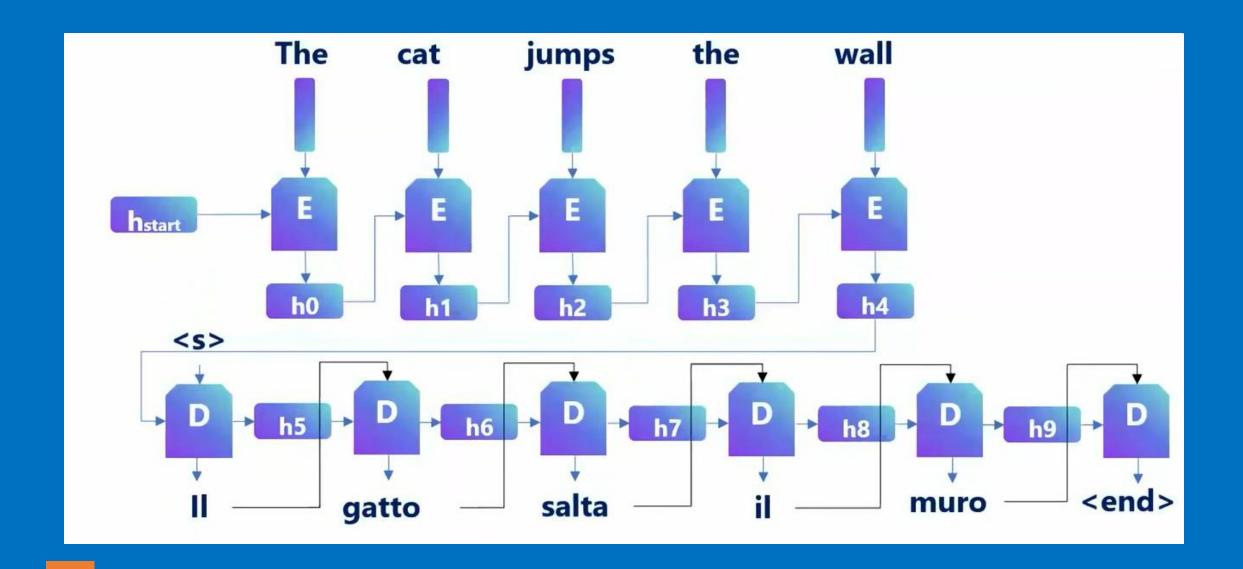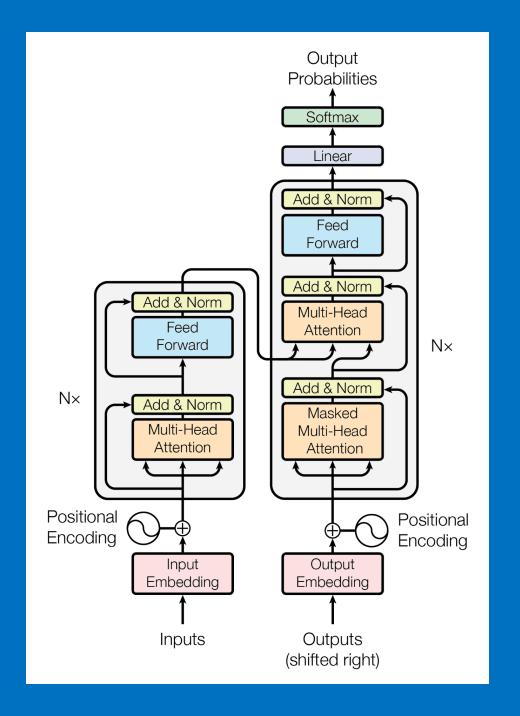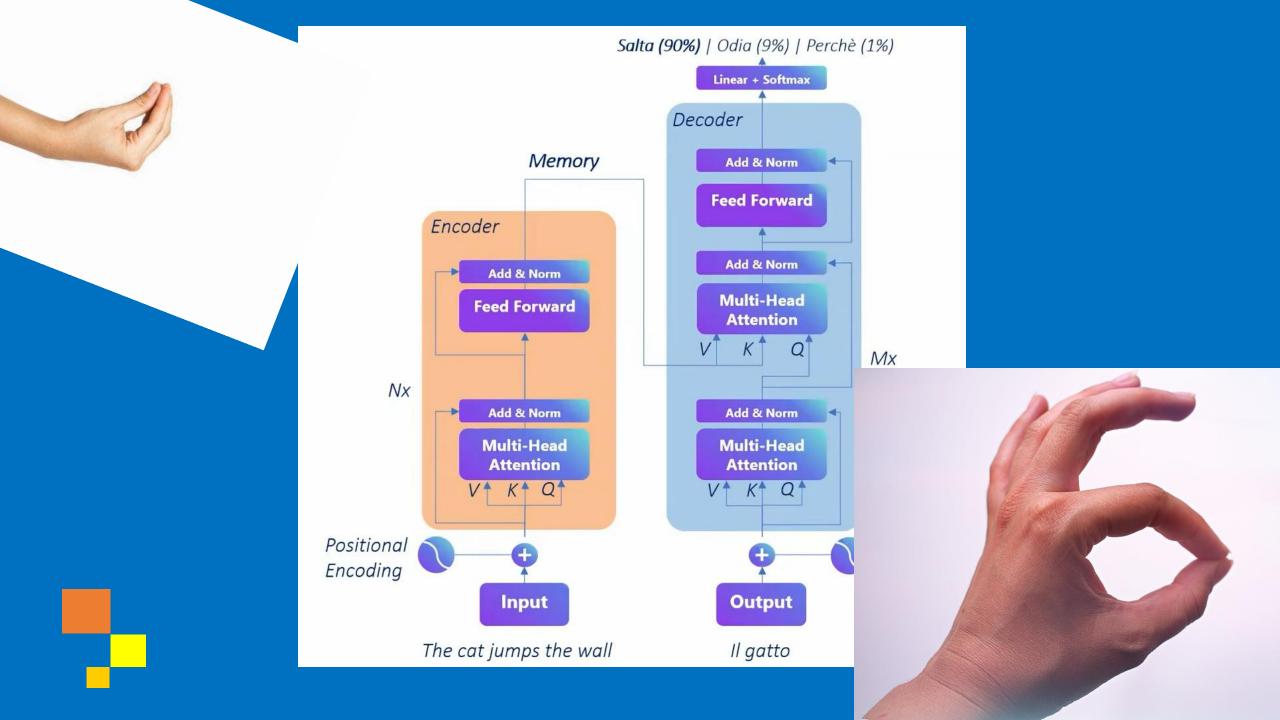# Computer Vision

_____

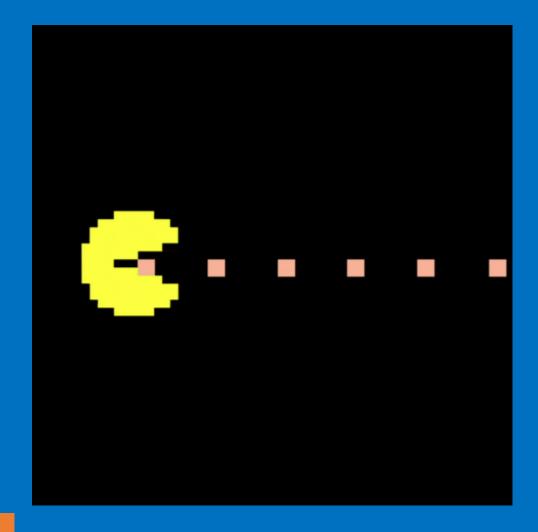by Piotr Krzemiński

# Agenda for Today:

- RNNs – the heart of all transformers
- A brief summary on Transformers
- Embedding layer, aka convert indices into dimensions
- Attention layer, single or multi-head
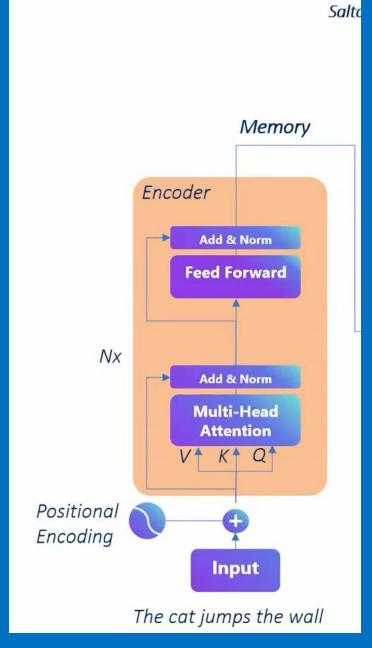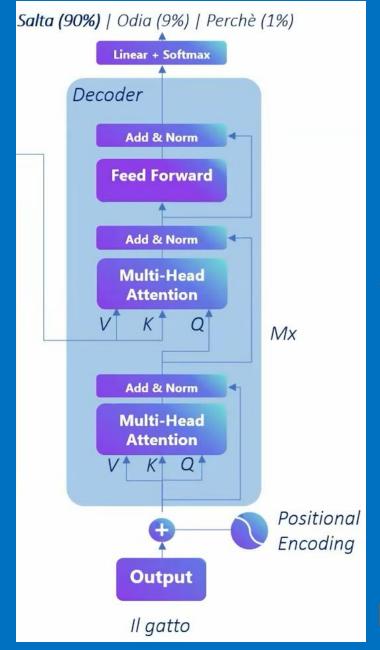- ViT - Transformers in Computer Vision
- DIY – How to code your own ViT
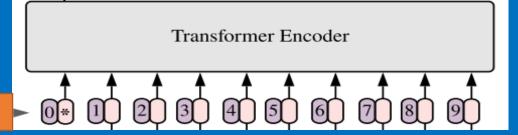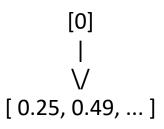- Should I code transformers?

Who has seen this picture?

# Positional Encoding

# Attention

# Attention... SIMPLIFIED!

# Attention... SIMPLIFIED!

## Attention

I **gave** my **dog** Charlie some **food**

**Ho dato** da **mangiare** al mio **cane** Charlie

- Source ≠ Target
  - Queries from **Source**
  - Key, Values from **Target**

**Multi-Head Attention**

$V$  $K$  $Q$

Target      Source

- Captures **inter-**sequence dependencies

## Self-Attention

I **gave** my dog Charlie some food

Who?        To whom?        What?

- Source = Target
  - Key, Queries, Values obtained from the same sentence

**Multi-Head Attention**

$V$  $K$  $Q$

- Captures **intra-**sequence dependencies

# Self - Attention

# Multi – Head Attention

# But I Can't Code Complex Stuffs Like These ☹

```python
tf.keras.layers.MultiHeadAttention(
    num_heads,
    key_dim,
    value_dim=None,
    dropout=0.0,
    use_bias=True,
    output_shape=None,
    attention_axes=None,
    kernel_initializer='glorot_uniform',
    bias_initializer='zeros',
    kernel_regularizer=None,
    bias_regularizer=None,
    activity_regularizer=None,
    kernel_constraint=None,
    bias_constraint=None,
    **kwargs
)
```
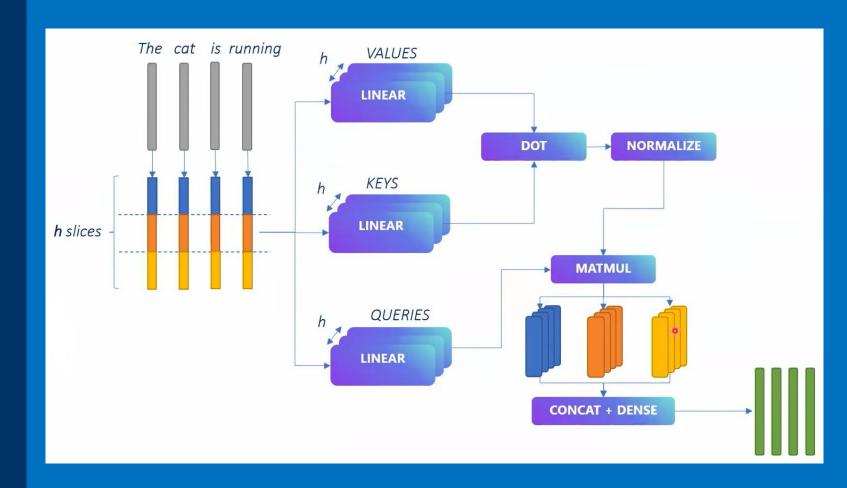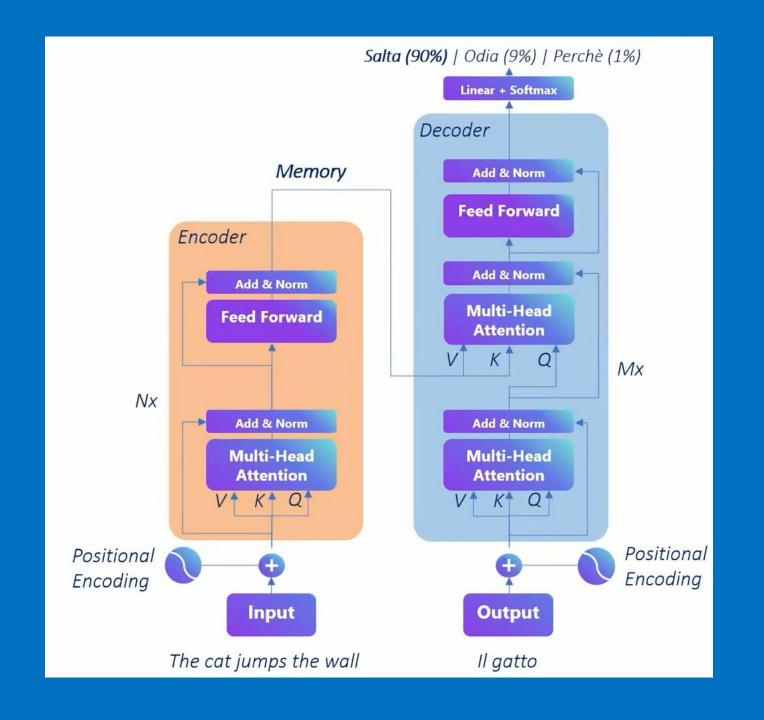
```python
torch.nn.MultiheadAttention(embed_dim, num_heads, dropout=0.0, bias=True,
add_bias_kv=False, add_zero_attn=False, kdim=None, vdim=None,
batch_first=False, device=None, dtype=None)  [SOURCE]
```
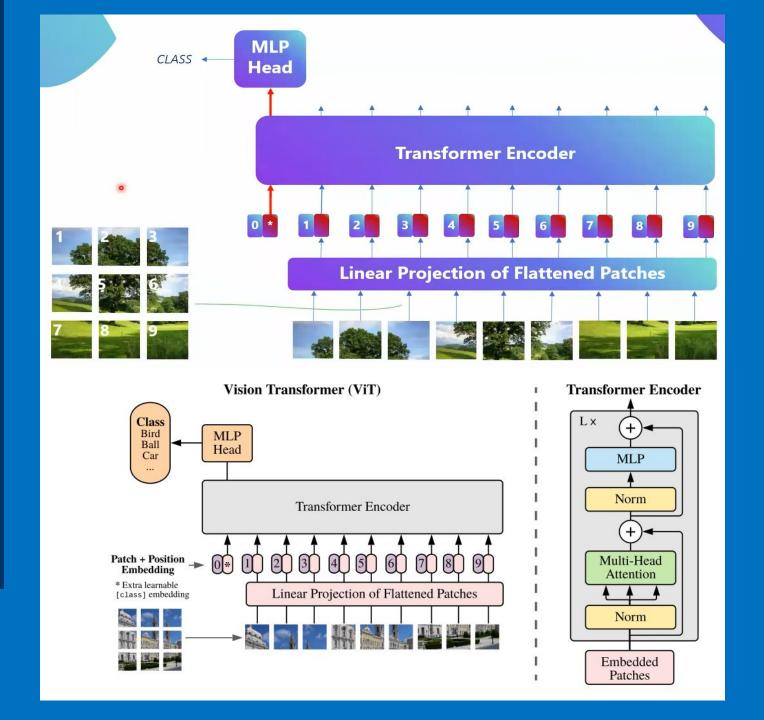
Allows the model to jointly attend to information from different representation subspaces as described in the paper: Attention Is All You Need.
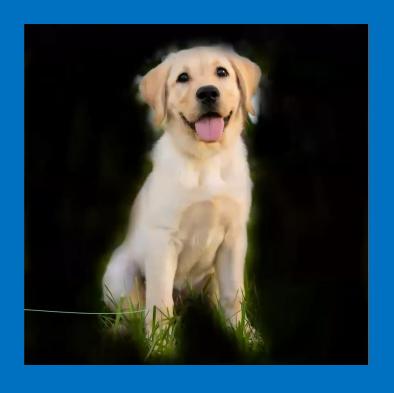
Multi-Head Attention is defined as:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(head_1, \ldots, head_h)W^O$$

where $head_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$.

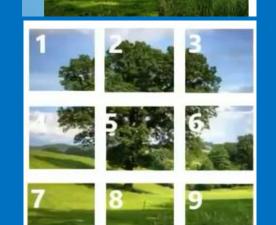# Transformers in Computer Vision

# Attention for Images

# Positional Encoding & Improving Attention Performance

Every Pixel

256^2 relations = (256^2)^2 operations

,Square' Attention

are_side//2)^2 – still a lot

Chunks

Best performance, but reduced Attention benefits



256px

256px

256px

image patch
1 layer

hidden layer 1
4 feature maps

hidden layer 2
8 feature maps

final layer
4 class units

36x36

28x28

14x14

10x10

5x5

convolution
(kernel: 9x9x1)

max
pooling

convolution
(kernel: 5x5x4)

max
pooling

convolution
(kernel: 5x5x8)

1  2  3
4  5  6
7  8  9

# DIY
# Time to
# Code It!

Dataset:
CIFAR100



ViT Model from
„Transformers for
Image
Recognition at
Scale" paper

# Should You Try It?



- Basis of ChatGPT – „GPT" is „Generative Pre-Trained Transformer"

- ChatGPT Chief Scientist Salary - ~2 000 000 $/yr

- ViTs are expected to be the next „big boom"

# References:

- Attention is All You Need
  https://arxiv.org/pdf/1706.03762.pdf

- An Image Is Worth 16x16 Words: Transformers For Image Rrecognition At Scale
  https://arxiv.org/pdf/2010.11929v2.pdf