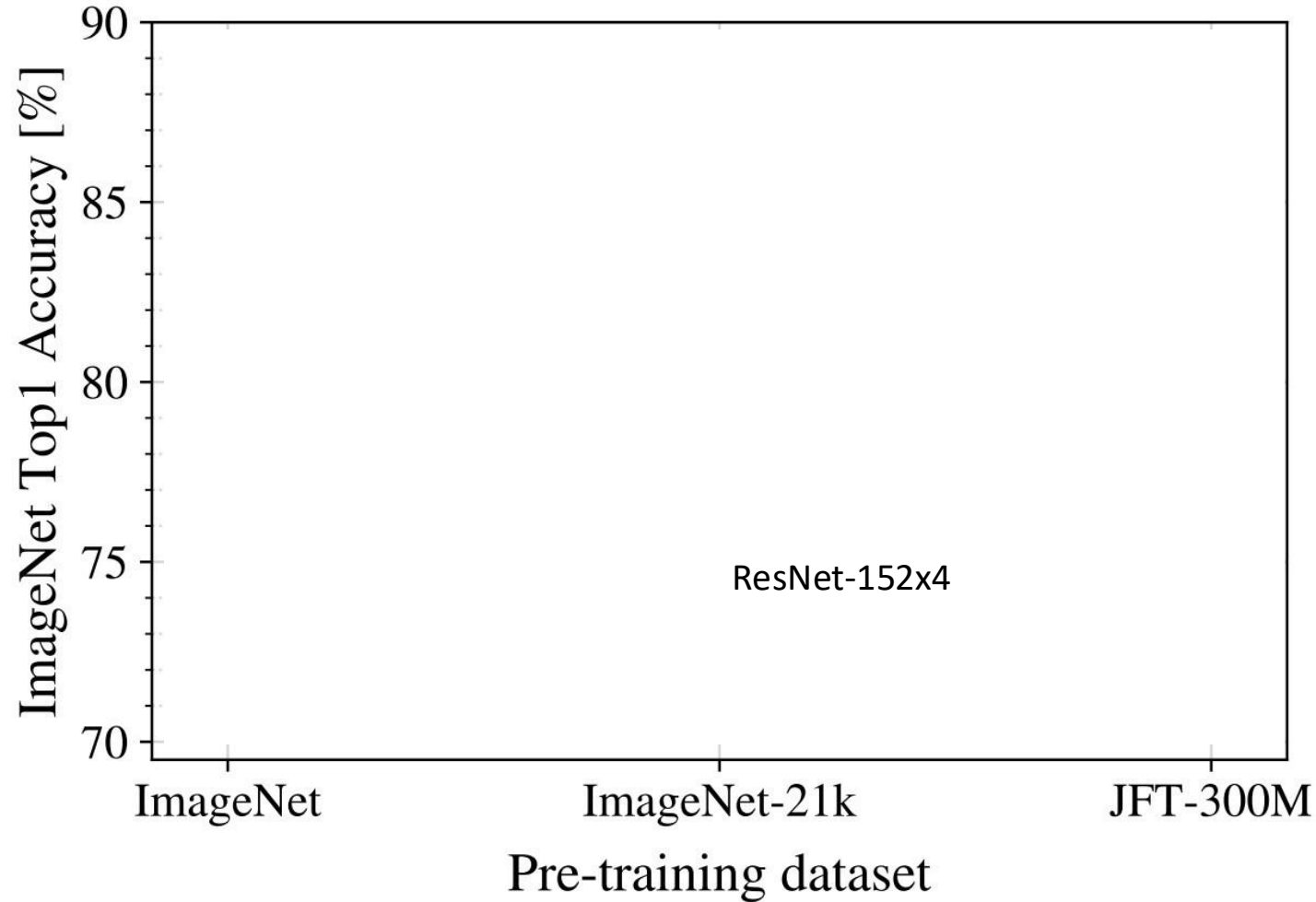
Special extra input: **classification token** (D dims, learned)



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformer (ViT) vs ResNets



B = Base  
L = Large  
H = Huge

/32, /16, /14 is patch size; smaller patch size is a bigger model (more patches)

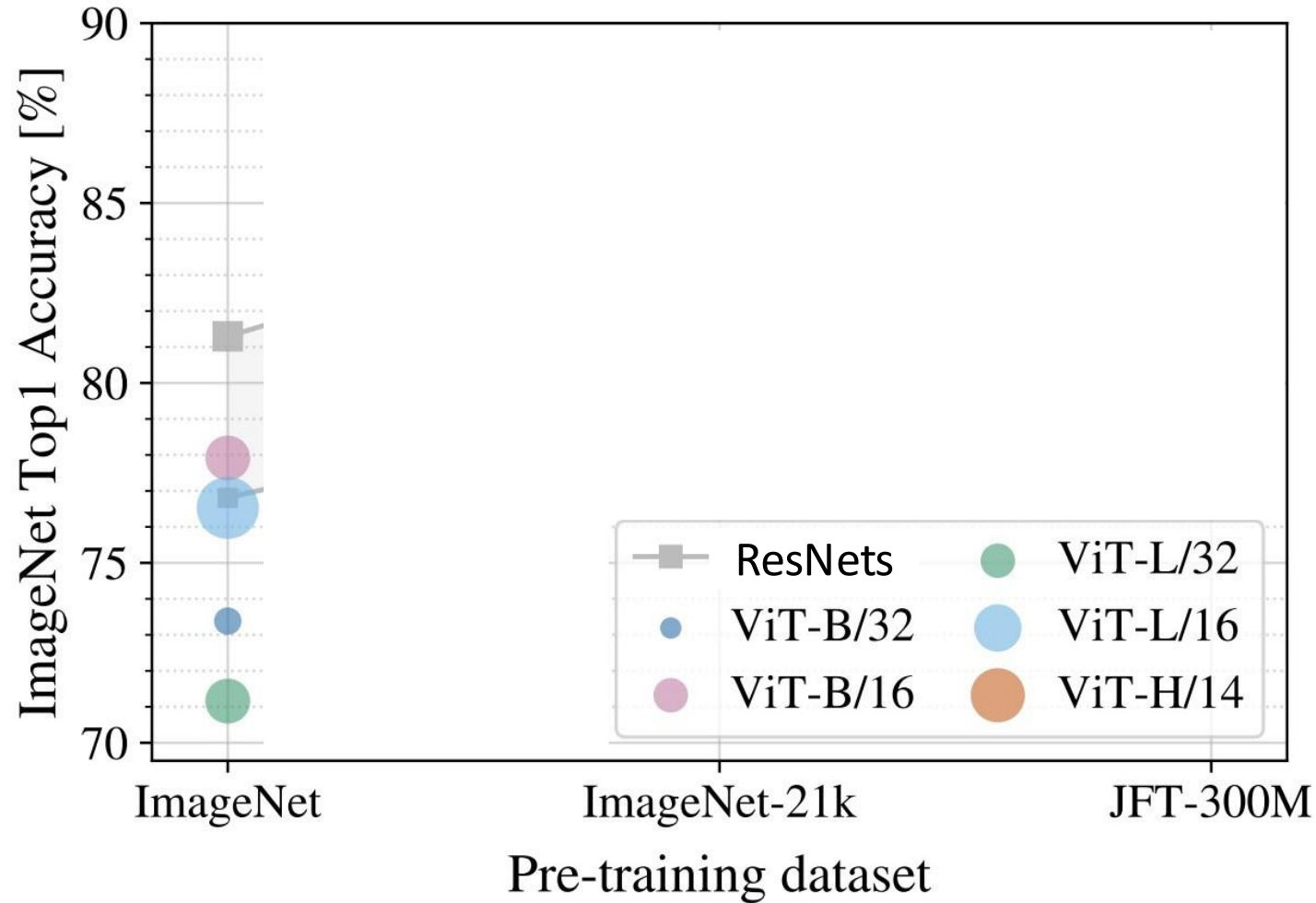
Imagenet 21k - 21k classes with 14M images  
JFT 300M - 18,291 classes

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

# Vision Transformer (ViT) vs ResNets

Recall: ImageNet dataset has 1k categories, 1.2M images

When trained on ImageNet, ViT models perform worse than ResNets



B = Base  
L = Large  
H = Huge

/32, /16, /14 is patch size; smaller patch size is a bigger model (more patches)

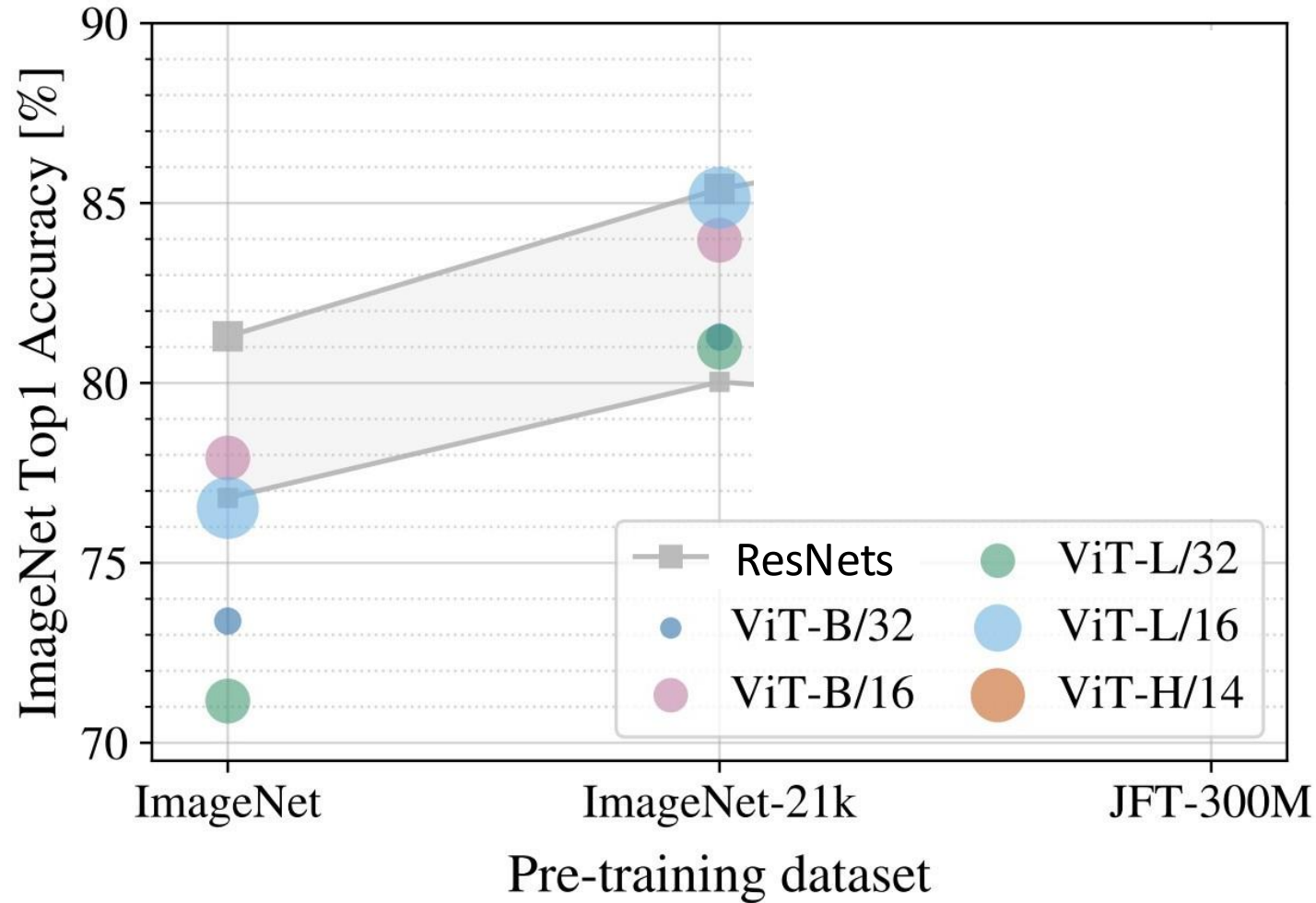
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021



# Vision Transformer (ViT) vs ResNets

ImageNet-21k has **14M** images with 21k categories

If you pretrain on ImageNet-21k and fine-tune on ImageNet, ViT does better: big ViTs match big ResNets



B = Base  
L = Large  
H = Huge

/32, /16, /14 is patch size; smaller patch size is a bigger model (more patches)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

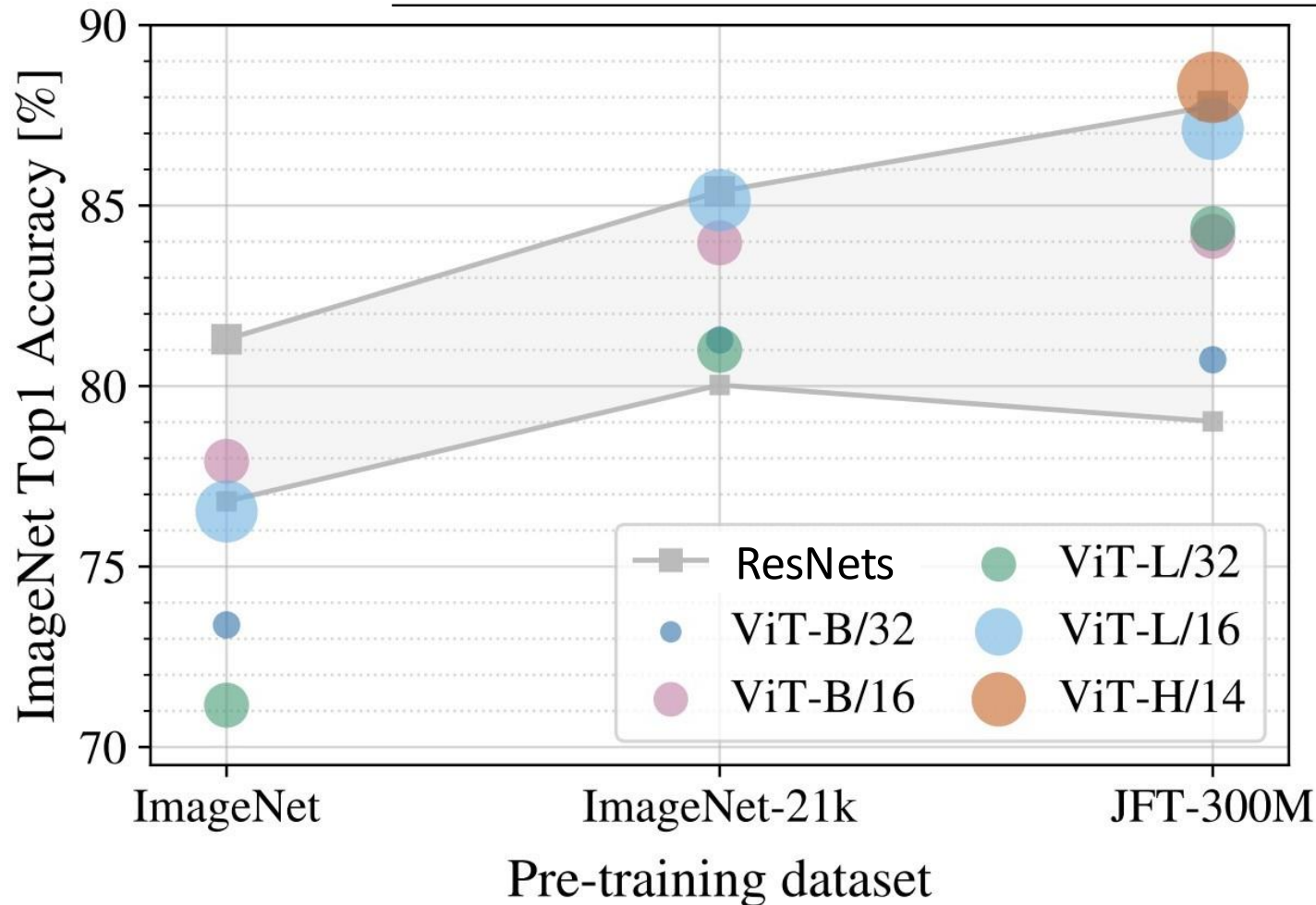
# Vision Transformer

## ResNets

JFT-300M is an internal Google dataset with 300M labeled images

If you pretrain on JFT and finetune on ImageNet, large ViTs outperform large ResNets

Model	Layers	Hidden size $D$	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M



Resnet-152 60M

B = Base  
L = Large  
H = Huge

/32, /16, /14 is patch size; smaller patch size is a bigger model (more patches)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

The "ViT-L/16" configuration of the  
**Patch Size:** 16x16 (The input image is divided into 16x16 patches)

•**Embedding Dimension:** 1024

•**Number of Transformer Blocks:** 24

•**Number of Heads in Multi-Head Attention:** 16

•**MLP Hidden Size:** 4096

•**Number of Classes:** Typically 1000 for ImageNet classification

So the "ViT-L/16" model has approximately 339 million parameters.

**1.Patch Embedding:**

$$16 \times 16 \times 3 \times 1024 = 786,432$$

**2.Positional Embedding:**

$$1024 \times 197 = 201,728$$

**3.Transformer Blocks:**

$$24 \times ((1024 \times 1024 \times 3) + (1024 \times 4096 \times 2)) \approx 337,002,496$$

**4.Classification Head:**

$$1024 \times 1000 = 1,024,000$$

The total is roughly:

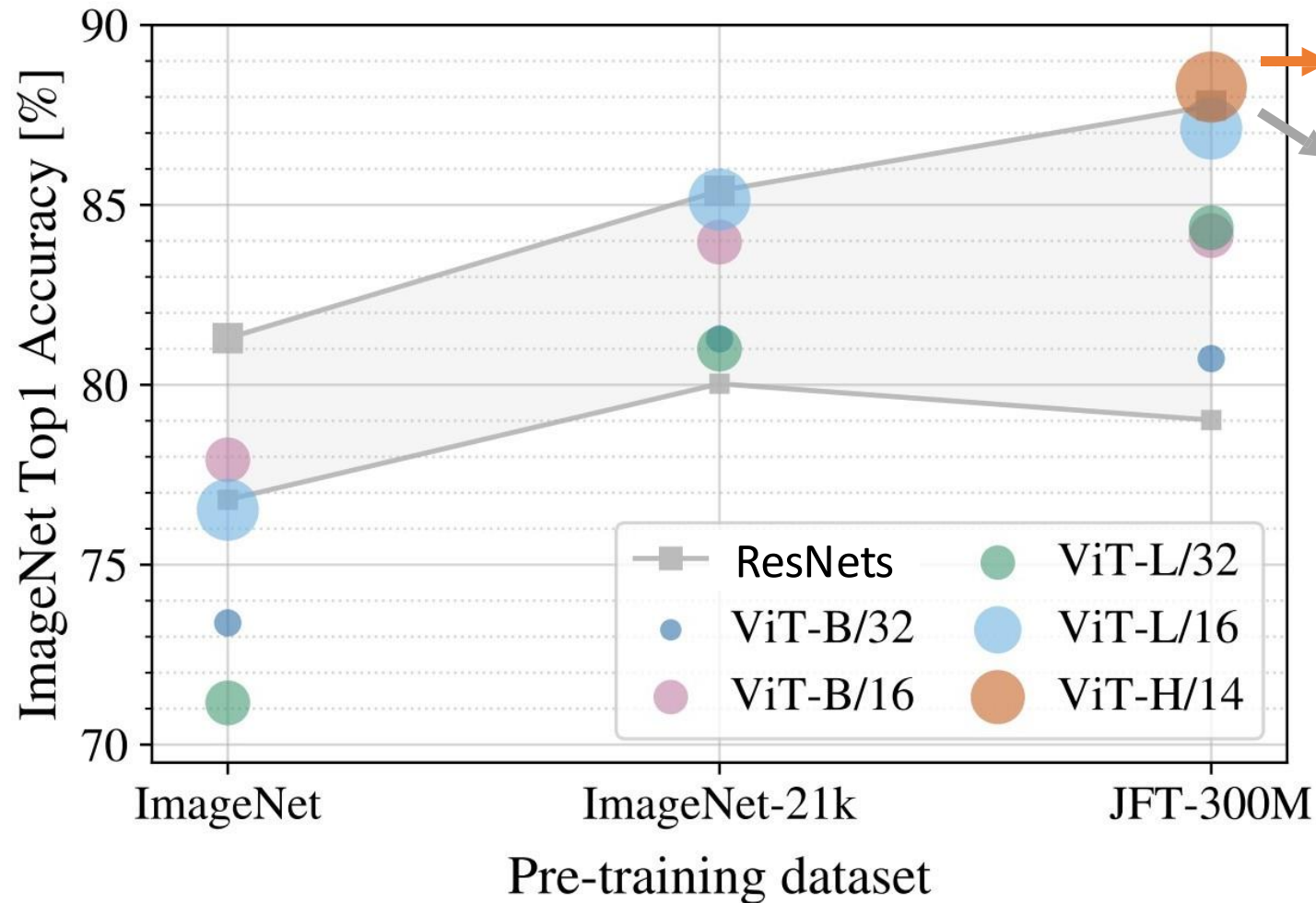
$$786,432 + 201,728 + 337,002,496 + 1,024,000 \approx 339,014,656$$

$$786,432 + 201,728 + 337,002,496 + 1,024,000 \approx 339,014,656$$

# Vision Transformer (ViT) vs ResNets

JFT-300M is an internal Google dataset with 300M labeled images

If you pretrain on JFT and finetune on ImageNet, large ViTs outperform large ResNets



ViT: 2.5k TPU-v3 core days of training

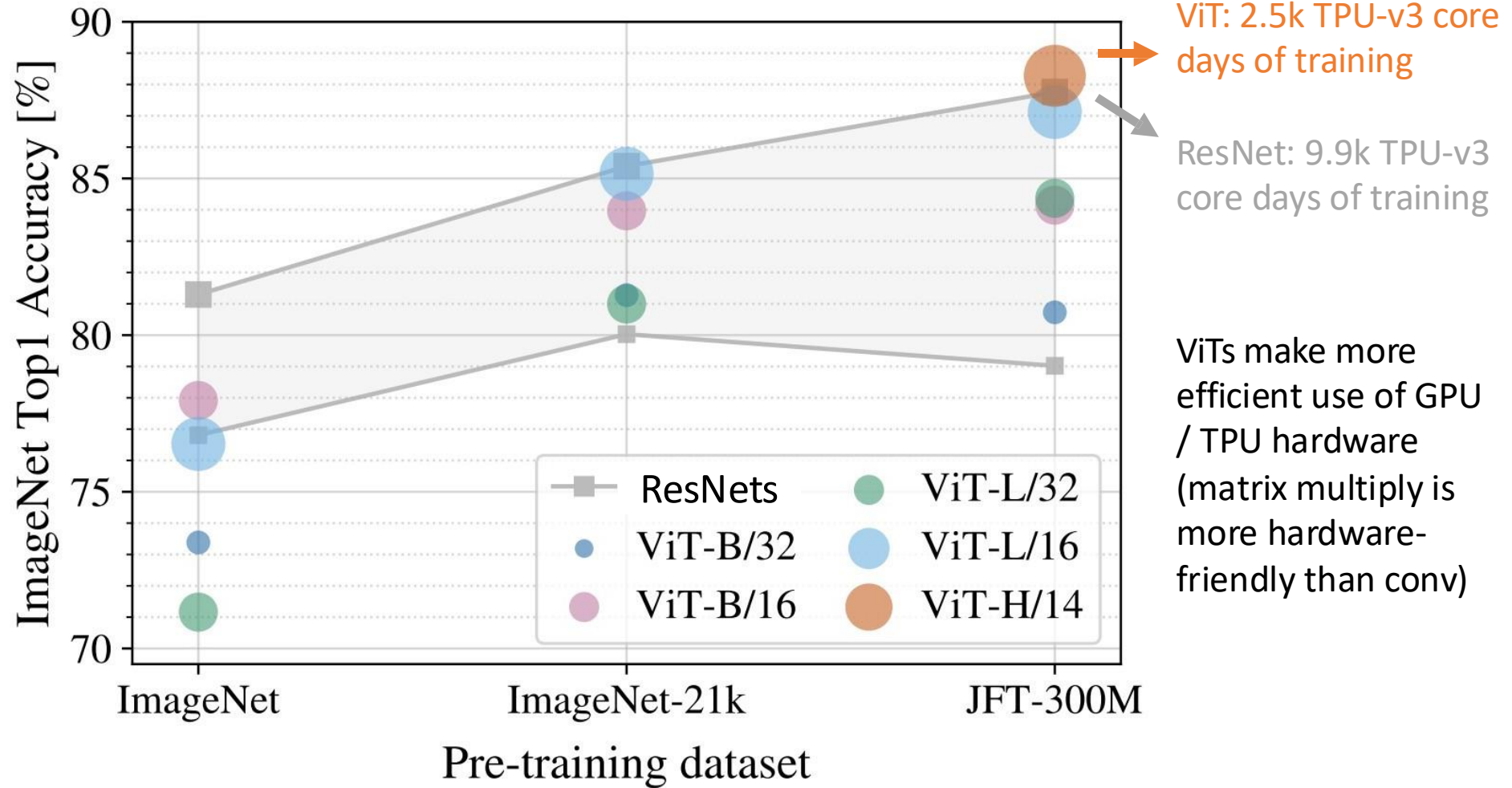
ResNet: 9.9k TPU-v3 core days of training

ViTs make more efficient use of GPU / TPU hardware (matrix multiply is more hardware-friendly than conv)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

# Vision Transformer (ViT) vs ResNets

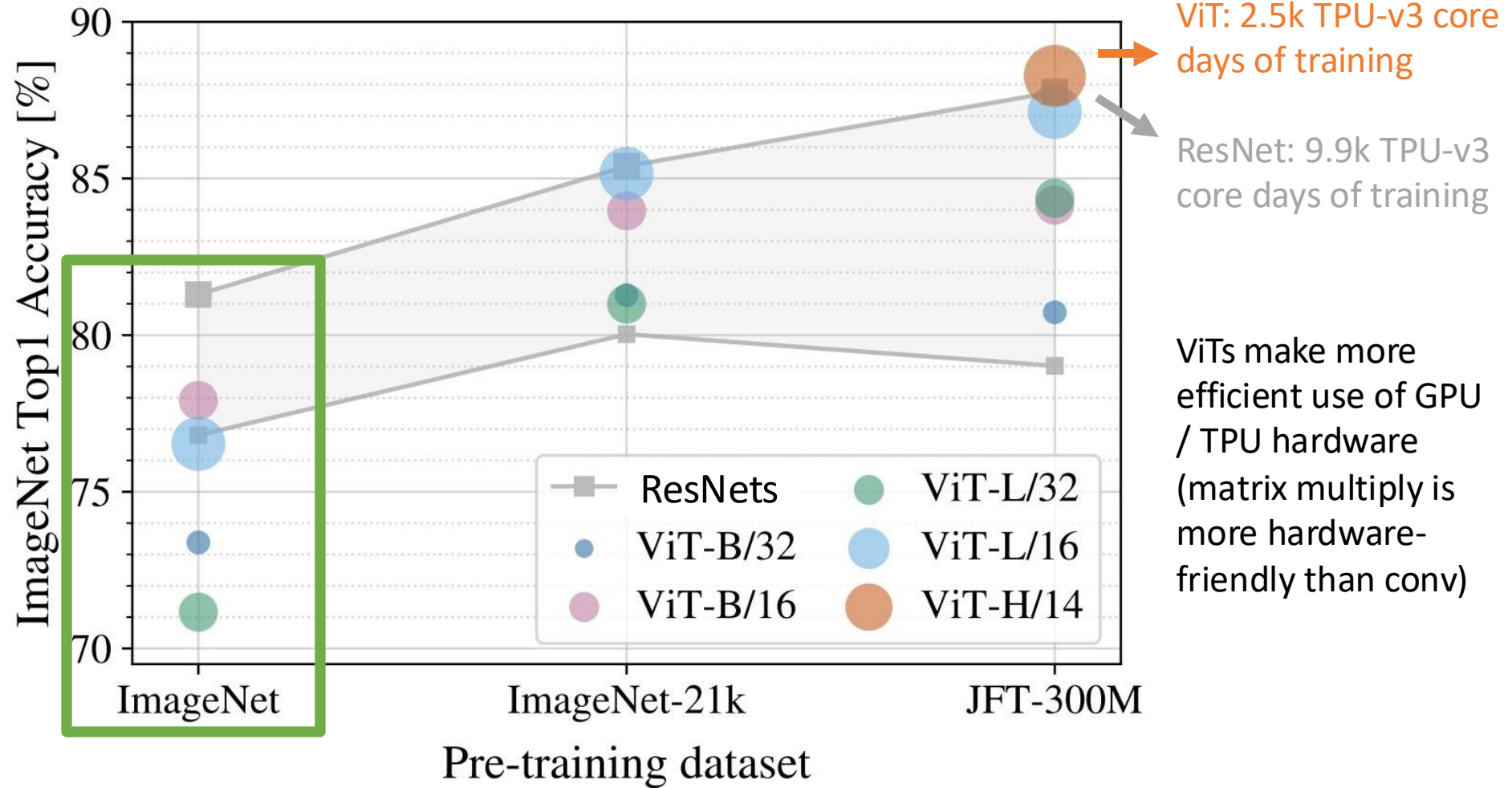
Claim: ViT models have “less inductive bias” than ResNets, so need more pretraining data to learn good features



Dosovitskiy et al, “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”, ICLR 2021

# Vision Transformer (ViT) vs ResNets

How can we improve the performance of ViT models on ImageNet?



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021



# Improving ViT: Augmentation and Regularization

Regularization for ViT models:

- Weight Decay
- Stochastic Depth
- Dropout (in FFN layers of Transformer)

Data Augmentation for ViT models:

- MixUp
- RandAugment

Stochastic Depth

**Layer Dropping:** During training, stochastic depth randomly drops certain layers in the neural network. Each mini-batch might use a different network depth due to this random dropping.

**Survival Probability:** Each layer has a predetermined probability, with which the layer will be active during a particular forward pass in training.

# Mix up

**Select Two Examples:** Randomly select two examples from the training data. Let's denote them as  $(x_i, y_i)$  and  $(x_j, y_j)$ , where  $x$  represents the input features (e.g., images) and  $y$  represents the labels.

**Mix Inputs and Labels:** Form a new example  $(x', y')$  by performing a weighted average of the two selected examples:

$$x' = \lambda x_i + (1 - \lambda) x_j$$

$$y' = \lambda y_i + (1 - \lambda) y_j$$



# RandAugment

**Predefined Transformations:** RandAugment uses a fixed set of image transformations, which might include operations like rotation, shearing, translation, color adjustment, contrast enhancement, and more.

## **Random Selection and Application:**

- For each image in the training dataset, a fixed number  $N$  of transformations are randomly selected from the pool.
- Each selected transformation is applied with a certain intensity or magnitude  $M$ .
- Both  $N$  and  $M$  are hyperparameters

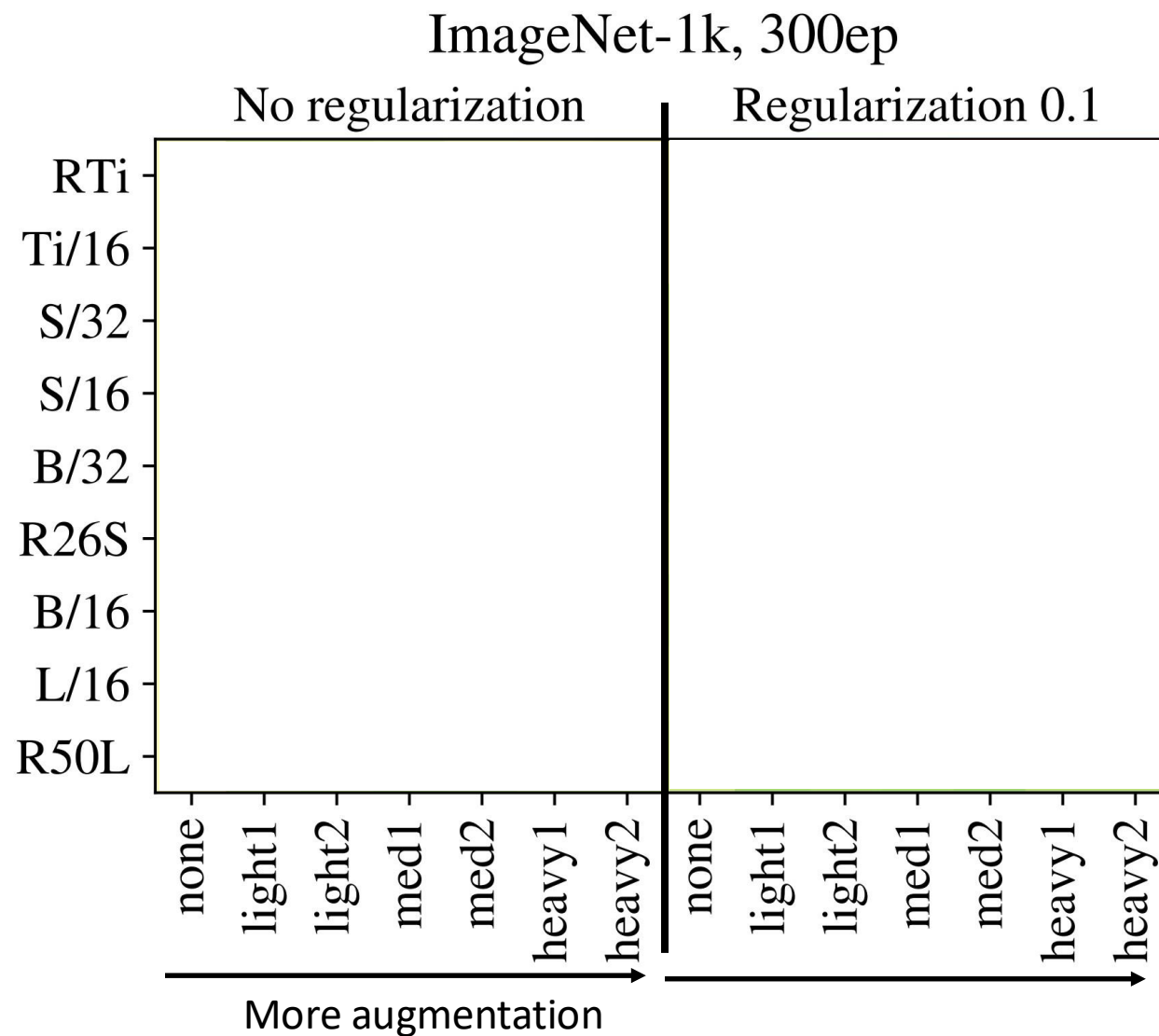
# Improving ViT: Augmentation and Regularization

Regularization for ViT models:

- Weight Decay
- Stochastic Depth
- Dropout (in FFN layers of Transformer)

Data Augmentation for ViT models:

- MixUp
- RandAugment



# Improving ViT: Augmentation and Regularization

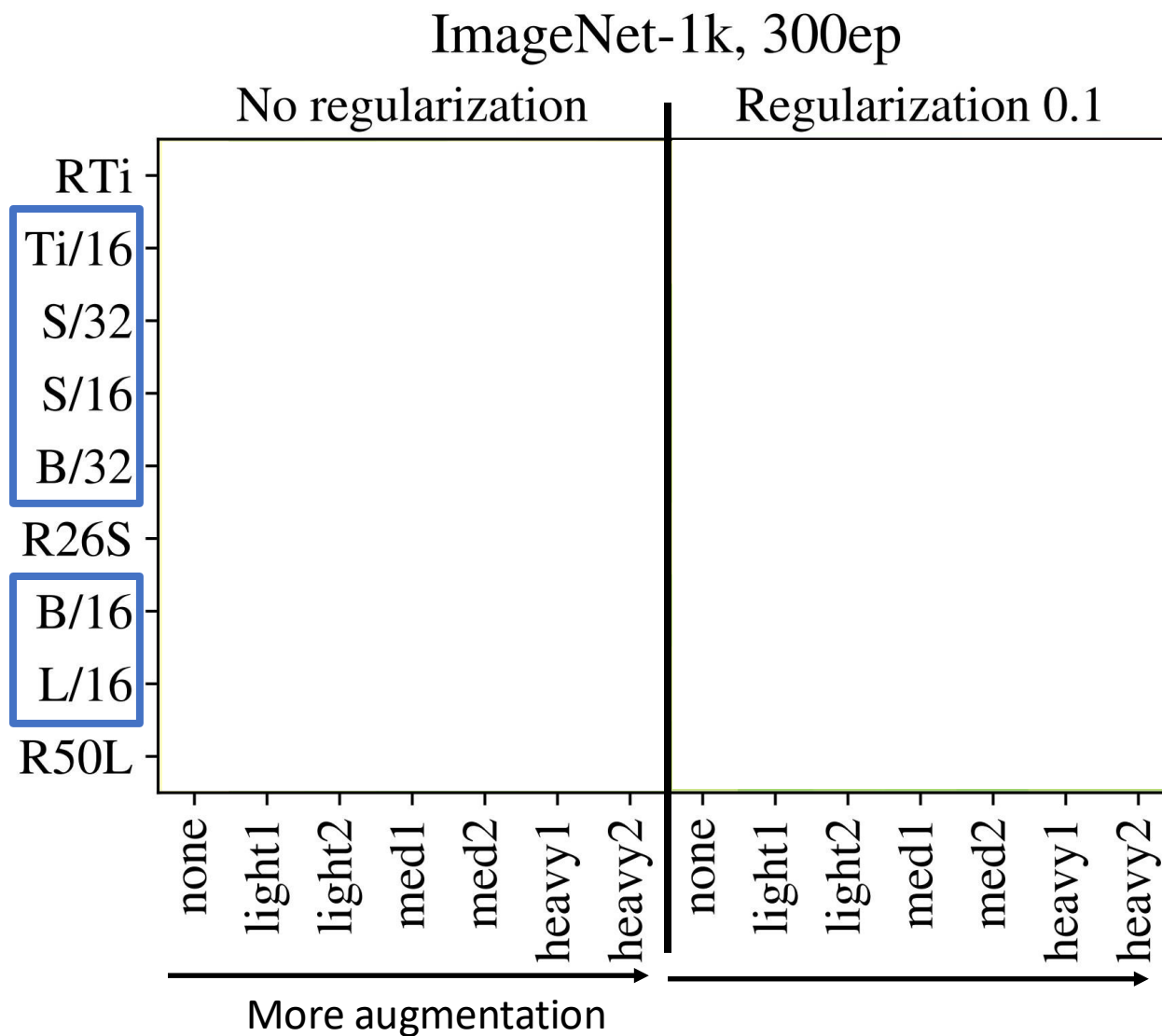
Regularization for ViT models:

- Weight Decay
- Stochastic Depth
- Dropout (in FFN layers of Transformer)

Data Augmentation for ViT models:

- MixUp
- RandAugment

ViT models:  
Ti = Tiny  
S = Small  
B = Base  
L = Large



# Improving ViT: Augmentation and Regularization

Regularization for ViT models:

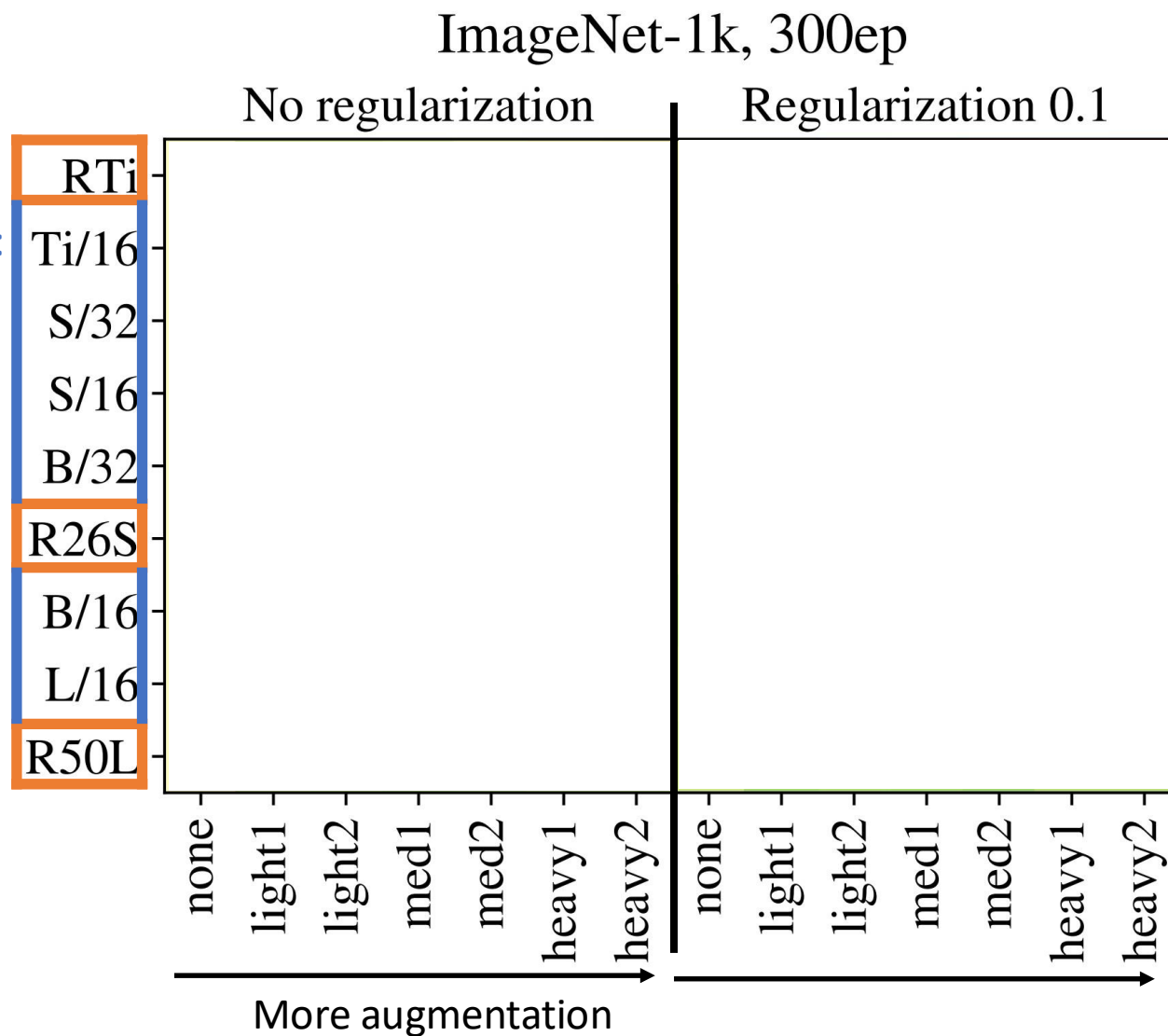
- Weight Decay
- Stochastic Depth
- Dropout (in FFN layers of Transformer)

Data Augmentation for ViT models:

- MixUp
- RandAugment

Hybrid models:  
ResNet blocks,  
then ViT blocks

ViT models:  
Ti = Tiny  
S = Small  
B = Base  
L = Large



# Improving ViT: Augmentation and Regularization

Regularization for ViT models:

- Weight Decay
- Stochastic Depth
- Dropout (in FFN layers of Transformer)

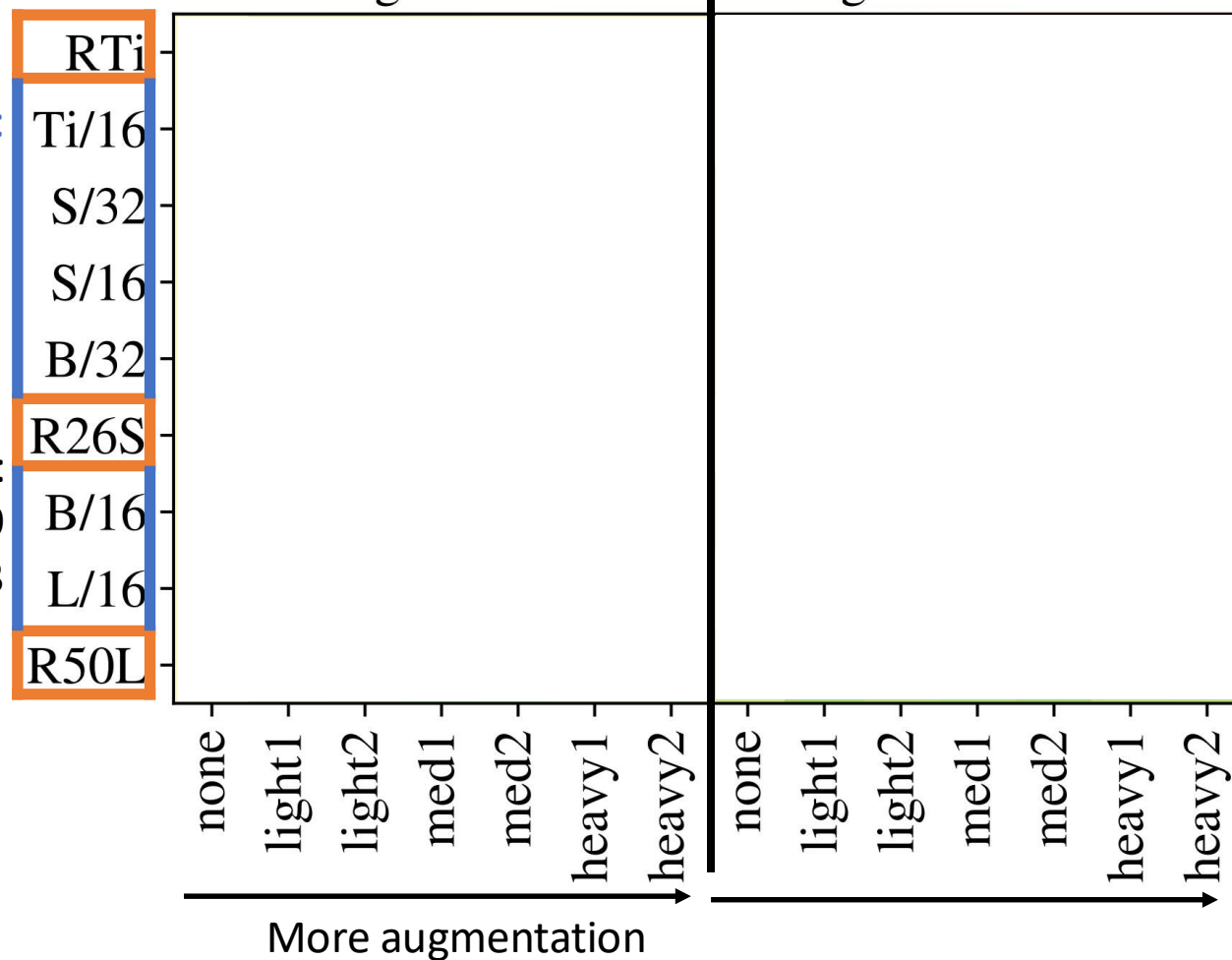
Data Augmentation for ViT models:

- MixUp
- RandAugment

Hybrid models:  
ResNet blocks,  
then ViT blocks

ViT models:  
Ti = Tiny  
S = Small  
B = Base  
L = Large

Original Paper:  
77.9  
76.53



# Improving ViT: Augmentation and Regularization

Regularization for ViT models:

- Weight Decay
- Stochastic Depth
- Dropout (in FFN layers of Transformer)

Data Augmentation for ViT models:

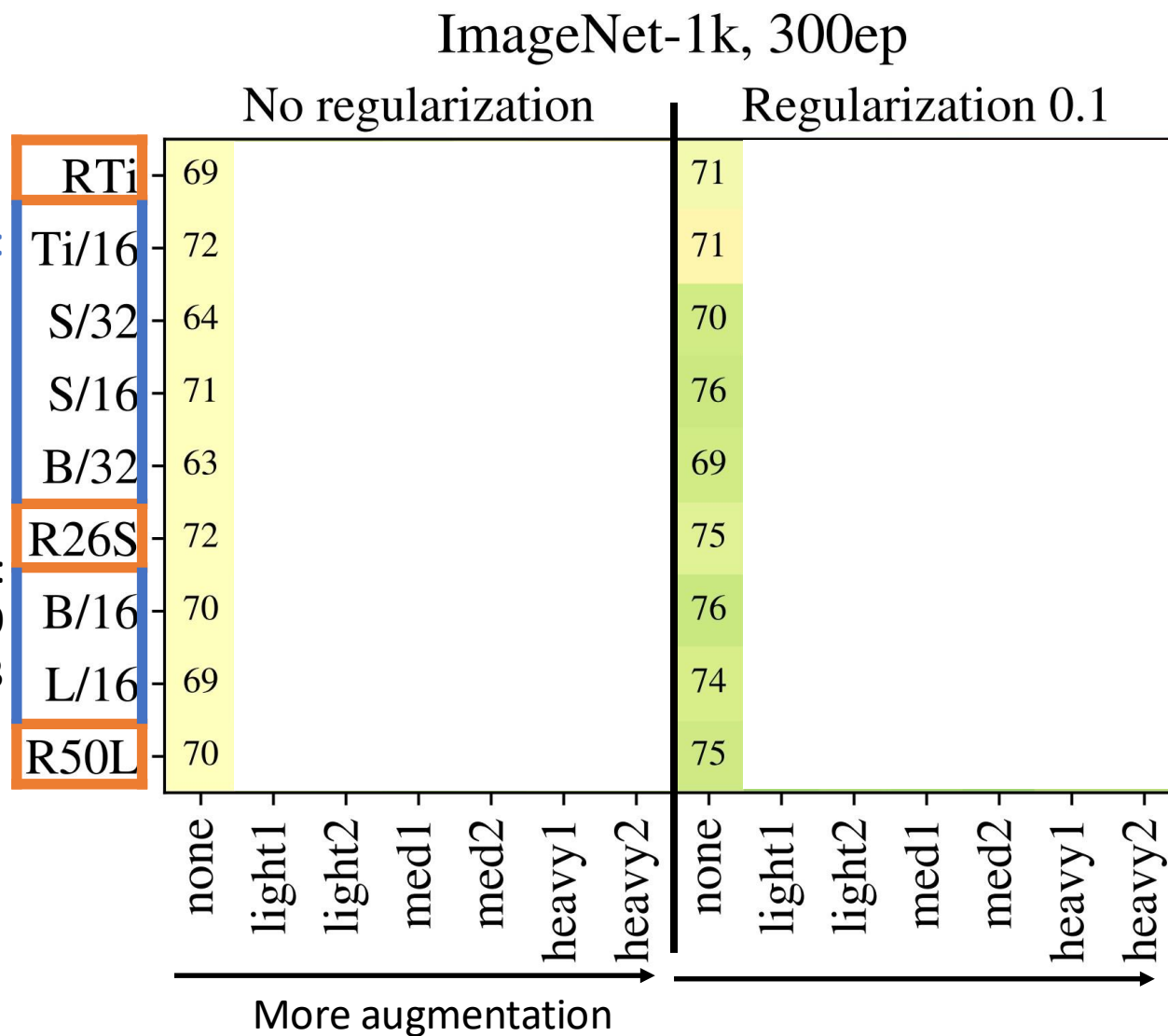
- MixUp
- RandAugment

Adding regularization is (almost) always helpful

Hybrid models:  
ResNet blocks,  
then ViT blocks

ViT models:  
Ti = Tiny  
S = Small  
B = Base  
L = Large

Original Paper:  
77.9  
76.53



# Improving ViT: Augmentation and Regularization

Regularization for ViT models:

- Weight Decay
- Stochastic Depth
- Dropout (in FFN layers of Transformer)

Data Augmentation for ViT models:

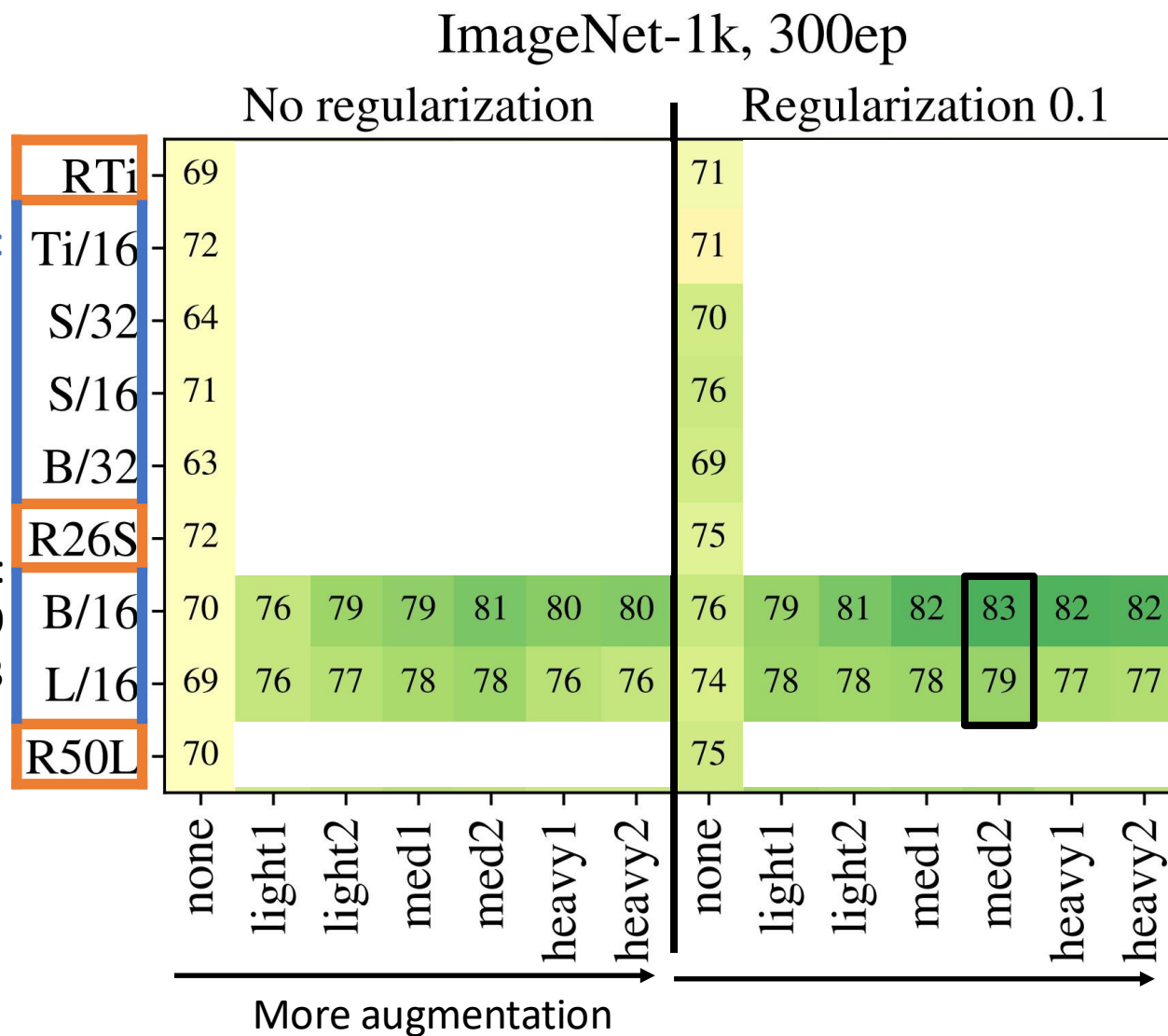
- MixUp
- RandAugment

Regularization +  
Augmentation gives  
big improvements  
over original results

Hybrid models:  
ResNet blocks,  
then ViT blocks

ViT models:  
Ti = Tiny  
S = Small  
B = Base  
L = Large

Original Paper:  
77.9  
76.53





# Improving ViT: Augmentation and Regularization

Regularization for ViT models:

- Weight Decay
- Stochastic Depth
- Dropout (in FFN layers of Transformer)

Data Augmentation for ViT models:

- MixUp
- RandAugment

Hybrid models:  
ResNet blocks,  
then ViT blocks

ViT models:  
Ti = Tiny  
S = Small  
B = Base  
L = Large

Original Paper:

Lots of other  
patterns in  
full results

ImageNet-1k, 300ep

	No regularization							Regularization 0.1						
	none	light1	light2	med1	med2	heavy1	heavy2	none	light1	light2	med1	med2	heavy1	heavy2
RTi	69	73	73	72	70	69	68	71	70	67	65	63	62	61
Ti/16	72	76	75	75	74	72	71	71	72	68	65	63	63	62
S/32	64	71	76	76	76	74	74	70	72	72	71	71	69	68
S/16	71	77	79	81	82	80	80	76	79	80	79	79	77	77
B/32	63	70	73	75	76	75	76	69	74	77	77	78	77	77
R26S	72	76	78	79	80	80	80	75	78	81	82	82	81	81
B/16	70	76	79	79	81	80	80	76	79	81	82	83	82	82
L/16	69	76	77	78	78	76	76	74	78	78	78	79	77	77
R50L	70	75	76	77	77	76	76	75	78	78	78	79	77	77

More augmentation



- Why right top is poor performer?

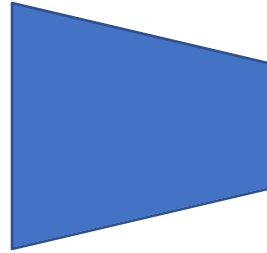
Given an image.

- Can you obtain its caption?
- Can you obtain its embedding?

Can you use text embedding and make image look like same as input but the manipulated image should give a different/irrelevant/harmful caption?  
Can you do vice-versa?

# Improving ViT: Distillation

Step 1: Train a **teacher model** on images and ground-truth labels



$P(\text{cat}) = 0.9$   
 $P(\text{dog}) = 0.1$



Cross  
Entropy  
Loss



GT label:  
Cat

Assumption:

$P(x_{\text{teacher}}) \neq P(x_{\text{student}})$ , for instance teacher sees day images, student is asked for night images

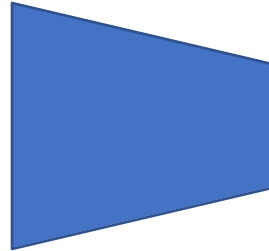
$P(y)$  is same for both teacher and student, that is, labels fully overlap

What we get      What comes in  
taught in school      the exam



# Improving ViT: Distillation

Step 1: Train a **teacher model** on images and ground-truth labels

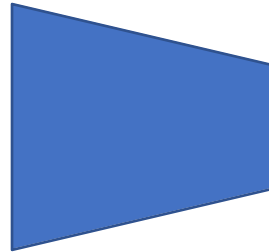
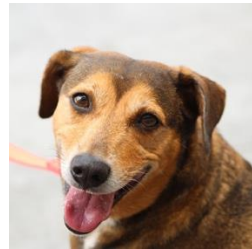


$P(\text{cat}) = 0.9$   
 $P(\text{dog}) = 0.1$

Cross  
Entropy  
Loss

GT label:  
Cat

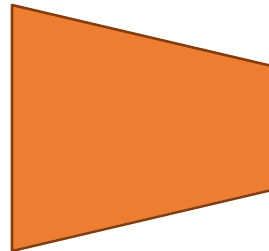
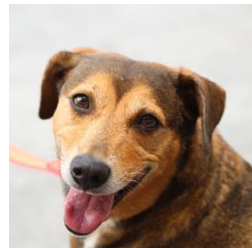
Step 2: Train a **student model** to match predictions from the **teacher**



$P(\text{cat}) = 0.1$   
 $P(\text{dog}) = 0.9$

$\lambda \tau^2 \text{KL}(\psi(Z_s/\tau), \psi(Z_t/\tau)).$

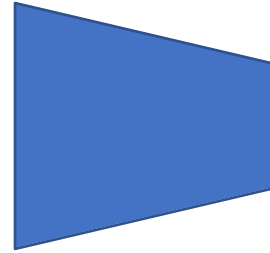
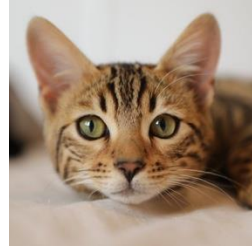
KL Divergence Loss



$P(\text{cat}) = 0.2$   
 $P(\text{dog}) = 0.8$

# Improving ViT: Distillation

Step 1: Train a **teacher model** on images and ground-truth labels

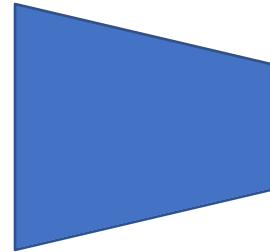
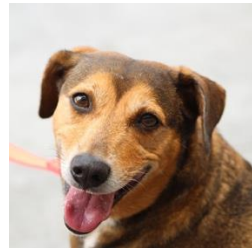


$P(\text{cat}) = 0.9$   
 $P(\text{dog}) = 0.1$

→ Cross Entropy Loss ←

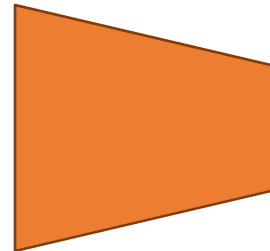
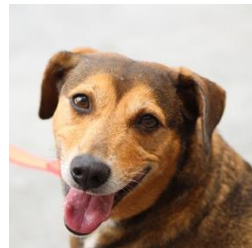
GT label: Cat

Step 2: Train a **student model** to match predictions from the **teacher** (sometimes also to match GT labels)



$P(\text{cat}) = 0.1$   
 $P(\text{dog}) = 0.9$

→ KL Divergence Loss



$P(\text{cat}) = 0.2$   
 $P(\text{dog}) = 0.8$

→ Cross Entropy Loss ←

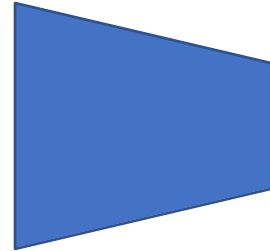
GT label: Dog

$$(1 - \lambda) \text{LCE}(\psi(Z_s), y)$$

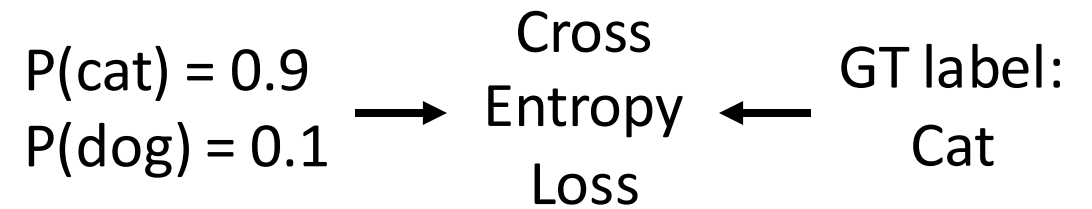
$$\lambda \tau^2 \text{KL}(\psi(Z_s/\tau), \psi(Z_t/\tau)).$$

# Improving ViT: Distillation

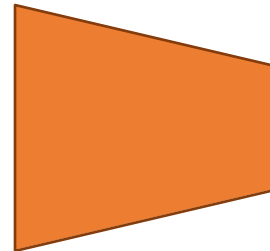
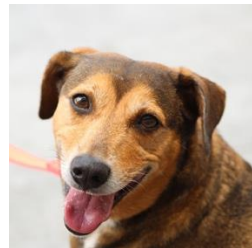
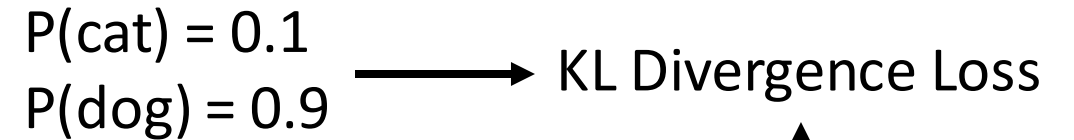
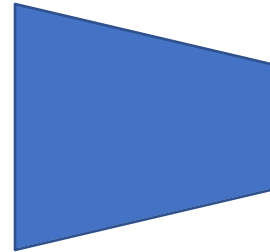
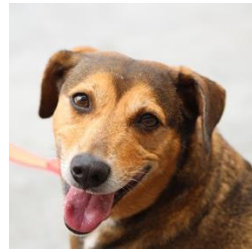
Step 1: Train a **teacher model** on images and ground-truth labels



Often works better than training student from scratch (especially if teacher is bigger than student)



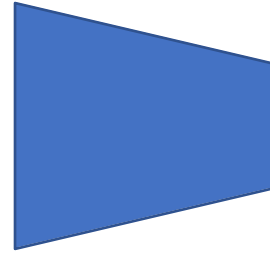
Step 2: Train a **student model** to match predictions from the **teacher** (sometimes also to match GT labels)



# Improving ViT: Distillation

Can also train student on **unlabeled** data! (Semi-supervised learning)

Step 1: Train a **teacher model** on images and ground-truth labels

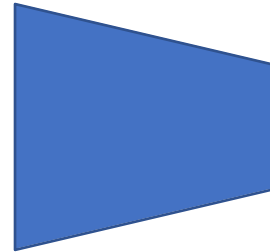
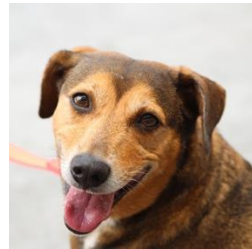


$$\begin{aligned} P(\text{cat}) &= 0.9 \\ P(\text{dog}) &= 0.1 \end{aligned}$$

Cross  
Entropy  
Loss

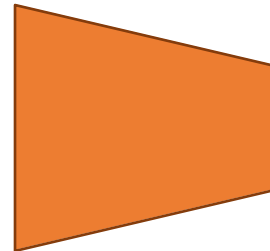
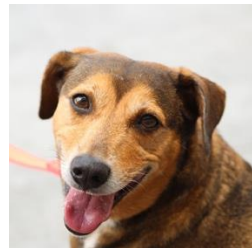
GT label:  
Cat

Step 2: Train a **student model** to match predictions from the **teacher** (sometimes also to match GT labels)



$$\begin{aligned} P(\text{cat}) &= 0.1 \\ P(\text{dog}) &= 0.9 \end{aligned}$$

KL Divergence Loss



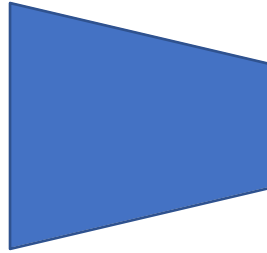
$$\begin{aligned} P(\text{cat}) &= 0.2 \\ P(\text{dog}) &= 0.8 \end{aligned}$$

Cross  
Entropy  
Loss

GT label:  
Dog

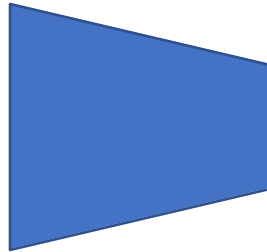
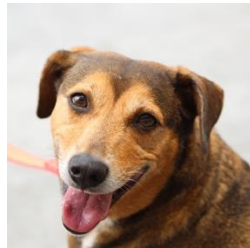
# Improving ViT: Distillation

Step 1: Train a teacher CNN on ImageNet

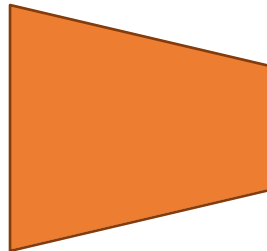
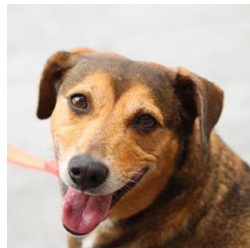


$P(\text{cat}) = 0.9$   
 $P(\text{dog}) = 0.1$  → Cross Entropy Loss ← GT label: Cat

Step 2: Train a student ViT to match ImageNet predictions from the teacher CNN (and match GT labels)



$P(\text{cat}) = 0.1$   
 $P(\text{dog}) = 0.9$  → KL Divergence Loss



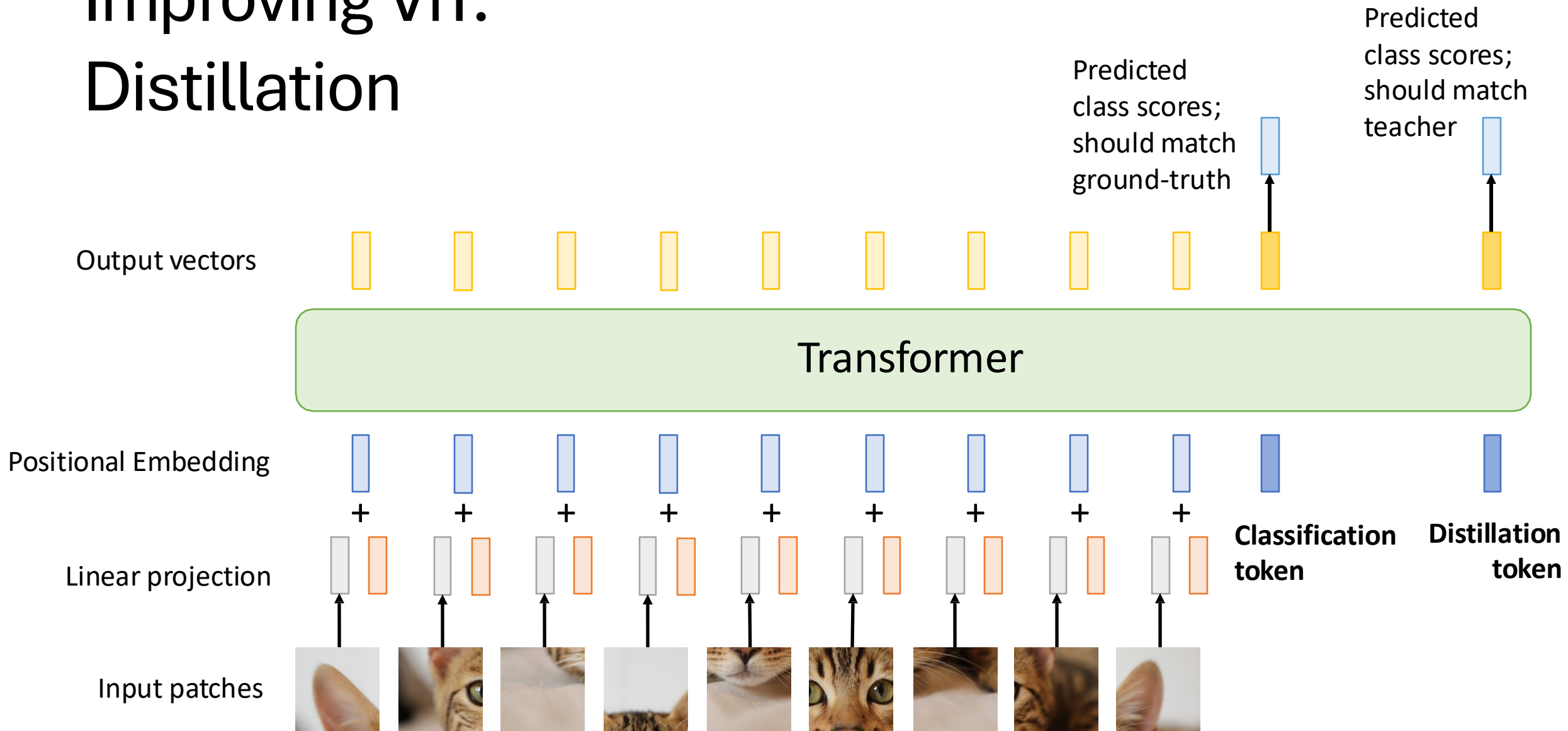
$P(\text{cat}) = 0.2$   
 $P(\text{dog}) = 0.8$  → Cross Entropy Loss ← GT label: Dog



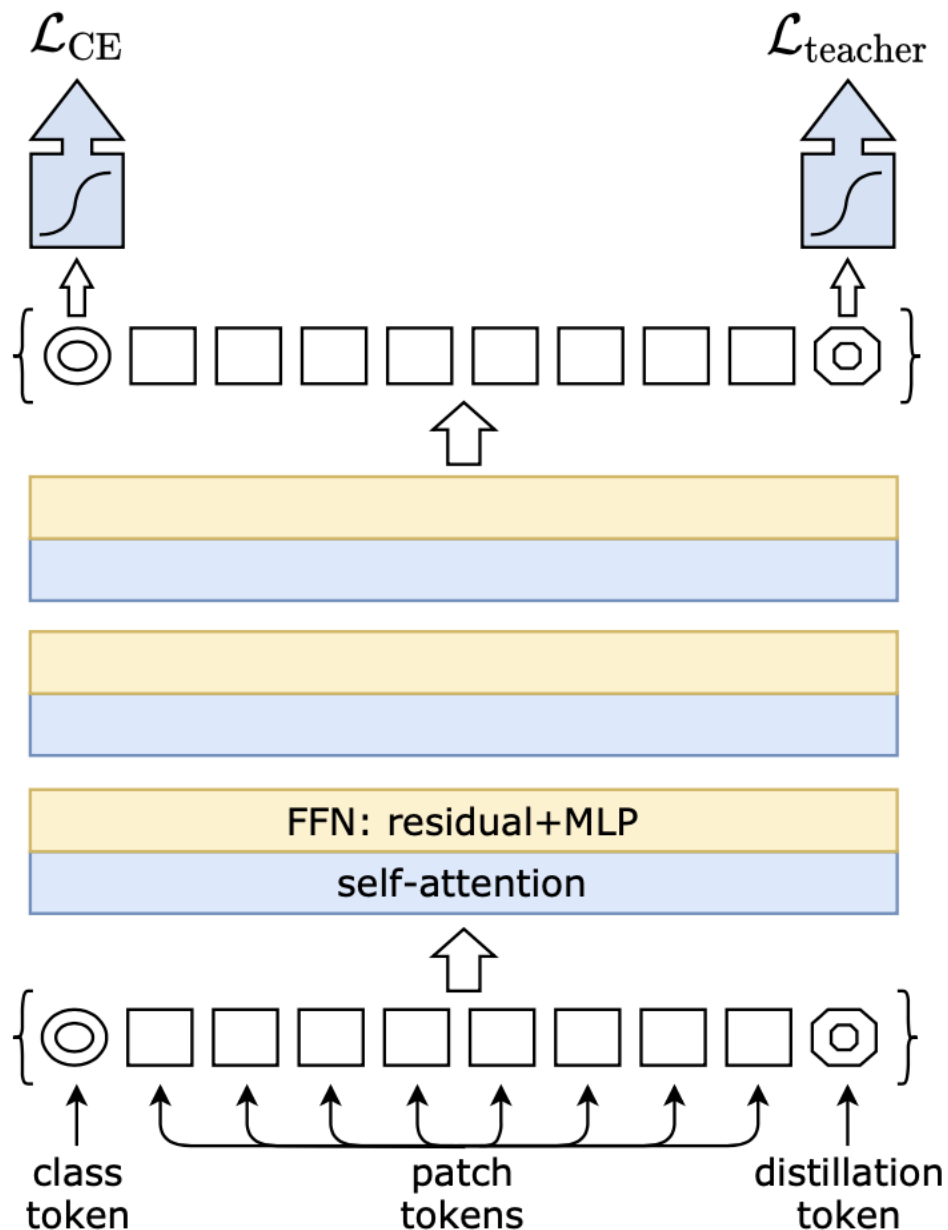
**Transfer of Knowledge:** In model distillation, the goal is to transfer knowledge from a larger, more complex teacher model to a smaller student model. By using a softer distribution, the teacher model conveys not just information about the correct class but also about the relationships between incorrect classes. This nuanced information helps the student model learn more about the underlying data structure.

**Avoiding Overconfidence:** Without temperature scaling, the teacher model's predictions might be overly confident, focusing only on the most likely class. By softening the distribution, the teacher model provides more balanced information, allowing the student model to learn from the teacher's uncertainty and make more robust predictions.

# Improving ViT: Distillation



Touvron et al, "Training data-efficient image transformers & distillation through attention", ICML 2021



$$\mathcal{L}_{global} = 0.5\mathcal{L}_{CE}(\psi(Z_s), y) + 0.5\mathcal{L}_{CE}(\psi(Z_s), y_t)$$

$$y_t = \operatorname{argmax}_c Z_t(c)$$

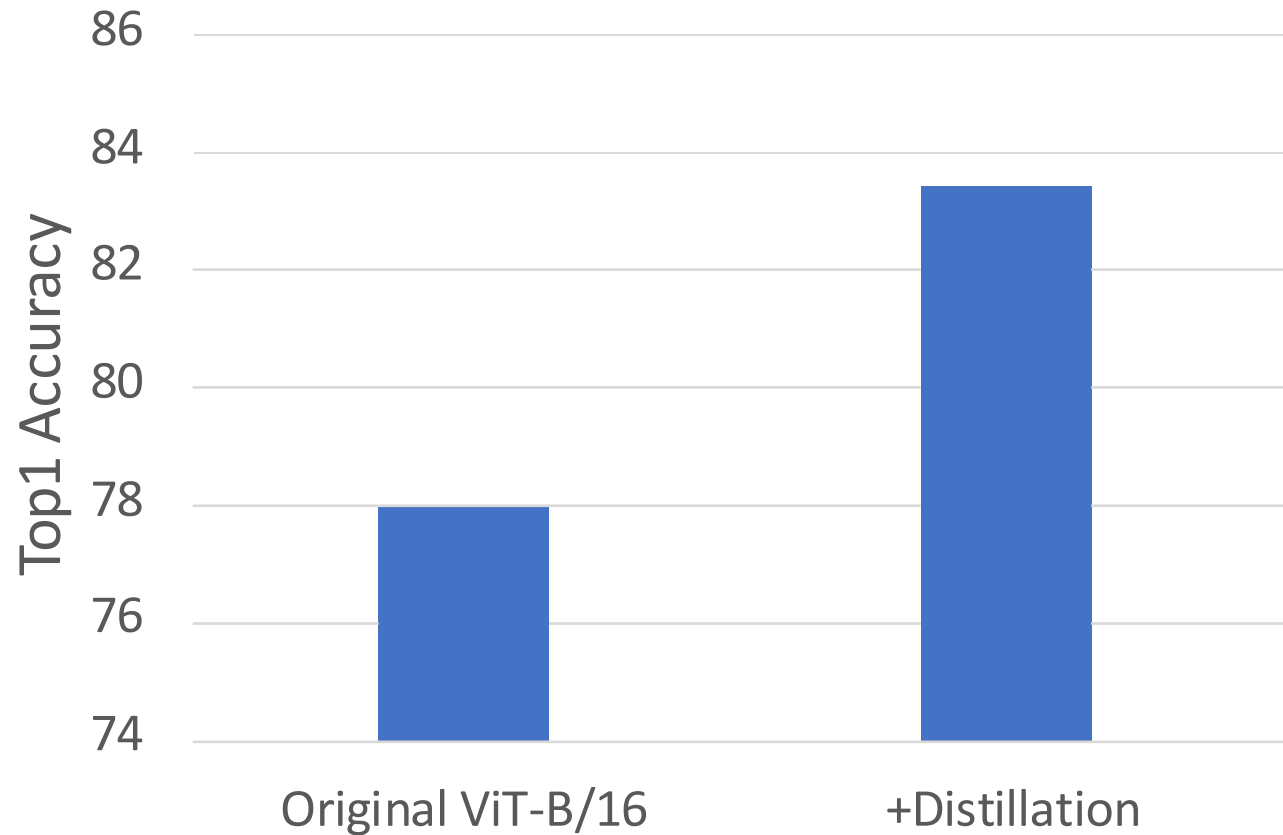
~~Seems fishy, directly taken from the paper~~  
 ~~$Z_t$  should be softmax probs.~~

- augmentation like mix-up can change the label, two different images are mixed, so are there labels.
- applied label smoothing, prediction gets prob of  $1-\epsilon$ ,  $\epsilon=0.1$ , applied uniformly to other classes

For prediction, softmax output is added

# Improving ViT: Distillation

ViT-B/16 on ImageNet

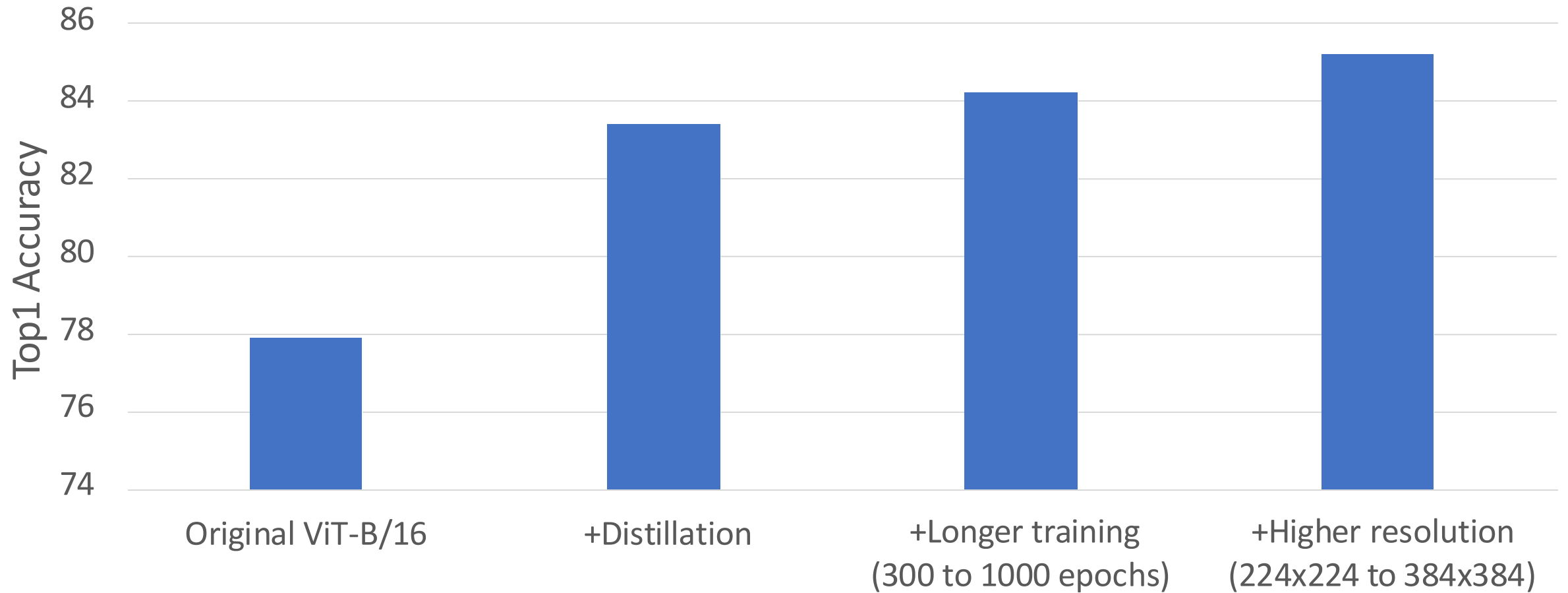


Teacher RegNetY-16GF ([Radosavovic et al., 2020](#)) with 84M parameters

Touvron et al, "Training data-efficient image transformers & distillation through attention", ICML 2021

# Improving ViT: Distillation

ViT-B/16 on ImageNet



Touvron et al, "Training data-efficient image transformers & distillation through attention", ICML 2021

- a. Suppose you want to distill from resnet trained on imagenet to cifar-10. The Setting is few-shot with 10-way-5-shot classification.
- b. Further to it, you need to deploy it on edge device.

Note, train data is less so you need to find a niche way to fine-tune resnet.

As edge will expect a small model and you have less train data so you cannot Train it in classical way.

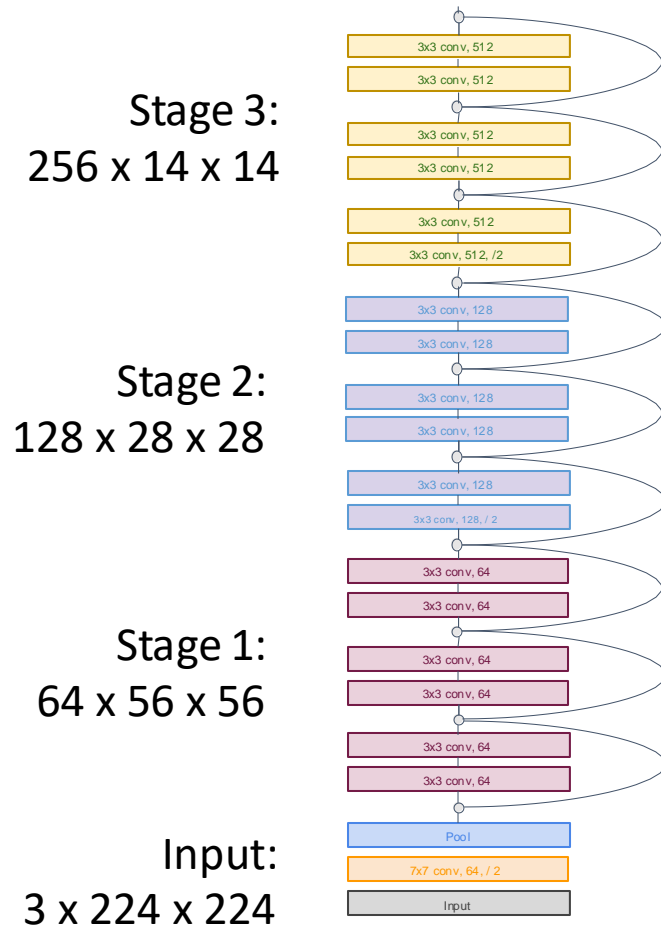
How will you address the two problems?

Can you find some difference between ViT and CNN?

# ViT vs CNN

In most CNNs (including ResNets), **decrease** resolution and **increase** channels as you go deeper in the network (Hierarchical architecture)

Useful since objects in images can occur at various scales



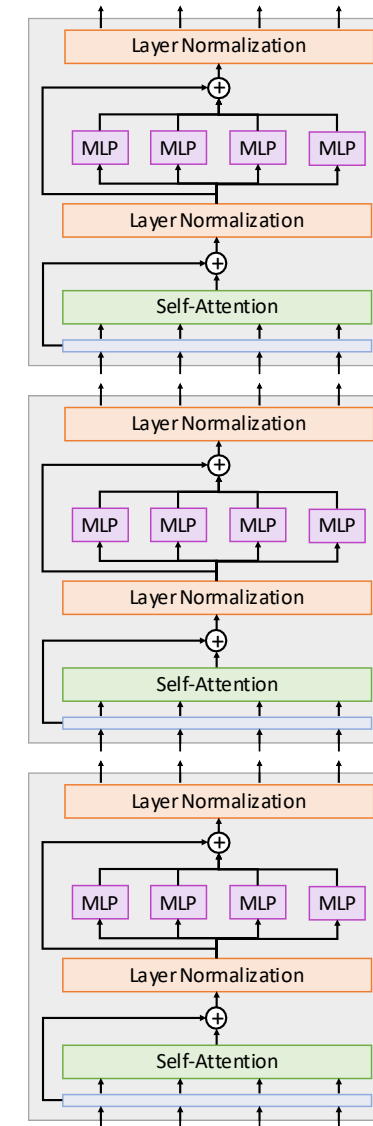
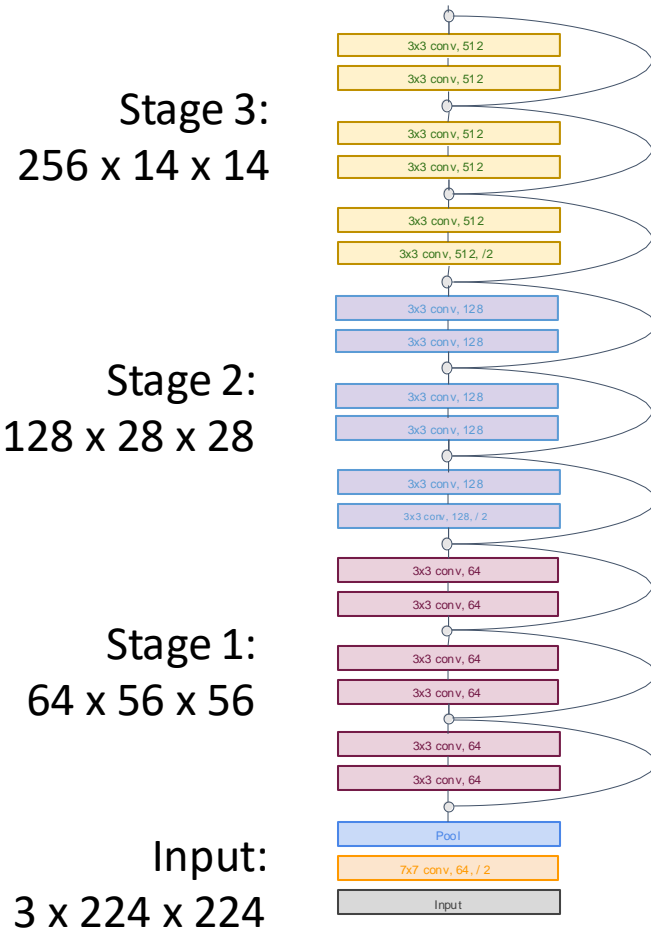


# ViT vs CNN

In most CNNs (including ResNets), **decrease** resolution and **increase** channels as you go deeper in the network (Hierarchical architecture)

Useful since objects in images can occur at various scales

In a ViT, all blocks have same resolution and number of channels (Isotropic architecture)



# ViT vs CNN

In most CNNs (including ResNets), **decrease** resolution and **increase** channels as you go deeper in the network (Hierarchical architecture)

Useful since objects in images can occur at various scales

In a ViT, all blocks have same resolution and number of channels (Isotropic architecture)

Can we build a **hierarchical** ViT model?

