# Image Classification

Vishal Reddy Mandadi
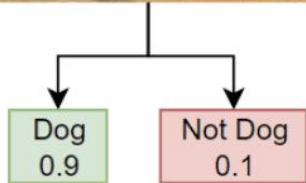
# Image Classification
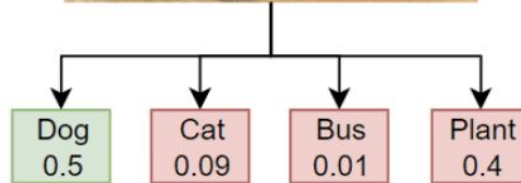


convolution + nonlinearity    max pooling    vec    bird $p_{bird}$    sunset $p_{sunset}$    dog $p_{dog}$    cat $p_{cat}$    ...

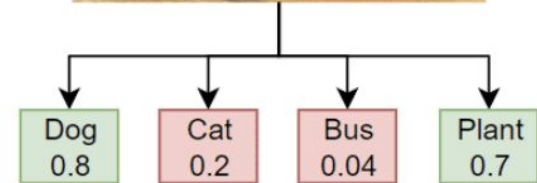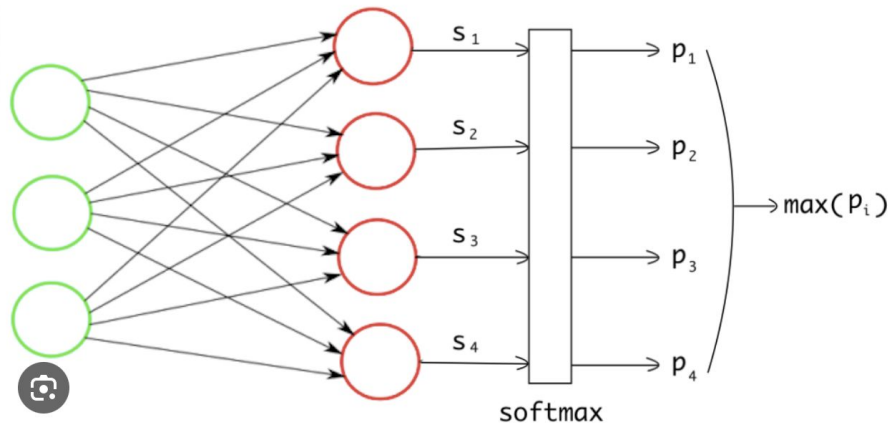# Binary vs Multi-class vs Multi-label Classification
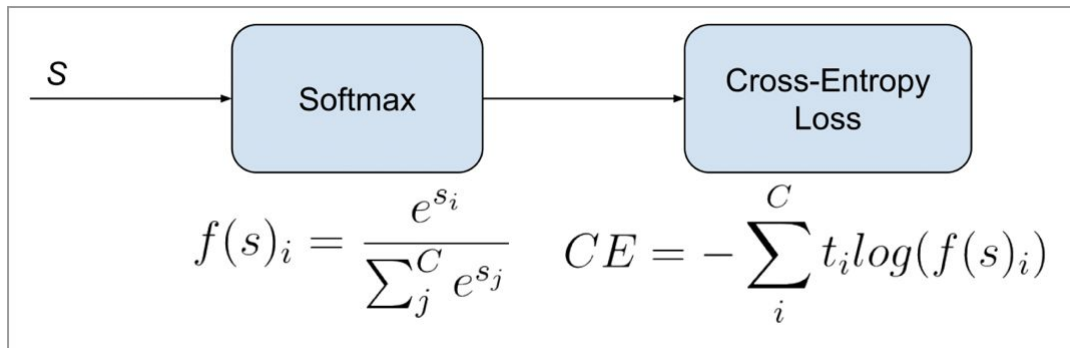
# Categorical Cross Entropy Loss



$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

$$f(s)_i = \frac{e^{s_i}}{\sum_{j}^{C} e^{s_j}} \qquad CE = -\sum_{i}^{C} t_i log(f(s)_i)$$

# Evolution of Image Classification

# LeNet (1998)



conv. layer     avg pool     conv. layer     avg pool

$f = 5$   $s = 1$     $f = 2$   $s = 2$     $f = 5$   $s = 1$     $f = 2$   $s = 2$

$\hat{y}$

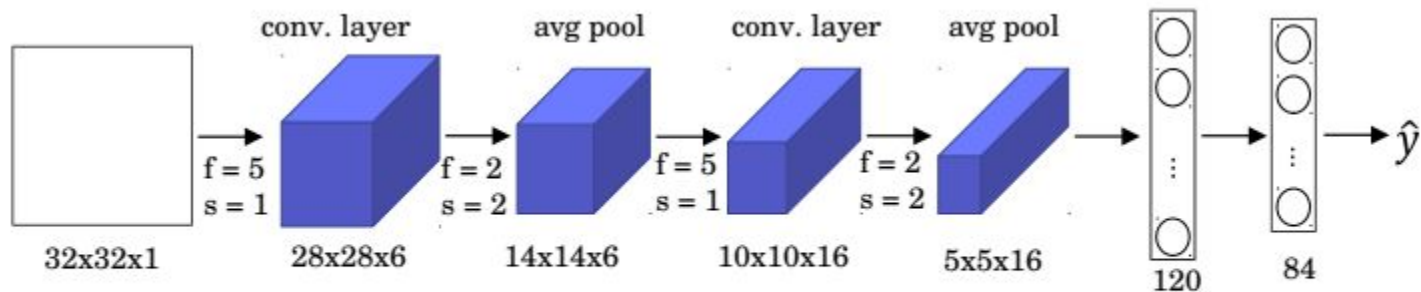32x32x1     28x28x6     14x14x6     10x10x16     5x5x16     120     84

Figure 3: LeNet-5 neural network. Around 60k parameters.

# AlexNet (2012)



Figure 4: AlexNet neural network. Around 60 million parameters.
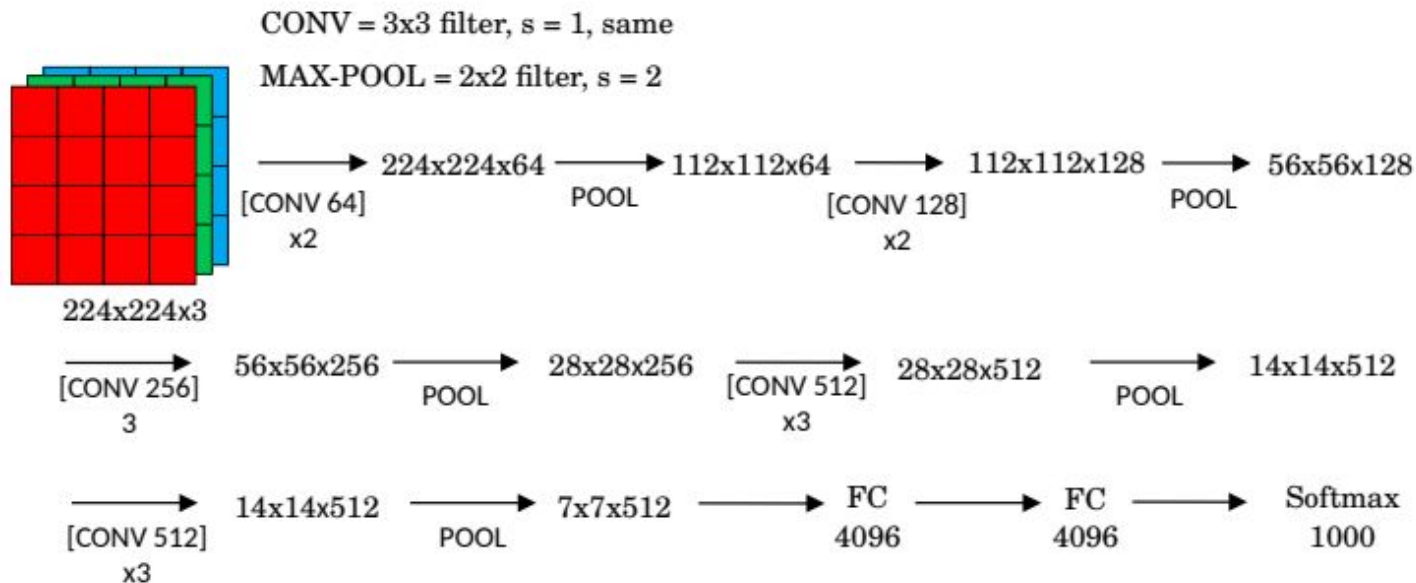
# VGG-16 and VGG-19 (2014)



CONV = 3x3 filter, s = 1, same

MAX-POOL = 2x2 filter, s = 2

224x224x3

[CONV 64] x2 → 224x224x64 → POOL → 112x112x64 → [CONV 128] x2 → 112x112x128 → POOL → 56x56x128

[CONV 256] 3 → 56x56x256 → POOL → 28x28x256 → [CONV 512] x3 → 28x28x512 → POOL → 14x14x512

[CONV 512] x3 → 14x14x512 → POOL → 7x7x512 → FC 4096 → FC 4096 → Softmax 1000

Figure 5: VGG-16. Around 138 million parameters.

# Inception Net (GoogLeNet) (2014)
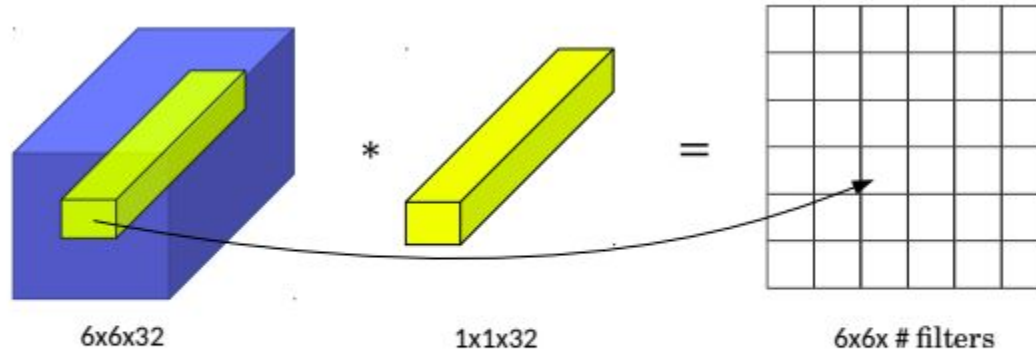


6x6x32      1x1x32      6x6x # filters

Figure 11: $1 \times 1$ convolution. The filter has size $1 \times 1 \times 32$ elements (weights). The number of filters correspond to the number of channels of the output.
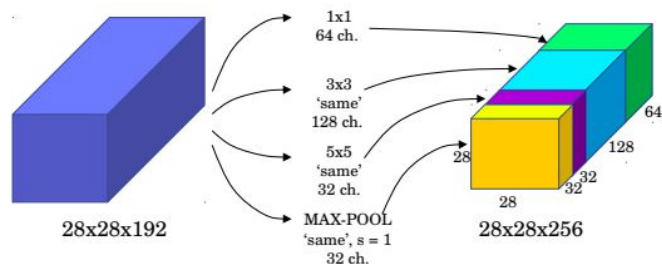
# Inception Net



Figure 13: Inception module with $1 \times 1$, $3 \times 3$, $5 \times 5$ convolutional layers, and max-pooling.



Figure 14: Inception module with $1 \times 1$, $3 \times 3$, $5 \times 5$ convolutional layers, and max-pooling with intermediate $1 \times 1$ convolutions.
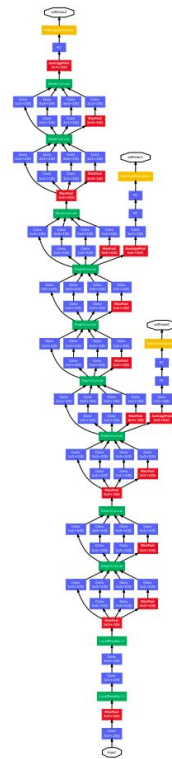


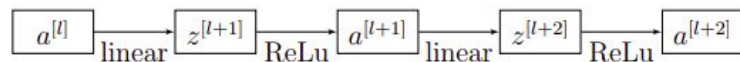Figure 15: GoogLeNet network with all the bells and whistles [7].

# ResNets (2015)


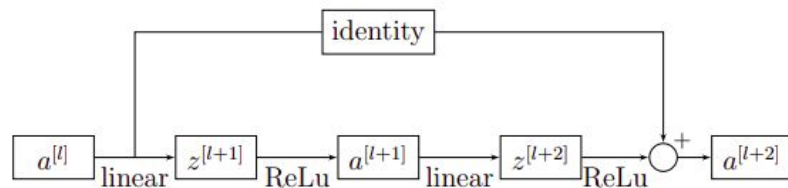
Figure 6: Plain network structure for layers $l$ to $l+2$.



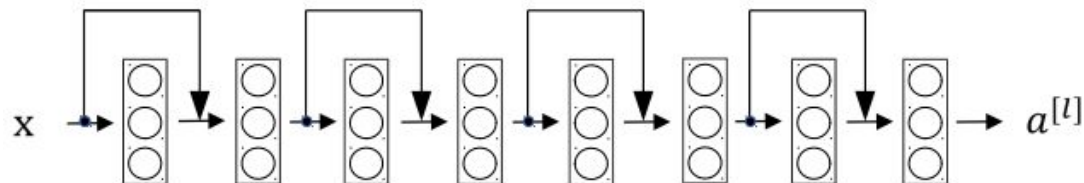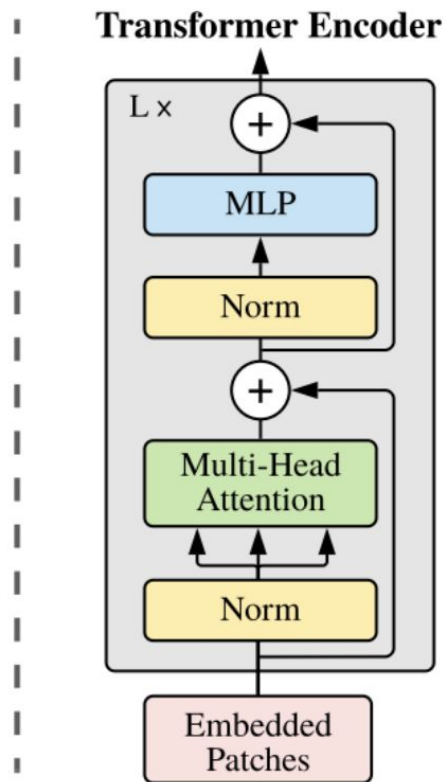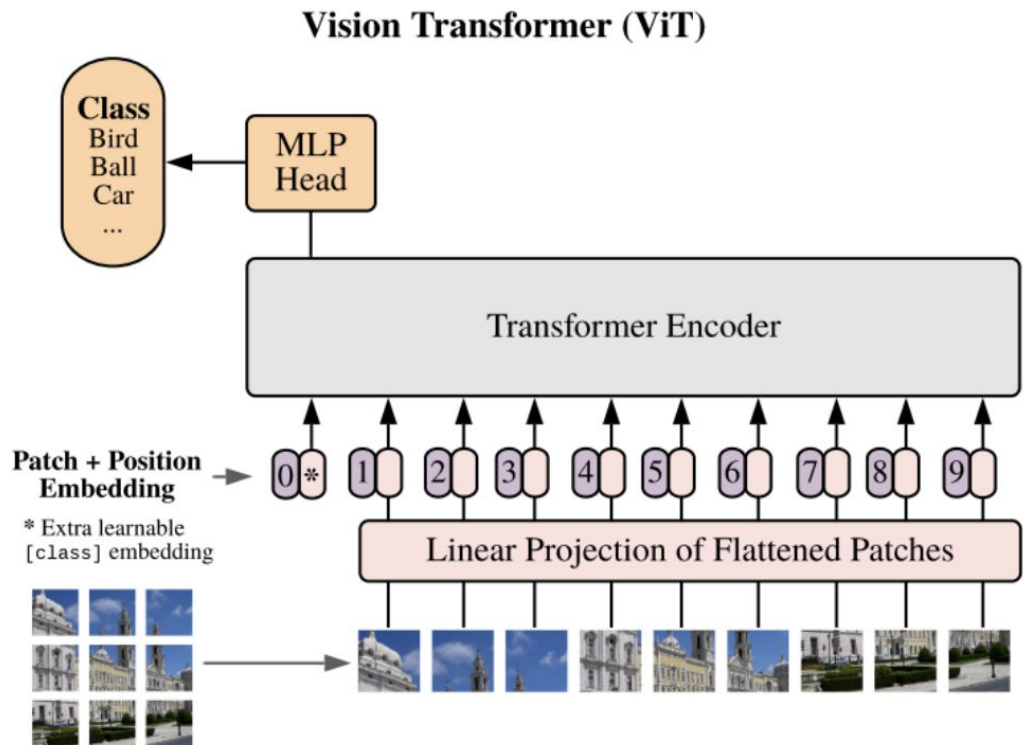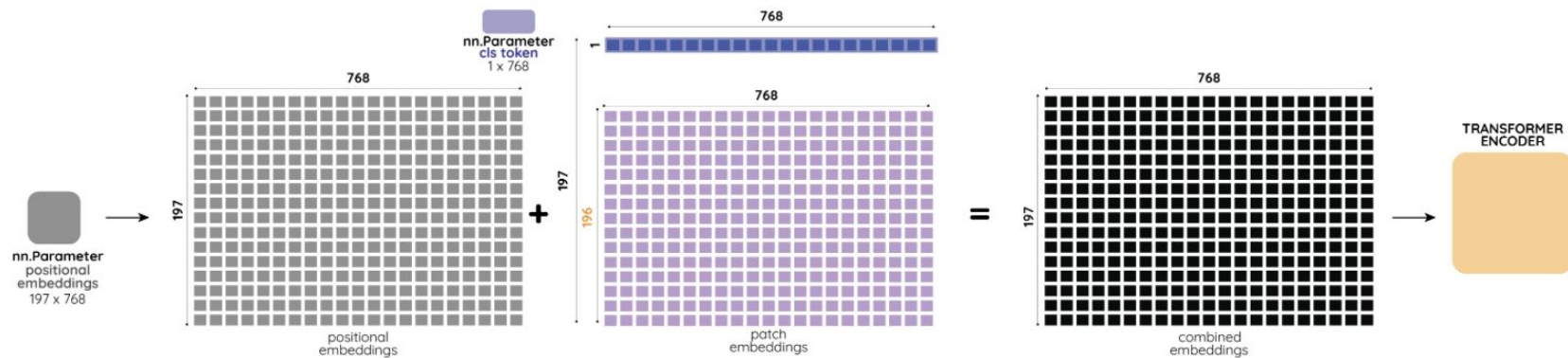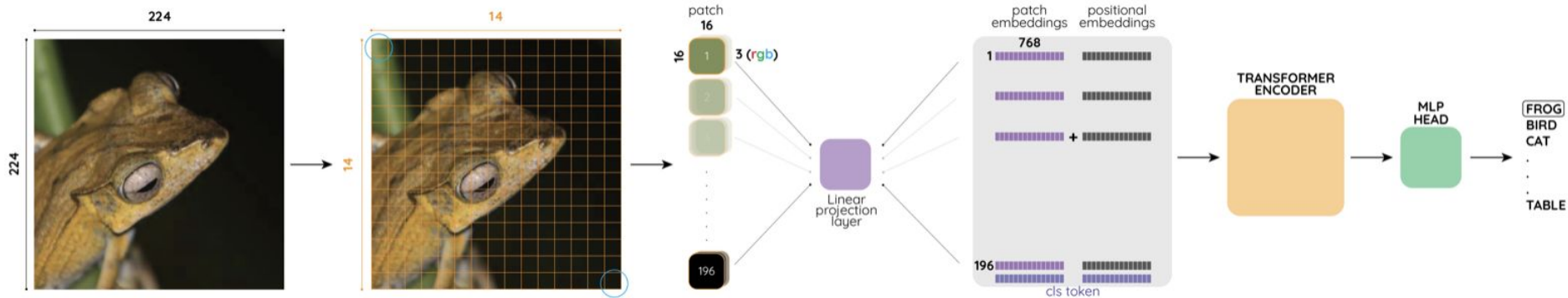Figure 7: Residual network structure for layers $l$ to $l+2$.



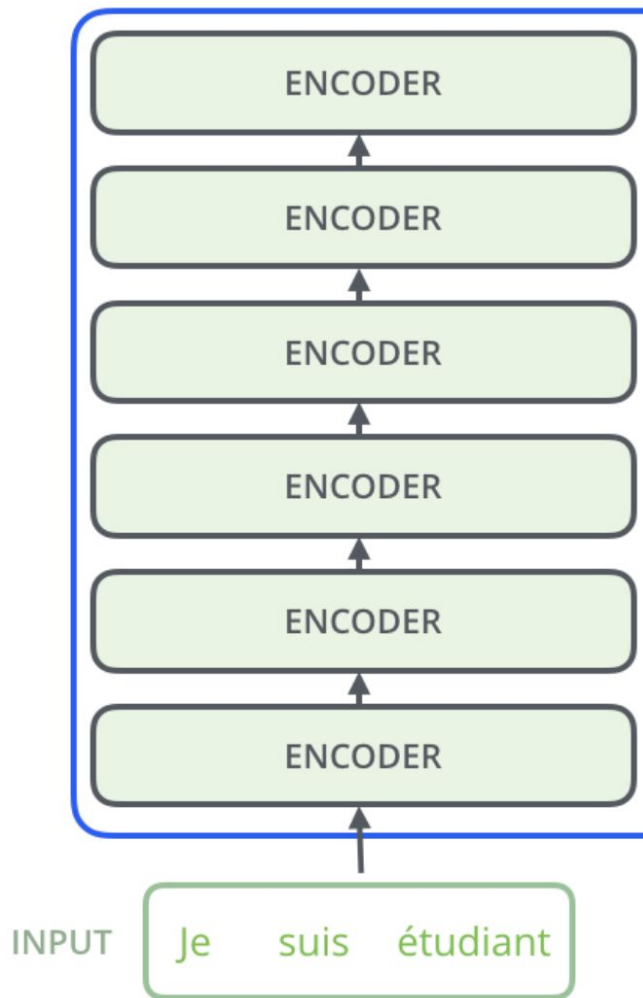Figure 8: Residual network structure for layers $l$ to $l+2$.

ViT (2020)

# Overall Architecture (ViT)



**Vision Transformer (ViT)**

Class
Bird
Ball
Car
...

MLP
Head

Transformer Encoder

**Patch + Position Embedding**

* Extra learnable
[class] embedding

0 * 1 2 3 4 5 6 7 8 9

Linear Projection of Flattened Patches

**Transformer Encoder**

L ×

MLP

Norm

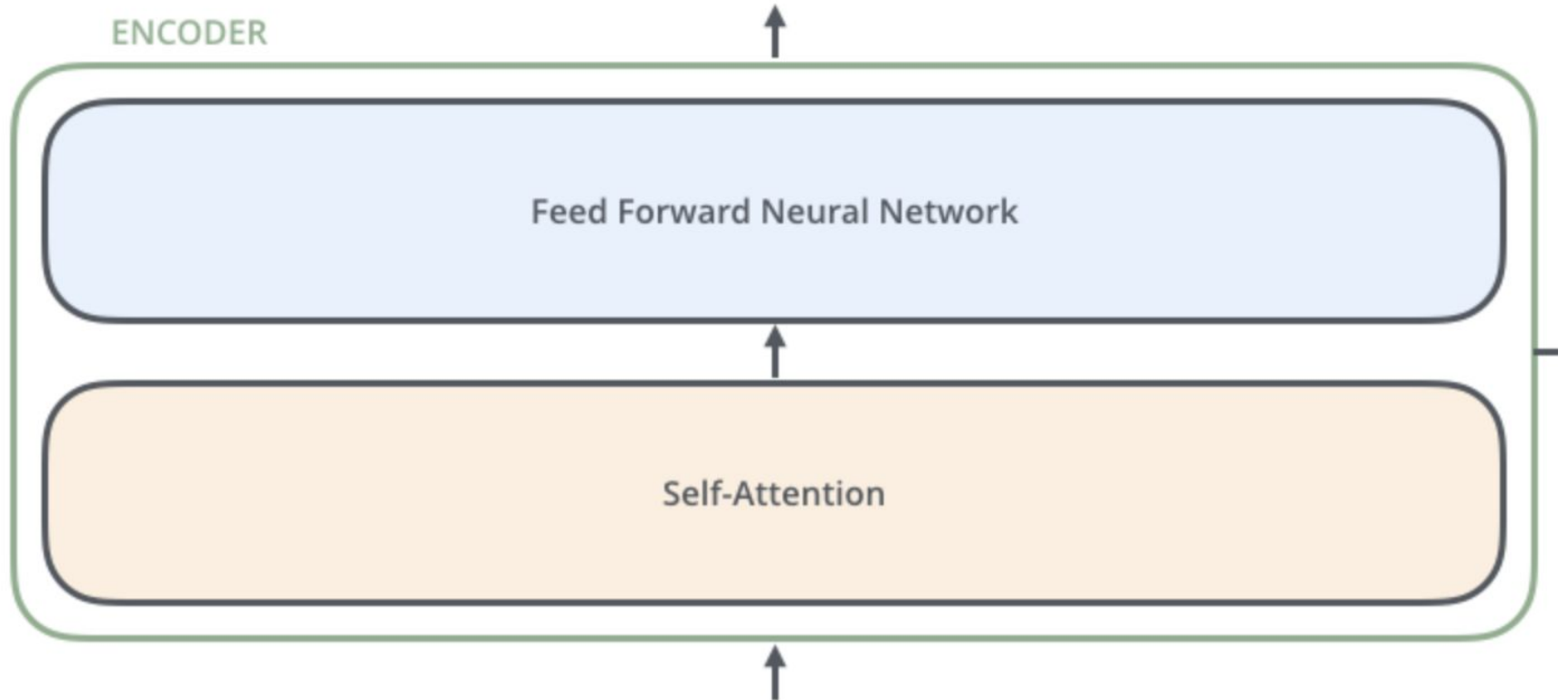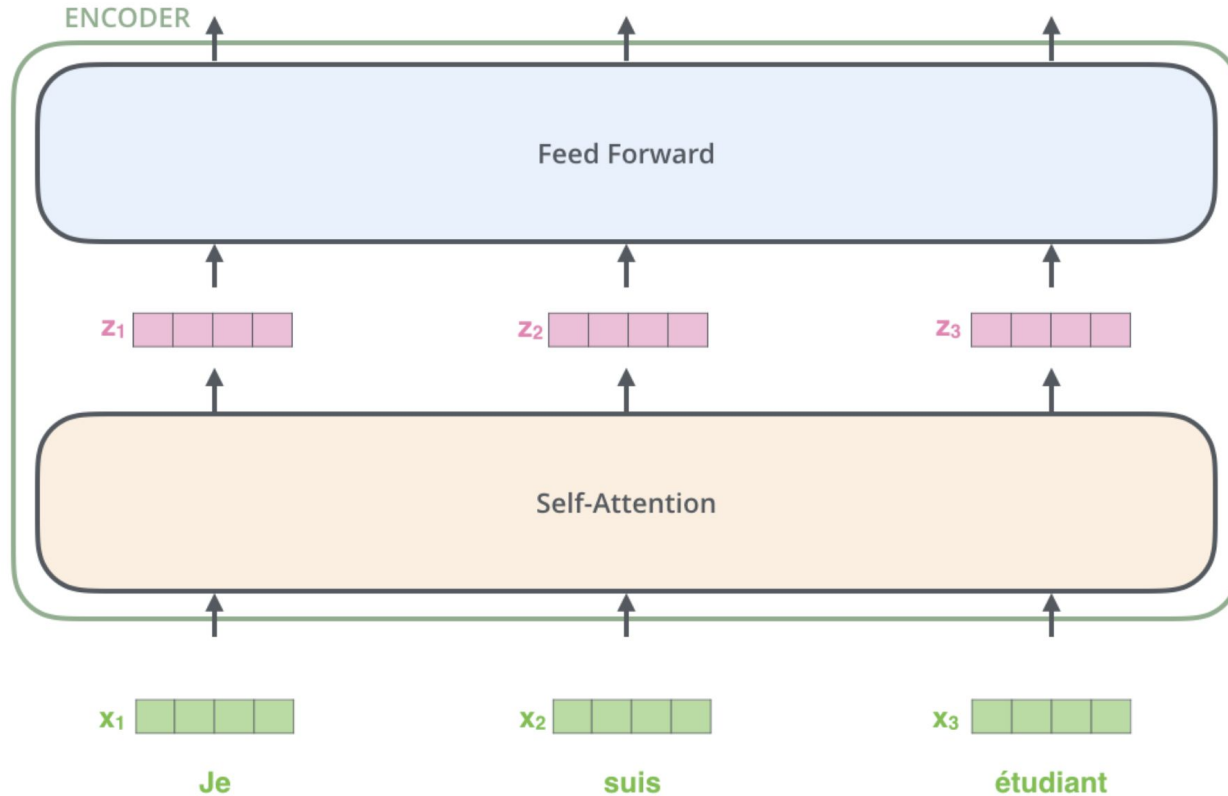Multi-Head Attention

Norm

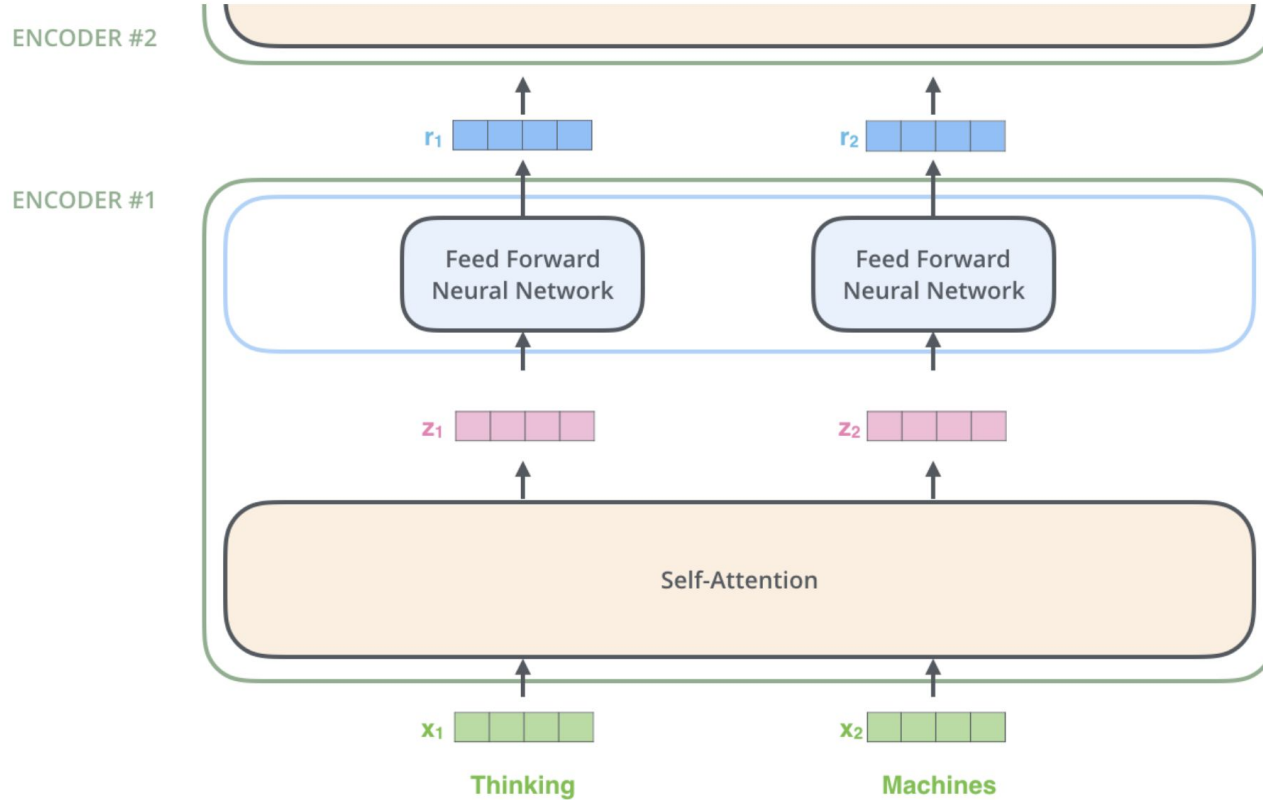Embedded Patches

# ViTs

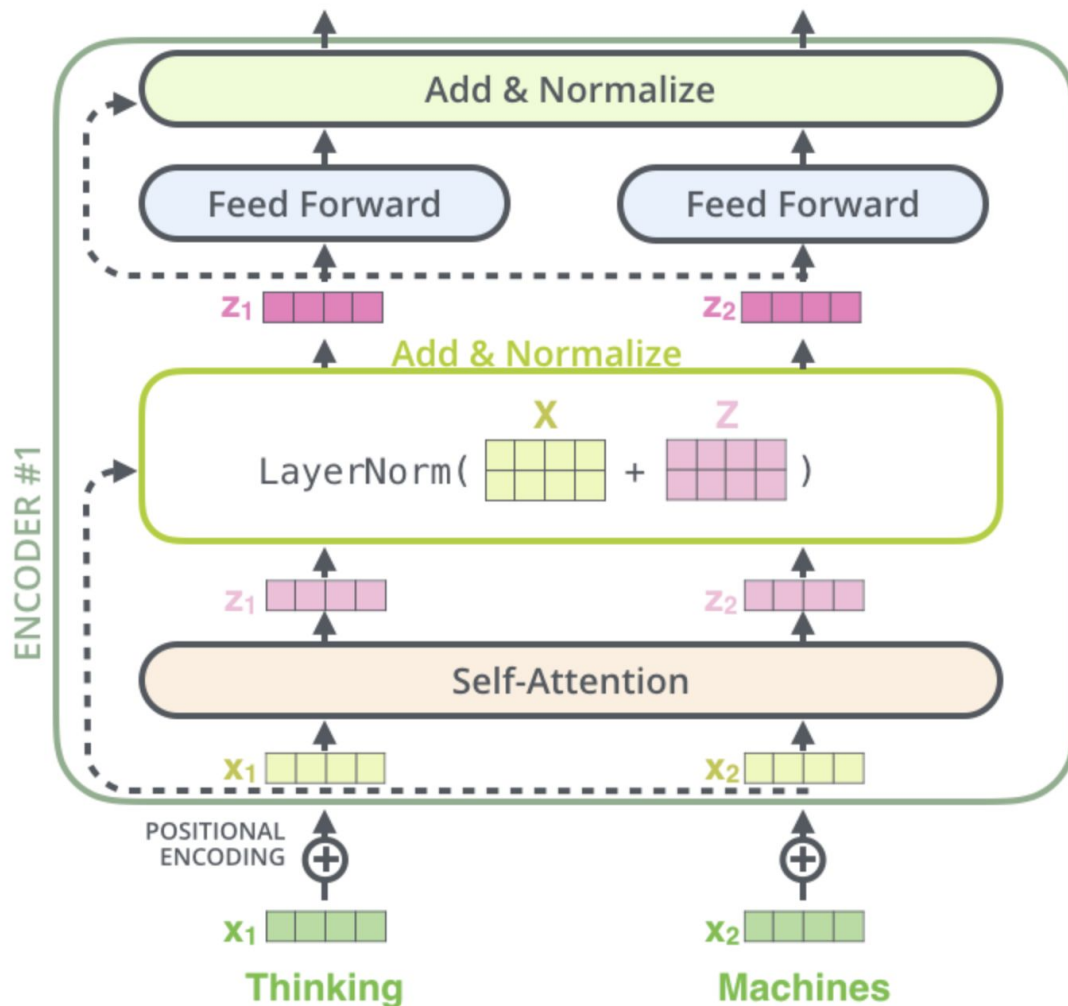# Transformer Encoder

# Transformer Encoder

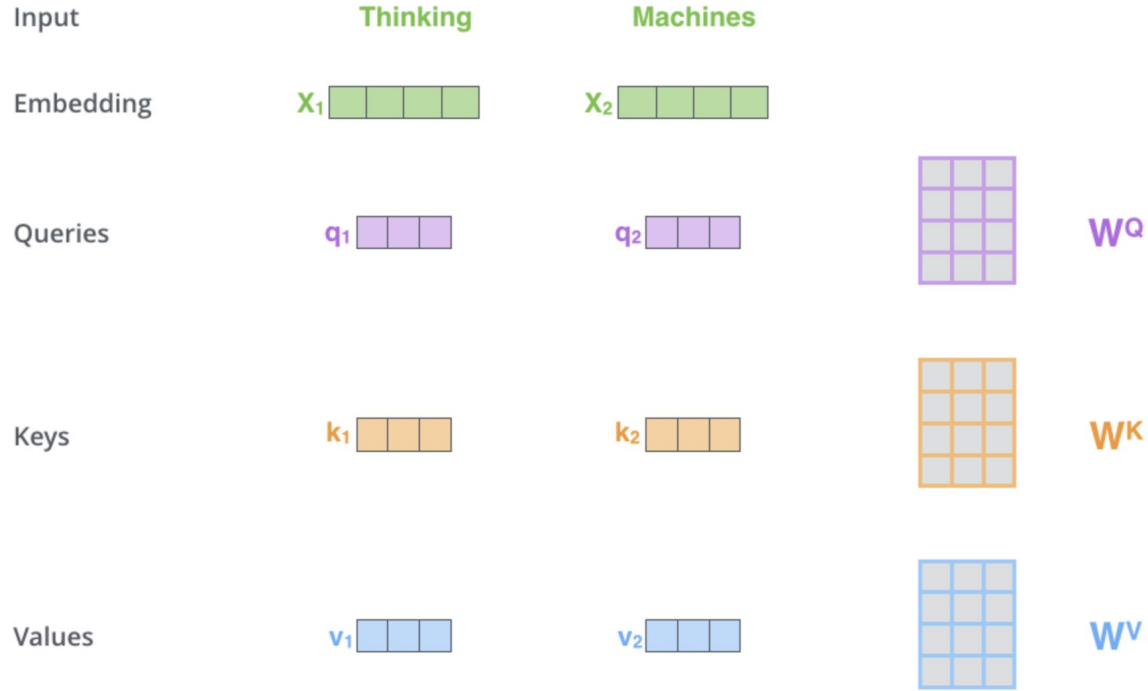# Transformer Encoder

# Transformer Encoder

# Transformer Encoder

# Transformer Encoder



Multiplying x1 by the WQ weight matrix produces q1, the "query" vector associated with that word. We end up creating a "query", a "key", and a "value" projection of each word in the input sentence.

# Transformer Encoder



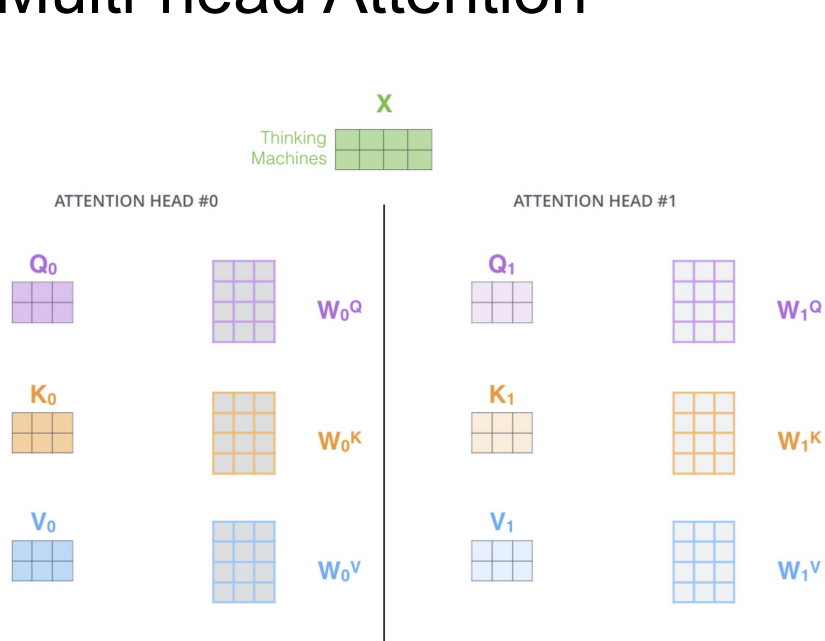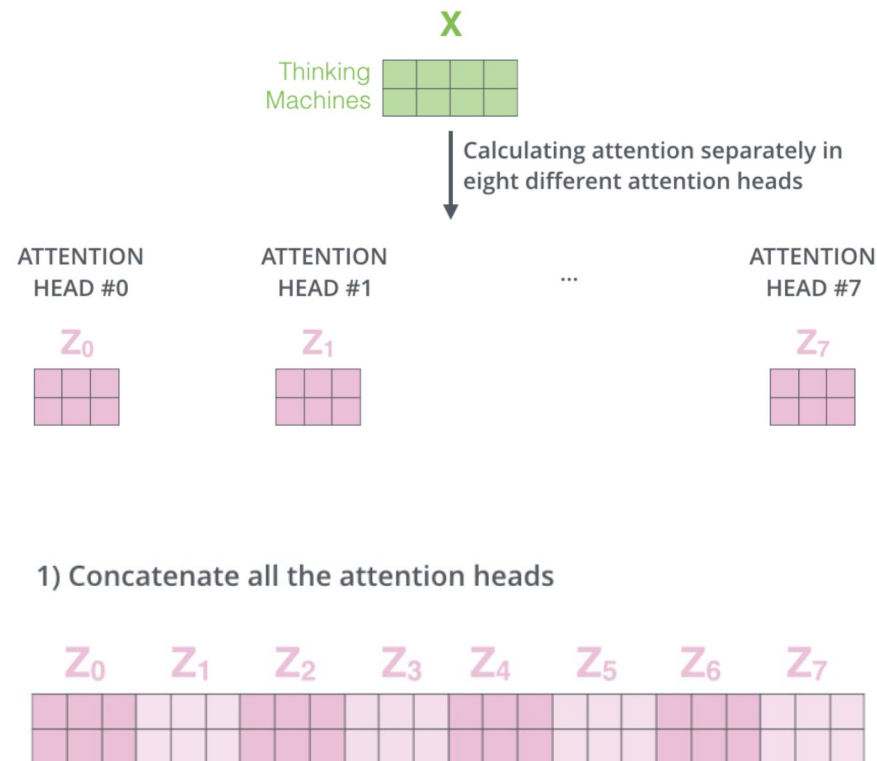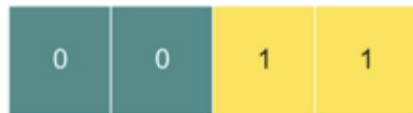| Input | Thinking | Machines |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

# Multi-head Attention



With multi-headed attention, we maintain separate Q/K/V weight matrices for each head resulting in different Q/K/V matrices. As we did before, we multiply X by the WQ/WK/WV matrices to produce Q/K/V matrices.
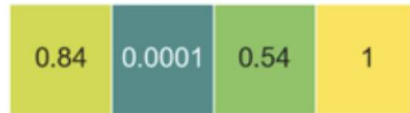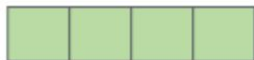
# Positional Encodings

# Resources

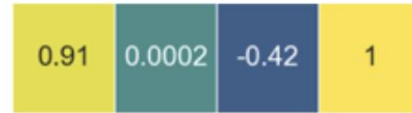1.  Blog on Attention - [Attention? Attention! | Lil'Log](Attention? Attention! | Lil'Log)
2.  A wonderful blog on Transformers - [The Illustrated Transformer – Jay Alammar](The Illustrated Transformer – Jay Alammar)
3.  Blog on Vision Transformers - [Vision Transformer](Vision Transformer)
4.  Original Paper on Transformers in NLP (2017) - [[1706.03762] Attention Is All You Need]([1706.03762] Attention Is All You Need)
5.  Original Paper on Vision Transformers(2020) - [[2010.11929] An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale]([2010.11929] An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale)

# The End