

# Global Context Vision Transformers (GC ViT)

Sofiia Ilnytska, Wiebke Teetz

# Contents

## 1. Transformer

- a. History
- b. Input and Output
- c. Attention
- d. Encoder and Decoder

## 2. Vision Transformer

- a. Encoding
- b. Overall structure
- c. Limitations

## 3. Global Context Vision Transformer

- a. Overall structure
- b. Downsampling
- c. Global Context (Query Generator)

# **Understanding the Foundations of Transformers**

# What Came Before Transformers?

Introduction of Recurrent Neural Networks (RNNs) for sequential data processing in machine learning.

**1986**



**1997**

Long Short-Term Memory (LSTM) networks developed to address RNNs' vanishing gradient problem.

Gated Recurrent Units (GRUs) introduced, simplifying LSTM architecture while maintaining performance.

**2014**



**2015**

Recognition of RNNs' limitations in parallelization, prompting the need for more efficient models.

Emergence of Transformers, utilizing self-attention mechanisms to eliminate recurrence entirely.

**2017**



**2018**

Transformers revolutionize deep learning, enabling effective processing of long-range dependencies in data.



# What Are Transformers?

A **sequence transduction model**: It turns one sequence into another (e.g., a sentence in English → a sentence in French).

**Built entirely on Attention mechanisms**: No Recurrent Neural Networks (RNNs), no Convolutions.

A **Transformer** model is built from two main parts:

- **Encoder** (on the left side)
- **Decoder** (on the right side)

Both of them are made up of **layers** that are **stacked** on top of each other—imagine building blocks stacked into a tower. Each layer does two important things:

1. **Self-Attention**
2. **Point-wise, Fully Connected Layers**

This is a simple neural network that looks at each part (word, image patch, etc.) **on its own**, and processes it further.

"Point-wise" just means it works on each piece individually, without looking at its neighbors.

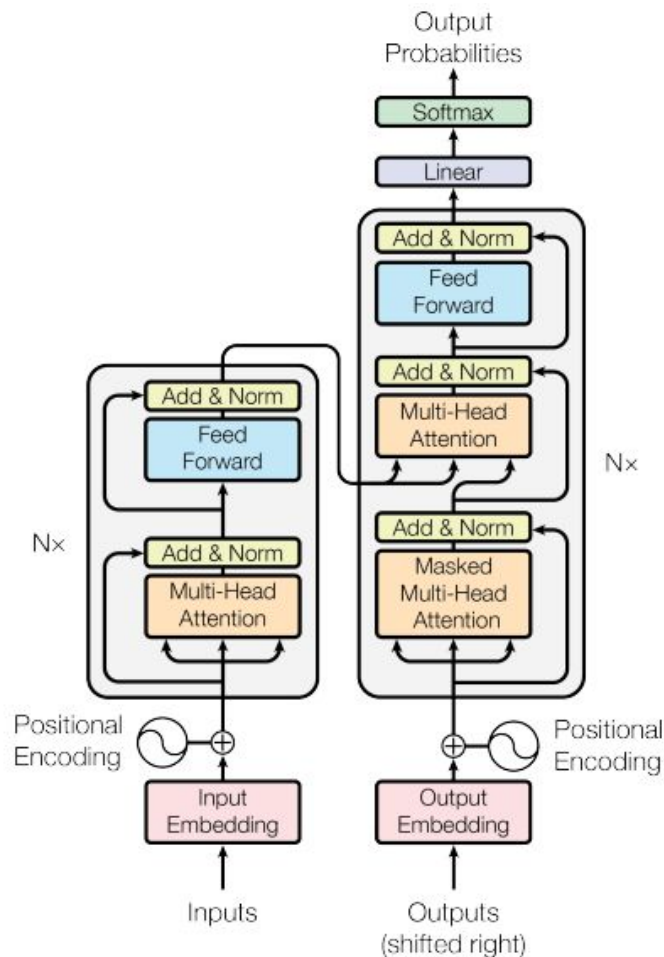


Figure 1: The Transformer - model architecture.

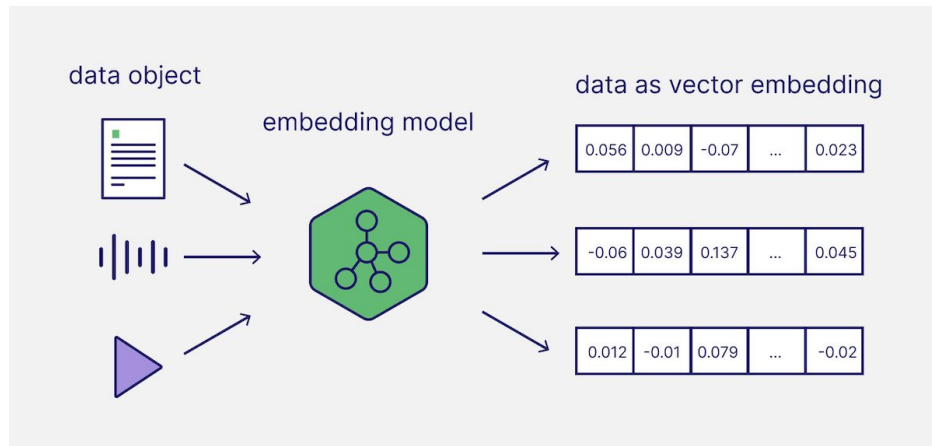
# Input and Output Embedding in Transformers

Transformers don't work directly with words — they work with **numbers**!

So, before processing, words are converted into vectors called embeddings.

These vectors capture:

- Semantic Meaning: Similar words get similar vectors.
- Contextual Clues: Embeddings help the model grasp relationships between words.



Input Embedding:

- Converts words into dense numerical vectors that capture their meaning.
- Adds Positional Encoding to show word order.

Output Embedding:

- Converts the final vectors back into words or tokens.
- Uses a softmax layer to pick the most likely word.

# Attention in transformers and its types

Attention is the mechanism that helps a model decide which parts of the input are important.  
It weighs the relevance of each word when making predictions.

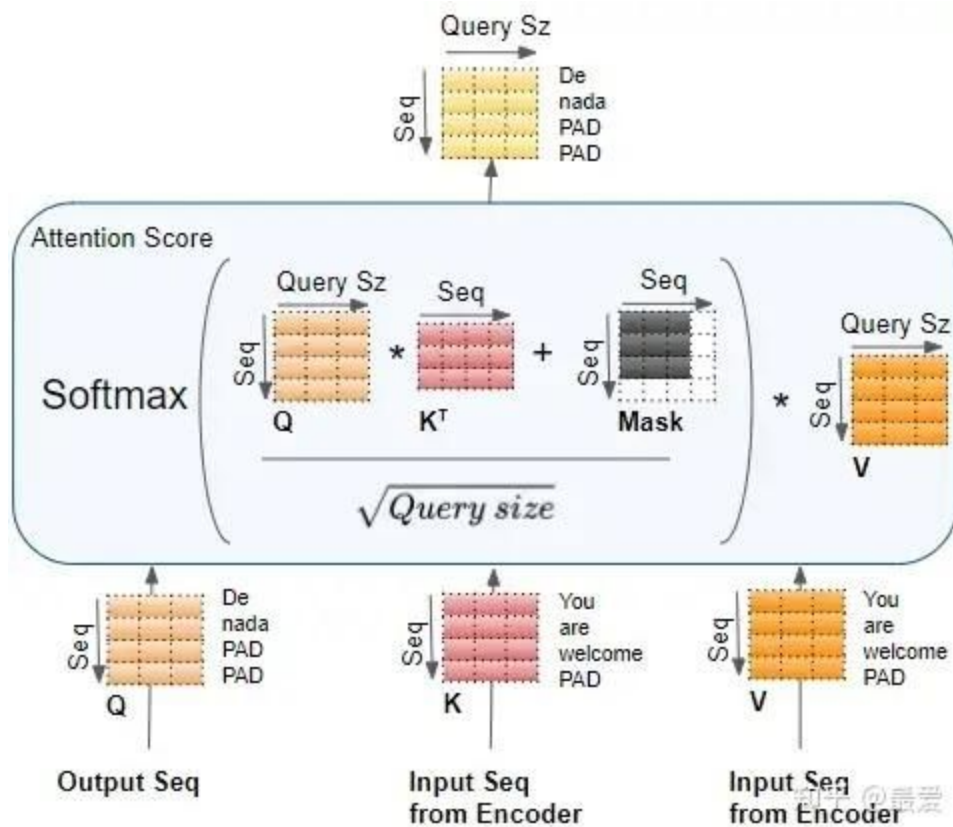
- Basic Attention
- Multi-Head Attention
- Masked Multi-Head Attention

Each word (or token) in the input is transformed into three vectors:

1. Query (Q): Represents what we're looking for.
2. Key (K): Represents what we have.
3. Value (V): Represents the actual information we use.

The model compares Queries with Keys to calculate attention scores, which determine how much Value to use from each word.

Multi-Head Attention helps the model look at different aspects of the input simultaneously, while Masked Attention makes sure the model doesn't cheat by looking ahead.





# Encoder in Transformers: How It Works + BERT

- 6 stacked layers
- two important components (sub-layers):
  - + Multi-Head Self-Attention
  - + Feed-Forward Network (FFN)
- Each layer also has:
  - + Residual Connections
  - + Layer Normalization

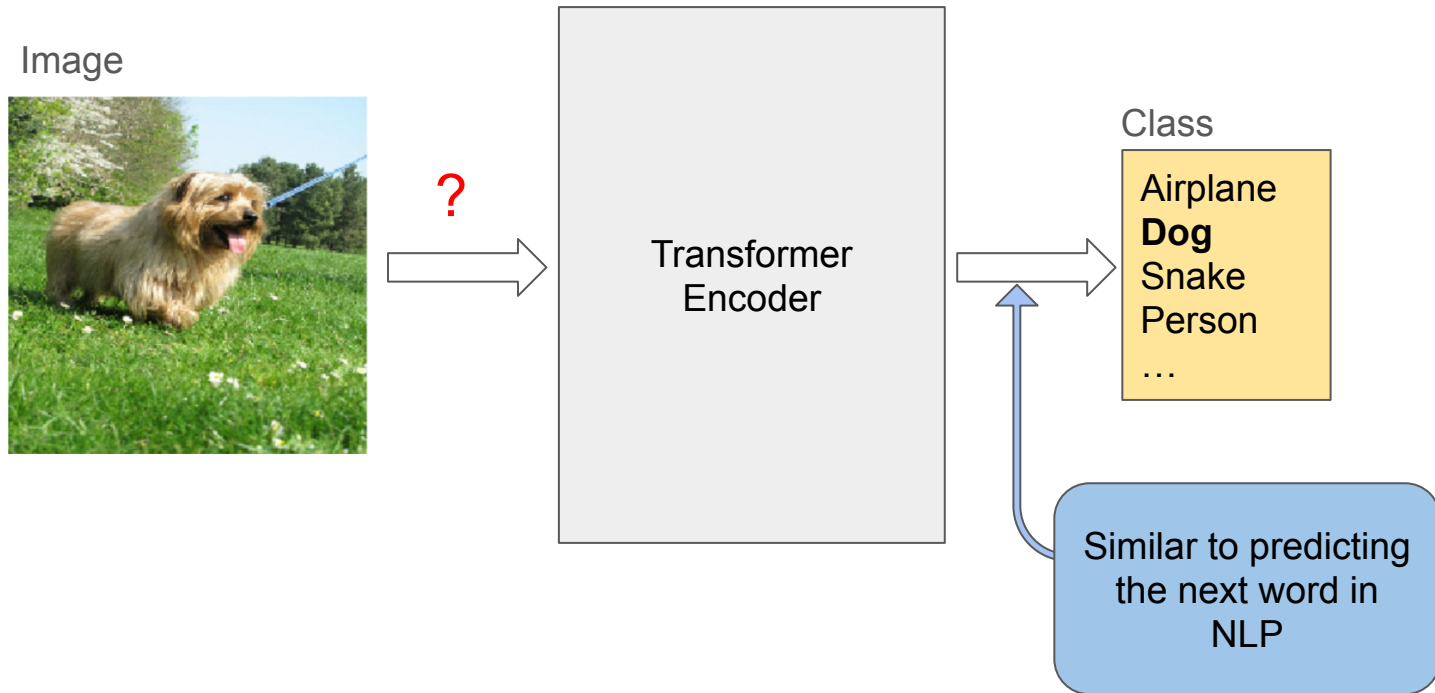
**BERT**, (Bidirectional Encoder Representations from Transformers) a popular language model, only uses the encoder because it focuses on **understanding text** rather than generating it.

# Decoder in Transformers: How It Works + GPT

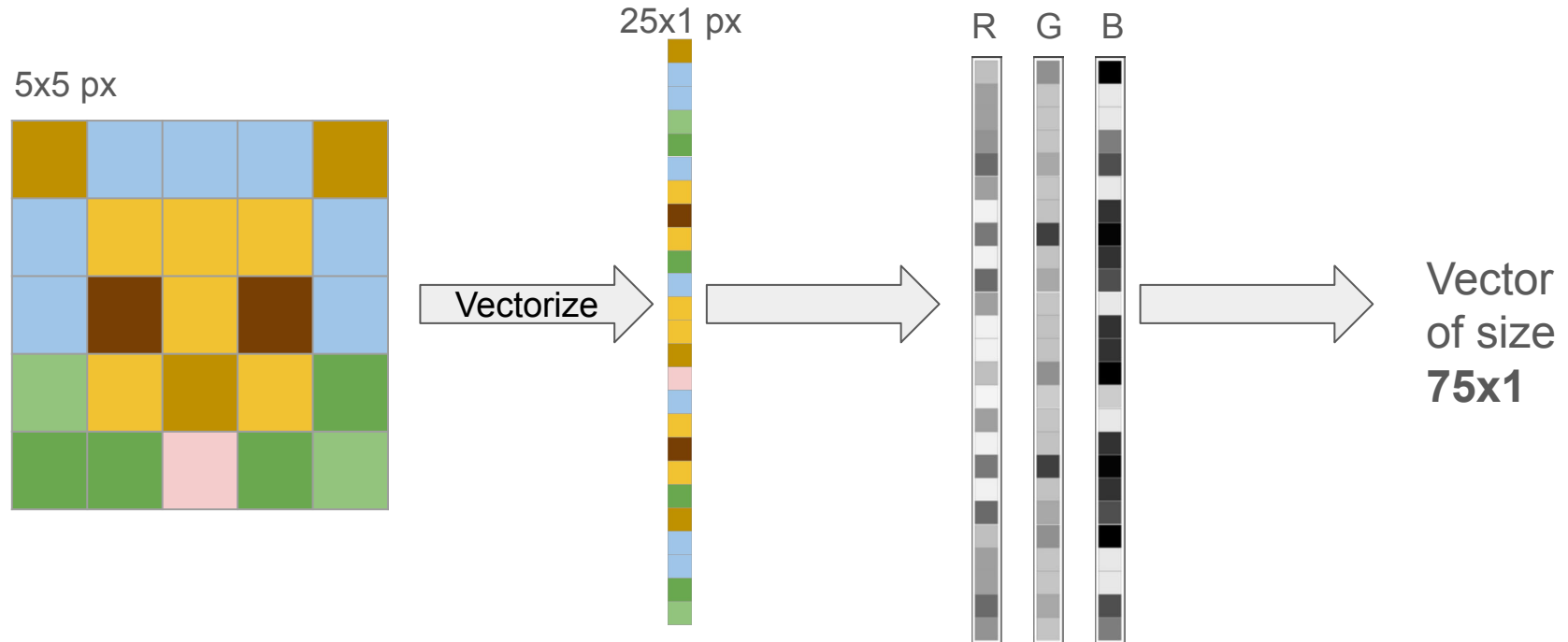
- 6 stacked layers
- two important components (sub-layers):
  - + Multi-Head Self-Attention
  - + Multi-Head Attention over Encoder Output
  - + Feed-Forward Network (FFN)
- Each layer also has:
  - + Residual Connections
  - + Layer Normalization

**GPT**, (generative pre-trained transformer) a popular language model, only uses the decoder because it focuses on **text generation** rather than understanding.

# Vision Transformers (ViT) for Image Classification



# How does Embedding an Image work?



$d1 \times d2 \times d3 \times d4$  tensor  $\rightarrow d1 * d2 * d3 * d4 \times 1$  tensor

# How does Embedding an Image work at scale?

- Self-Attention is **expensive!**

e.g. **224 x 224 px**  $\rightarrow 224 \times 224 \times 3 \times 1$  vector  $\rightarrow 150'528 \times 150'528$  Self-Attention Matrix  
 $= 150'528 \times 1 \approx$  **22.6 Billion**

- Transformers usually work with a **sequence of tokens** as input

# How does Embedding an Image work at scale?

Published as a conference paper at ICLR 2021

---

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

**Alexey Dosovitskiy<sup>\*,†</sup>, Lucas Beyer<sup>\*</sup>, Alexander Kolesnikov<sup>\*</sup>, Dirk Weissenborn<sup>\*</sup>,  
Xiaohua Zhai<sup>\*</sup>, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby<sup>\*,†</sup>**

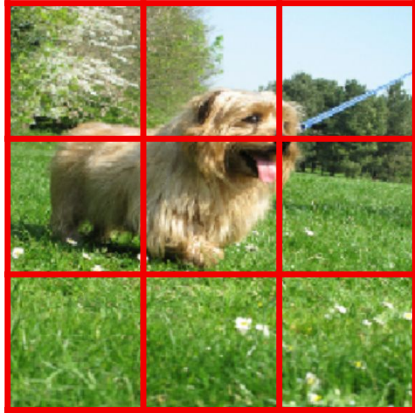
<sup>\*</sup>equal technical contribution, <sup>†</sup>equal advising

Google Research, Brain Team

{adosovitskiy, neilhoulby}@google.com

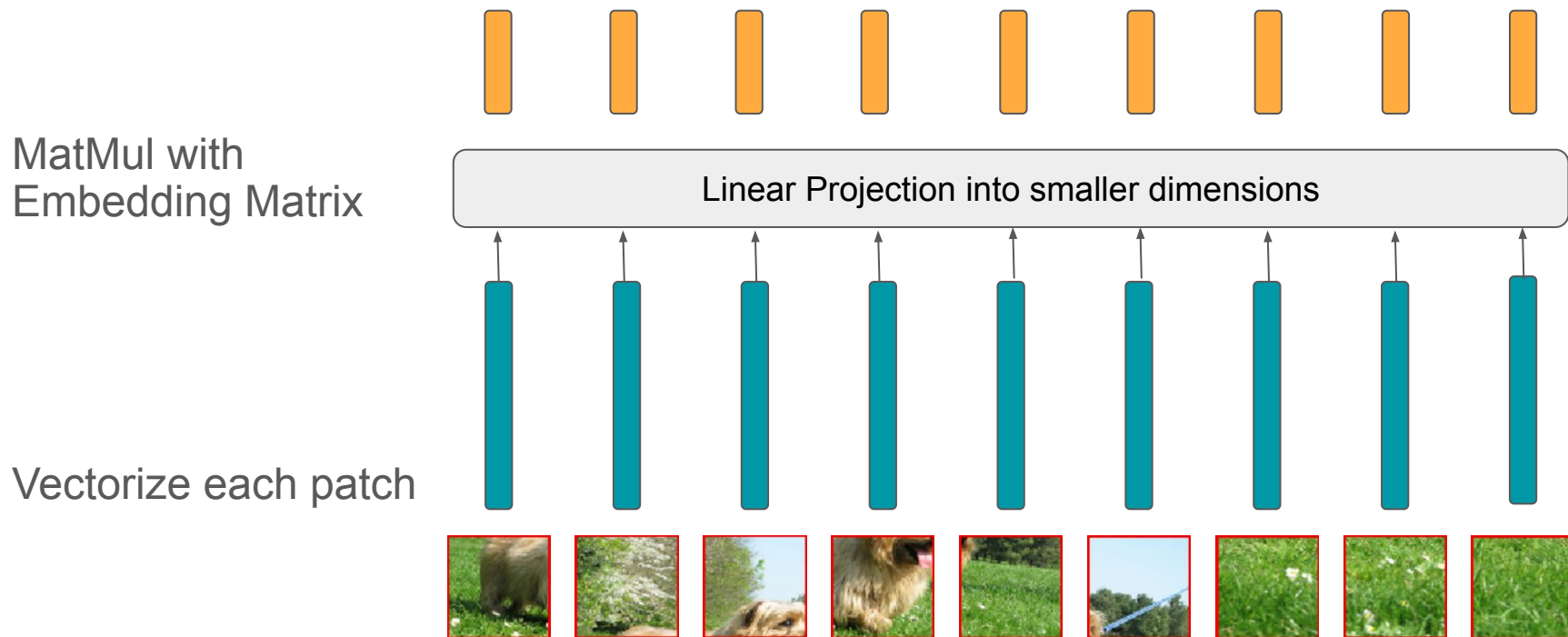
# How does Embedding an Image work at scale?

Split the Input into **patches** of size 16 x 16 px



1 Patch (16x16) = 1 Token ( = 1 Word )

# How does Embedding an Image work at scale?

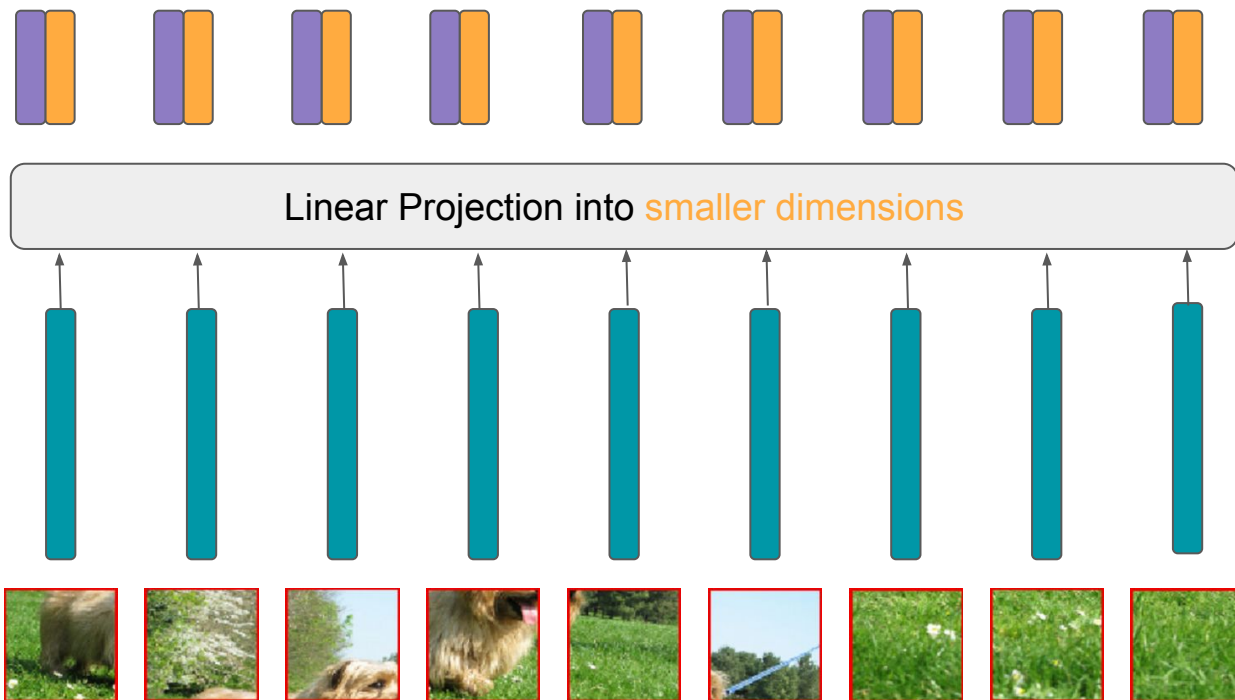




# How does Embedding an Image work at scale?

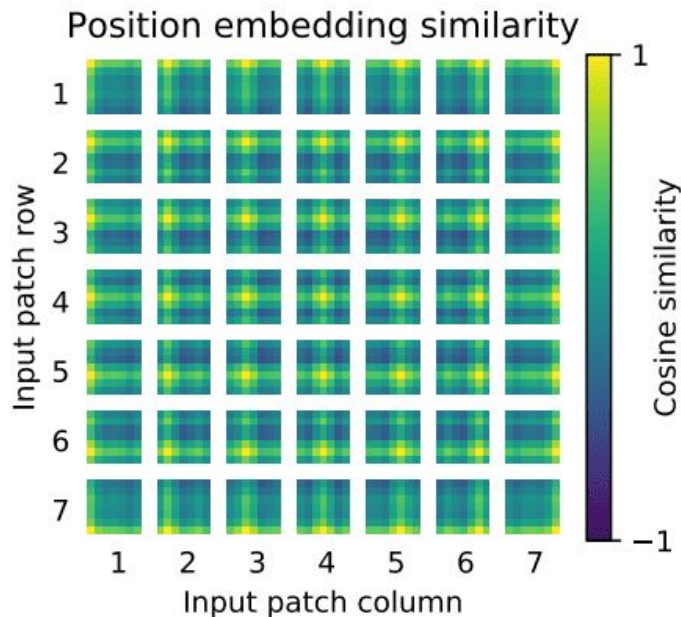
Add **Positional Embeddings**  
(trainable)

Vectorize each patch

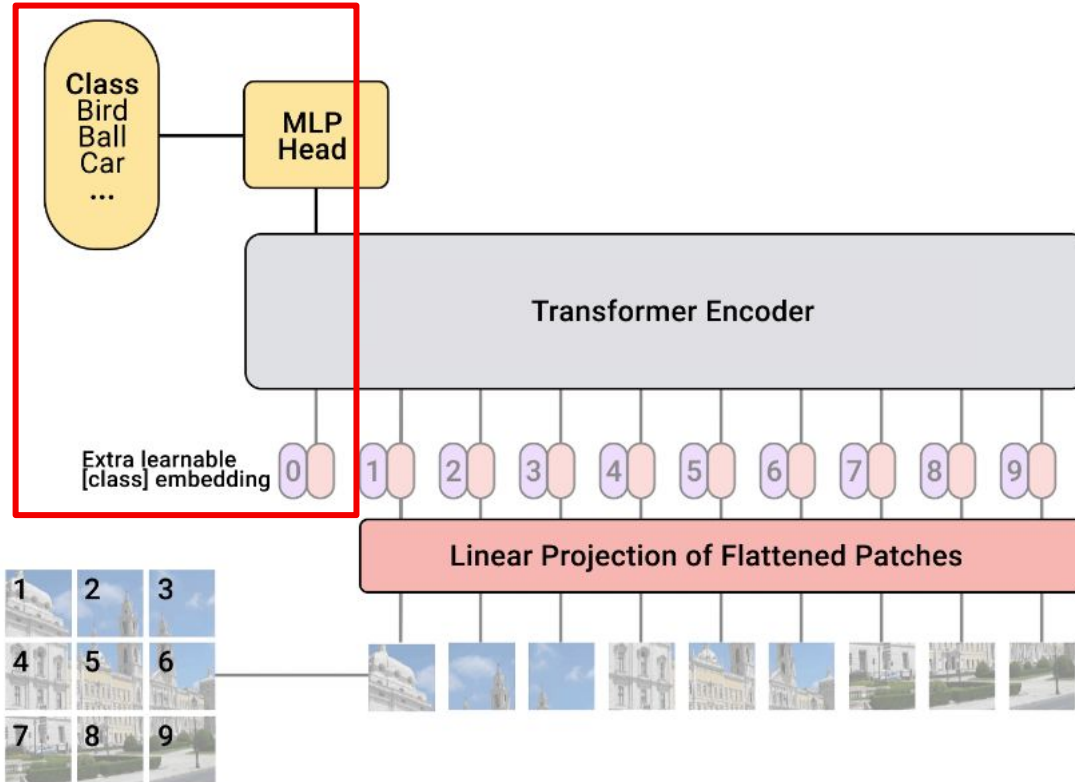


# Positional Embeddings

The transformer has no intrinsic understanding of where each patch is positioned in the original image



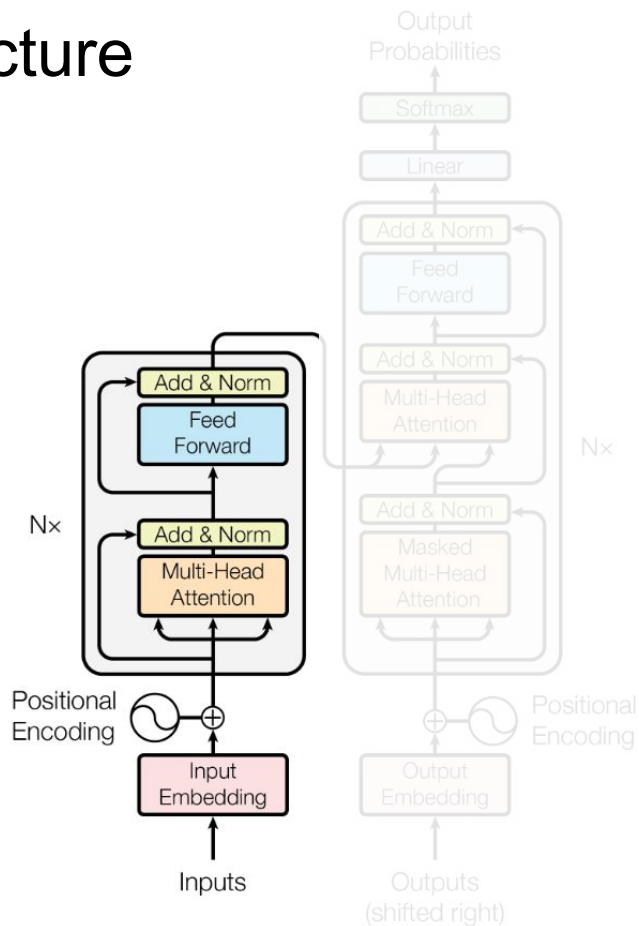
# CLS Token



# Overview of ViT Architecture

Self Attention on each individual patch

Combine infos of the patches



Original Transformer Architecture

# Limitations of ViT

- relative computational inefficiency (scales quadratically with image size)
- ViT lack the inductive bias of convolutional neural networks (CNN)
  - locally restricted inductive fields
  - translation invariance

ViTs treat **all patches equally**

→lack of local feature emphasis (e.g. edges, texture)

→large-scale training data needed to outperform older models (e.g. CNNs)

# What is convolution?

Source layer

5	2	6	8	2	0	1	2
4	3	4	5	1	9	6	3
3	9	2	4	7	7	6	9
1	3	4	6	8	2	2	1
8	4	6	2	3	1	8	8
5	8	9	0	1	0	2	3
9	2	6	6	3	6	2	1
9	8	8	2	6	3	4	5

Convolutional  
kernel

-1	0	1
2	1	2
1	-2	0

Destination layer

	5						

$$\begin{aligned} &(-1 \times 5) + (0 \times 2) + (1 \times 6) + \\ &(2 \times 4) + (1 \times 3) + (2 \times 4) + \\ &(1 \times 3) + (-2 \times 9) + (0 \times 2) = 5 \end{aligned}$$

# Limitations of ViT

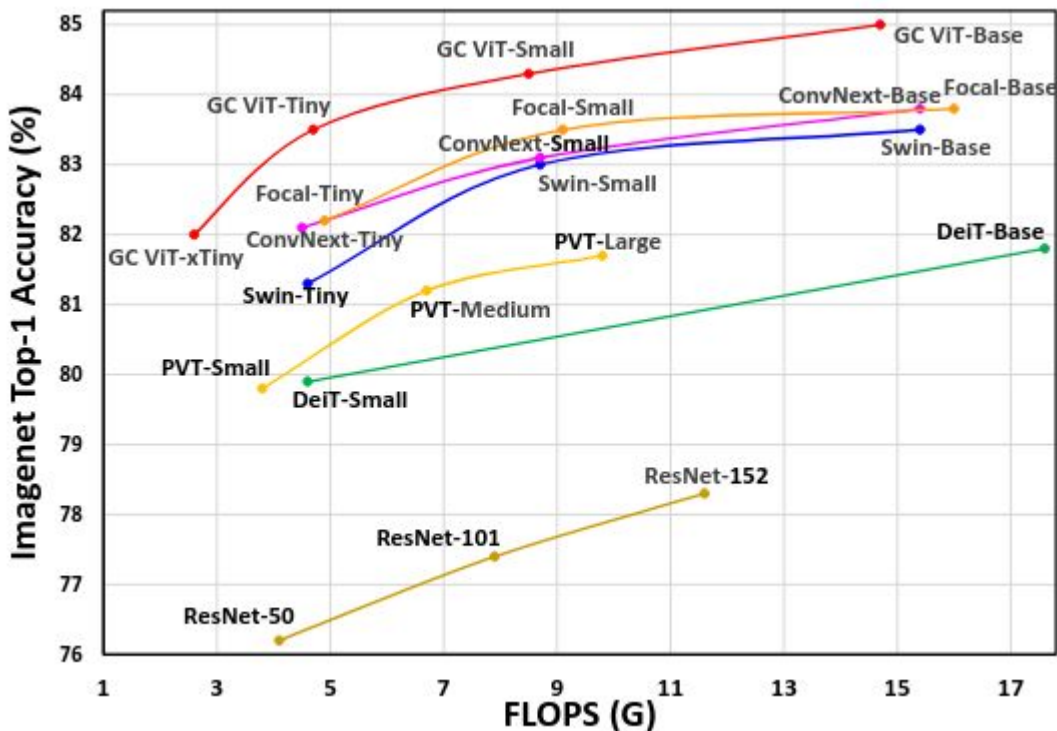
- relative computational inefficiency (scales quadratically with image size)
- ViT lack the inductive bias of convolutional neural networks (CNN)
  - locally restricted inductive fields
  - translation invariance

ViTs treat **all patches equally**

→lack of local feature emphasis (e.g. edges, texture)

→large-scale training data needed to outperform older models (e.g. CNNs)

# Global Context Visual Transformer (GC ViT)

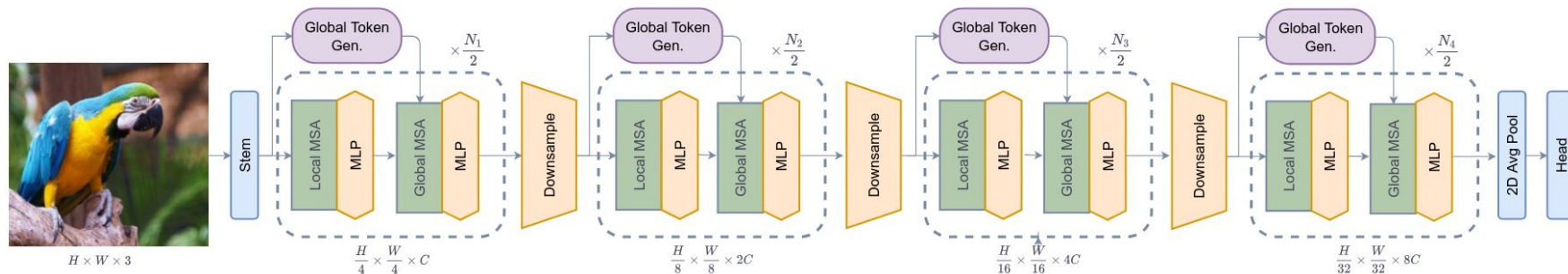




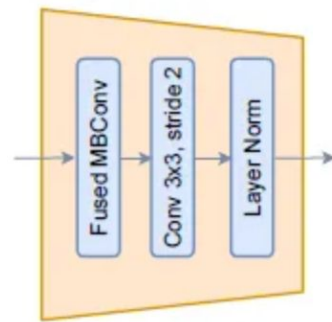
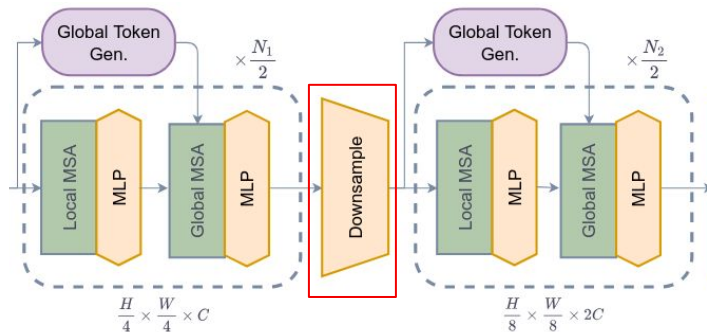
# Global Context Visual Transformer (GC ViT)

## Main Adjustments:

- optimized number of parameters
- CNN like token generator for global queries
- downsampling module that integrates inductive bias



# GC ViT - Downsampler

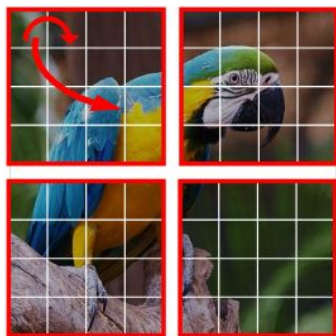


$$\begin{aligned}\hat{x} &= \text{DW-Conv}_{3 \times 3}(x), \\ \hat{x} &= \text{GELU}(\hat{x}), \\ \hat{x} &= \text{SE}(\hat{x}), \\ x &= \text{Conv}_{1 \times 1}(\hat{x}) + x,\end{aligned}$$

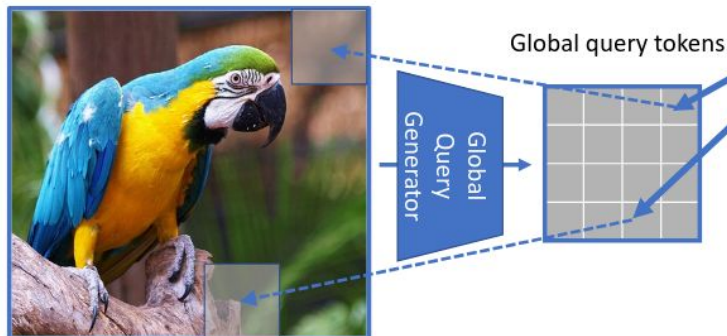
- applies convolution to efficiently capture local features
  - enforcing inductive biases
- reduces spatial resolution

# GC ViT - Global and Local Attention

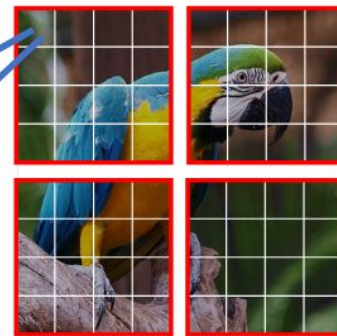
Local attention



Global Attention

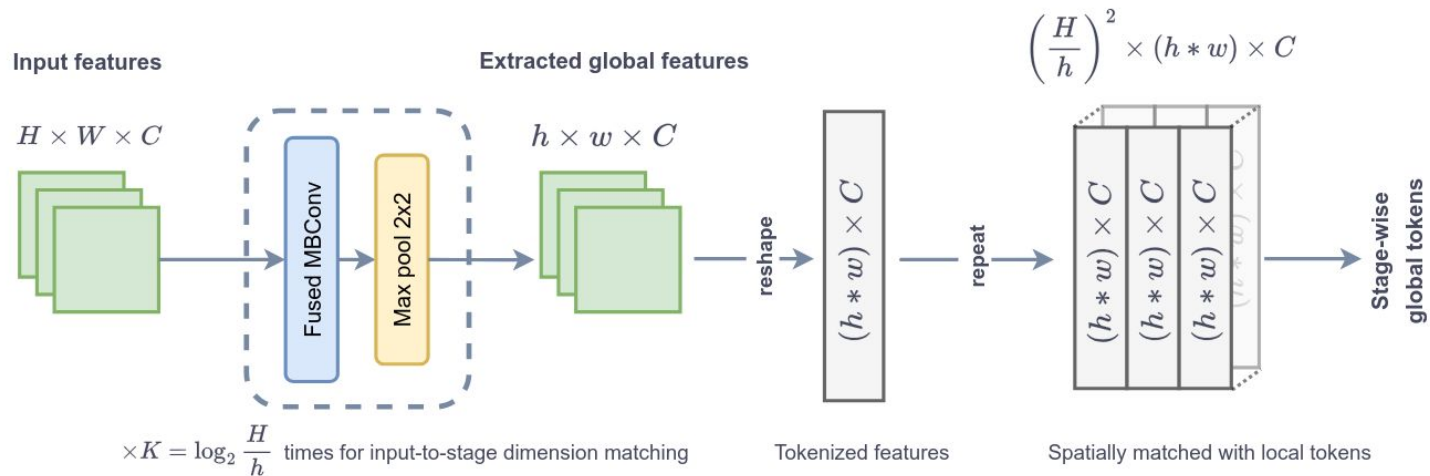


Local attention with global query

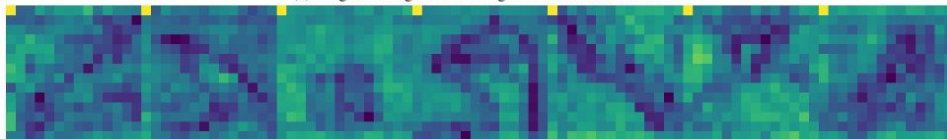


Computes global attention by multiplying the **global query**  
with each **local key** matrix

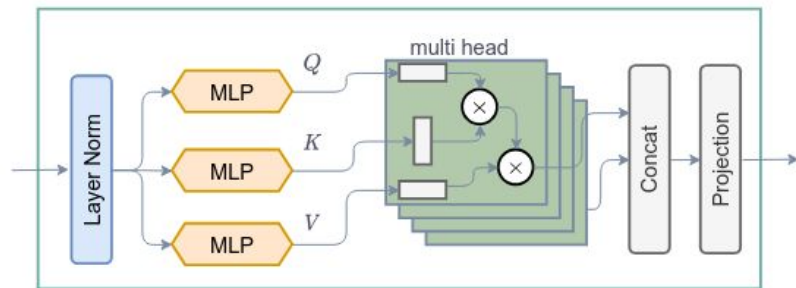
# GC ViT - Global Query Generator



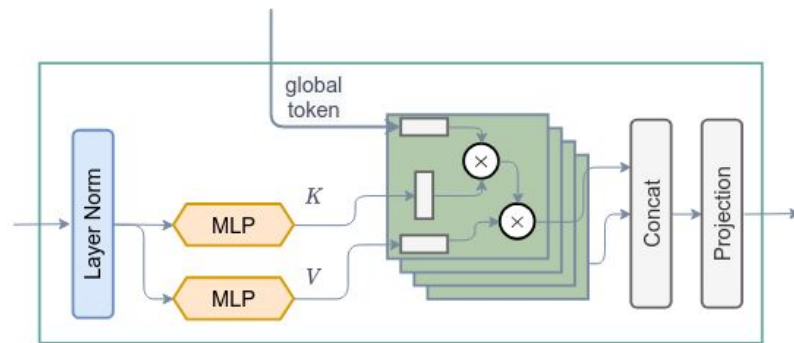
(a) Original images from ImageNet-1K validation set.



# GC ViT - Local and Global Self-Attention



Local MSA



Global MSA



(a) Original images from ImageNet-1K validation set.



(b) **Global attention** maps from GC ViT model (ours).

# Thank you

Any questions?