



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, neilhoulsby}@google.com

ABSTRACT

While the Transformer architecture has become the de-facto standard for natural language processing tasks, its applications to computer vision remain limited. In vision, attention is either applied in conjunction with convolutional networks, or used to replace certain components of convolutional networks while keeping their overall structure in place. We show that this reliance on CNNs is not necessary and a pure transformer applied directly to sequences of image patches can perform very well on image classification tasks. When pre-trained on large amounts of data and transferred to multiple mid-sized or small image recognition benchmarks (ImageNet, CIFAR-100, VTAB, etc.), Vision Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train.¹

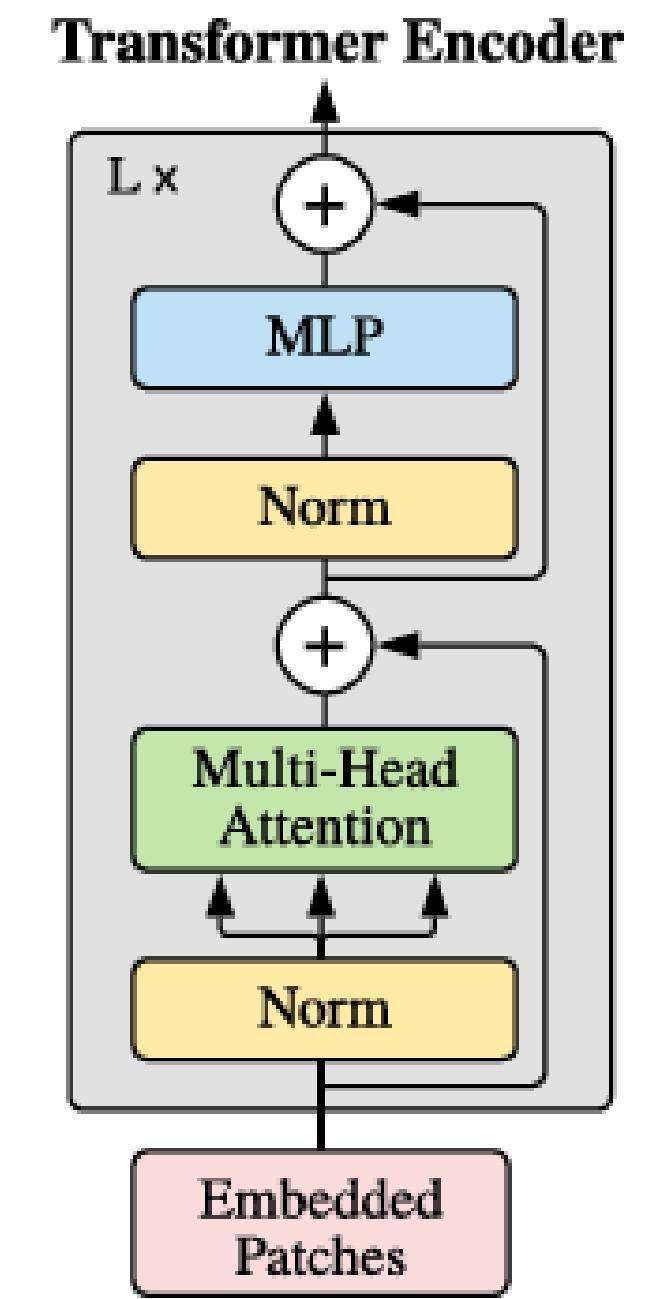
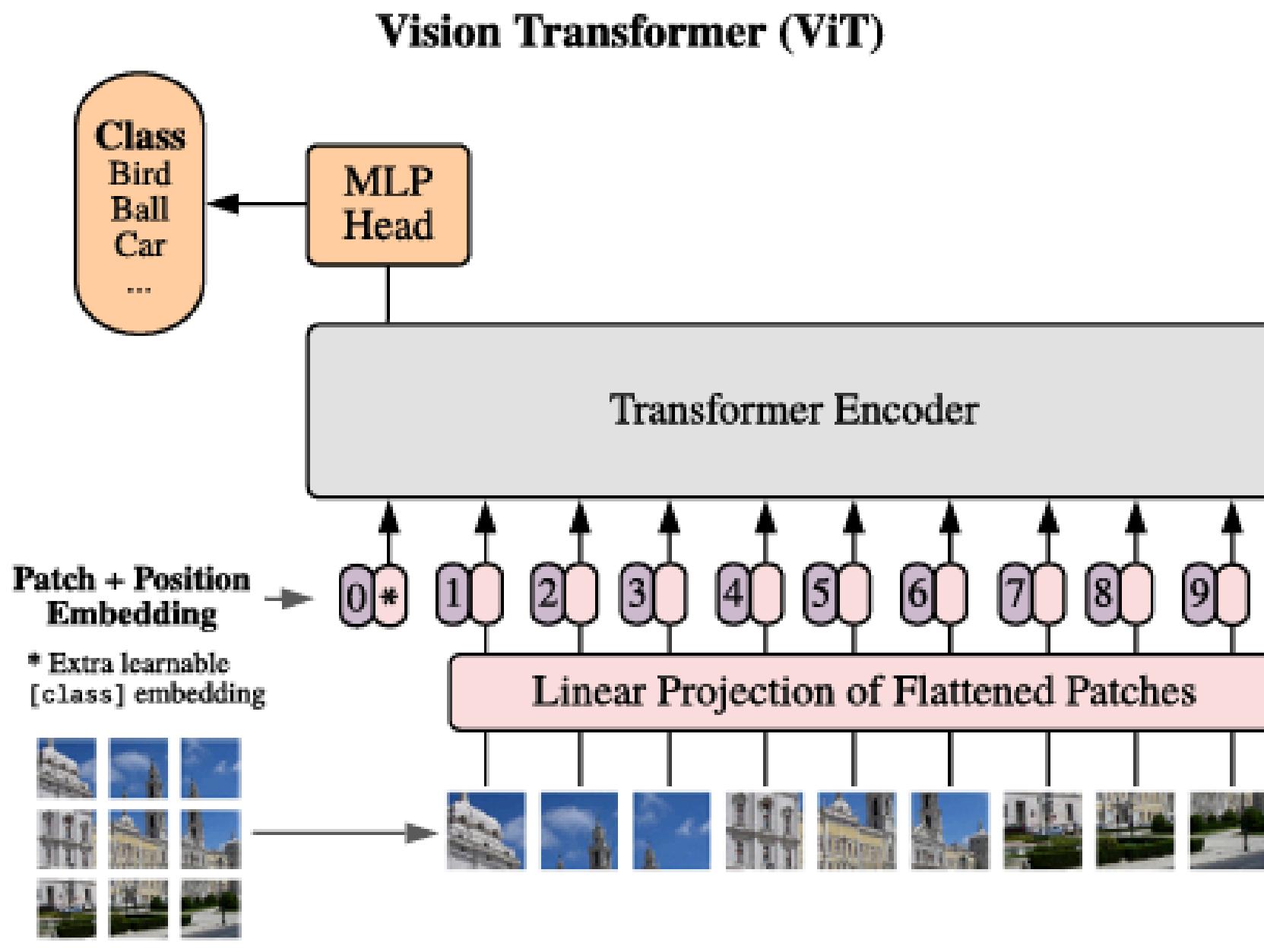
1 INTRODUCTION

Self-attention-based architectures, in particular Transformers (Vaswani et al., 2017), have become the model of choice in natural language processing (NLP). The dominant approach is to pre-train on a large text corpus and then fine-tune on a smaller task-specific dataset (Devlin et al., 2019). Thanks

What is ViT?

Inspired by Transformers, we **split an image into patches** and **provide the sequence of linear embeddings of these patches** as an input to a **Transformer**.

Image patches are treated the same way as **tokens** (words) in an NLP application



224



224

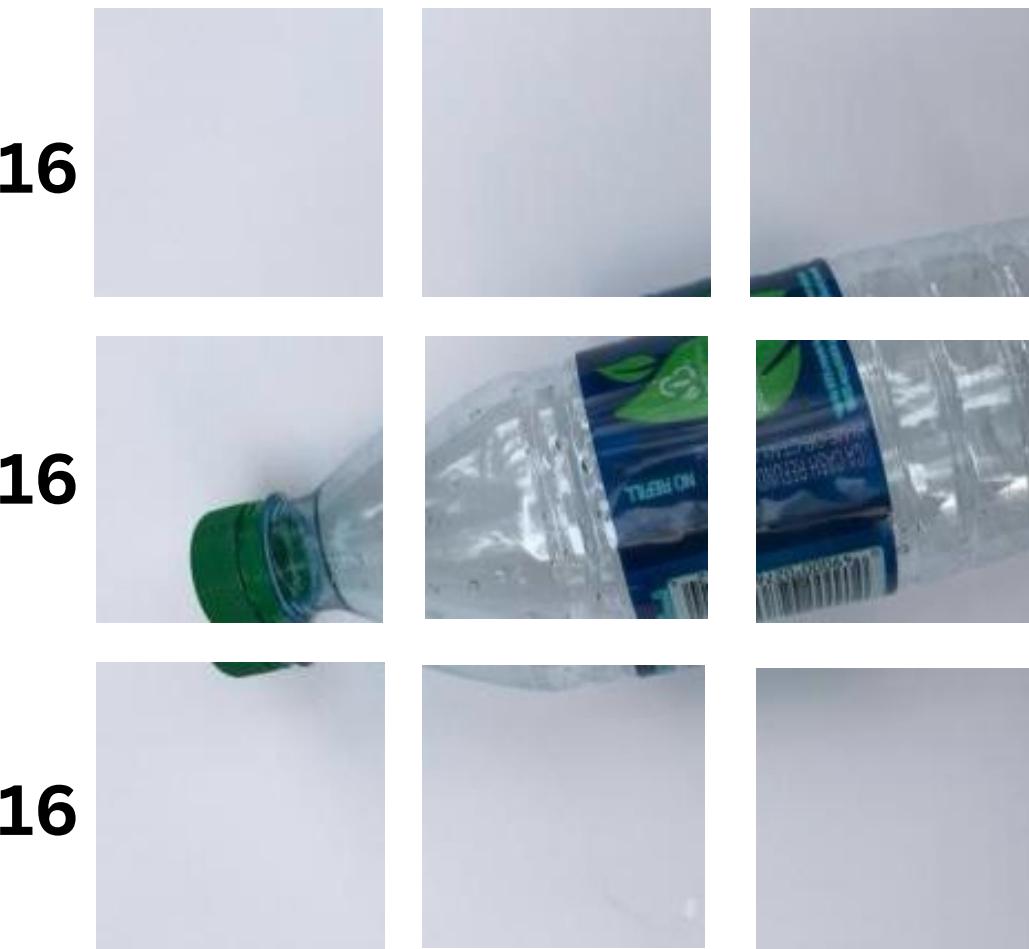
1

Input Image

$$N = \frac{HW}{P^2}$$

14 x 14

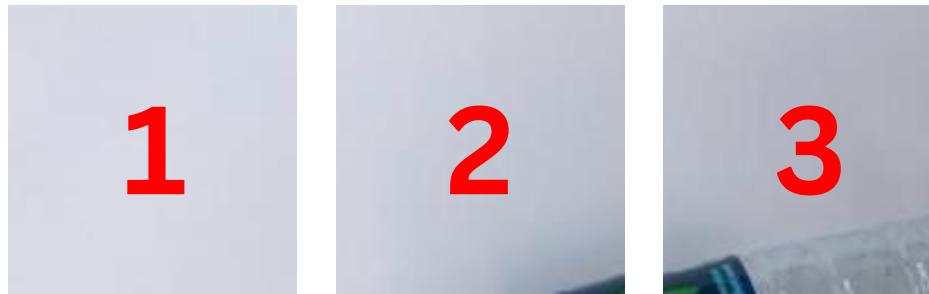
Image Patches



16

16

16

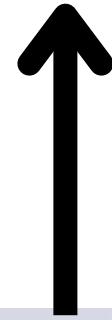


flattened patch

x_1



256 values, each representing 1 pixel



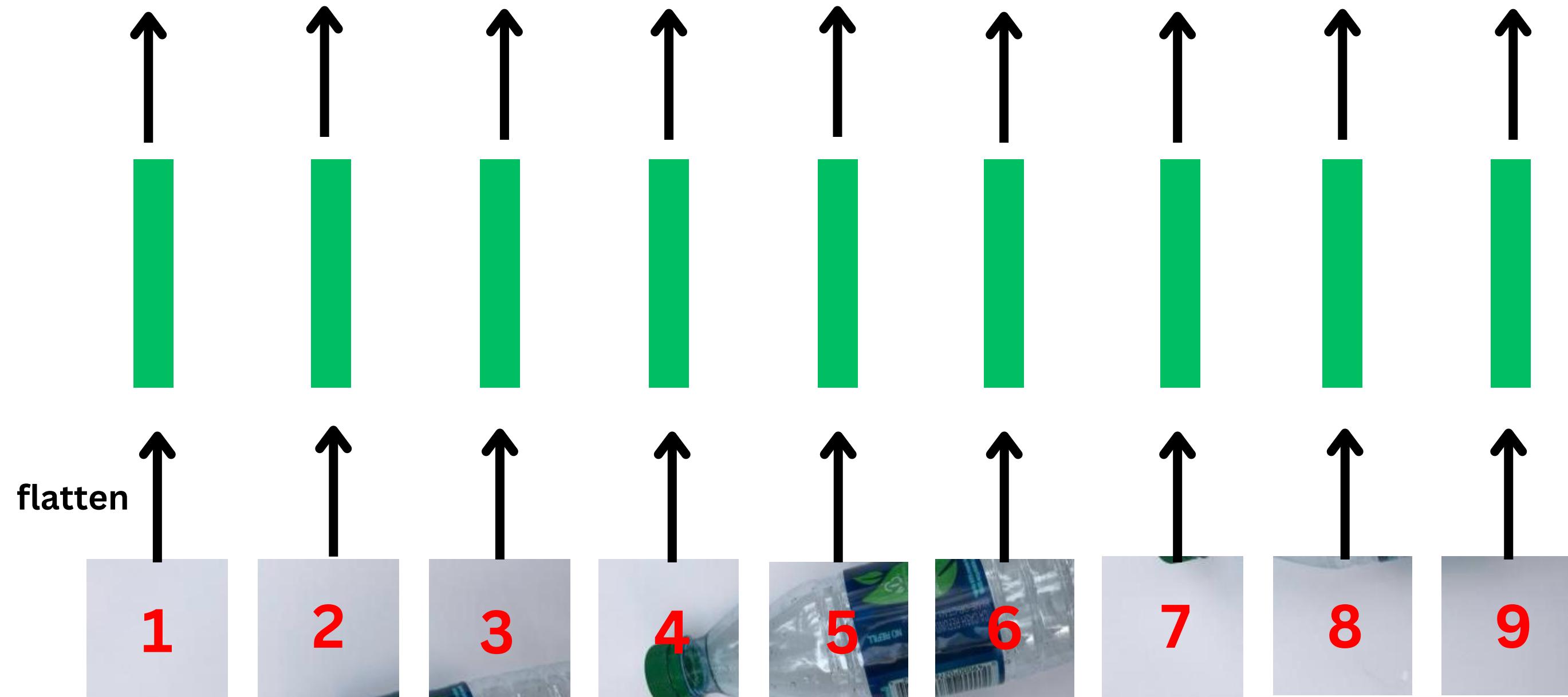
16

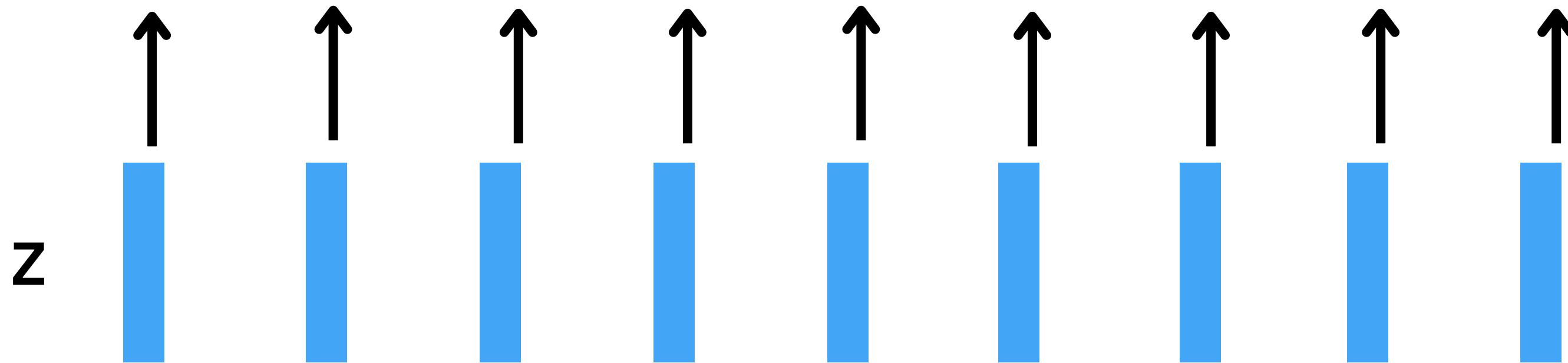
1

1 patch

16

Linear Projection Layer





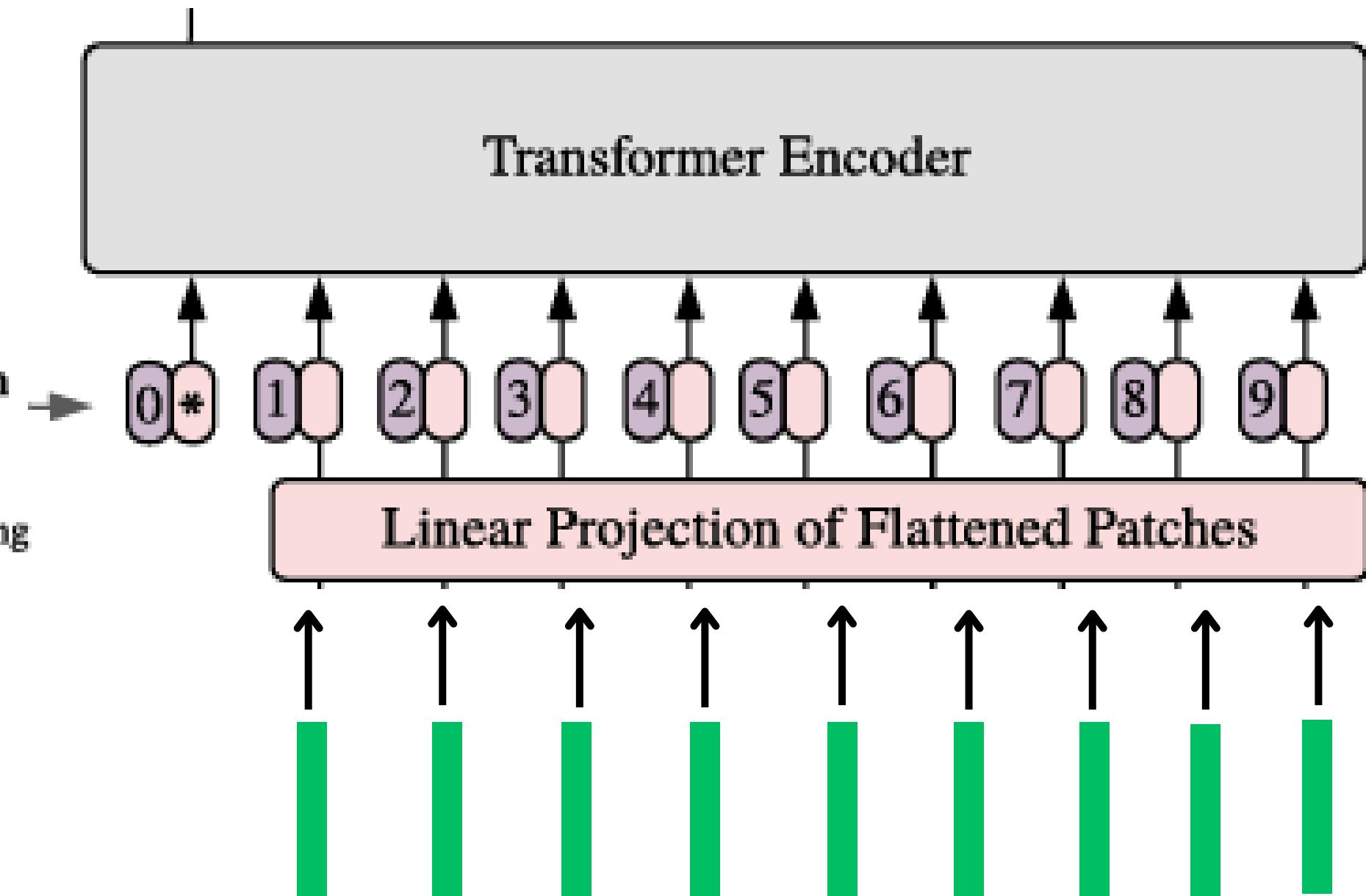
Linear Projection Layer

Positional Embeddings and Class Embedding



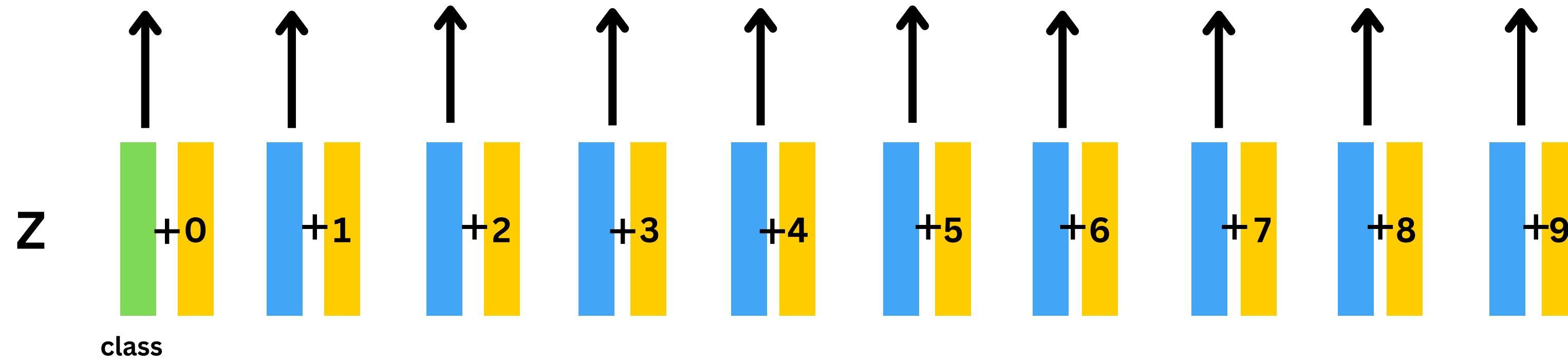
Patch + Position Embedding

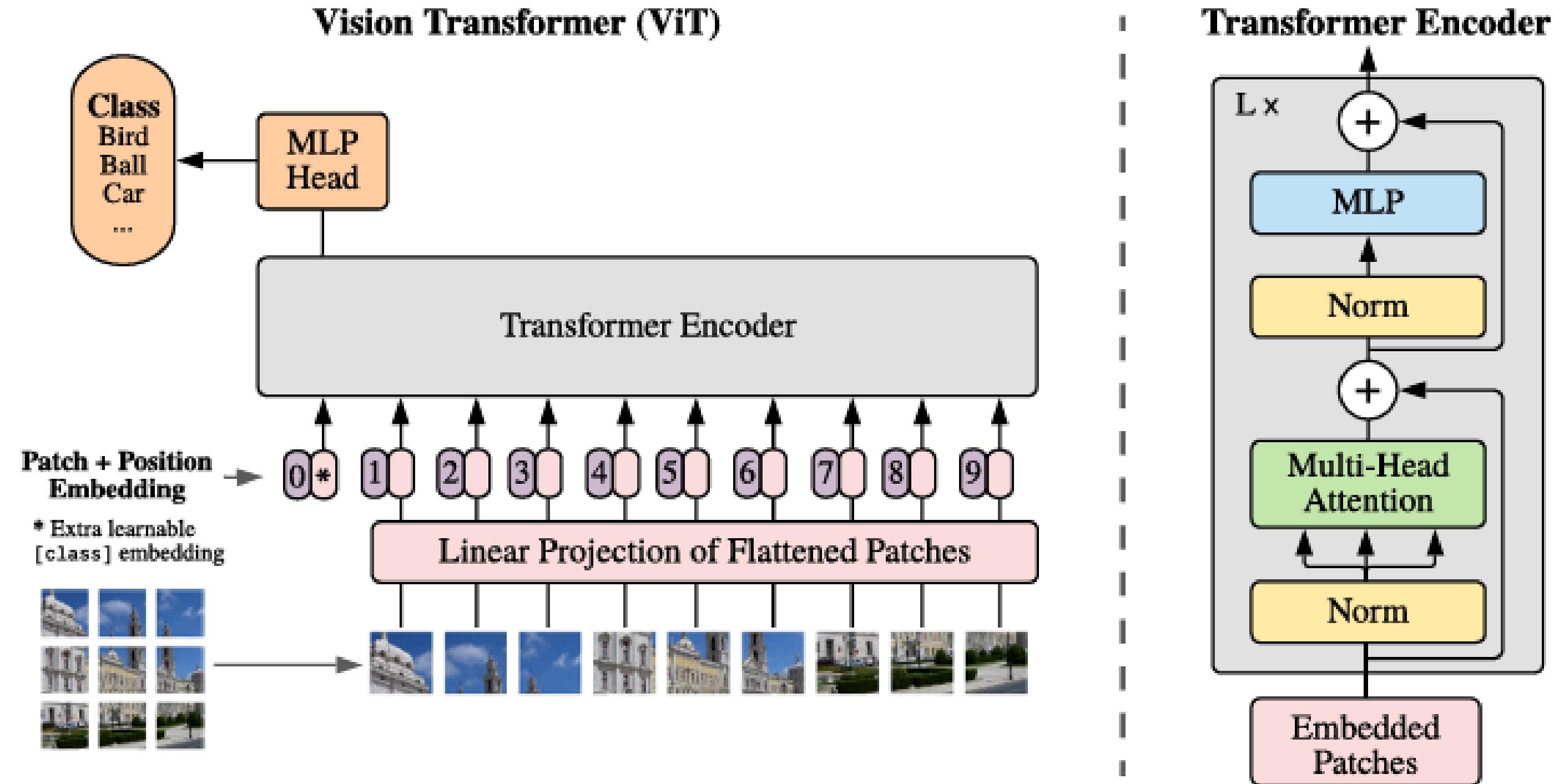
* Extra learnable [class] embedding



Pos emb

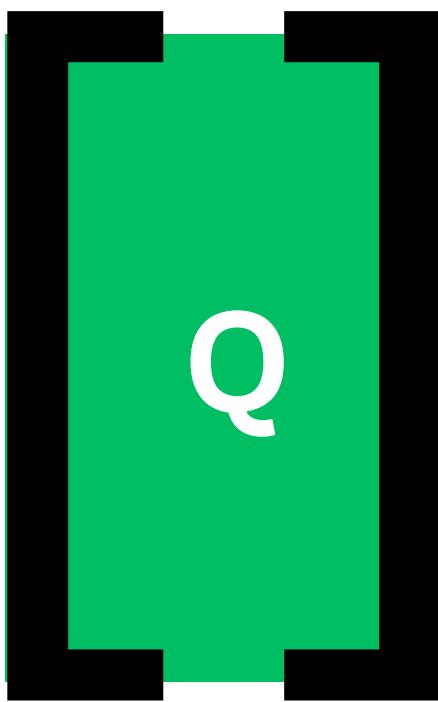
Transformer Encoder



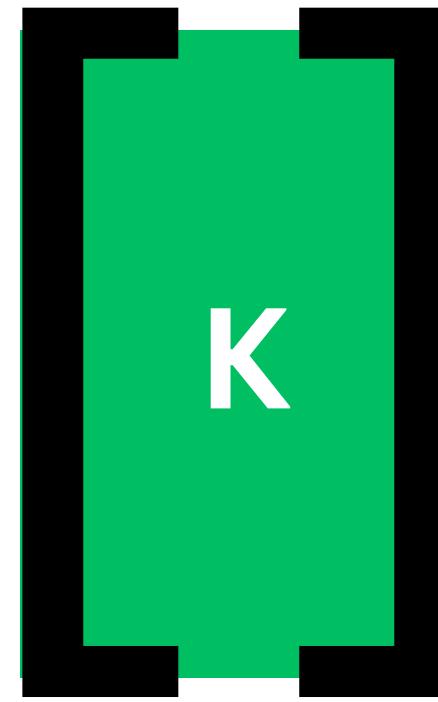


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Query



Key



Value



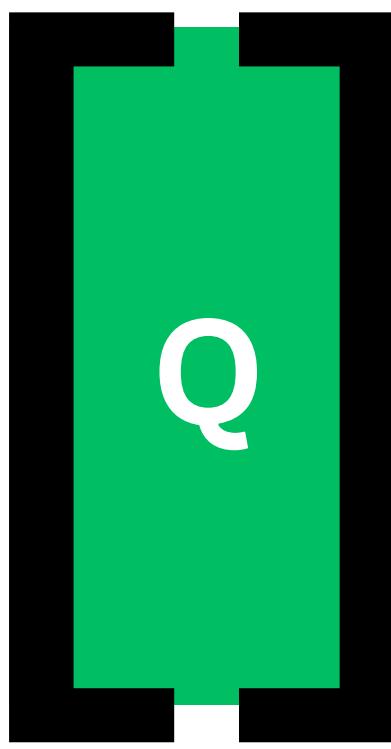
Weights Matrix



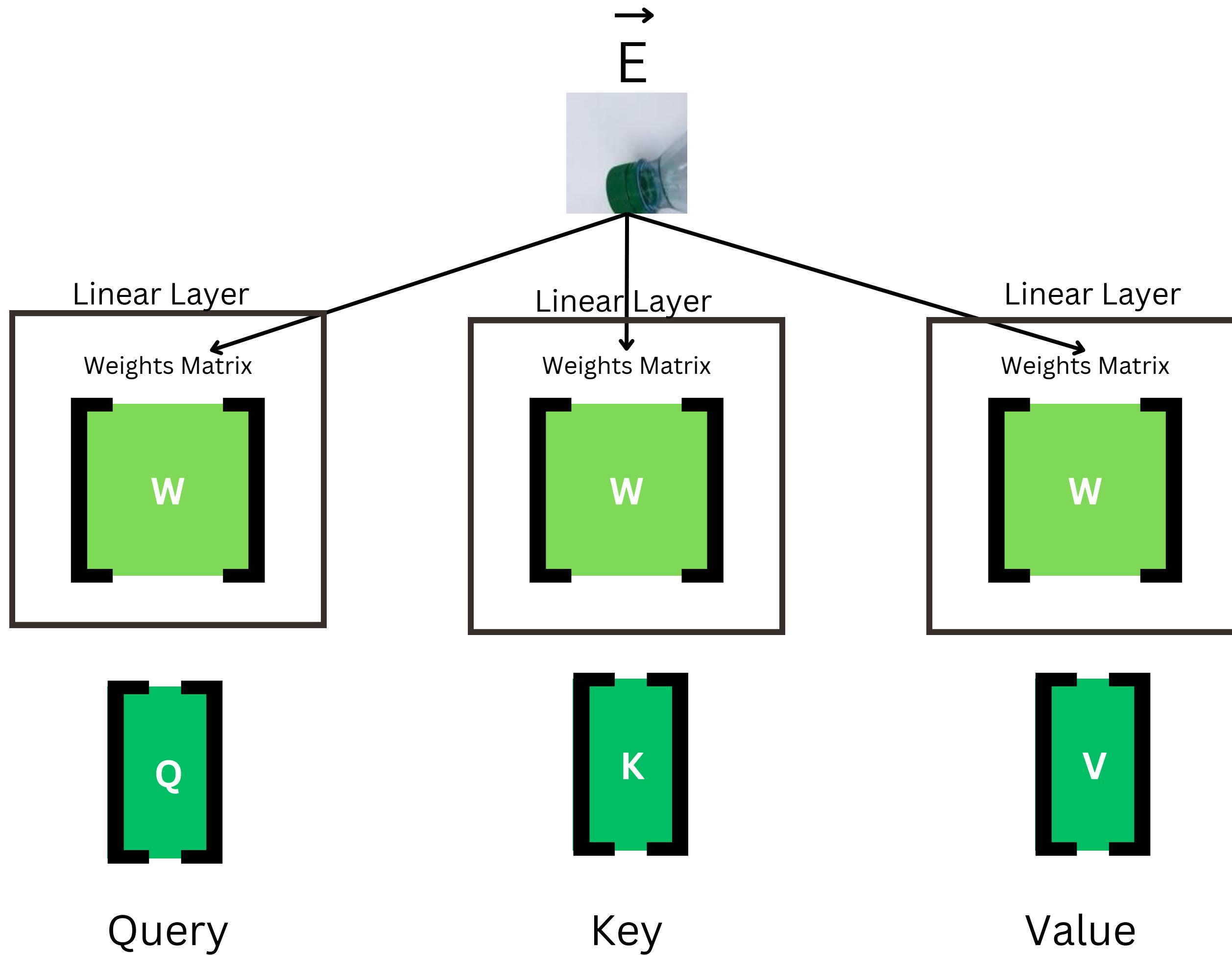
\rightarrow
E



Query Vector



=



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \longrightarrow \frac{Q \cdot K^T}{\sqrt{d_k}} = A$$

The diagram illustrates the computation of the attention matrix A . It shows three green rectangular matrices with black borders. The first matrix on the left is labeled Q and has a white 'Q' in its center. The second matrix in the middle is labeled K^T and has a white ' K^T ' in its center. The third matrix on the right is labeled A and has a white 'A' in its center. Between the first two matrices is a black asterisk (*) symbol, indicating multiplication. To the right of the second matrix is an equals sign (=). Below the first matrix is the word 'Query'. Below the second matrix is the word 'Key^T'. Below the third matrix is the word 'Attn'.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Sentence:

The fluffy rabbit
quickly hops past the
green field

Keys

Q

Queries

The fluffy rabbit quickly hops past the green field

K^*Q $\frac{1}{\sqrt{dk}}$	K^*Q $\frac{1}{\sqrt{dk}}$
...		+92.0						...
		-37.5						
...								...
K^*Q $\frac{1}{\sqrt{dk}}$...	K^*Q $\frac{1}{\sqrt{dk}}$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\boxed{\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V}$$

Image:

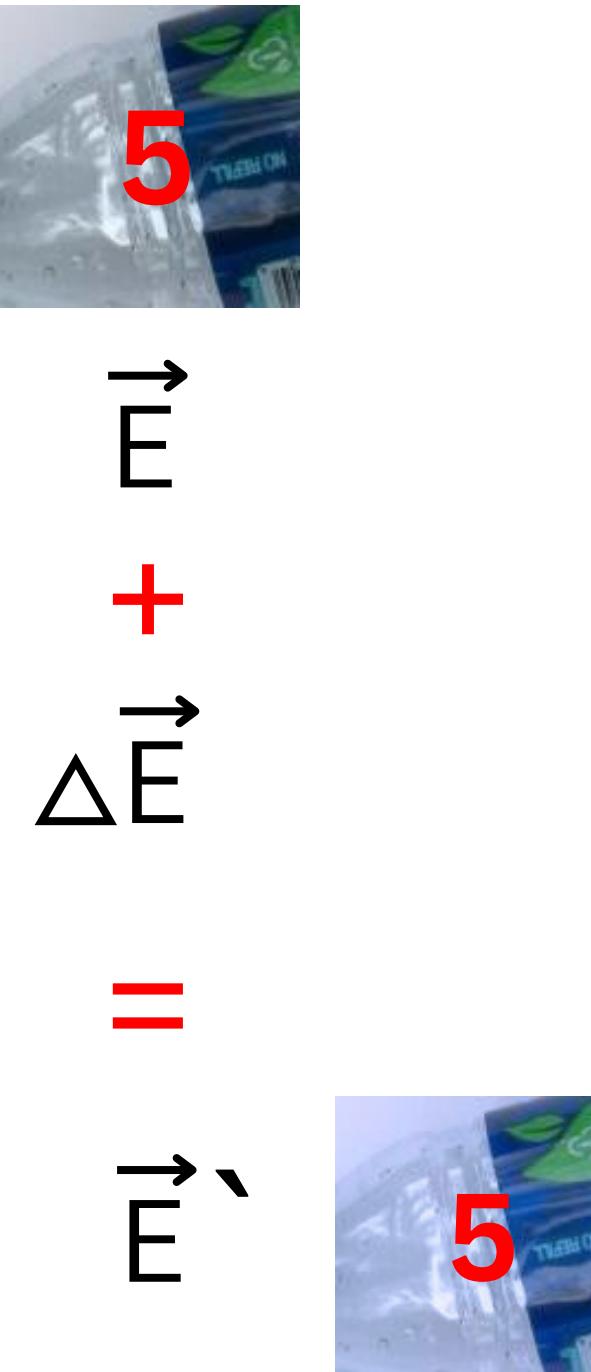
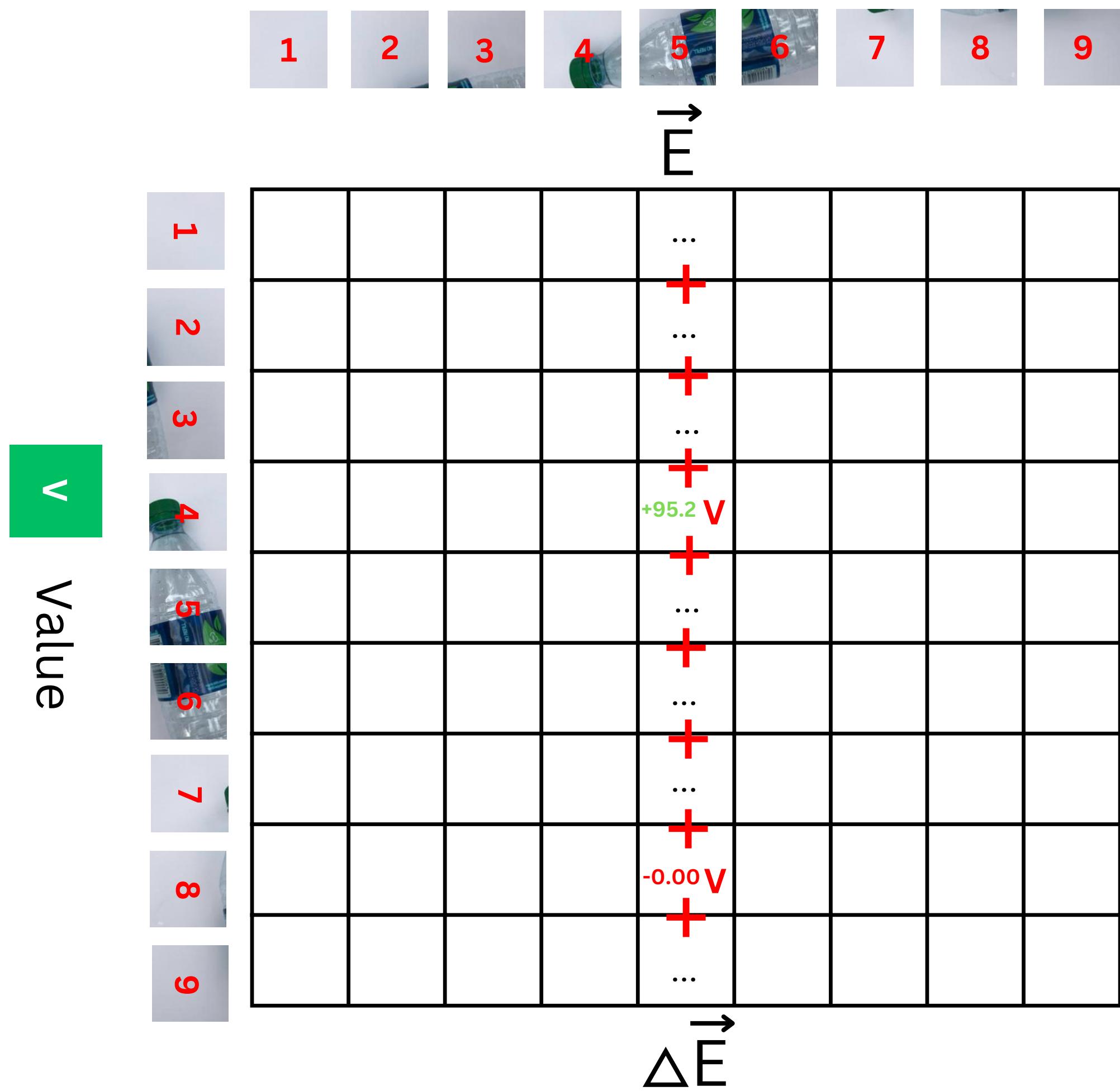


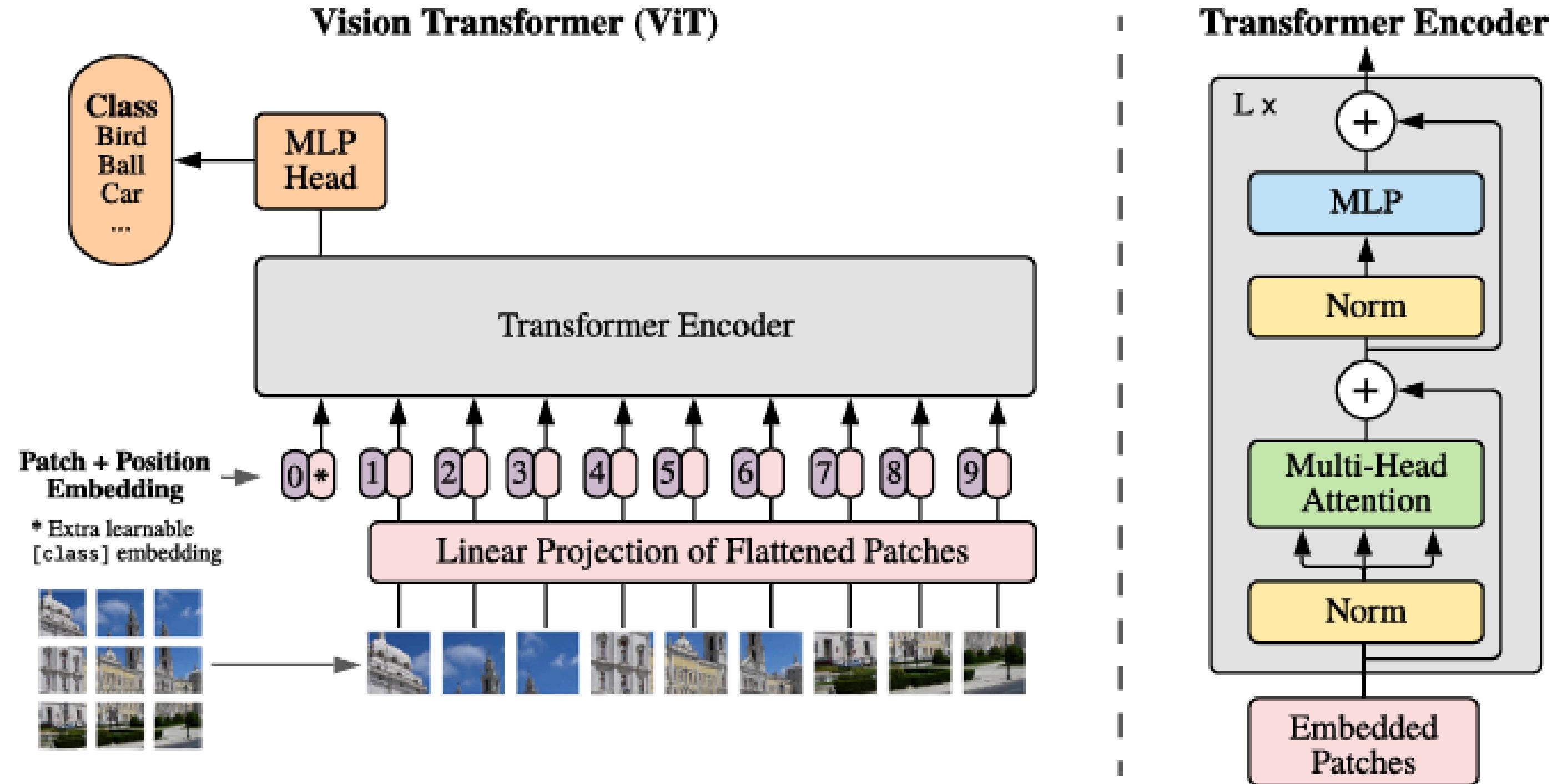
K
Keys^T



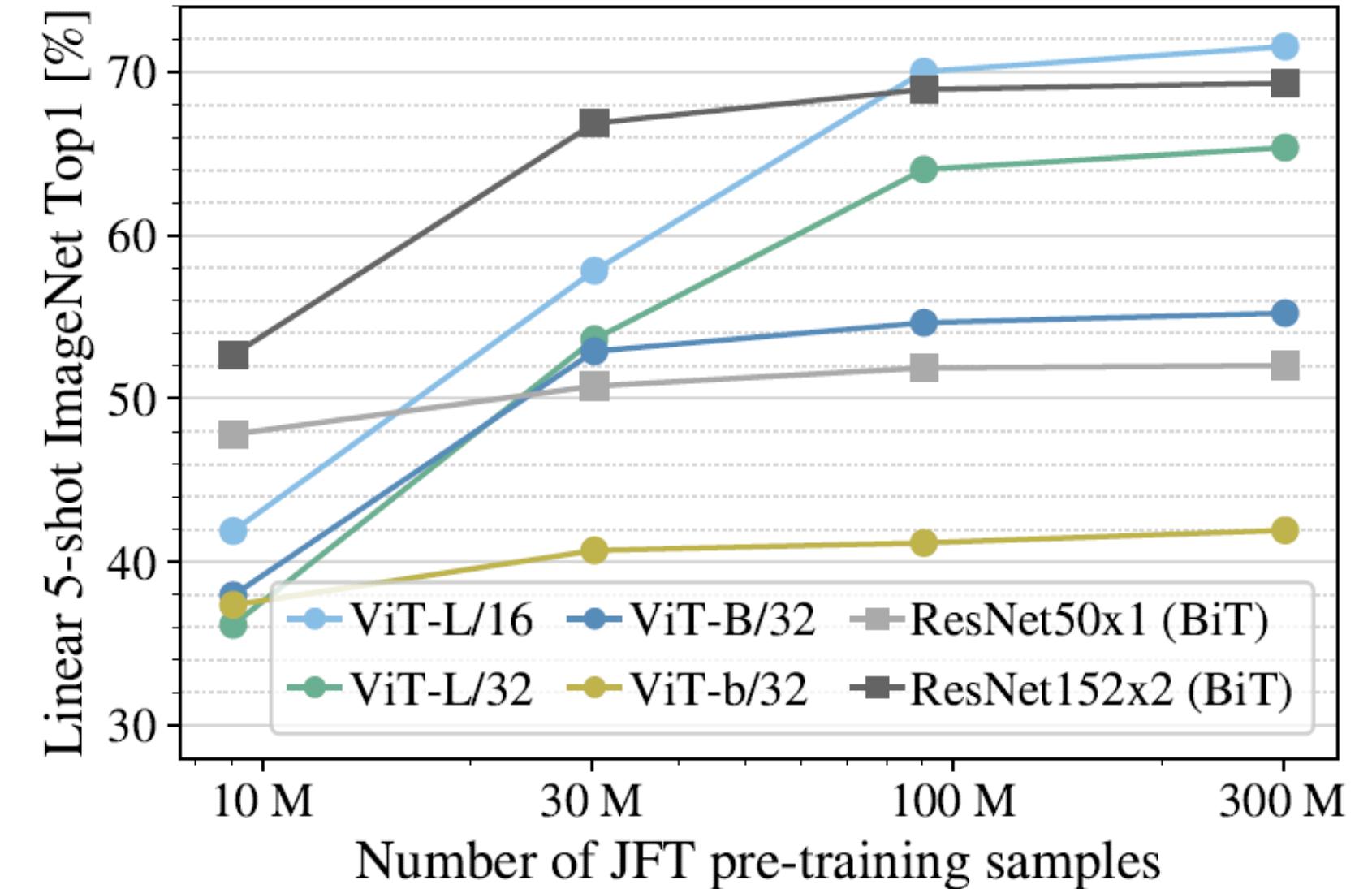
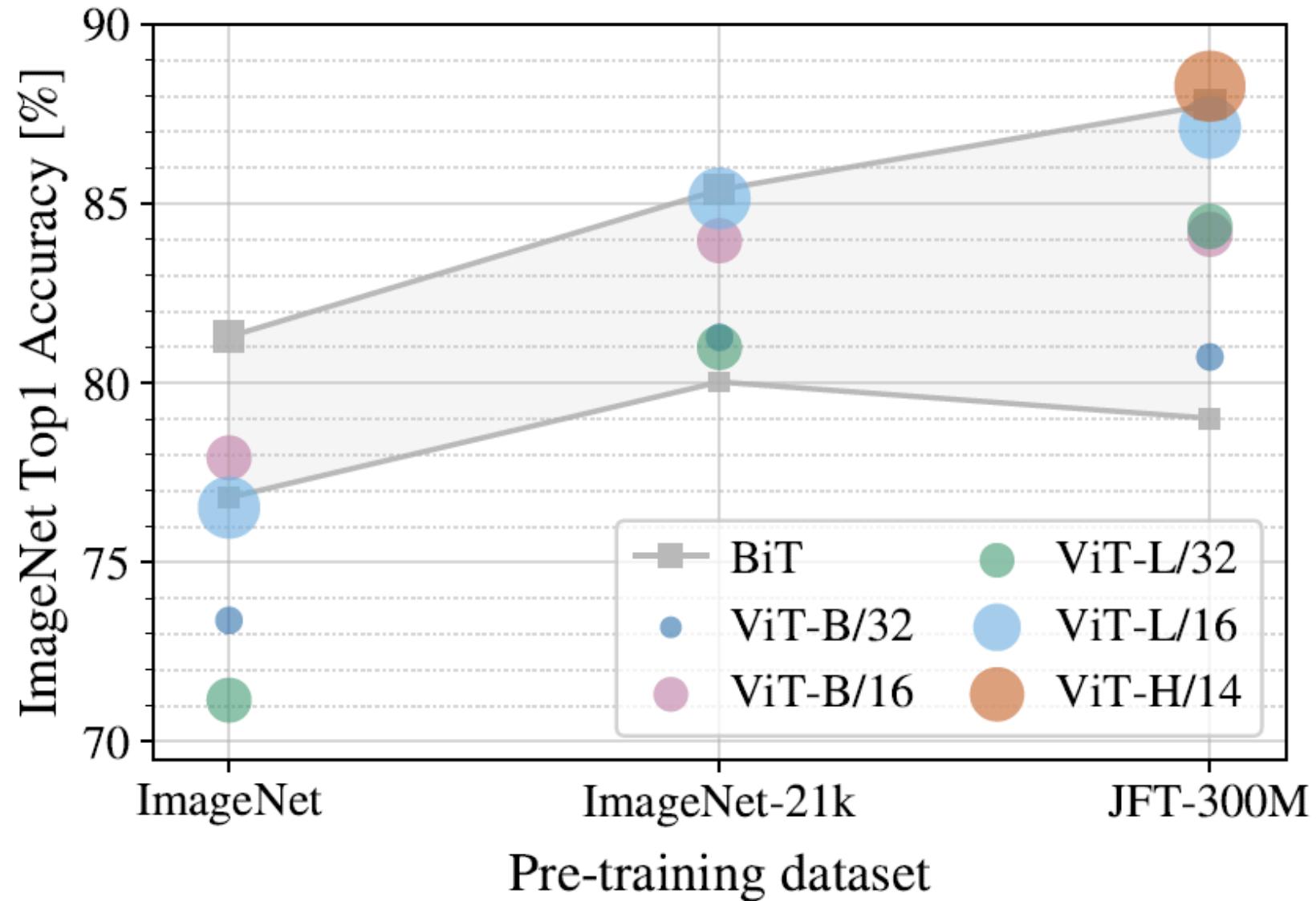
Q		Queries								
$\frac{K^*Q}{\sqrt{dk}}$...	SM	SM	SM	SM	SM	SM	SM	SM	$\frac{K^*Q}{\sqrt{dk}}$
1
2
3										
4										
5	+95.2									
6										
7										
8	-22.4									...
9										$\frac{K^*Q}{\sqrt{dk}}$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

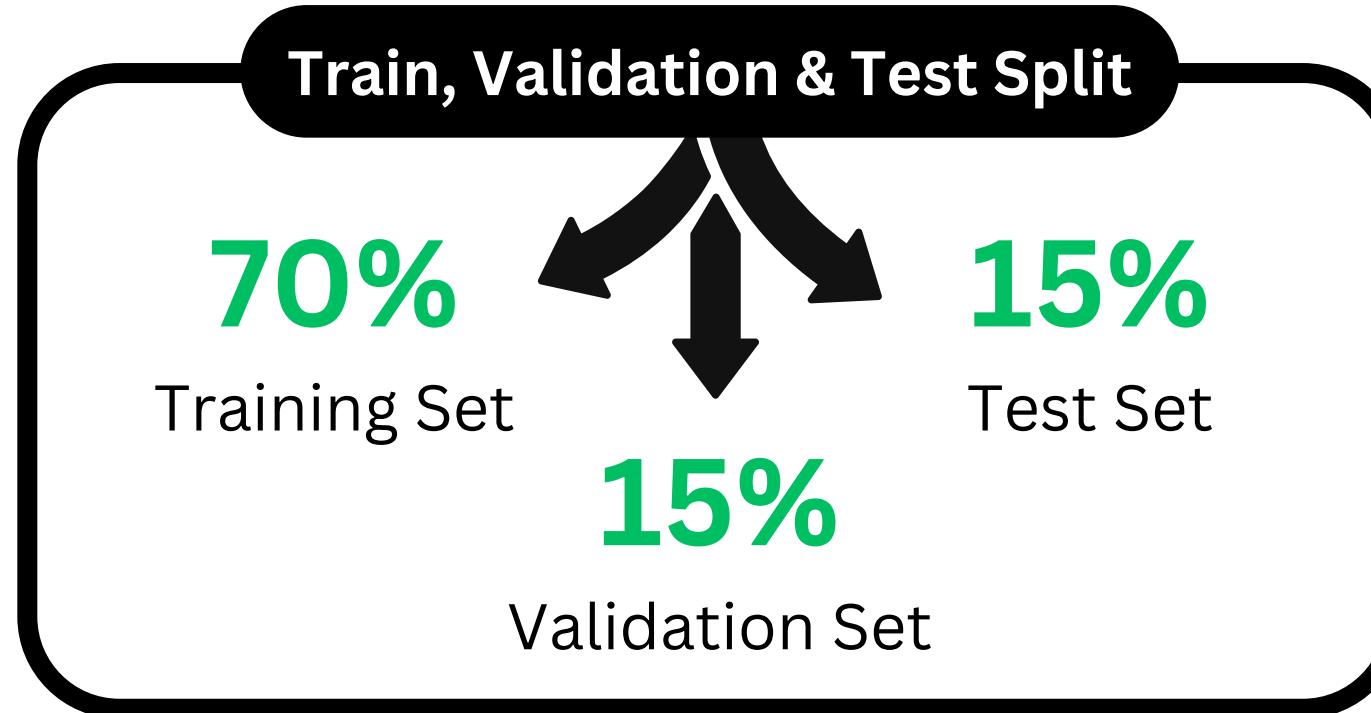




Model Performance



ViT Model Experimentation



Test Accuracy

97%

Computational Time

~1 Hrs
T4 GPU

Fine-Tuning

- Learning Rate = 0.0002
- Batch Size = 16
- Epochs = 4

Advantages

- Ease of Parallelization
- Scalability
- Transfer Learning Capabilities