

Vision Transformer (ViT)

기초심화CV팀

김수란

Table of contents

01

**Background
of ViT**

02

**ViT Architecture
and Workflow**

03




**ViT Performance
Analysis**

01

Background of ViT

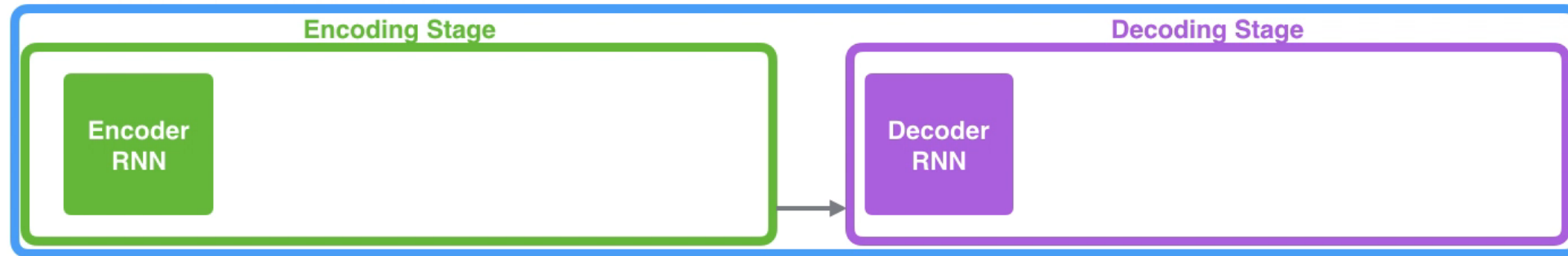
Sequential Data

Sequential data is common in various fields.

Speech recognition		→	"The quick brown fox jumped over the lazy dog."
Music generation	∅	→	
Sentiment classification	"There is nothing to like in this movie."	→	★☆☆☆☆
DNA sequence analysis	AGCCCCTGTGAGGAACTAG	→	AGCCCCTGTGAGGAACTAG
Machine translation	Voulez-vous chanter avec moi?	→	Do you want to sing with me?
Video activity recognition		→	Running
Name entity recognition	Yesterday, Harry Potter met Hermione Granger.	→	Yesterday, Harry Potter met Hermione Granger .

Seq2Seq Model with RNN

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL

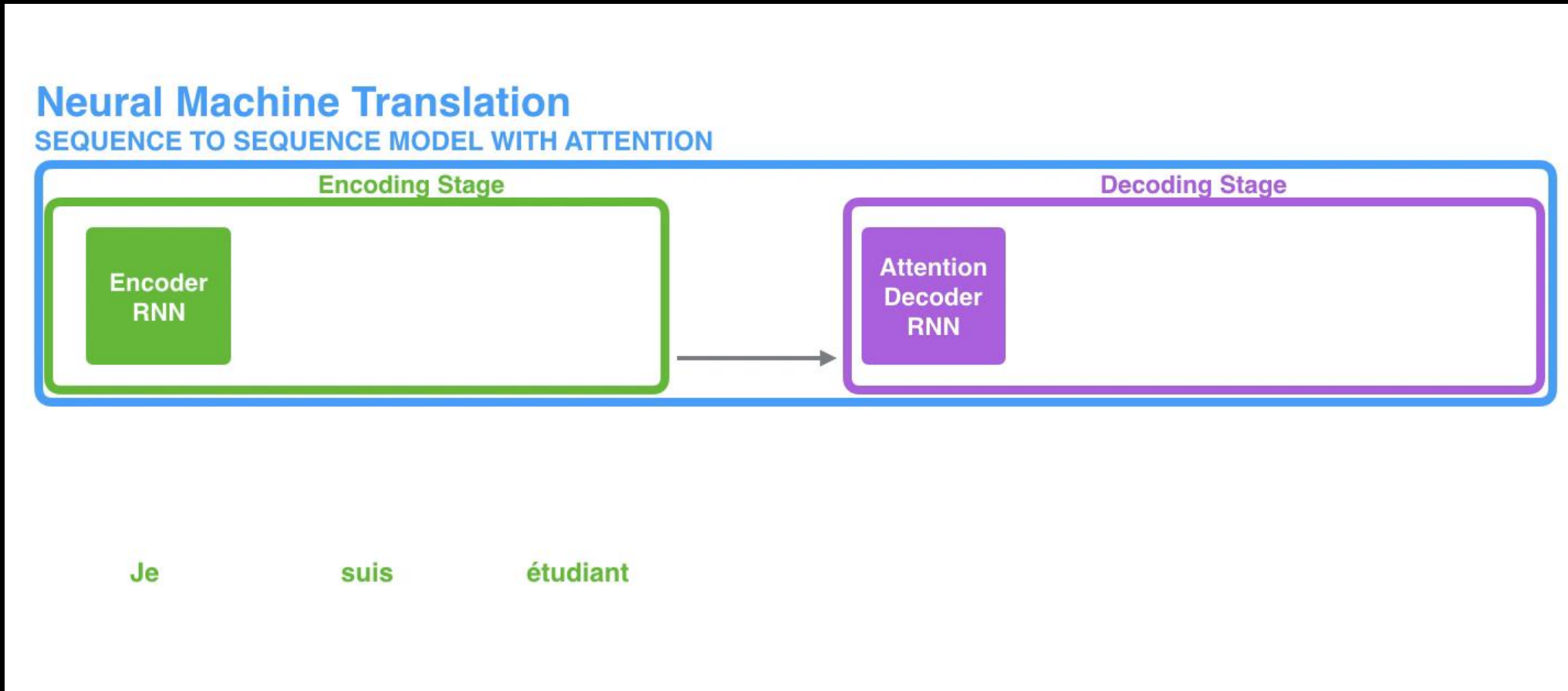


Je

suis

étudiant

Seq2Seq Model with Attention



Seq2Seq Model with Attention



Attention is great

- Attention significantly **improves performance** (in many applications)
 - It's very useful to allow decoder to focus on certain parts of the source
- Attention **solves the bottleneck problem**
 - Attention allows decoder to look directly at source; bypass bottleneck
- Attention **helps with vanishing gradient problem**
 - Provides shortcut to faraway states
- Attention provides **some interpretability**
 - By inspecting attention distribution, we can see what the decoder was focusing on

Transformer

Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez*[†]

University of Toronto

aidan@cs.toronto.edu

Łukasz Kaiser*

Google Brain

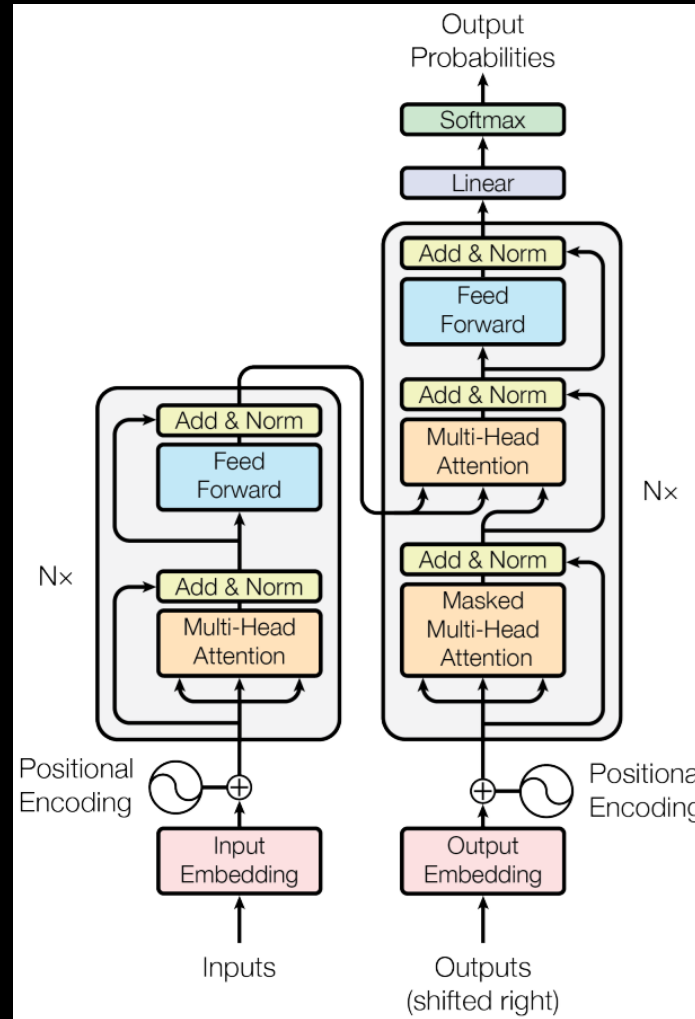
lukaszkaizer@google.com

Illia Polosukhin*[‡]

illia.polosukhin@gmail.com

Transformer

Encoder ->



< - Decoder

Transformer

Transformers have become the model of choice in NLP.

In CV, however, classic CNN architectures are still SOTA.



Chat GPT

Vision Transformer!

Published as a conference paper at ICLR 2021

AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

**Alexey Dosovitskiy^{*,†}, Lucas Beyer^{*}, Alexander Kolesnikov^{*}, Dirk Weissenborn^{*},
Xiaohua Zhai^{*}, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby^{*,†}**

^{*}equal technical contribution, [†]equal advising

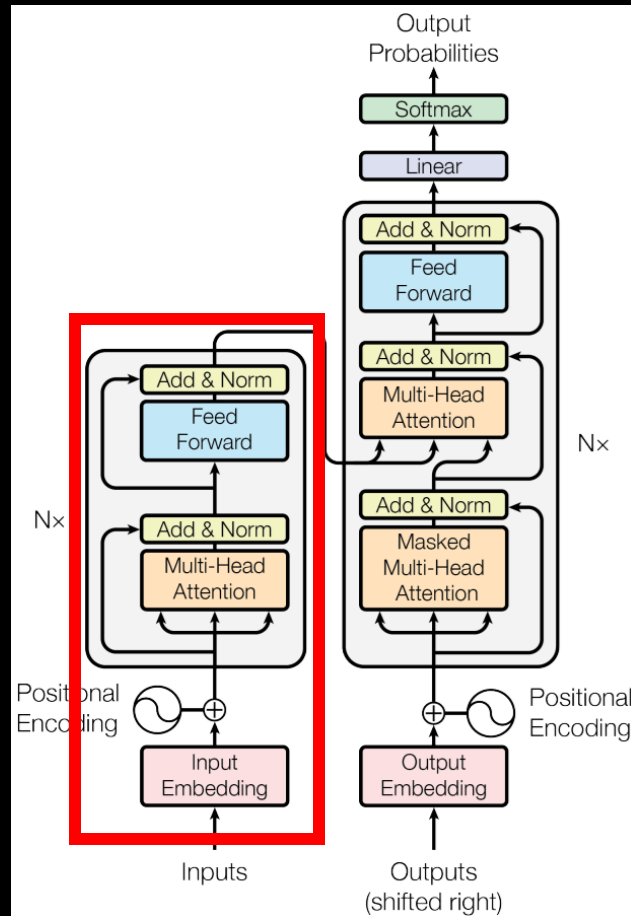
Google Research, Brain Team

{adosovitskiy, neilhoulby}@google.com

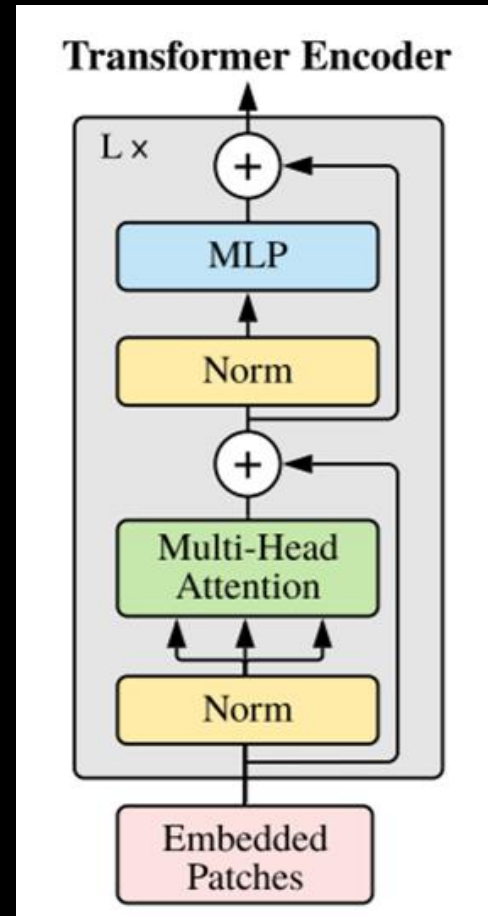
02

ViT Architecture and Workflow

Transformer -> Vision Transformer (ViT)



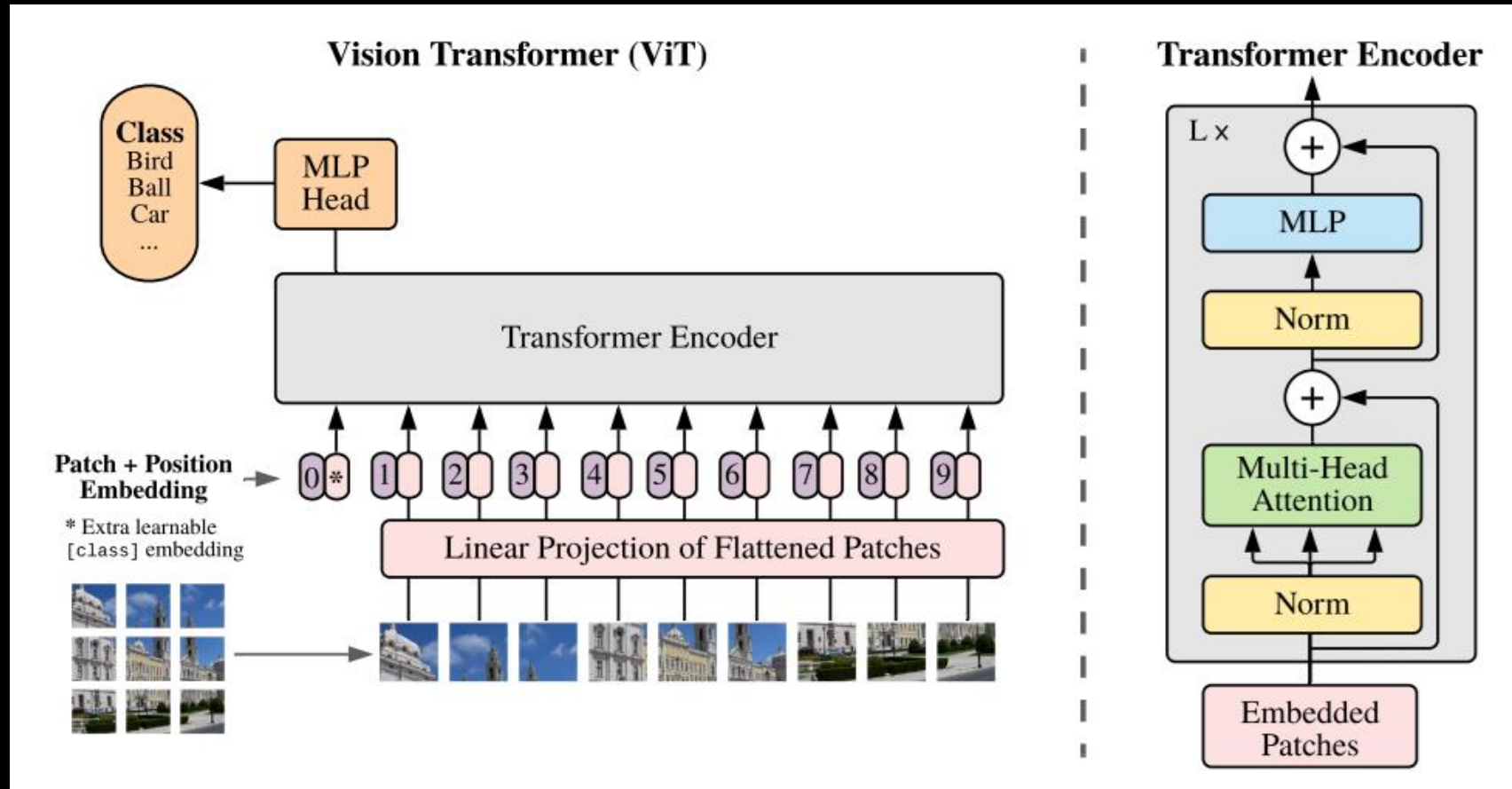
Transformer



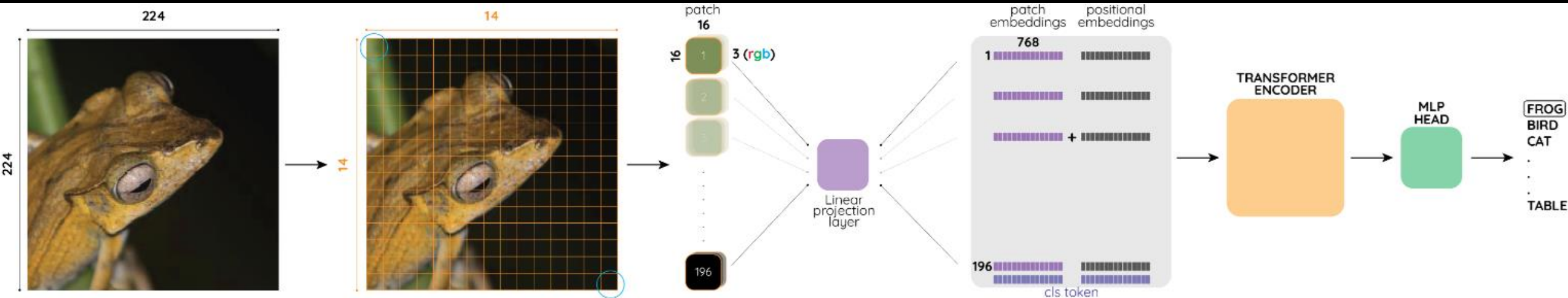
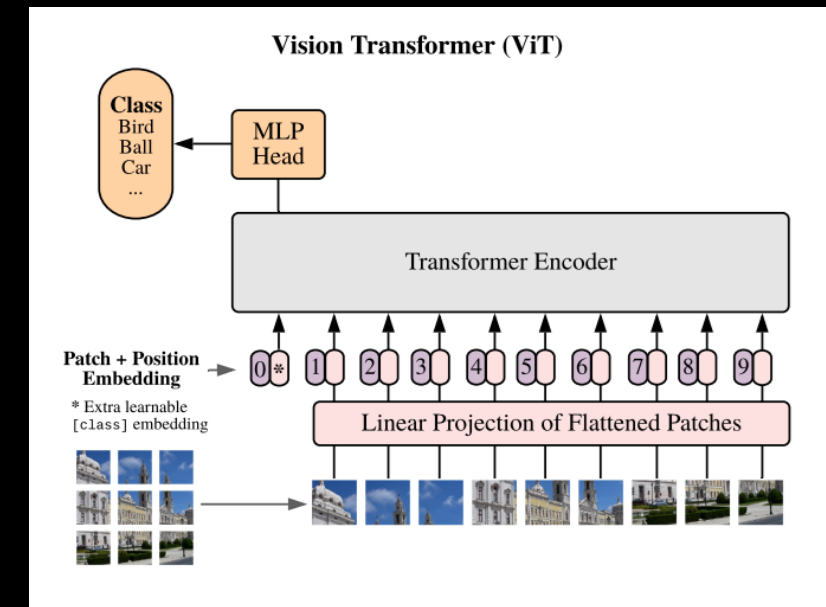
ViT

- Note that small difference compared to the transformer encoder.
- Norm comes before the MHA and MLP.

Vision Transformer (ViT)

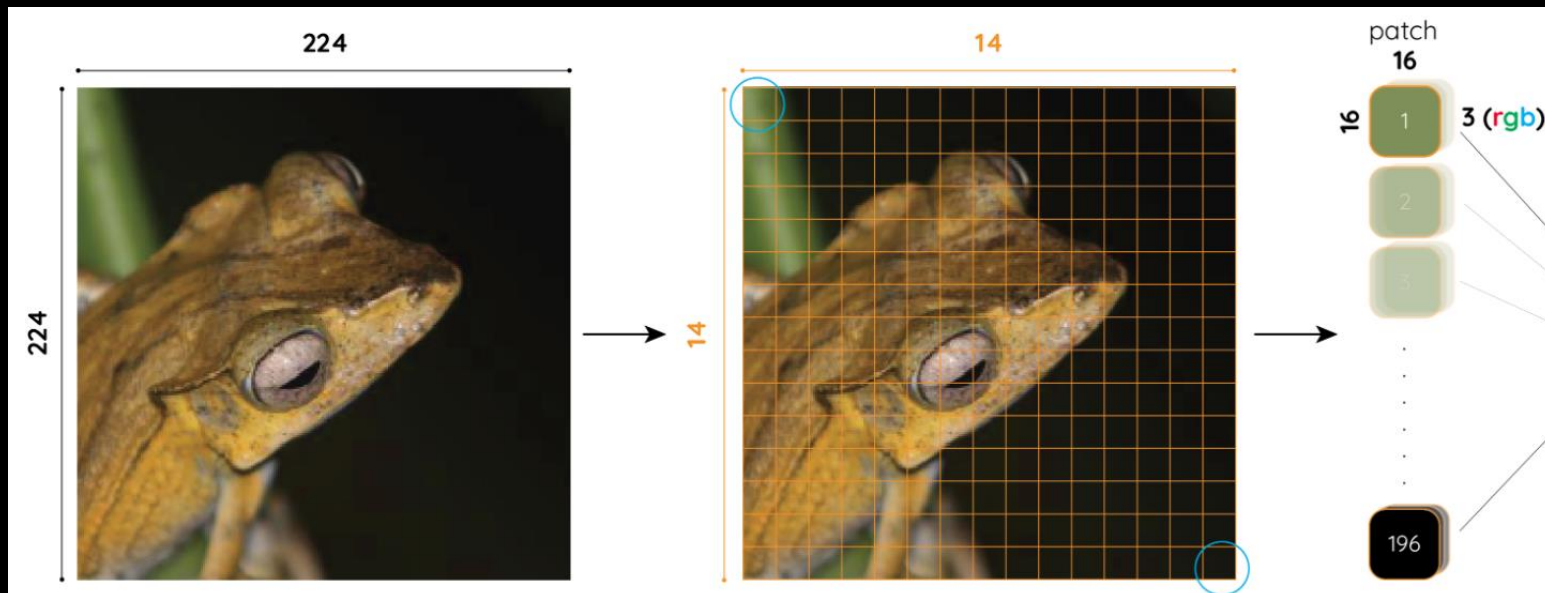
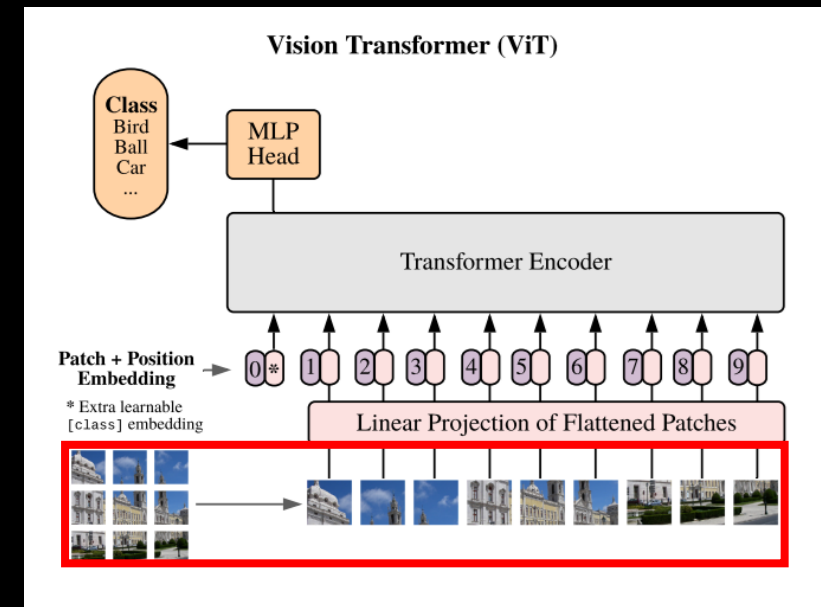


ViT: Step-by-Step Example



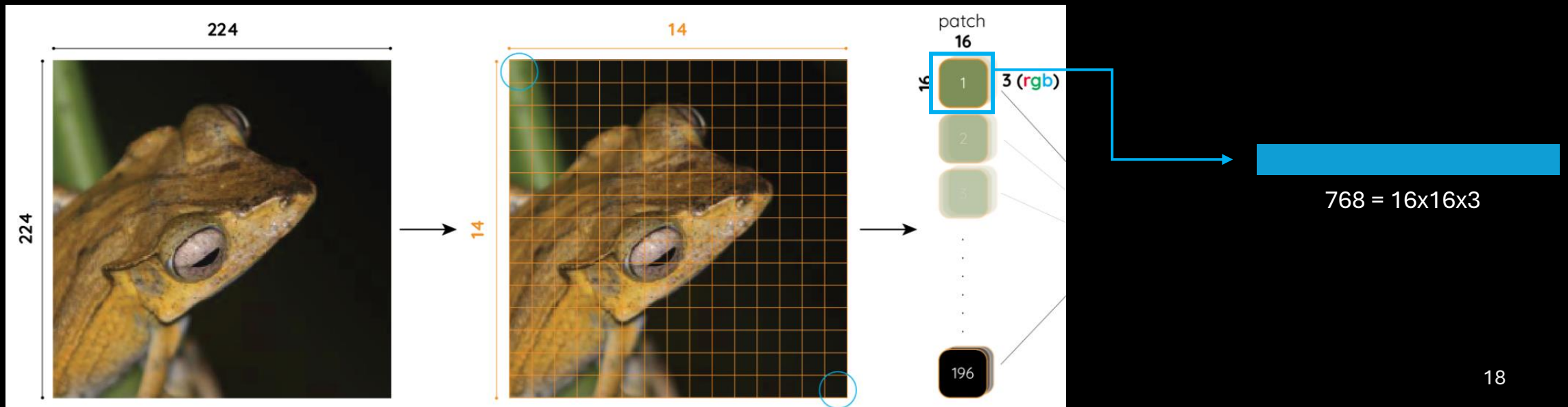
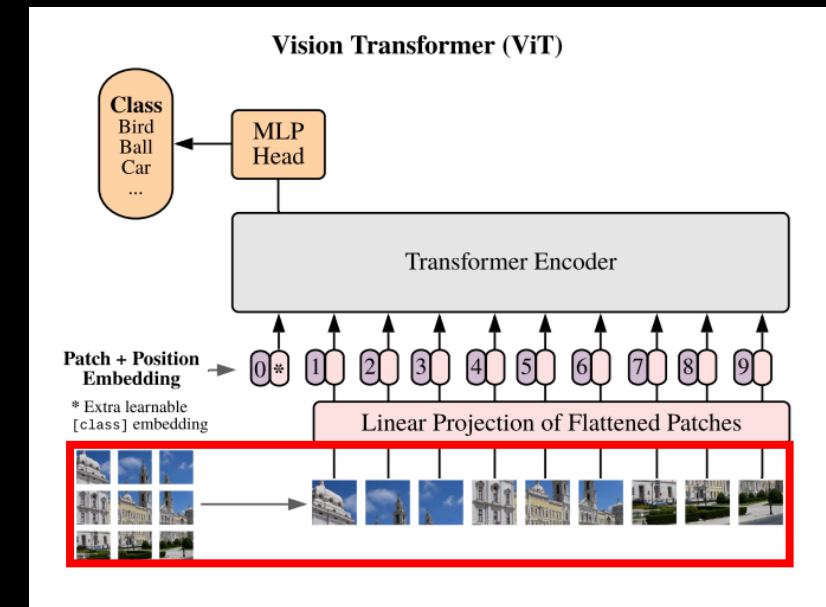
Patches

- Using an image as a sequence of patches.
- Partition an image into equal sized patches.
224*224*3 image to 14 patches in each dimension.
Thus, patch size = 16*16*3
Total of 196 = 14*14 patches



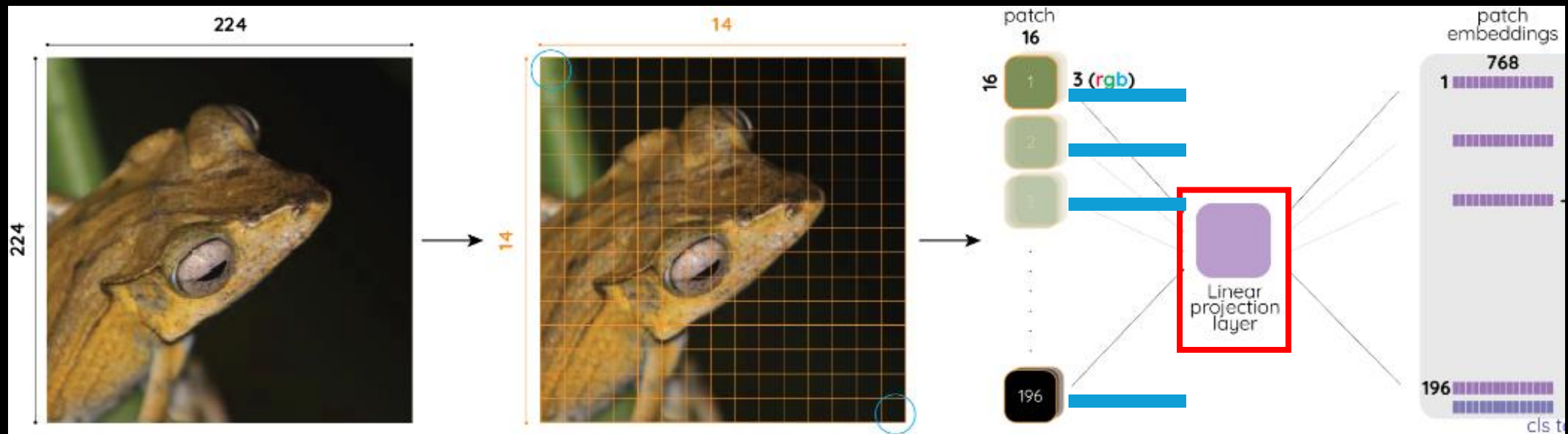
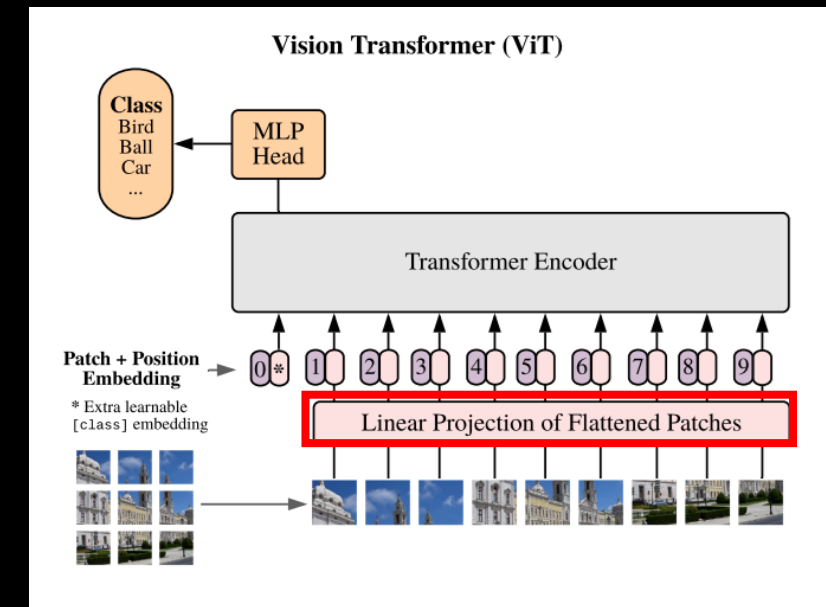
Flattening

- Each patch is technically $16 \times 16 \times 3$ dimension.
- Flattening (a.k.a. vectorization):
Tensor becomes a 1D-vector



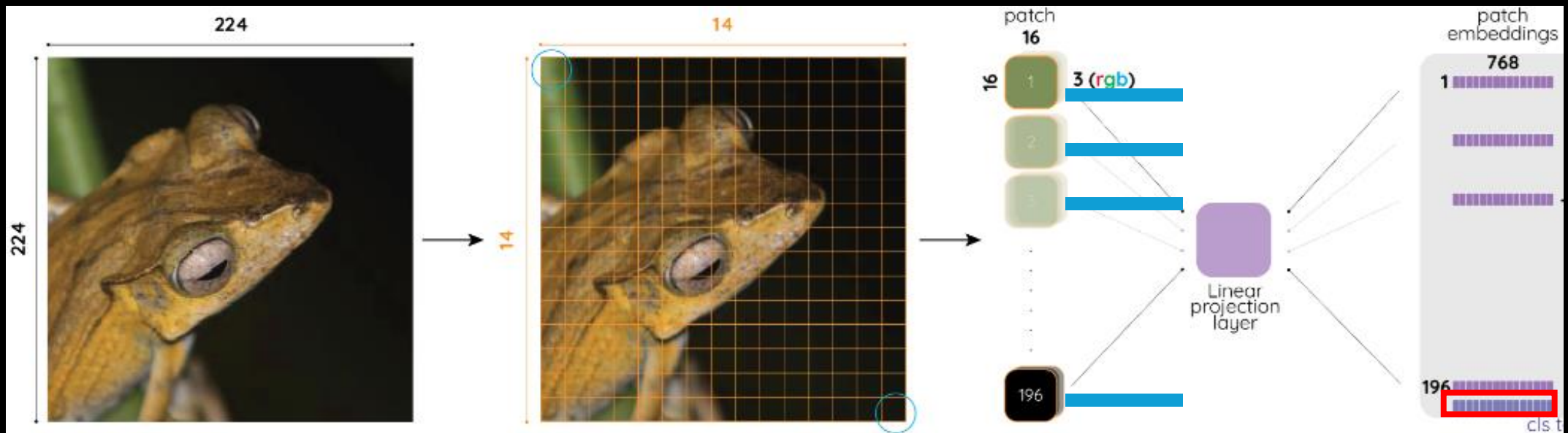
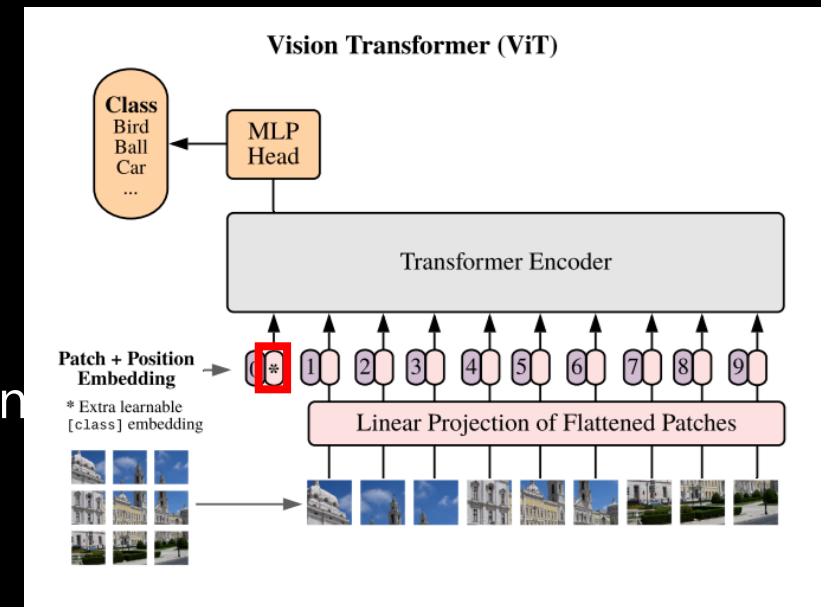
Patch Embedding

- Each 768×1 flattened patch goes through a fully connected layer to become a new 768D embedding.



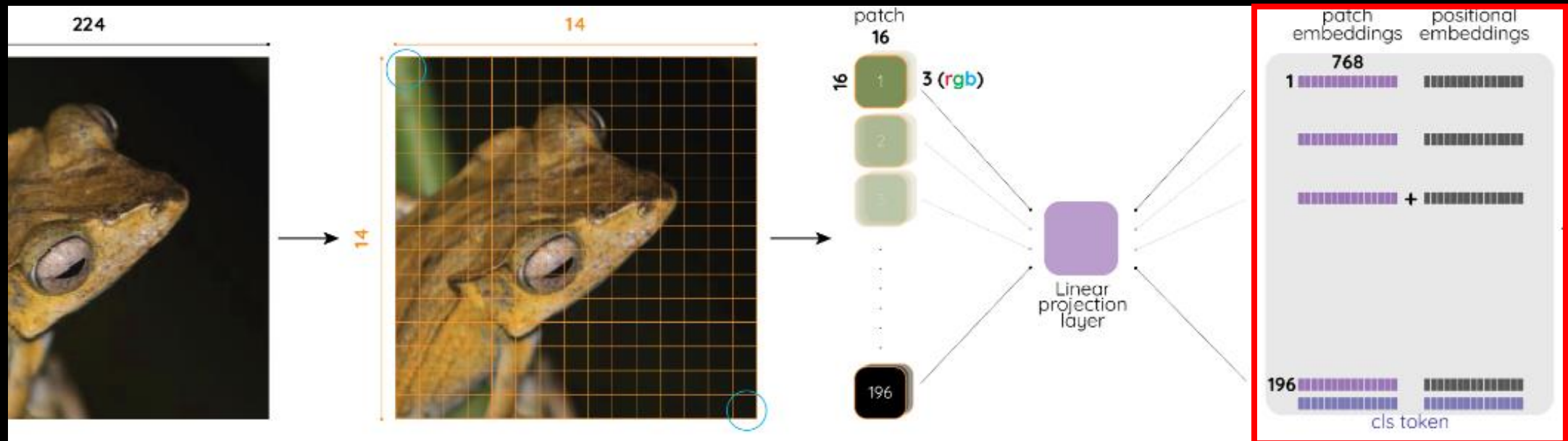
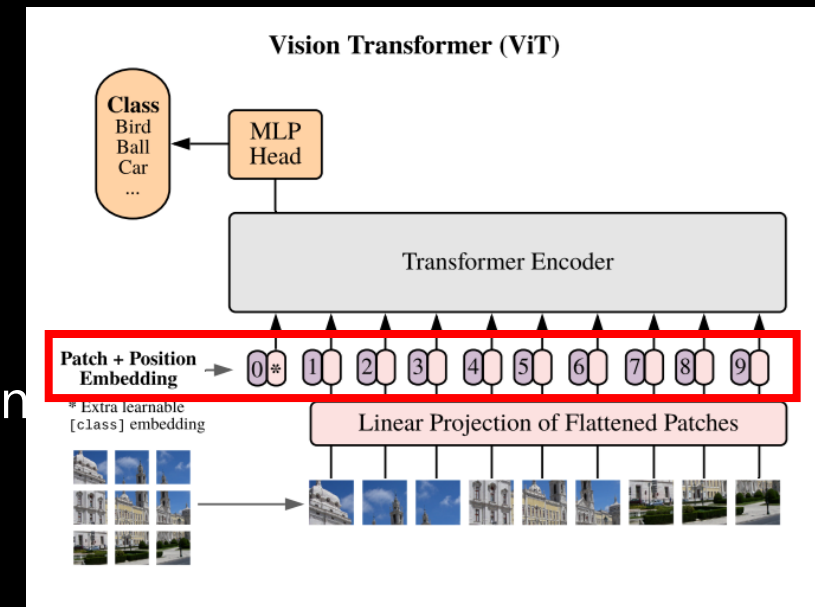
Patch Embedding

- Extra learnable class embedding with non-patch information
- $197 = 196(\text{patch embeddings}) + 1(\text{class embedding})$

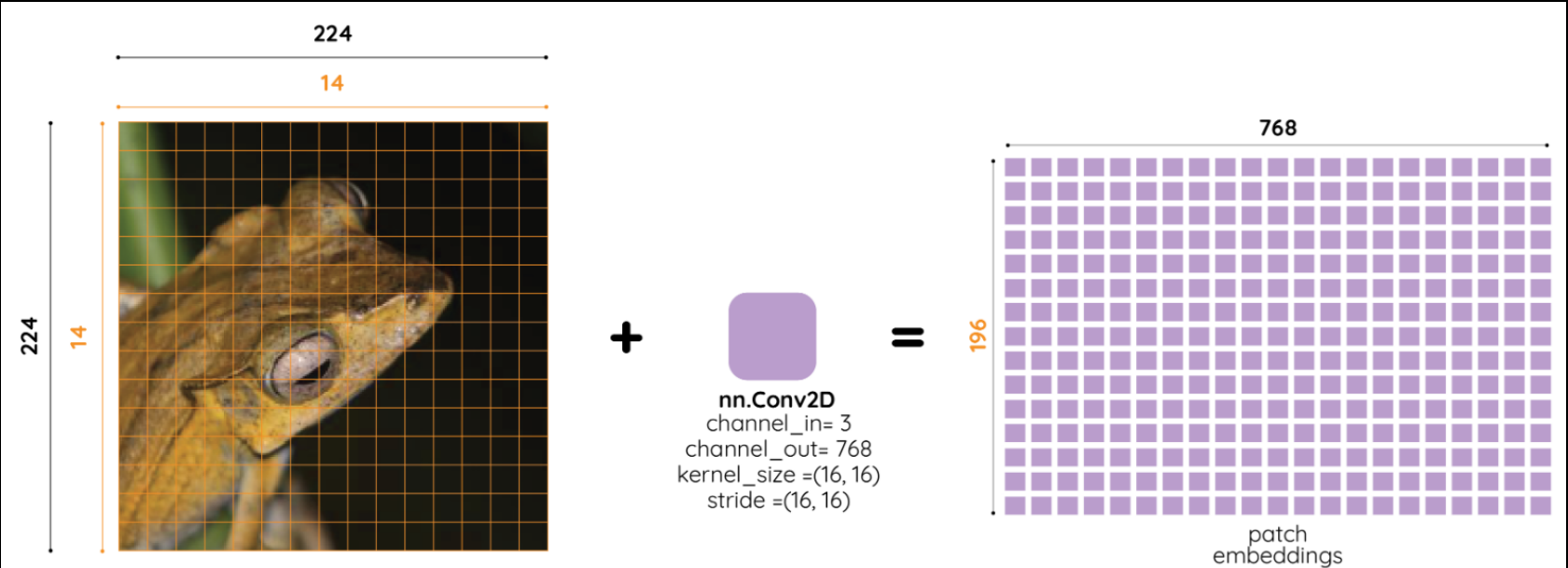
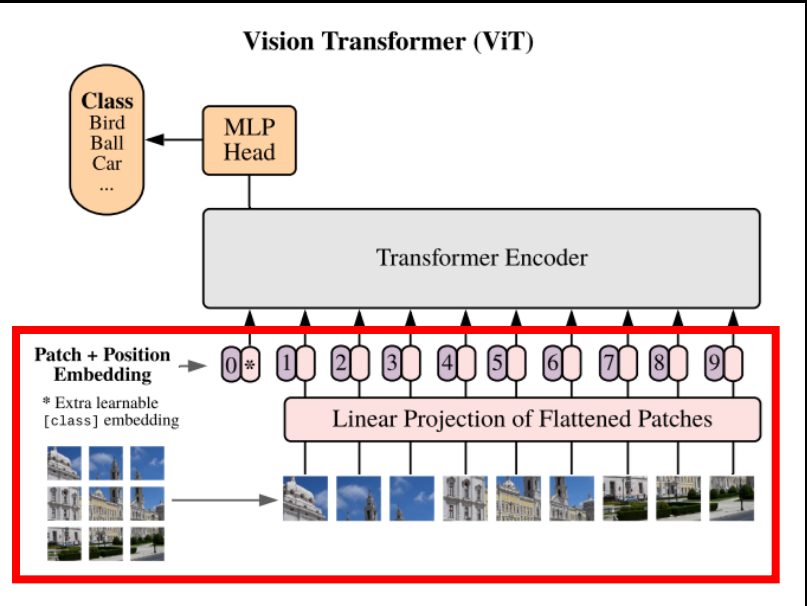


Positional Embedding

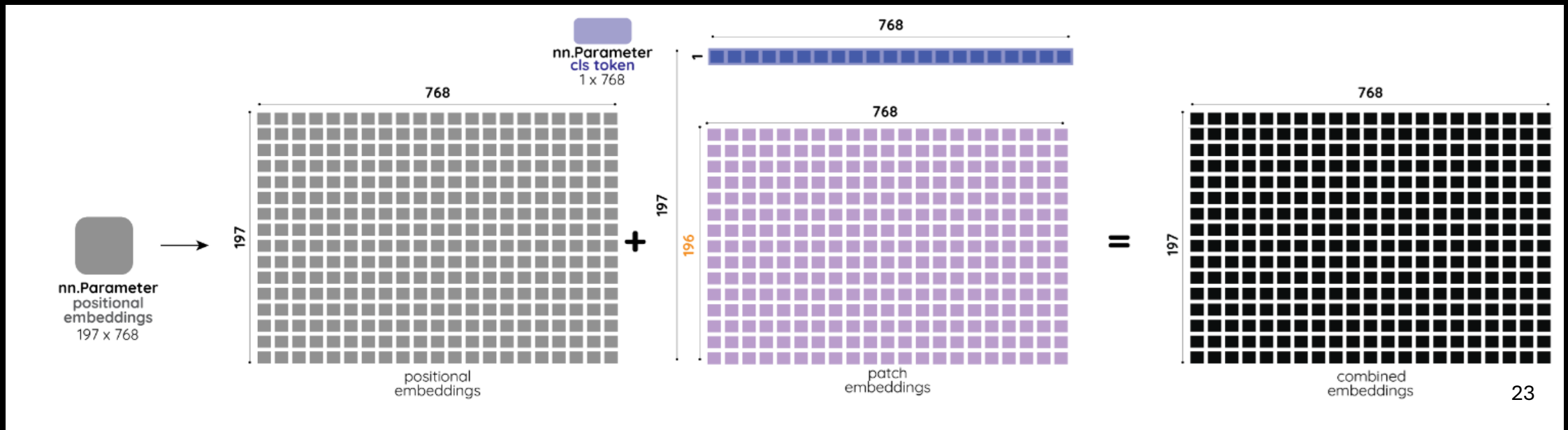
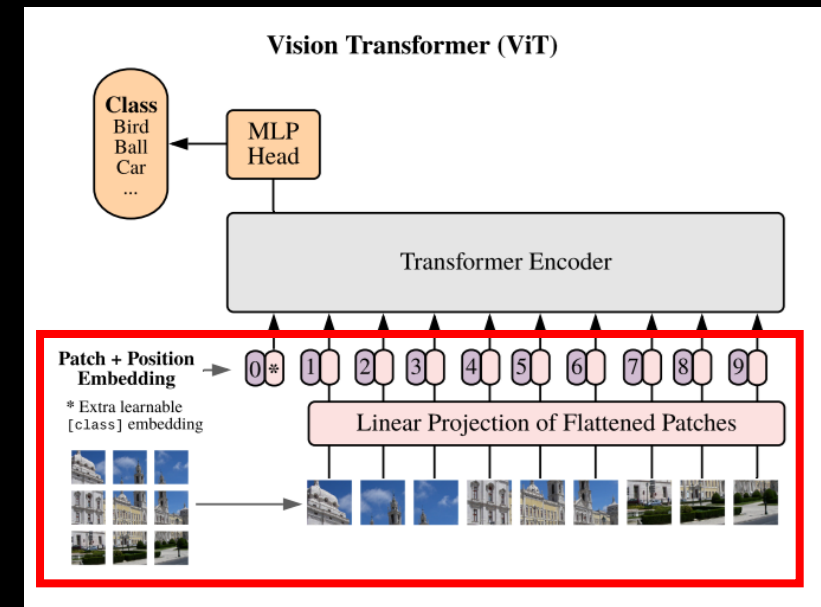
- Extra learnable class embedding with non-patch information
- $197 = 196(\text{patch embeddings}) + 1(\text{class embedding})$



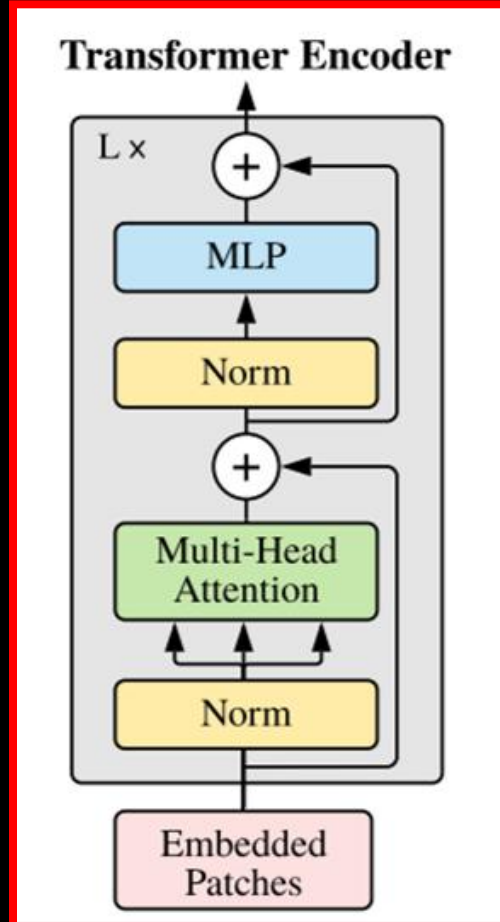
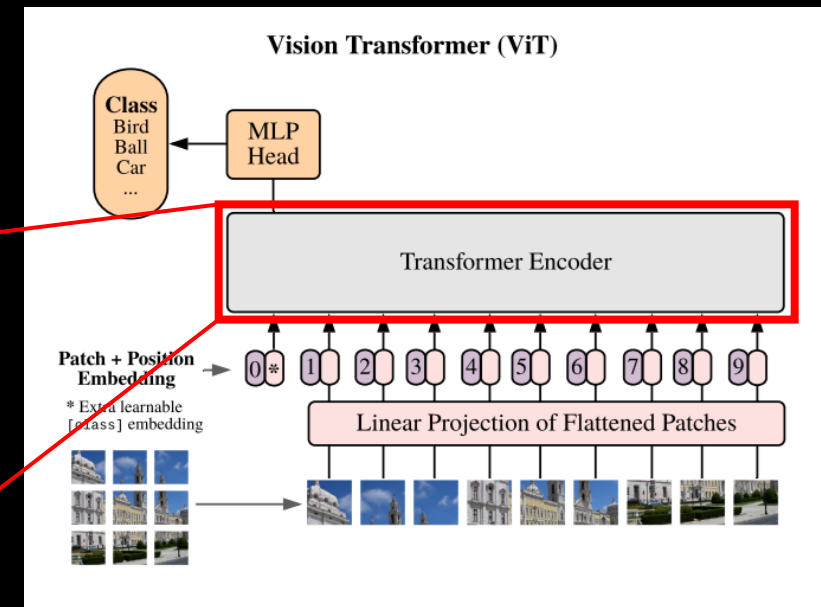
Summary So Far



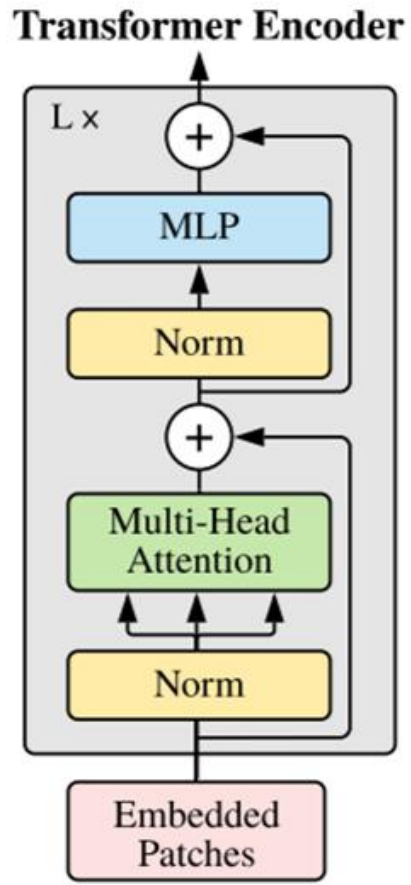
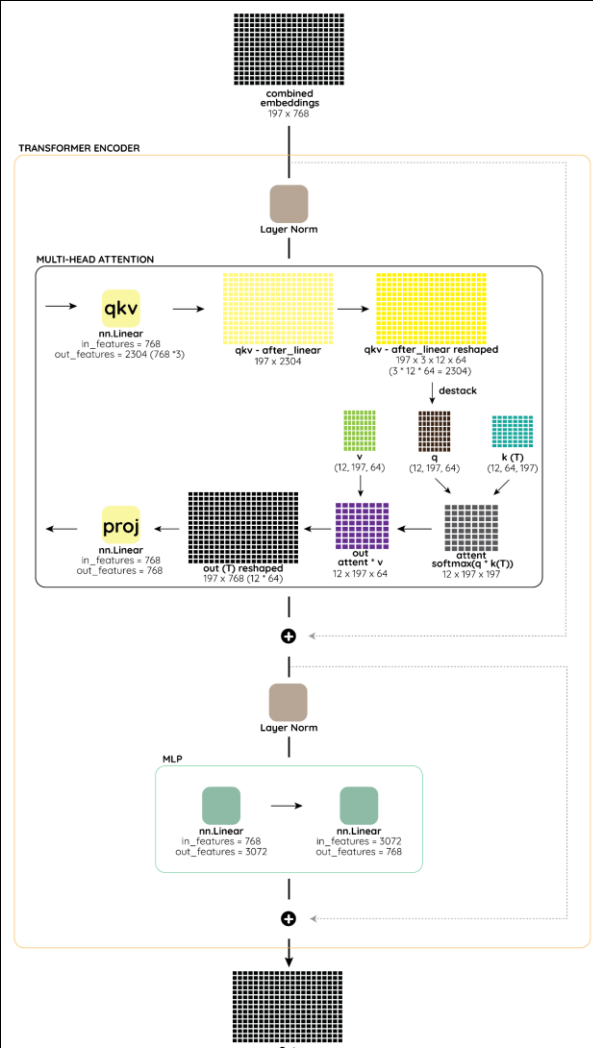
Summary So Far



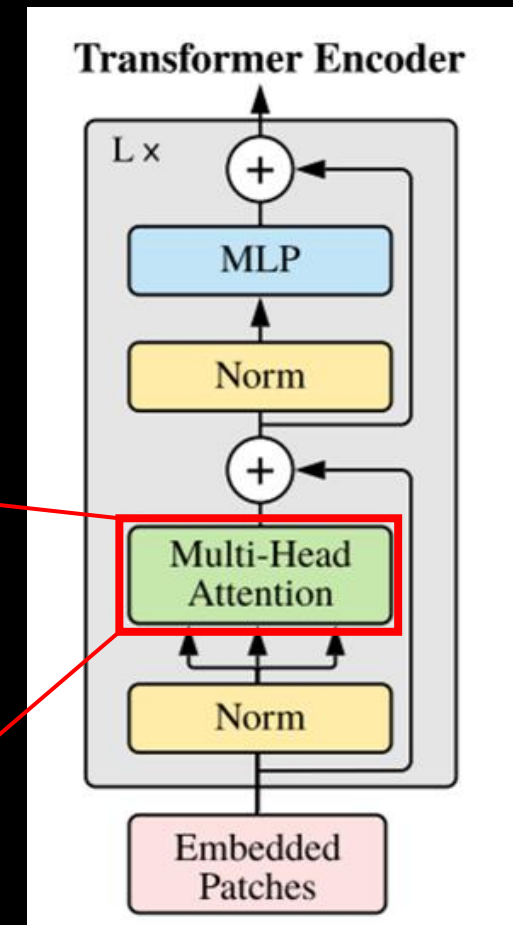
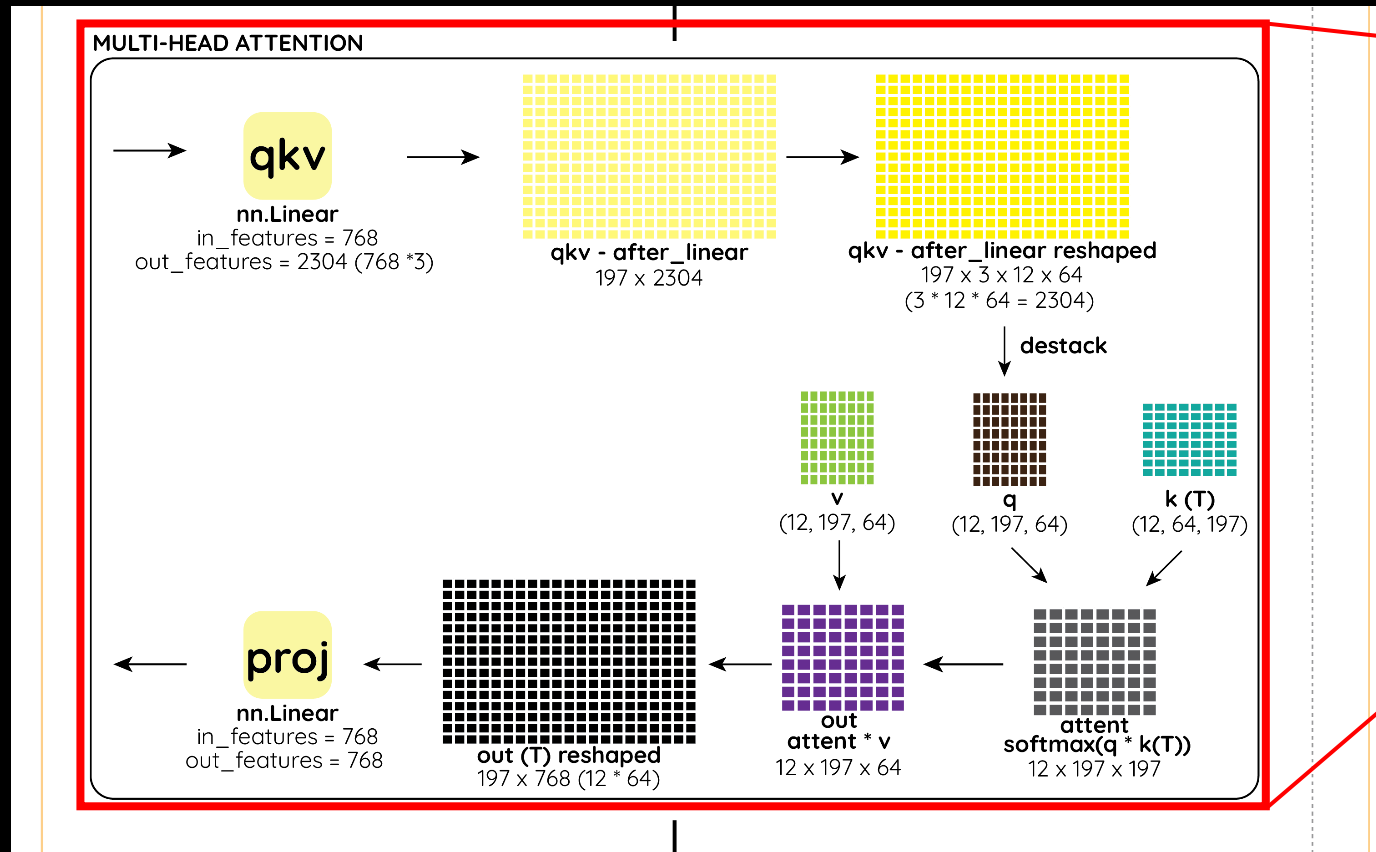
Transformer Encoder



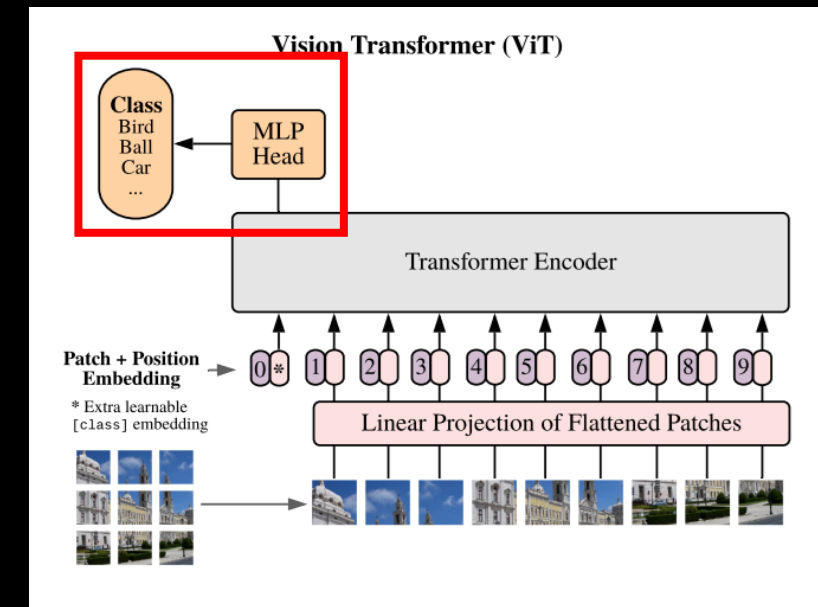
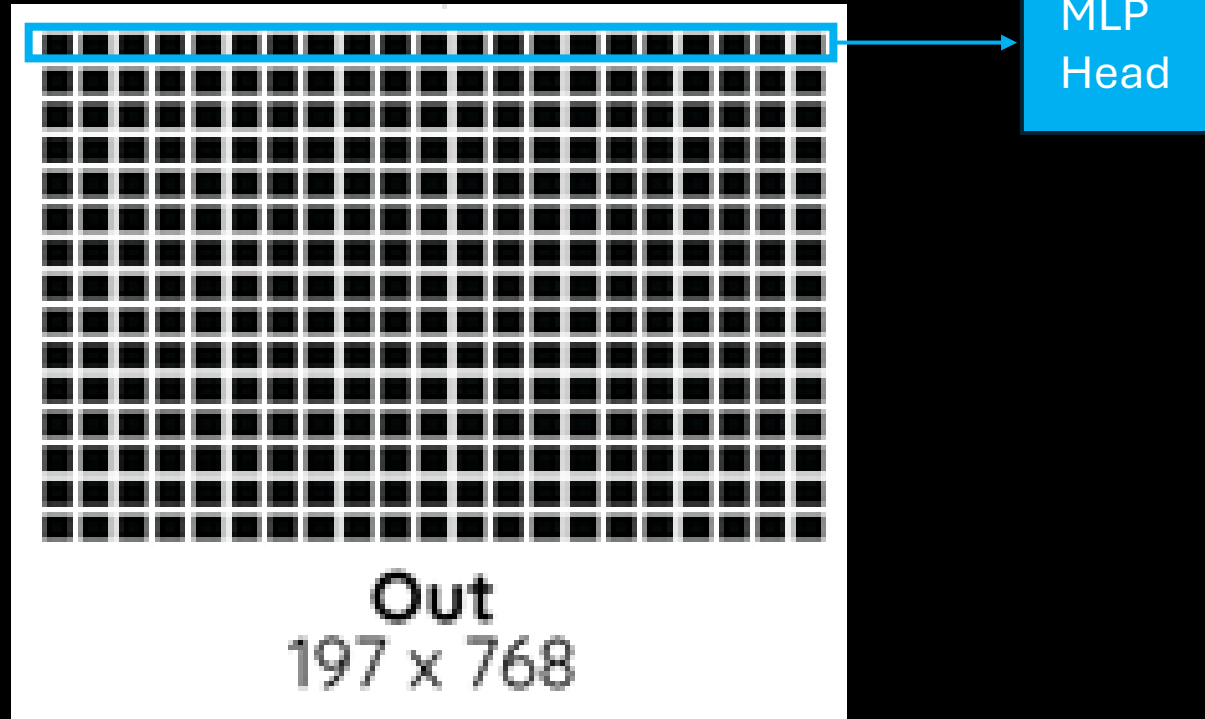
Transformer Encoder



Transformer Encoder

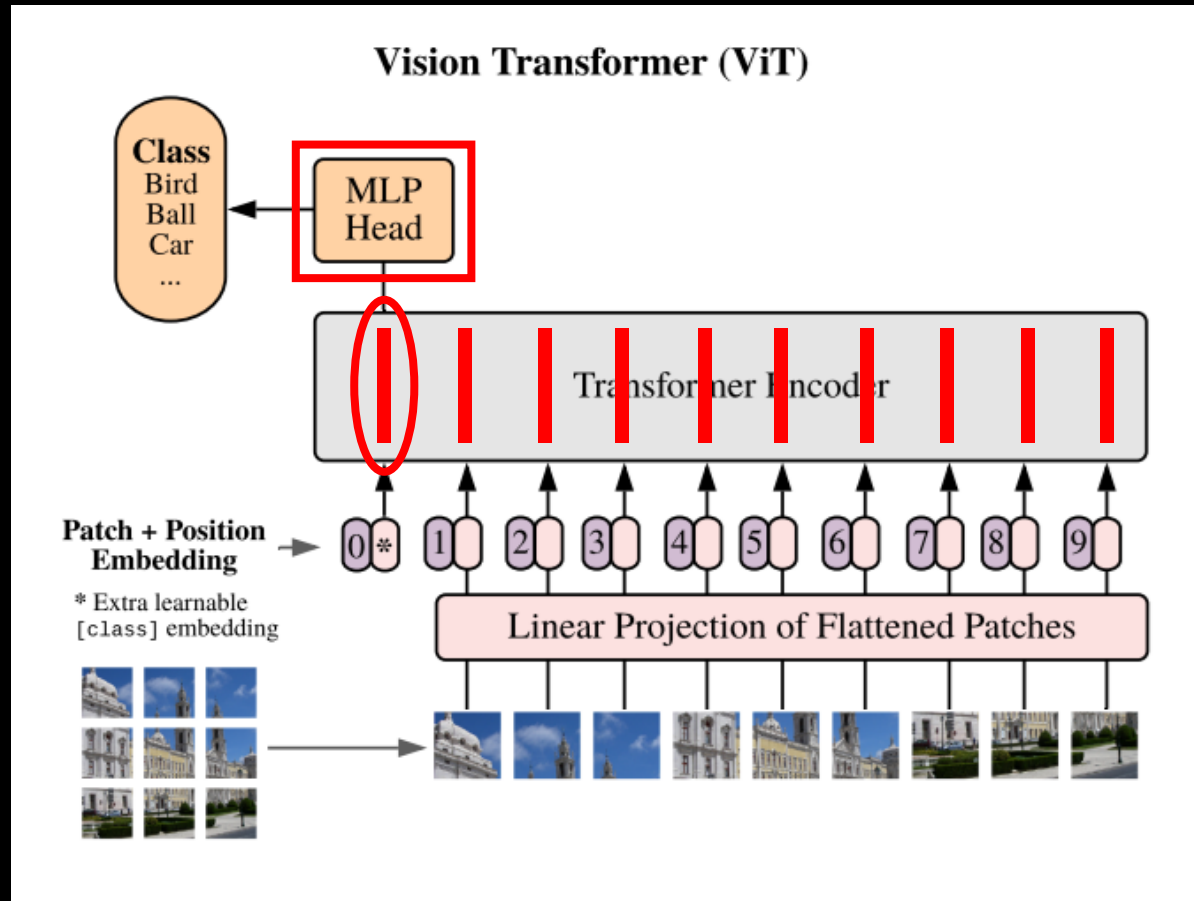


Classification



Only the [class token]'s output from the Transformer Encoder layer is used as the image representation for class prediction.

Classification



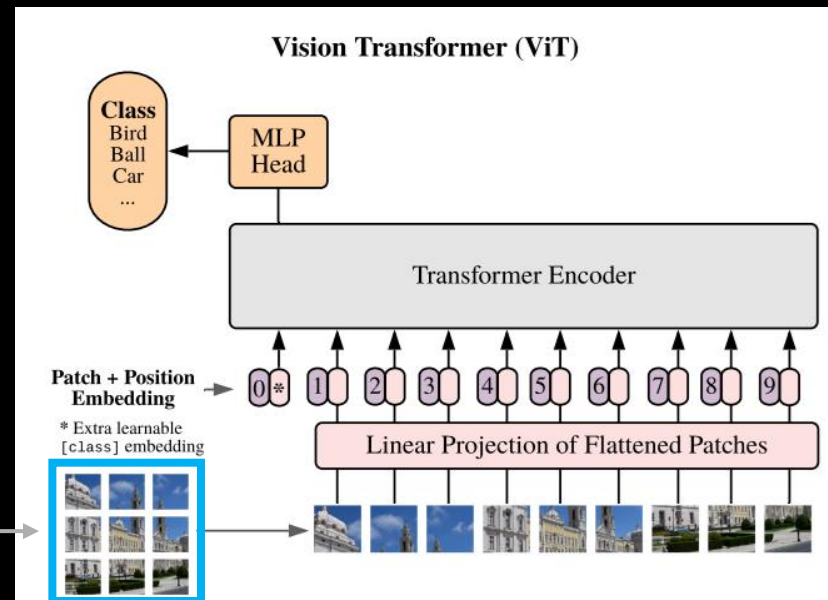
03

ViT Performance Analysis

Experiments

We evaluate the representation learning capabilities of

- ResNet (BiT)
- ViT
- Hybrid



Feature map of CNN

Model Variants

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

Table 1: Details of Vision Transformer model variants.

For instance, ViT-L/16 means the “Large” variant with 16x16 input patch size.

Datasets

- ImageNet: 1k classes with 1.3M images
- ImageNet-21k: 21k classes with 14M images
- JFT: 18k classes and 303M high-resolution images
- 19-task VTAB: 1000 training examples per task

Metrics

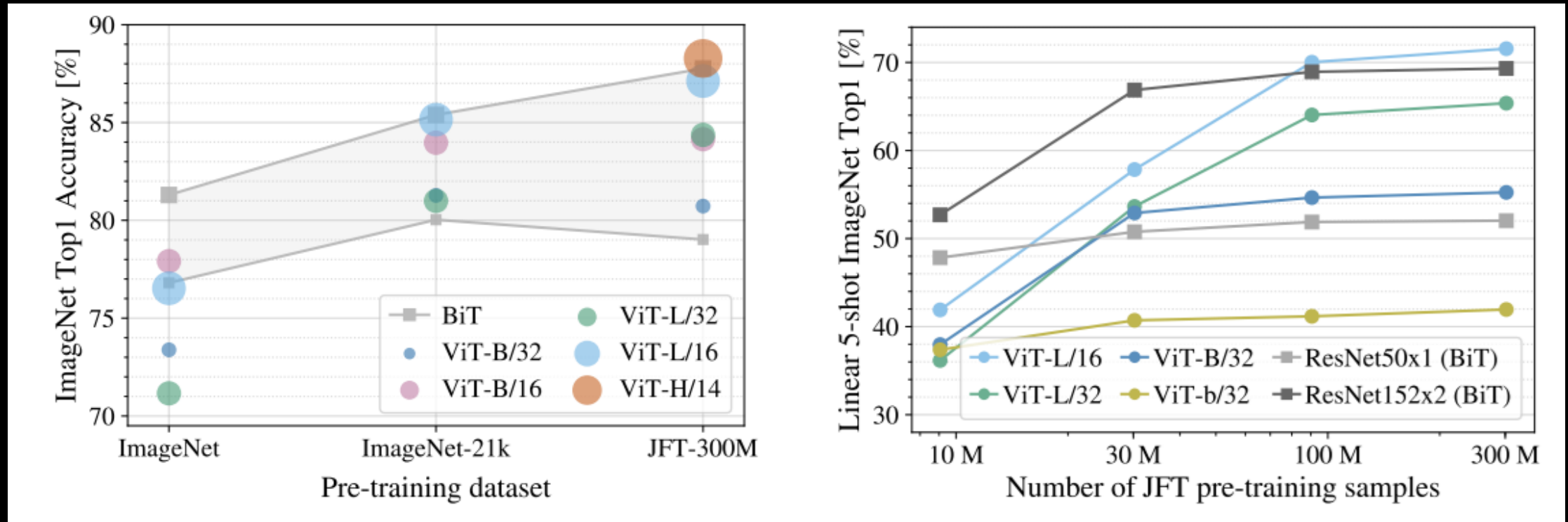
We report results on downstream datasets either through **few-shot** or **finetuning accuracy**.

Comparison to SOTA

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	88.55 ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4/88.5*
ImageNet ReaL	90.72 ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	99.50 ± 0.06	99.42 ± 0.03	99.15 ± 0.03	99.37 ± 0.06	—
CIFAR-100	94.55 ± 0.04	93.90 ± 0.05	93.25 ± 0.05	93.51 ± 0.08	—
Oxford-IIIT Pets	97.56 ± 0.03	97.32 ± 0.11	94.07 ± 0.15	96.62 ± 0.23	—
Oxford Flowers-102	99.68 ± 0.02	99.74 ± 0.00	99.61 ± 0.02	99.63 ± 0.03	—
VTAB (19 tasks)	77.63 ± 0.23	76.28 ± 0.46	72.72 ± 0.21	76.29 ± 1.70	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

Pretraining Data Requirements

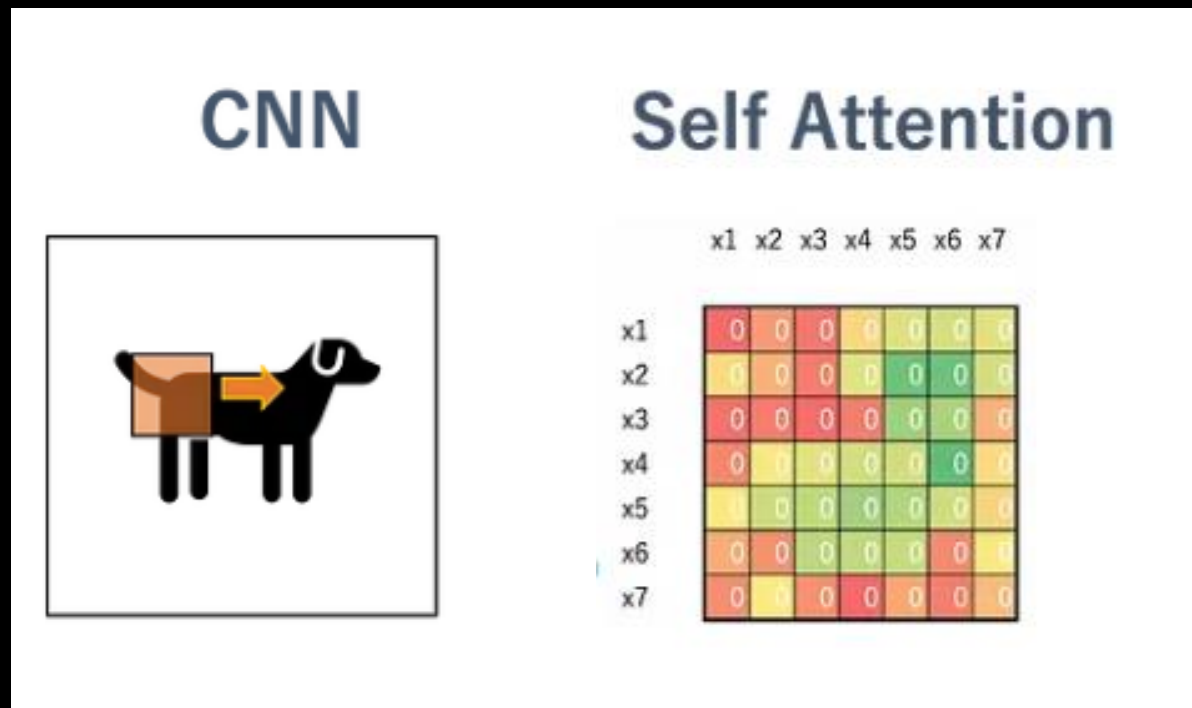
Good scalability: more data, higher performance



Inductive Bias: CNN vs. ViT

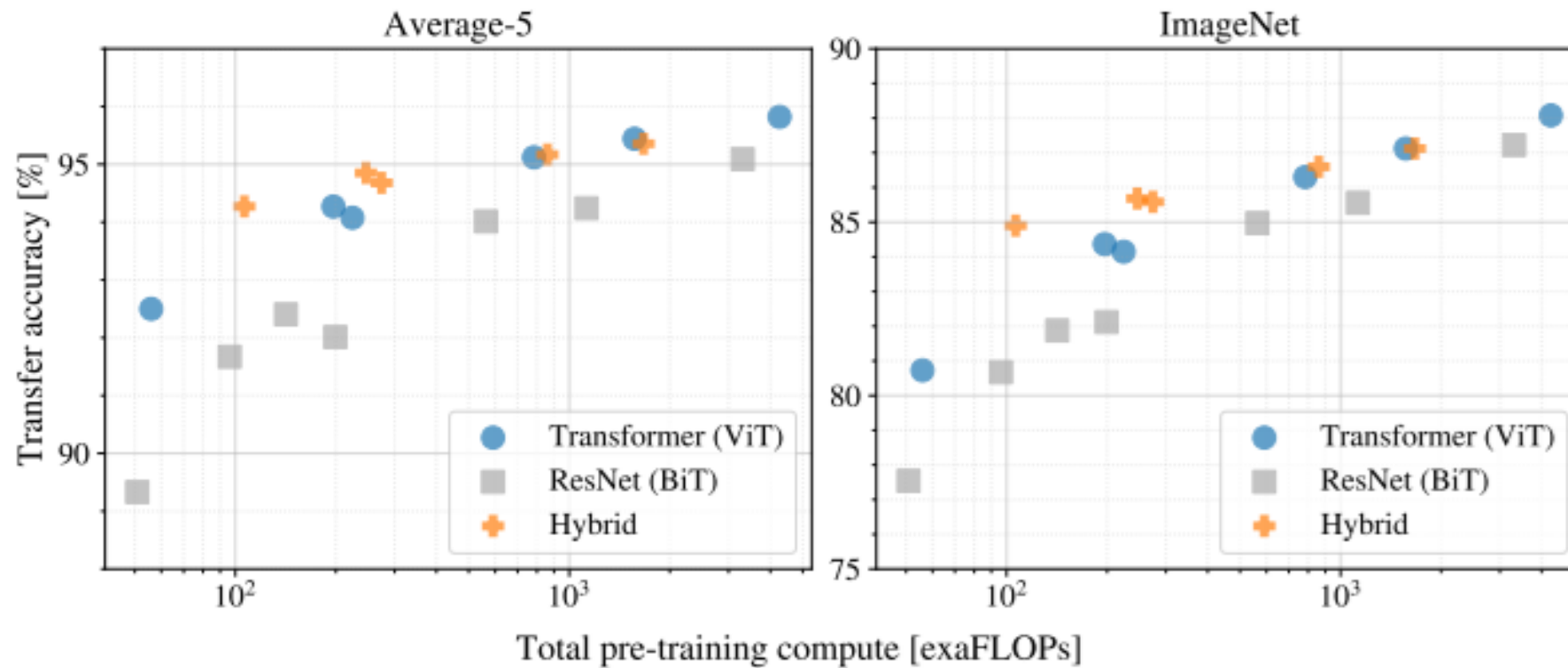
CNN, which has a **strong inductive bias** that the information is locally aggregated.

ViT, which has a relatively **weak inductive bias** because it only correlates all features.

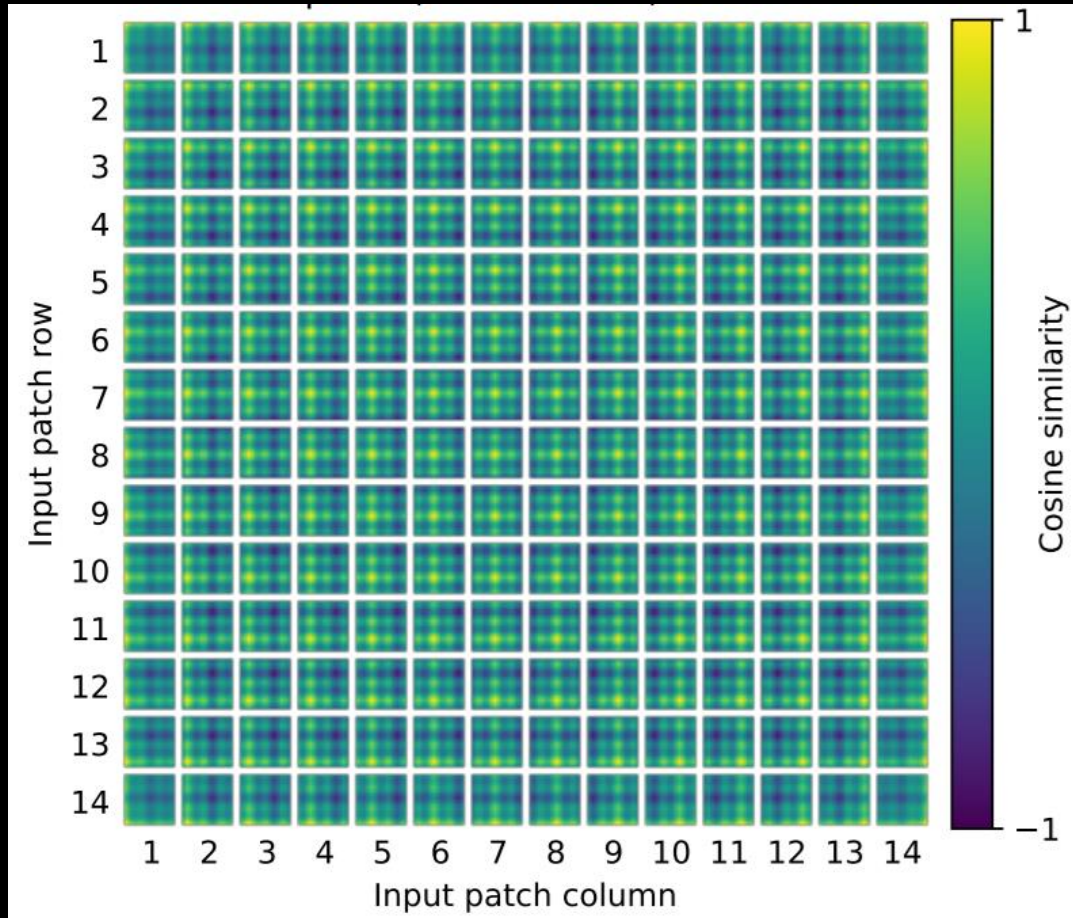


ViTs vs. ResNets: Efficiency Edge

ViTs dominate ResNets on the performance/compute tradeoff!

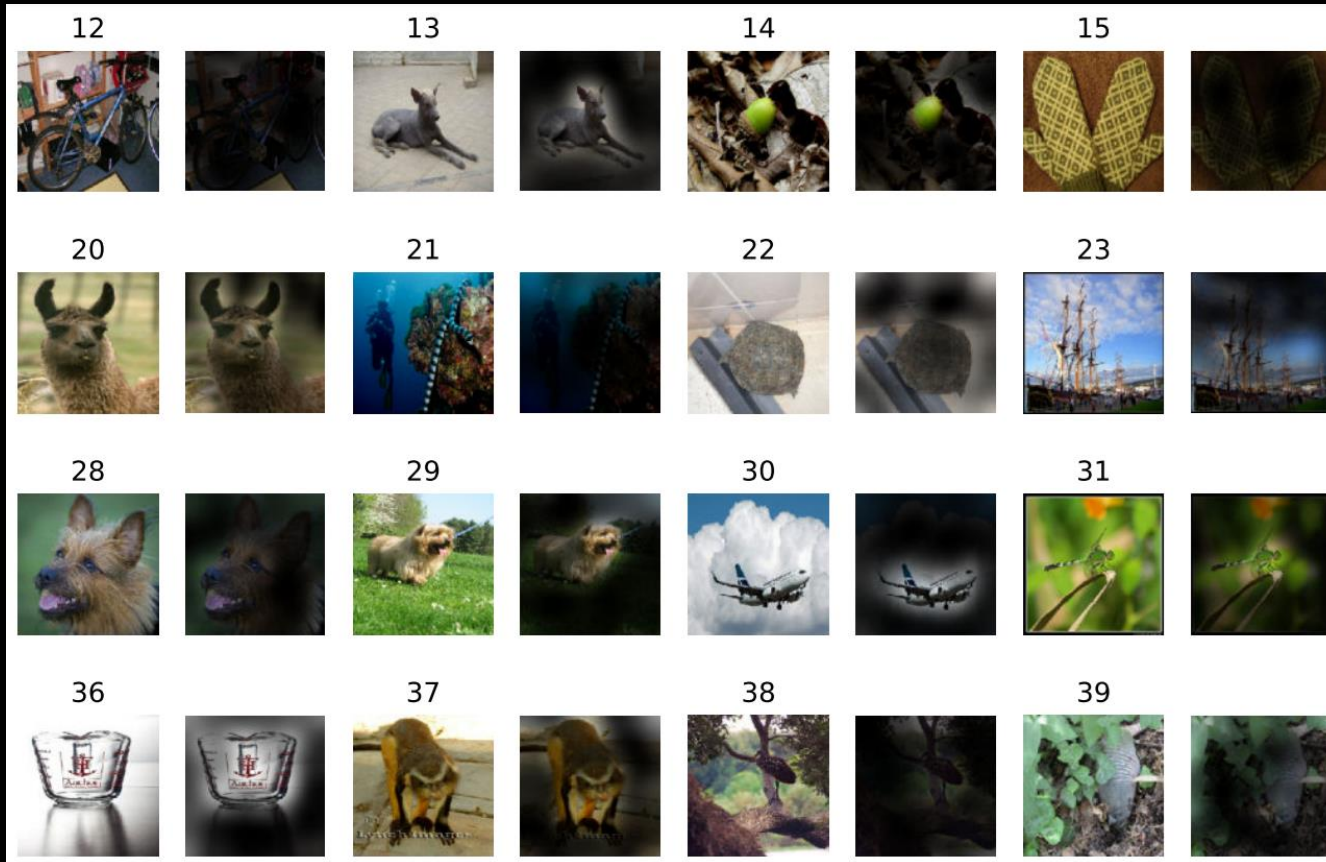


Position Embedding Similarity



The learned positional embedding happen to distinguish their positions!

Attention map



Examples of attention from the output token to the input space.

Significance of ViT

- Direct application of Transformers to image recognition.
- ViT matches or exceeds the SOTA of many image classification datasets and is cost-effective to pre-train.
- Many challenges remain...

Thank You!
Any Questions?