# Deep learning

Transformers
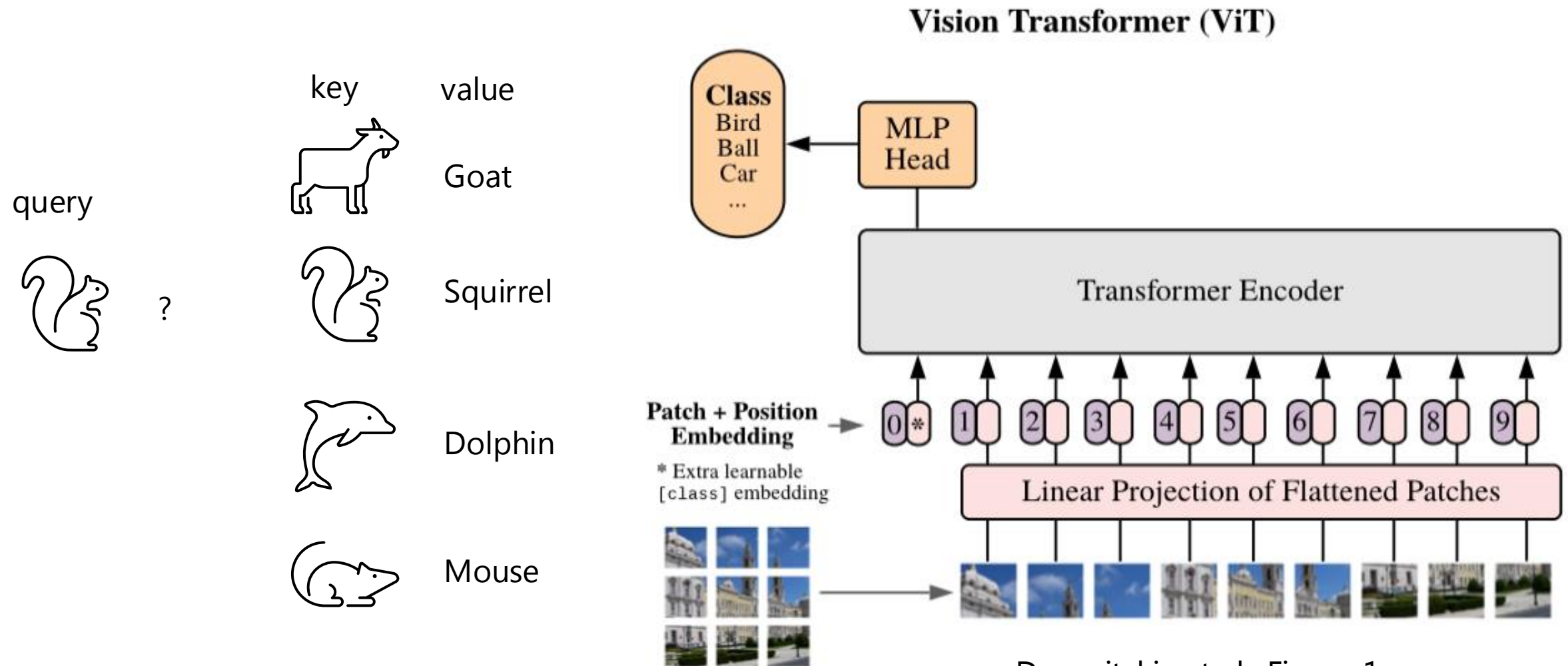
2024/12/18

**Jon Sporring**,
Department of Computer Science

UNIVERSITY OF COPENHAGEN
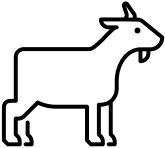
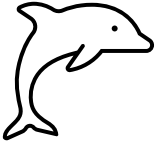# Transformer networks – Vision Transformer ViT: Vaswani ea (2017) & Dosovitskiy ea (2020)



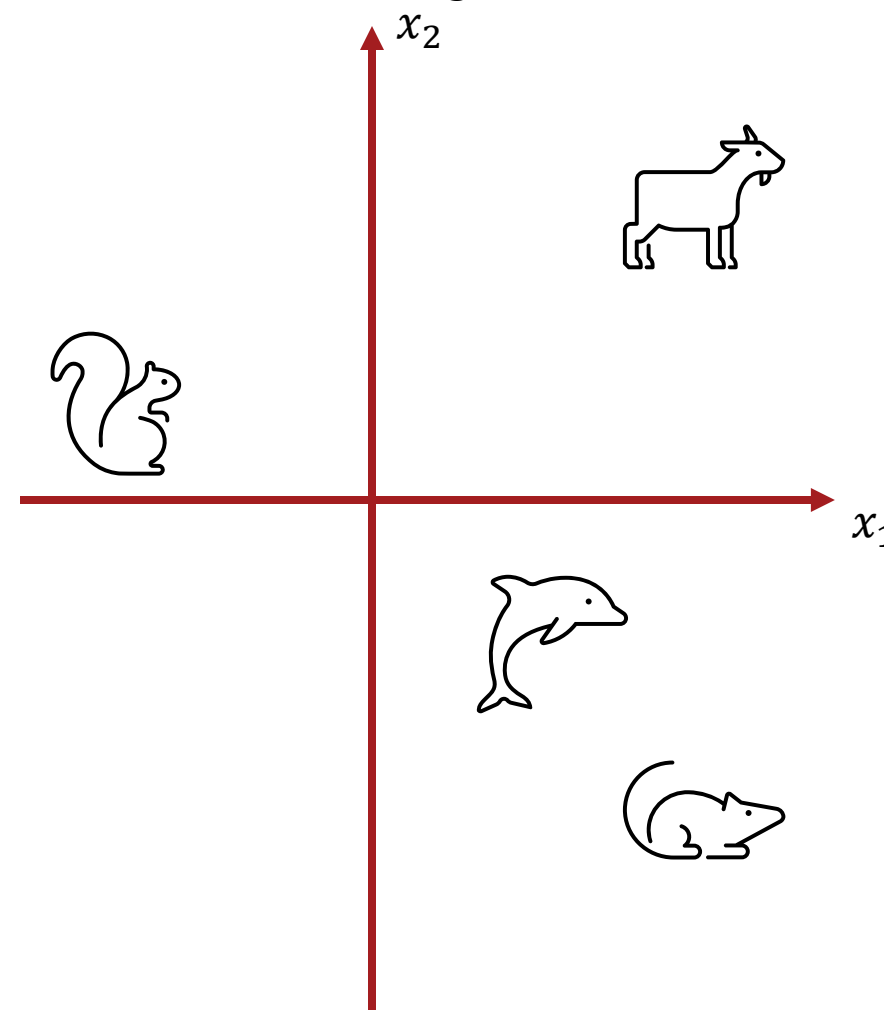key    value

Goat

query

?    Squirrel

Dolphin

Mouse

**Vision Transformer (ViT)**

Dosovitskiy et al., Figure 1

# Tokens: Cut into patches, flatten and linearly embed in lower dimensional space

Tokens as embedding

key    value

Goat

query

?    Squirrel

Dolphin

Mouse

$x_2$

$x_1$

# Embeddings++



Tokens as embedding

$$I \in \mathbb{R}^{H \times W \times C}, I_p \in \mathbb{R}^{M \times M \times C}, \text{Flatten}: \mathbb{R}^{M \times M \times C} \to \mathbb{R}^{M^2 C}, \text{Token}: \mathbb{R}^{M^2 C} \to \mathbb{R}^{D}$$

Convention: row vectors

$$x_j = \text{Flatten}(I_j) \in \mathbb{R}^{M^2 C}$$

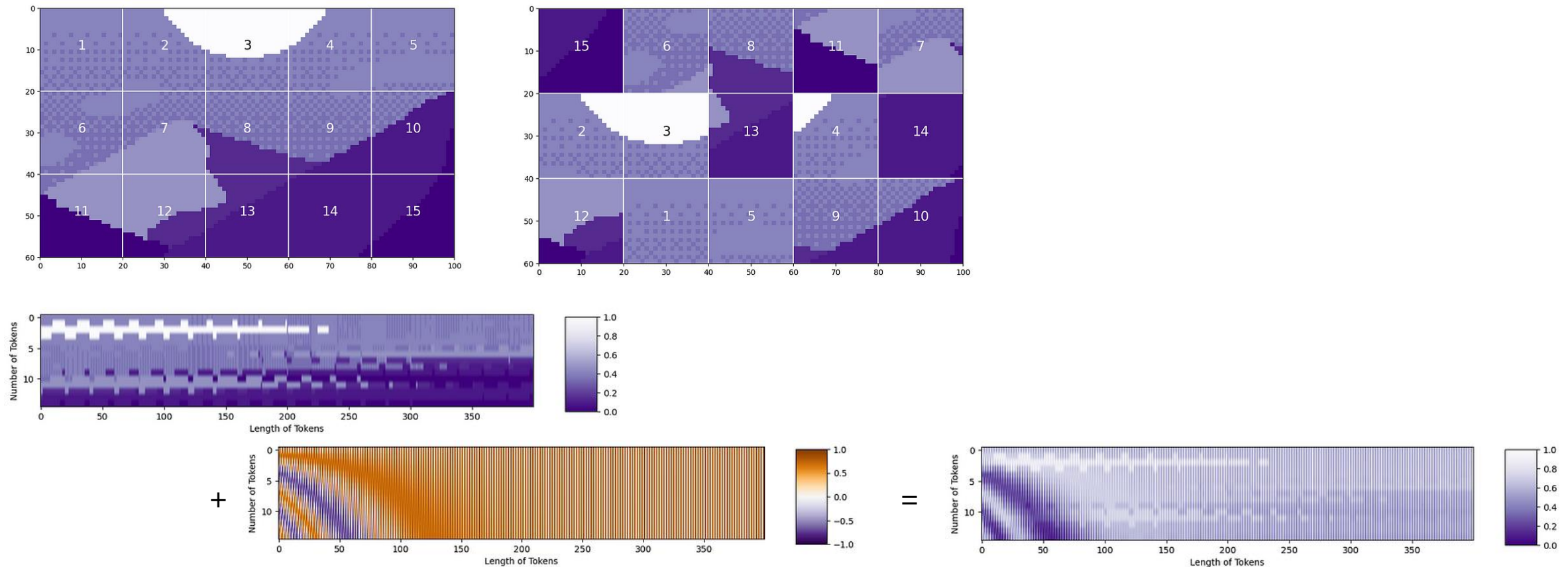$$t_j = \text{Token}(x_j) = x_j \mathbf{E} \in \mathbb{R}^{D}$$

$$z^0 = [x_{\text{class}}; t_1; t_2; \ldots; t_N] + \text{PositionEmbedding}() \in \mathbb{R}^{(N+1) \times D}$$
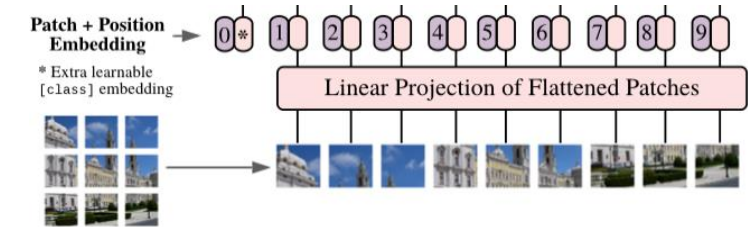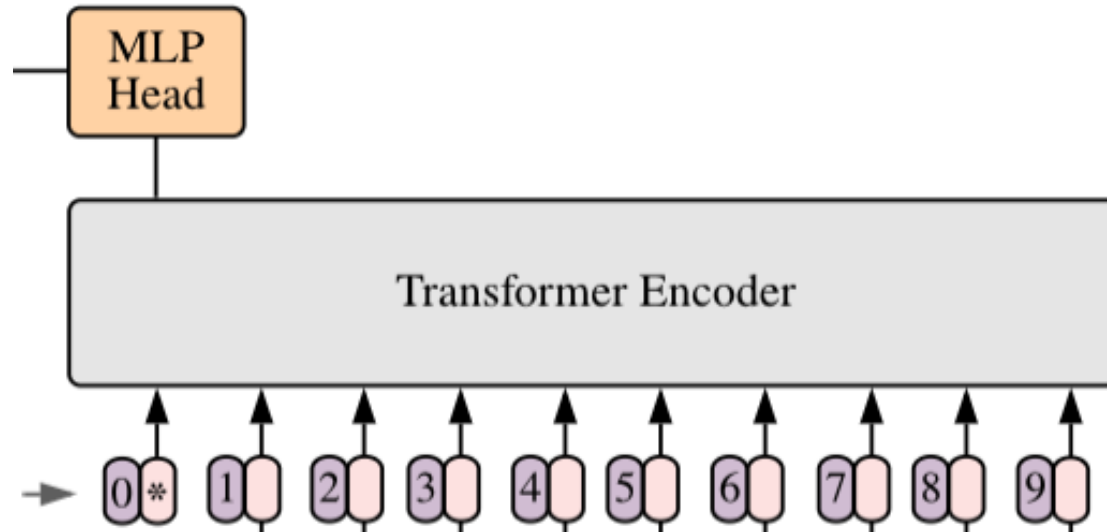
# Positional Embedding: Vision depends on position

Yuan et al (2021). *Tokens-to-Token ViT: Training Vision Transformers from Scratch on ImageNet.*
Figures from
https://towardsdatascience.com/position-embeddings-for-vision-transformers-explained-a6f9add341d5

# Encoding



A sequence of attention layers: $l = 1 \ldots L$

$$\boldsymbol{z}_0 = [\boldsymbol{x}_{\text{class}}; \boldsymbol{t}_1; \boldsymbol{t}_2; \ldots; \boldsymbol{t}_N] \in \mathbb{R}^{(N+1) \times D}$$

$$\boldsymbol{z}'_l = \text{MSA}\big(\text{LN}(\boldsymbol{z}'_{l-1})\big) + \boldsymbol{z}'_{l-1} \in \mathbb{R}^{(N+1) \times D}$$

$$\boldsymbol{z}_l = \text{MLP}_2\big(\text{LN}(\boldsymbol{z}_l)\big) + \boldsymbol{z}'_l \in \mathbb{R}^{(N+1) \times D}$$
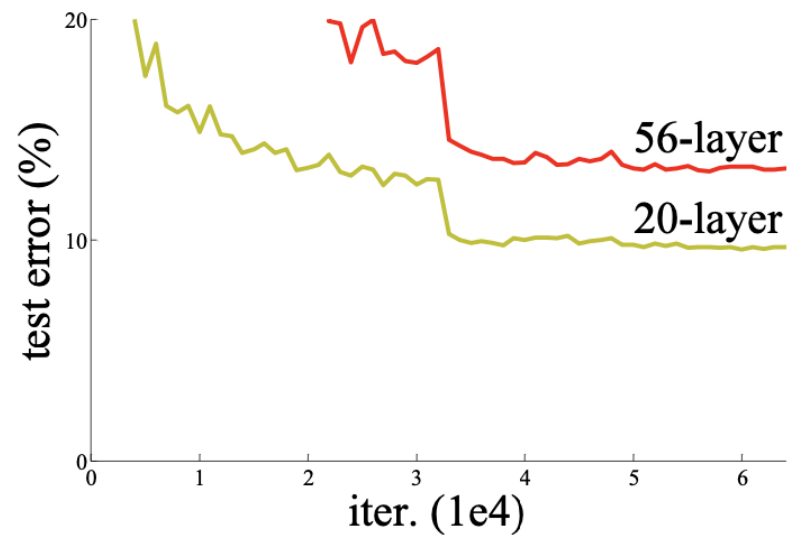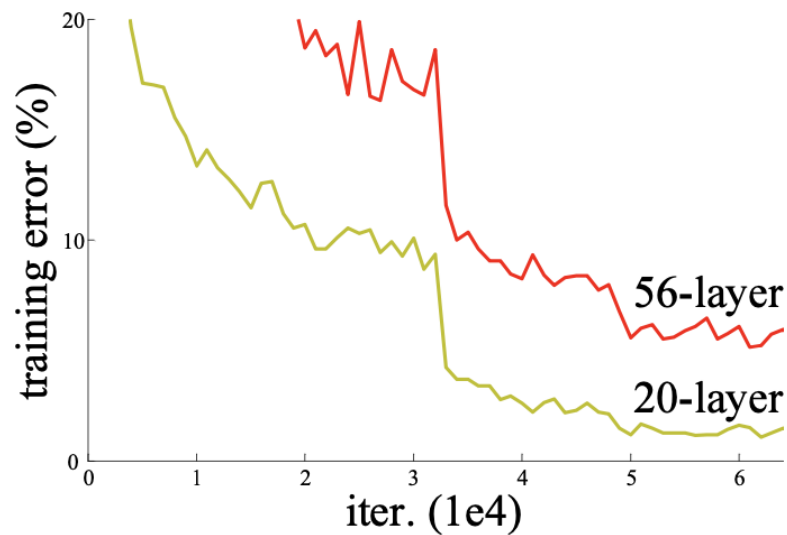
MSA: Multiheaded self-attention
$\text{MLP}_n$: n-layer Multilayer perceptron
LN: Layer normalization

Residual learning

# Residual learning for optimizing deeper networks

He et al, "Deep Residual Learning for Image Recognition, CVPR 2016

Observation: Deeper networks = poorer convergence



He et al. Figure 1

# Residual learning for optimizing deeper networks

He et al, "Deep Residual Learning for Image Recognition, CVPR 2016

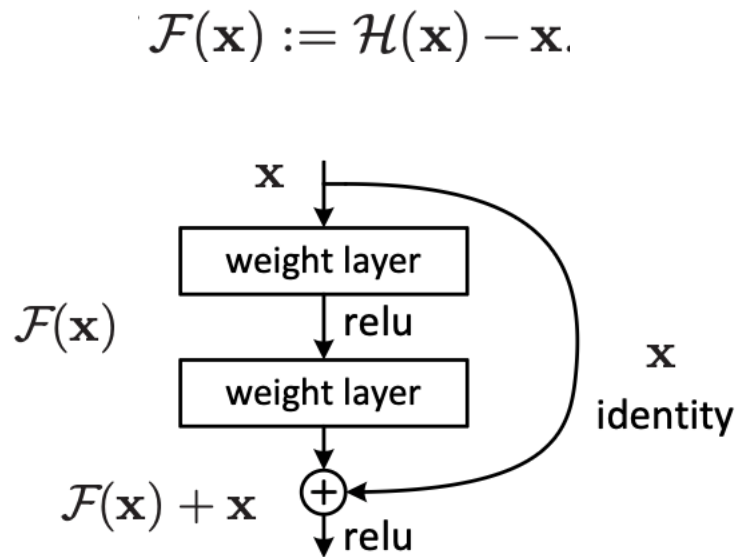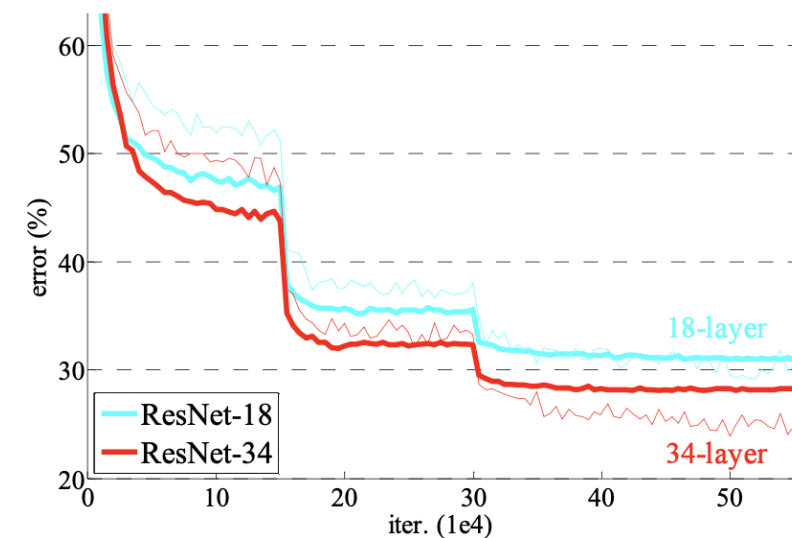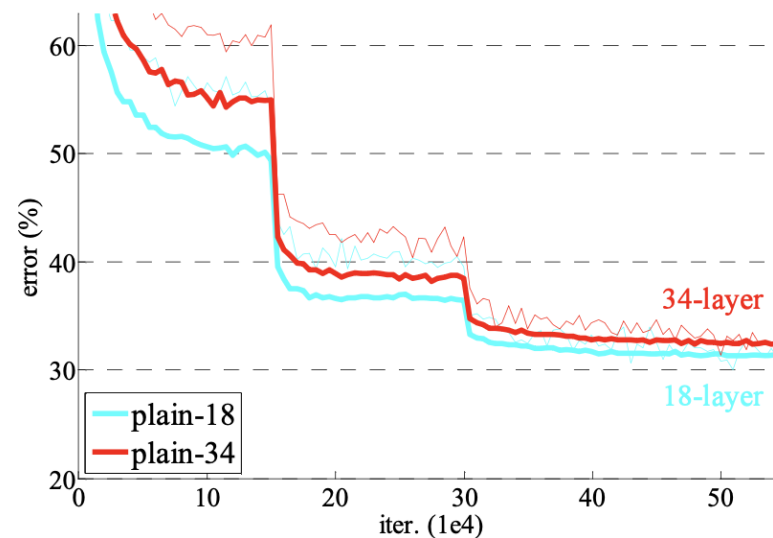Observation: Deeper networks = poorer convergence
Solution: Train on residuals

$$\mathcal{F}(\mathbf{x}) := \mathcal{H}(\mathbf{x}) - \mathbf{x}.$$

He et al. Figure 2

He et al. Figure 4

# Multilayer perceptron (MLP) and Layer Normalization (LN)

Ba, Kiros, and Hinton, Layer normalization, 2016

$$\text{LN}(\boldsymbol{z}) = \frac{\boldsymbol{z} - \bar{\boldsymbol{z}}}{\sqrt{\text{Var}(\boldsymbol{z}) + \epsilon}}$$

A sequence of attention layers: $l = 1 \dots L$

$$\boldsymbol{z}_0 = [\boldsymbol{x}_{\text{class}}; \boldsymbol{t}_1; \boldsymbol{t}_2; \dots,; \boldsymbol{t}_N] \in \mathbb{R}^{(N+1)\times D}$$
$$\boldsymbol{z}'_l = \text{MSA}\big(\text{LN}(\boldsymbol{z}'_{l-1})\big) + \boldsymbol{z}'_{l-1} \in \mathbb{R}^{(N+1)\times D}$$
$$\boldsymbol{z}_l = \text{MLP}\big(\text{LN}(\boldsymbol{z}'_l)\big) + \boldsymbol{z}'_l \in \mathbb{R}^{(N+1)\times D}$$

MLP: Single layer with two GELU
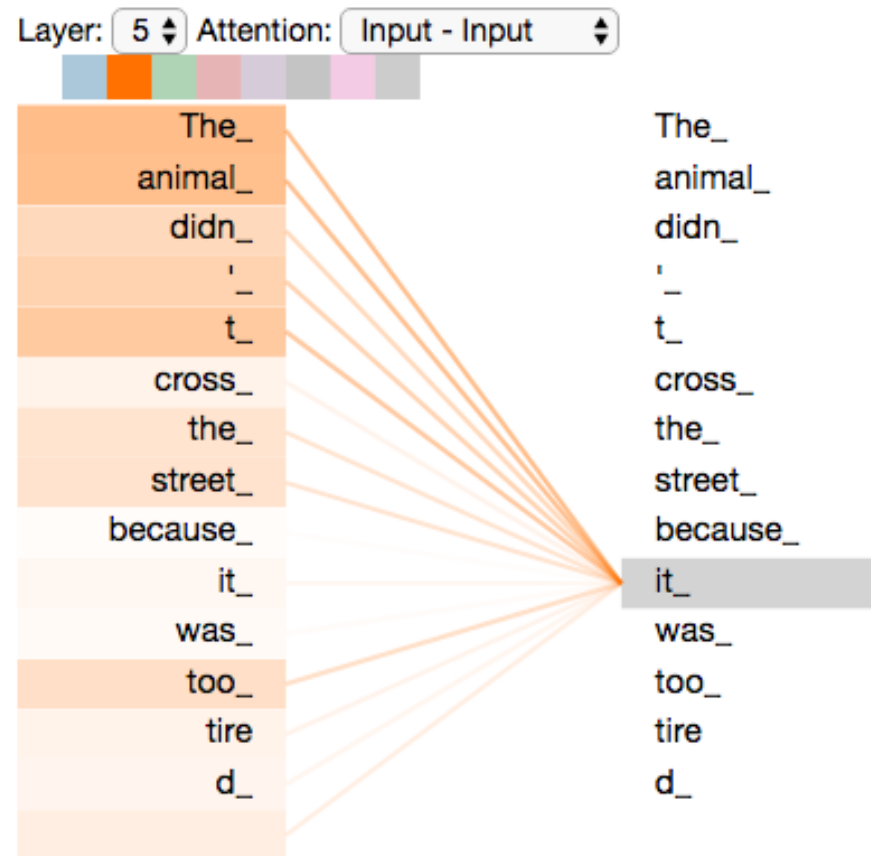GELU: Hendryks & Gimpel, Gaussian Error Linear Units (GELU)



Hendryks & Gimpel, Figure 1

$$\text{GELU}(x) = \frac{x}{2}\left[1 + \text{erf}\left(\frac{x}{\sqrt{2}}\right)\right]$$

$$\simeq x/2\left(1 + \tanh\left((x + 0.044715x^3)\sqrt{\frac{2}{\pi}}\right)\right)$$

# Attention:
# The animal didn't cross the street because it was too tired

# Self-attention head

$I \in \mathbb{R}^{H \times W \times C}, I_p \in \mathbb{R}^{M \times M \times C}, \text{Flatten: } \mathbb{R}^{M \times M \times C} \to \mathbb{R}^{M^2 C}, \text{Token: } \mathbb{R}^{M^2 C} \to \mathbb{R}^D$

$\mathbf{z} \in \mathbb{R}^{(N+1) \times D}$



Vaswani, Fig 2

Query $\quad \mathbf{q}_j = \mathbf{z}_{j*} U_q \in \mathbb{R}^{D_h}$

Key $\quad \mathbf{k}_j = \mathbf{z}_{j*} U_k \in \mathbb{R}^{D_h}$

Value $\quad \mathbf{v}_j = \mathbf{z}_{j*} U_v \in \mathbb{R}^{D_h}$

$U_* \in \mathbb{R}^{D \times D_h}$ learnable

$\mathbf{a}_{ij} = \text{softmax}\left(\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{D_h}}\right)$
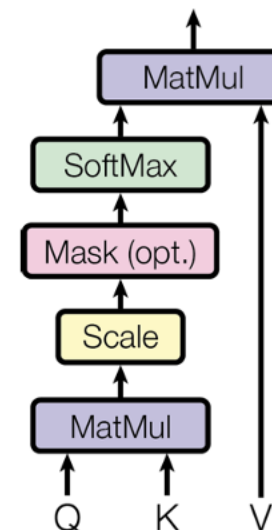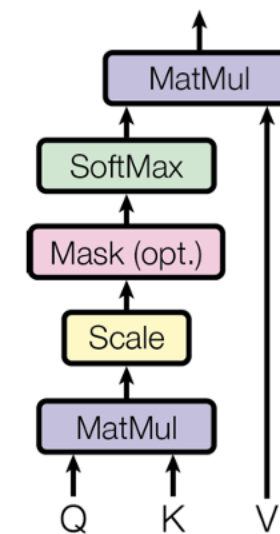
$\text{softmax}(\{s_j\}) = \left\{\frac{e^{s_j}}{\sum_{j=1}^n e^{s_j}}\right\}$

$\text{SA}(\mathbf{z}_{i*}) = \sum_j \mathbf{a}_{ij} \mathbf{v}_j \in \mathbb{R}^{D_h}$

Query how one embedding ($\mathbf{q}$) matches keys ($\mathbf{k}$), which are all other embeddings.

Softmax select the best maching querys to their cosine difference

Output is now computed as a weigted sum of values ($\mathbf{v}$).

# Multihead self-attention

$I \in \mathbb{R}^{H \times W \times C}, I_p \in \mathbb{R}^{M \times M \times C}, \text{Flatten: } \mathbb{R}^{M \times M \times C} \to \mathbb{R}^{M^2 C}, \text{Token: } \mathbb{R}^{M^2 C} \to \mathbb{R}^D$

$\boldsymbol{z} \in \mathbb{R}^{(N+1) \times D}$

Query     $\boldsymbol{Q} = \{\boldsymbol{q}_j\} \in \mathbb{R}^{(N+1) \times D_h}$

Key         $\boldsymbol{K} = \{\boldsymbol{k}_j\} \in \mathbb{R}^{(N+1) \times D_h}$         $U_* \in \mathbb{R}^{D \times D_h} \text{ learnable}$

Value      $\boldsymbol{V} = \{\boldsymbol{v}_j\} \in \mathbb{R}^{(N+1) \times D_h}$

$$\boldsymbol{A} = \text{softmax}\left(\frac{\boldsymbol{Q}\boldsymbol{K}^T}{\sqrt{D_h}}\right) \in \mathbb{R}^{(N+1) \times (N+1)}$$

$\text{SA}(\boldsymbol{z}) = \boldsymbol{A}\boldsymbol{V} \in \mathbb{R}^{(N+1) \times D_h}$

$\boldsymbol{U}_{msa} \in \mathbb{R}^{kD_h \times D}$

$\text{MSA}(\boldsymbol{z}) = [\text{SA}_1(\boldsymbol{z}), \text{SA}_2(\boldsymbol{z}), \dots, \text{SA}_k(\boldsymbol{z})]\boldsymbol{U}_{msa} \in \mathbb{R}^{(N+1) \times D}$



Vaswani, Fig 2

# Encoding



A sequence of attention layers: $l = 1 \dots L$

$$\boldsymbol{z}_0 = [\boldsymbol{x}_{\text{class}}; \ \boldsymbol{t}_1; \boldsymbol{t}_2; \dots; \boldsymbol{t}_N] \in \mathbb{R}^{2(N+1) \times D}$$

$$\boldsymbol{z}'_l = \text{MSA}\big(\text{LN}(\boldsymbol{z}'_{l-1})\big) + \boldsymbol{z}'_{l-1} \in \mathbb{R}^{2(N+1) \times D}$$

$$\boldsymbol{z}_l = \text{MLP}_2\big(\text{LN}(\boldsymbol{z}'_l)\big) + \boldsymbol{z}'_l \in \mathbb{R}^{2(N+1) \times D}$$

$$\boxed{\text{Class} = \ \text{MLP}_1(\boldsymbol{z}_L)}$$

# Inductive bias [Dosoviskiy]

**Inductive bias.** We note that Vision Transformer has much less image-specific inductive bias than CNNs. In CNNs, locality, two-dimensional neighborhood structure, and translation equivariance are baked into each layer throughout the whole model. In ViT, only MLP layers are local and translationally equivariant, while the self-attention layers are global. The two-dimensional neighborhood

# ImageNet: https://image-net.org/
14*10^6 images, 2*10^4 categories, 10^6 images with bounding boxes



© 2020 Stanford Vision Lab, Stanford University, Princeton University    imagenet.help.desk@gmail.com    Copyright infringement

# State-of-the-art 2024

**Swin Transformer: Hierarchical Vision Transformer using Shifted Windows**

xuan Wei[†]

o



(a) Swin Transformer (ours)

(b) ViT

# ConvNeXt

**A ConvNet for the 2020s**

Zhuang Liu[1,2]*   Hanzi Mao[1]   Chao-Yuan Wu[1]   Christoph Feichtenhofer[1]   Trevor Darrell[2]   Saining Xie[1†]

[1]Facebook AI Research (FAIR)    [2]UC Berkeley

Code: https://github.com/facebookresearch/ConvNeXt

# Papers with code https://paperswithcode.com/