# Deep Learning Essentials

## 8. Backbone Architectures

ResNet, EfficientNet, Self-attention, Transformers

Lukáš Neumann

# Image Classification

# ImageNet



Deng, Jia, et al. "ImageNet: A large-scale hierarchical image database." CVPR 2009

- 1000 image classes
- 1.2M training images, 100k validation

# ImageNet

Easiest classes

red fox (100)  hen-of-the-woods (100)  ibex (100)  goldfinch (100)  flat-coated retriever (100)

tiger (100)  hamster (100)  porcupine (100)  stingray (100)  Blenheim spaniel (100)

Hardest classes

muzzle (71)  hatchet (68)  water bottle (68)  velvet (68)  loupe (66)
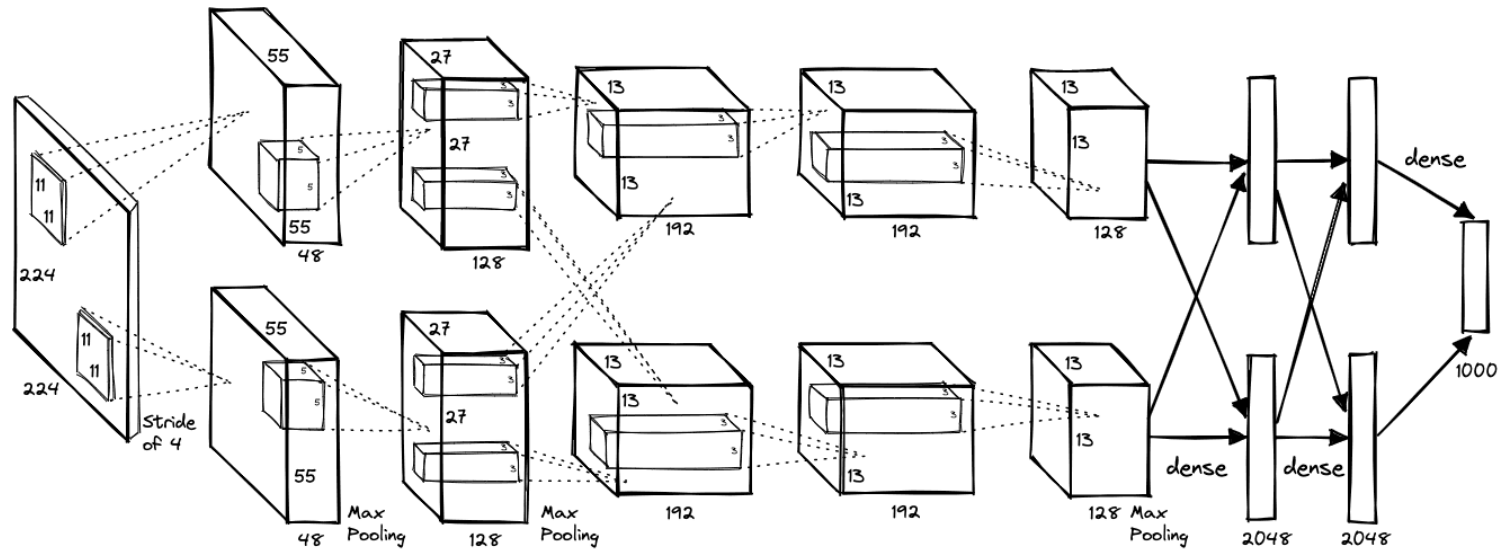
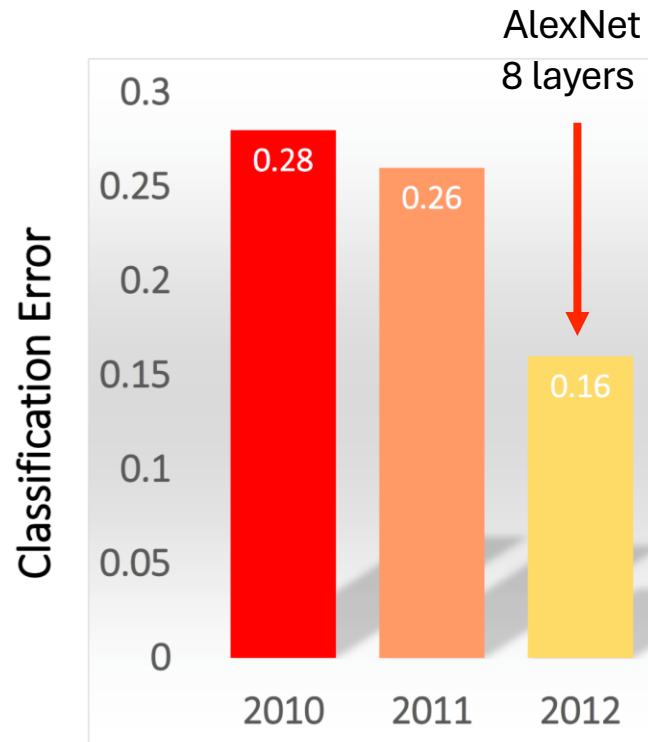hook (66)  spotlight (66)  ladle (65)  restaurant (64)  letter opener (59)

# AlexNet [2012]



- Input = 224x224
- First Layer = 11x11 Conv
- 8 layers
- 60M parameters
- Used 2 Nvidia GTX 580 3GB



**Alex Krizhevsky et al, ImageNet classification with deep convolutional neural networks, NIPS, 2012**

# ImageNet

# VGGNet [2013]



large filters
shallow (8 layers)

small filters
deeper (19 layers)

- 19 layers
- 138M parameters

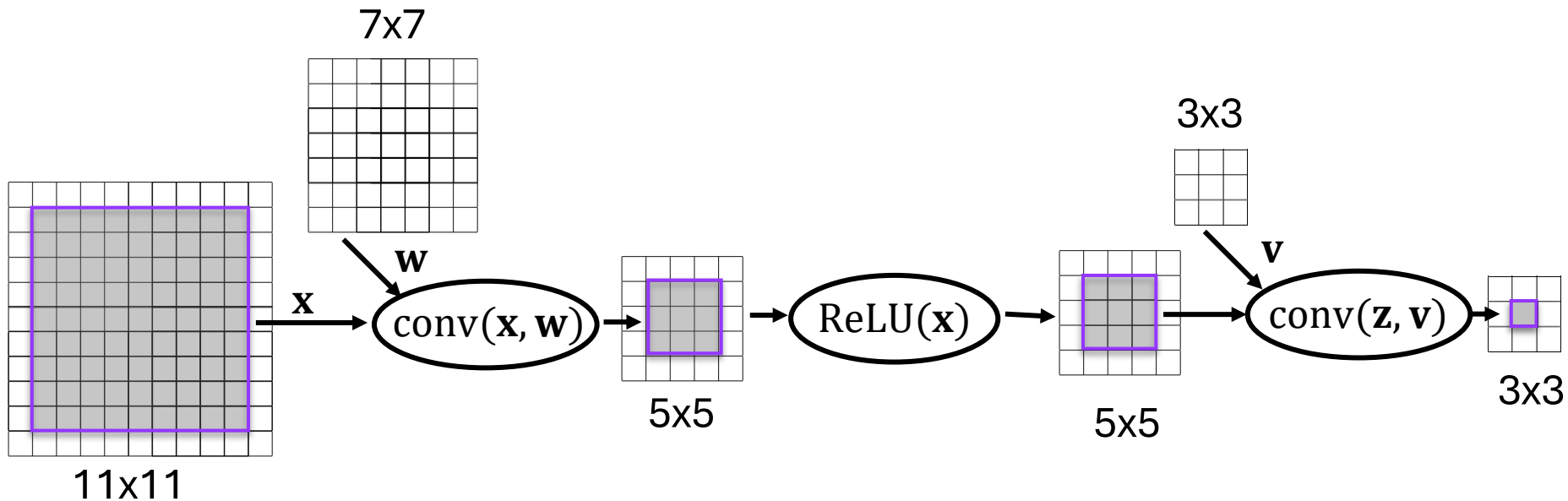Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition.", arXiv 2014

# Receptive field

- Receptive field = area in the image whose values affect given cell (neuron)



receptive field

7x7

$\mathbf{w}$

$\mathbf{x}$

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$
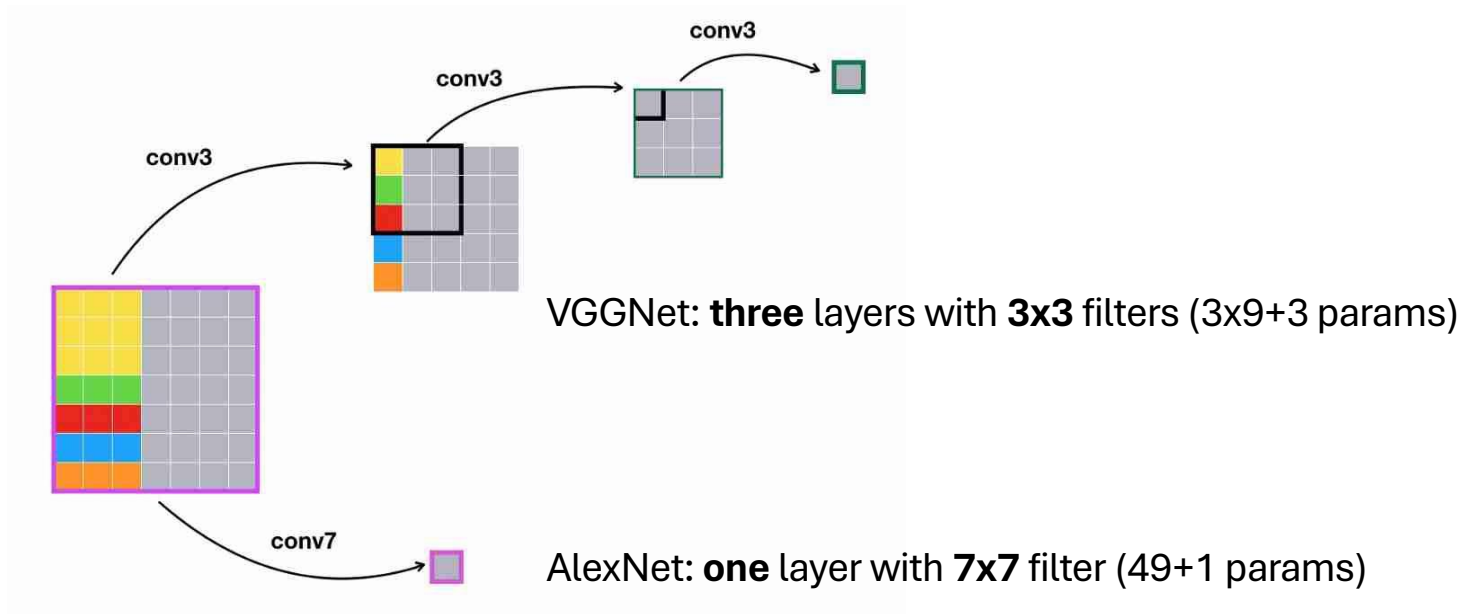
$\mathbf{y}$

5x5

11x11

# Receptive field

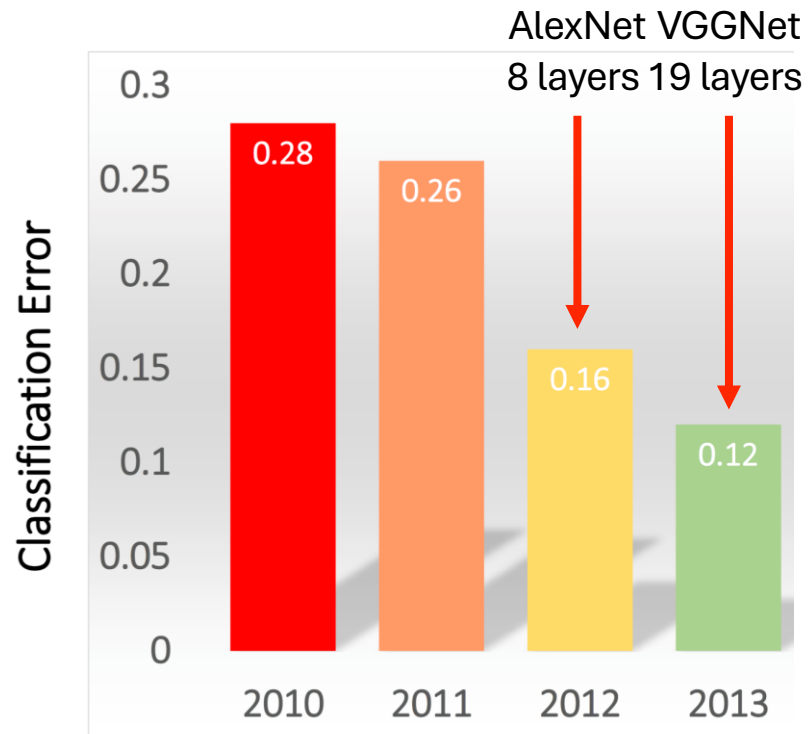- Receptive field = area in the image whose values affected selected cell (neuron)

# VGGNet [2013]

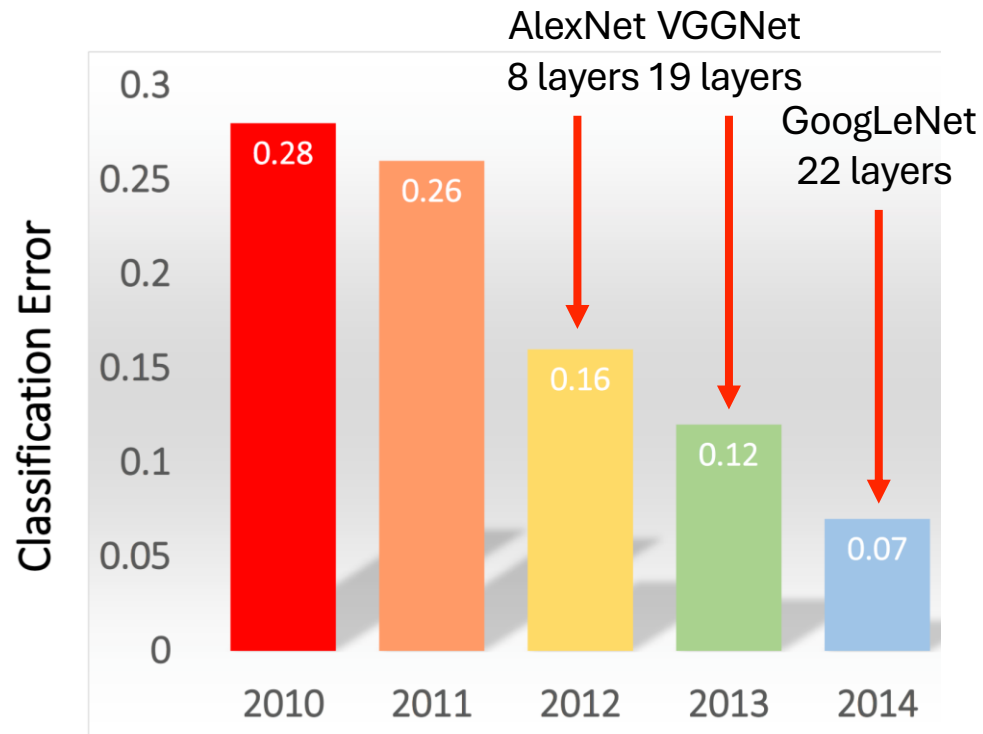- VGGNet has **the same receptive field** with **less parameters**



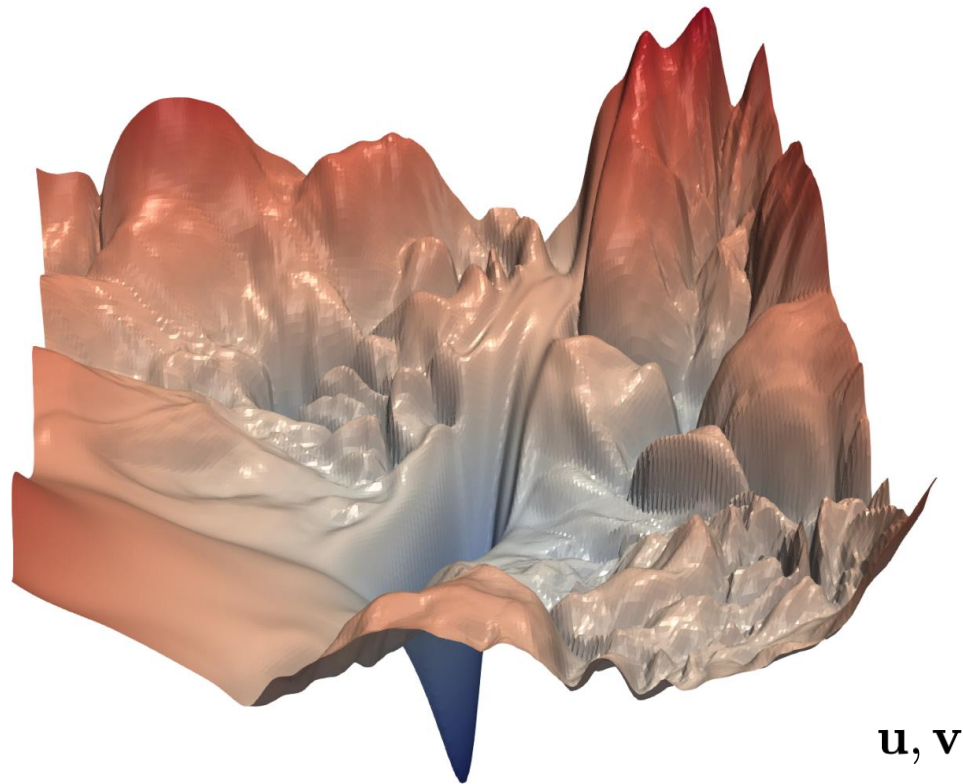VGGNet: **three** layers with **3x3** filters (3x9+3 params)

AlexNet: **one** layer with **7x7** filter (49+1 params)

# ImageNet

# ImageNet

# ResNet [2015]
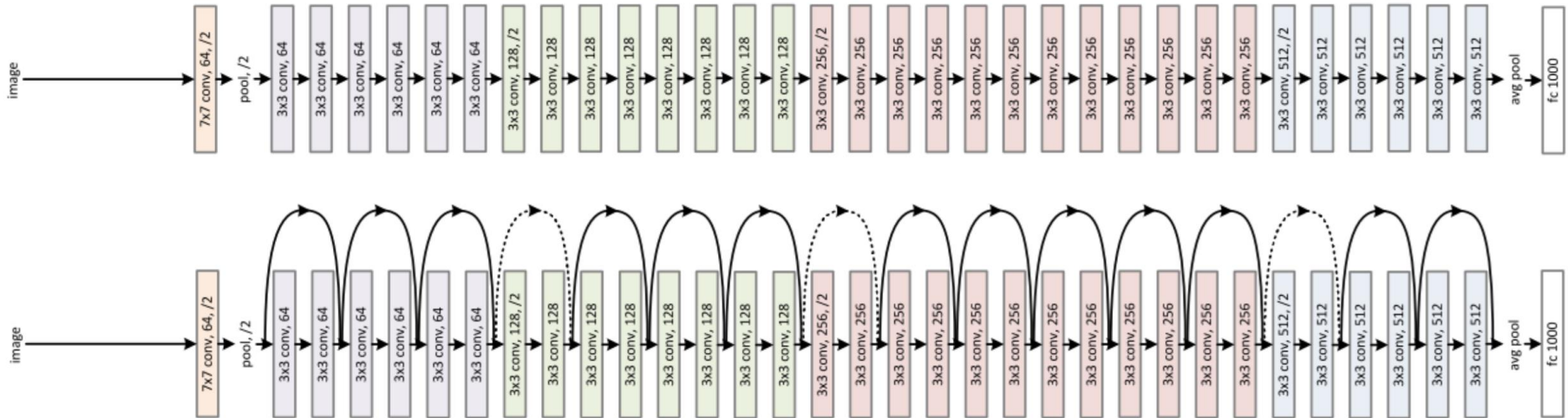


- Deeper architectures had higher training errors

- Is it overfitting?
- No overfitting, but vanishing gradients !

He et al. "Going Deeper with Convolutions", CVPR 2015

# ResNet [2015]

$f(\alpha, \beta) = \mathcal{L}(\mathbf{w}^* + \alpha\mathbf{u} + \beta\mathbf{v})$ for randomly chosen (and normalized) directions
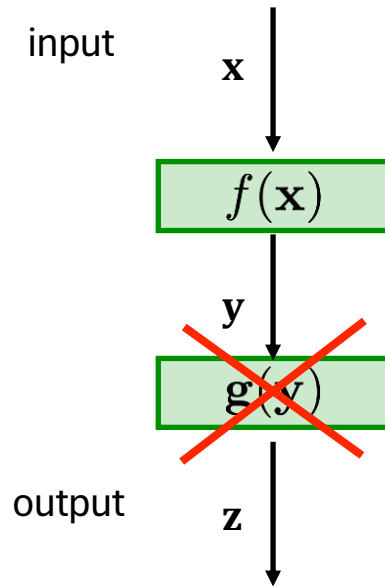


$\mathbf{u}, \mathbf{v}$

Li, Hao, et al. "Visualizing the loss landscape of neural nets.", NIPS 2018

# ResNet [2015]



- ResNet **adds skip connections** to prevent vanishing gradients
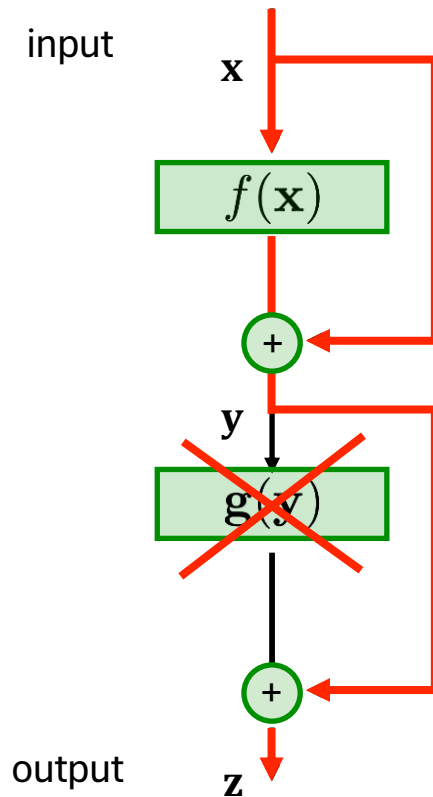- Allows training of very deep networks (e.g. ResNet-152)

He et al. "Going Deeper with Convolutions", CVPR 2015

# ResNet [2015]

**forward pass**

input

$\mathbf{x}$

$f(\mathbf{x})$

$\mathbf{y}$
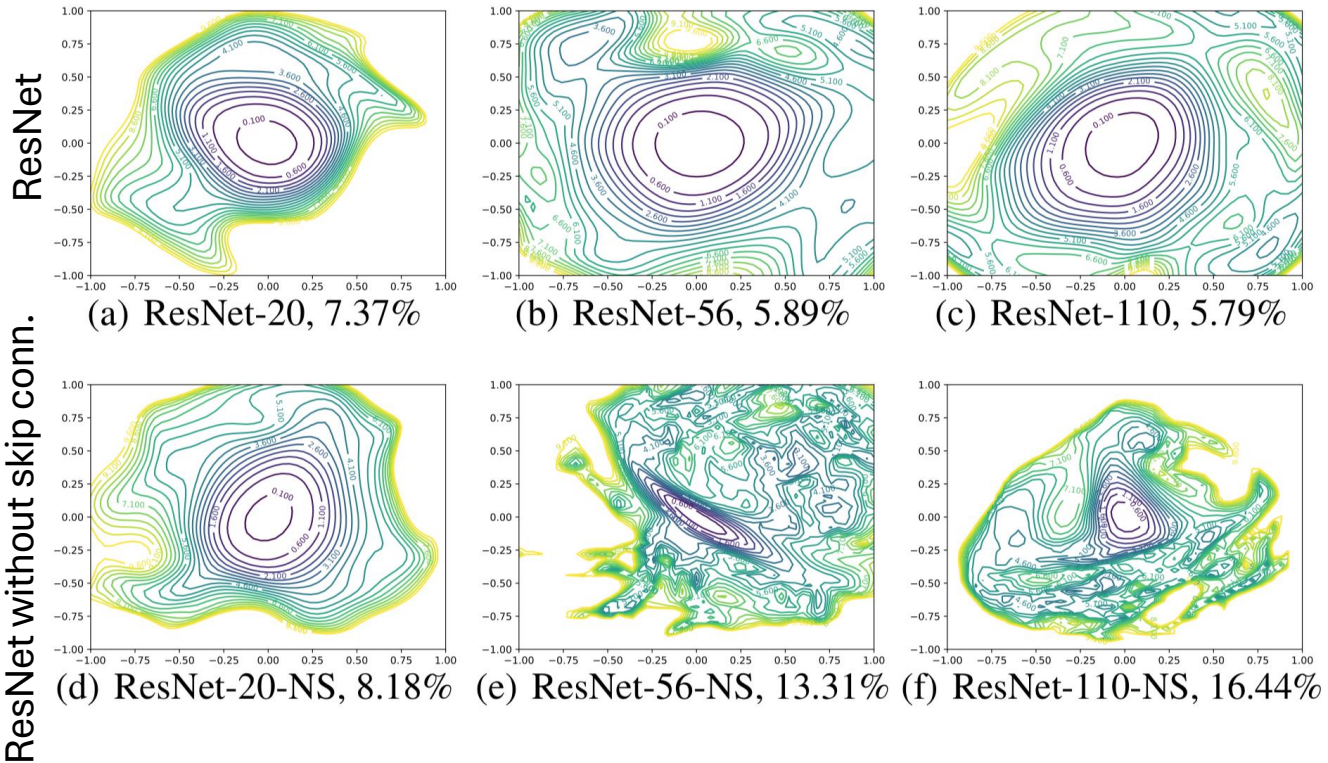
$\mathbf{g}(y)$

output

$\mathbf{z}$

**backward pass**

gradient

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \boxed{\frac{\partial \mathbf{z}}{\partial \mathbf{y}} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}}} \approx \mathbf{0}$$

$\approx \mathbf{0}$ $\qquad \approx \mathbf{0}$

if any local gradient is zero

# ResNet [2015]

input

**x**

$$f(\mathbf{x})$$

+

**y**

~~$$\mathbf{g}(\mathbf{y})$$~~

+

output

**z**

gradient

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = (\frac{\partial \mathbf{z}}{\partial \mathbf{y}} + 1) \quad (\frac{\partial \mathbf{y}}{\partial \mathbf{x}} + 1) \quad \neq \mathbf{0}$$

$$\approx \mathbf{0} \qquad \approx \mathbf{0}$$

if any local gradient is zero

the gradient can still flow through another path

# ResNet [2015]



(a) ResNet-20, 7.37%  (b) ResNet-56, 5.89%  (c) ResNet-110, 5.79%

(d) ResNet-20-NS, 8.18%  (e) ResNet-56-NS, 13.31%  (f) ResNet-110-NS, 16.44%

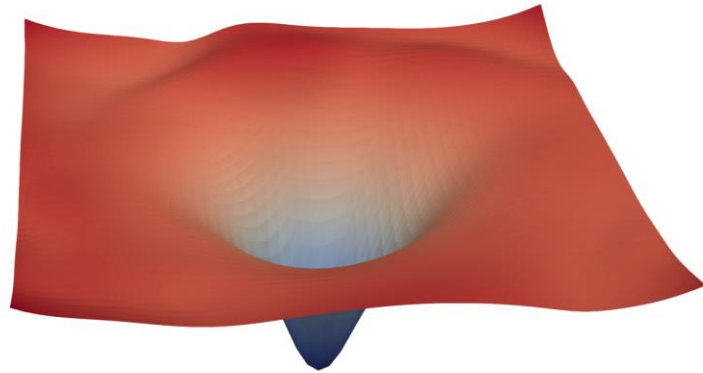He et al. "Going Deeper with Convolutions", CVPR 2015

# ResNet [2015]

$$f(\alpha, \beta) = \mathcal{L}(\mathbf{w}^* + \alpha\mathbf{u} + \beta\mathbf{v})$$ for randomly chosen (and normalized) directions
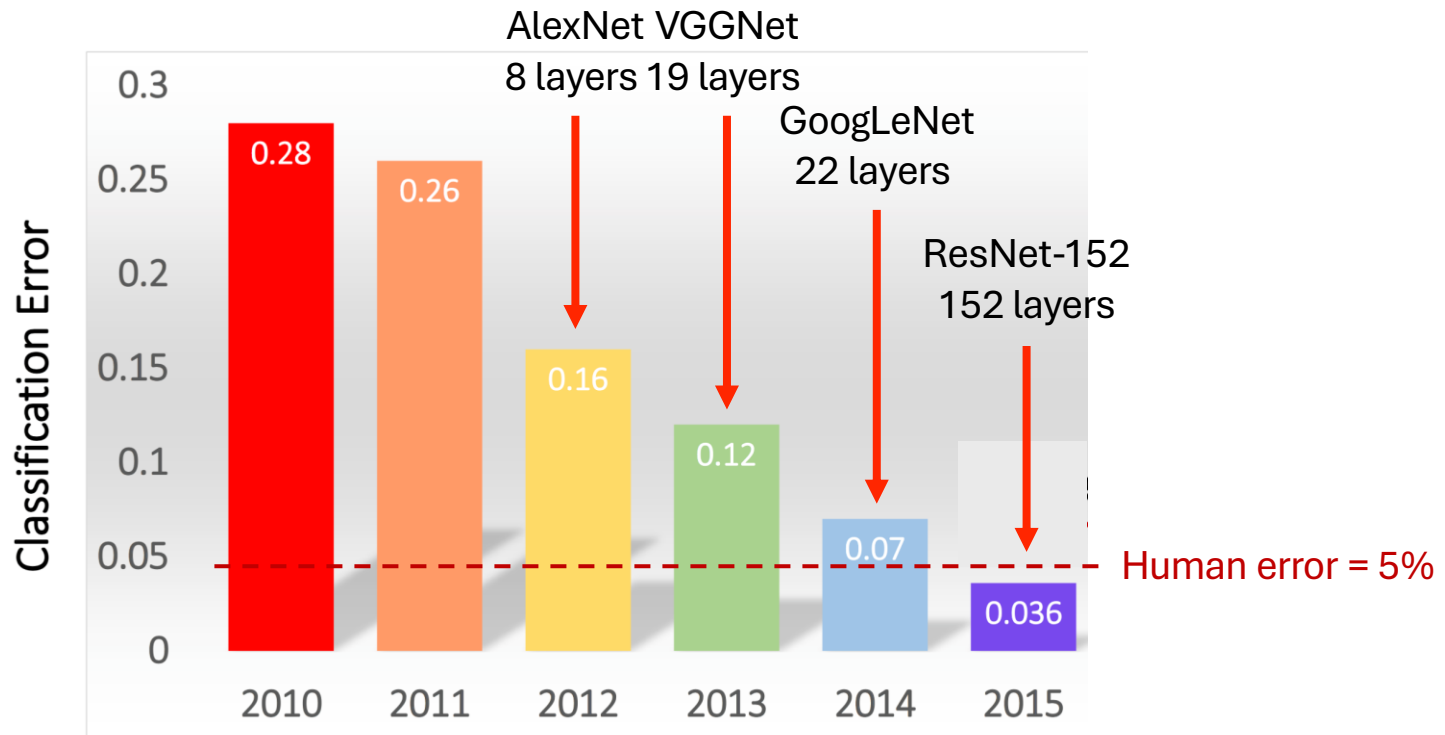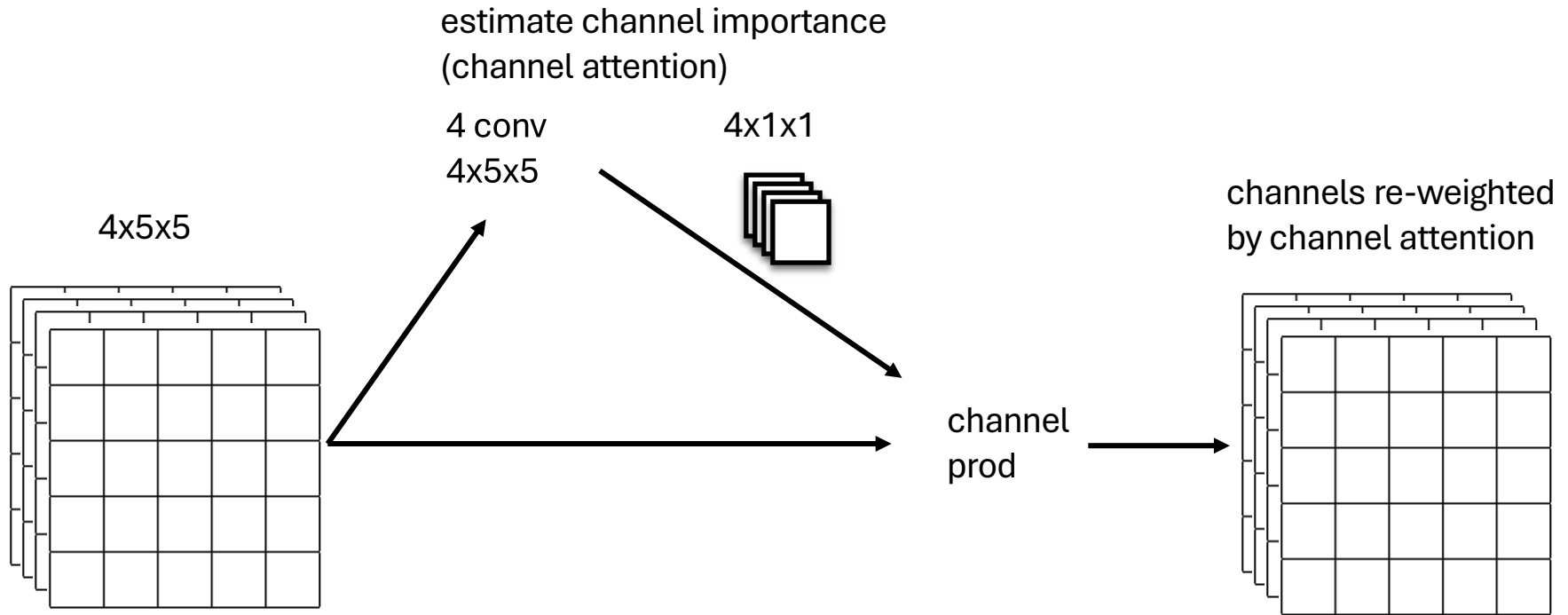


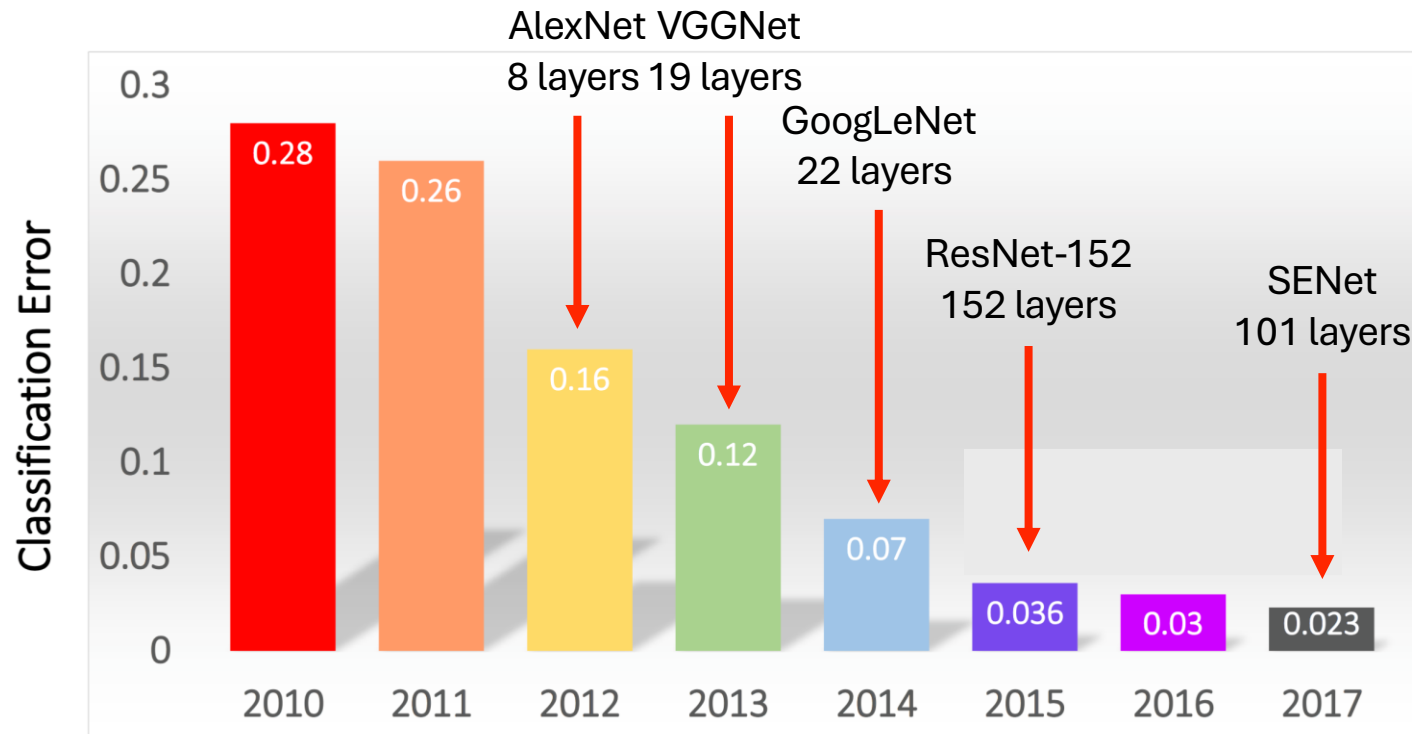(a) without skip connections

(b) with skip connections

Li, Hao, et al. "Visualizing the loss landscape of neural nets.", NIPS 2018

# ImageNet

# **Squeeze and Excitation Networks (SEN)**

- Channel attention

estimate channel importance
(channel attention)

4 conv
4x5x5

4x1x1

4x5x5

channels re-weighted
by channel attention

channel
prod

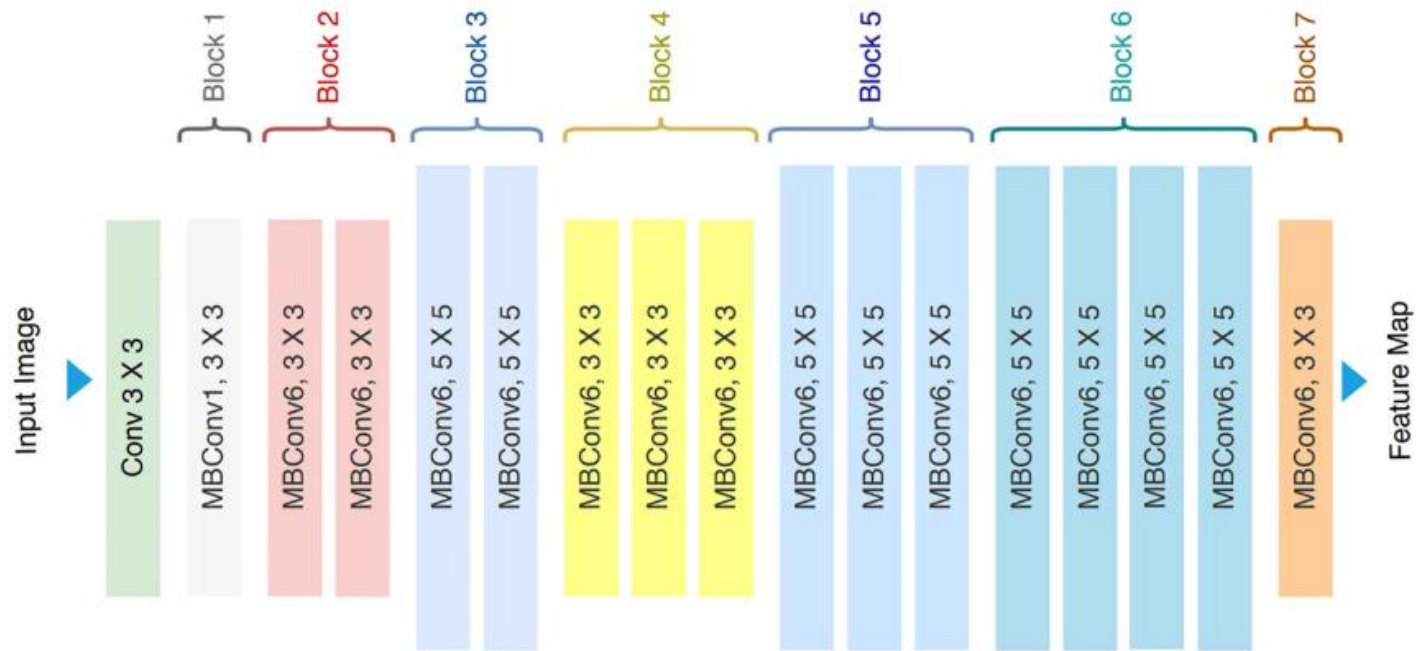**Hu, Jie, Li Shen, and Gang Sun. "Squeeze-and-excitation networks.", CVPR 2018**
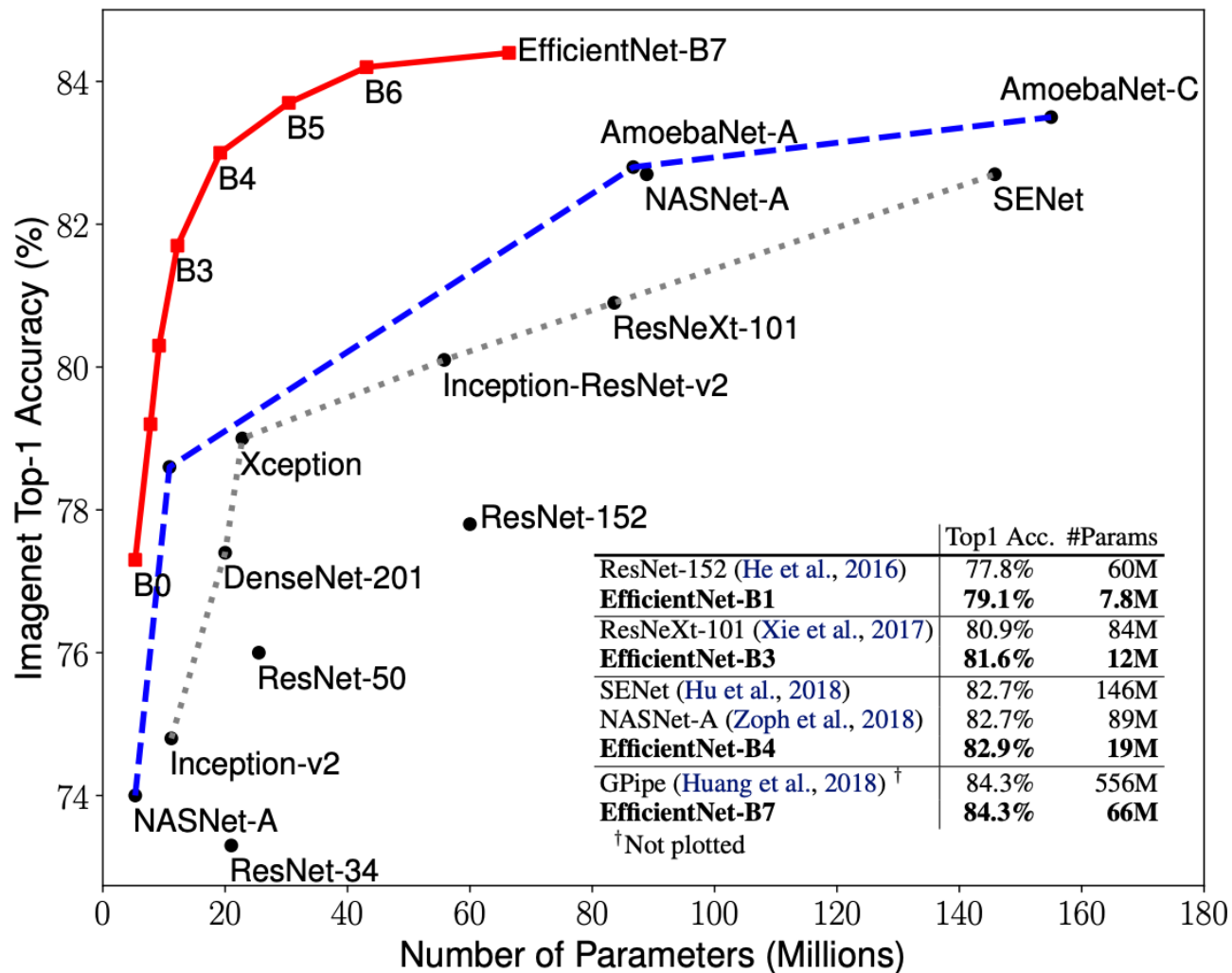
# ImageNet

# EfficientNet [2019]

- Good balance between number of channels, depth and resolution
- Good accuracy with fewer parameters → faster to train, less overfitting
- **In practice, good first choice as a backbone (EfficientNetV2 [2021])**
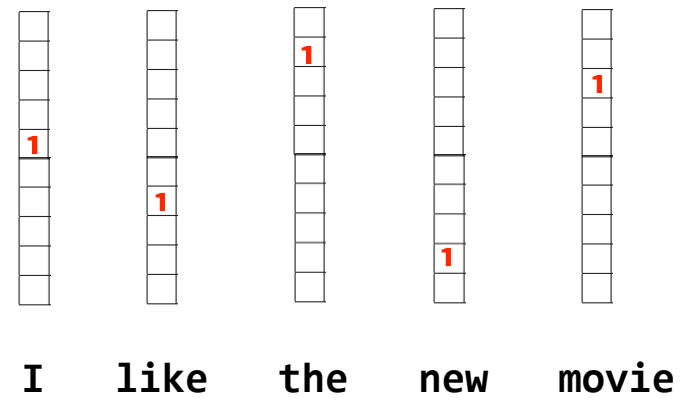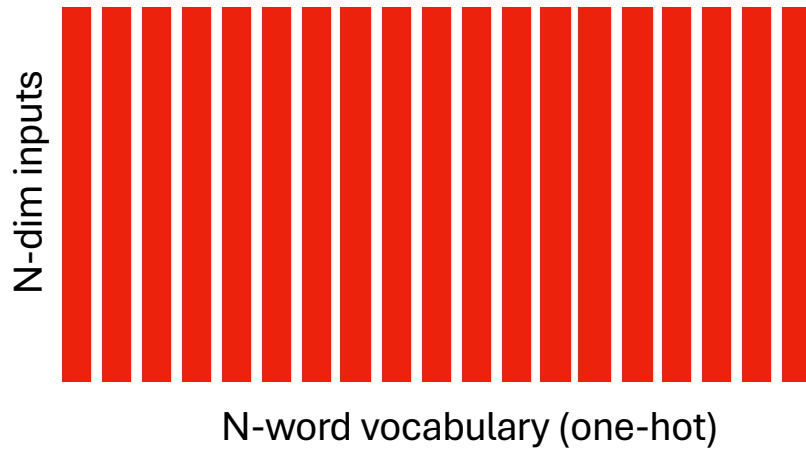
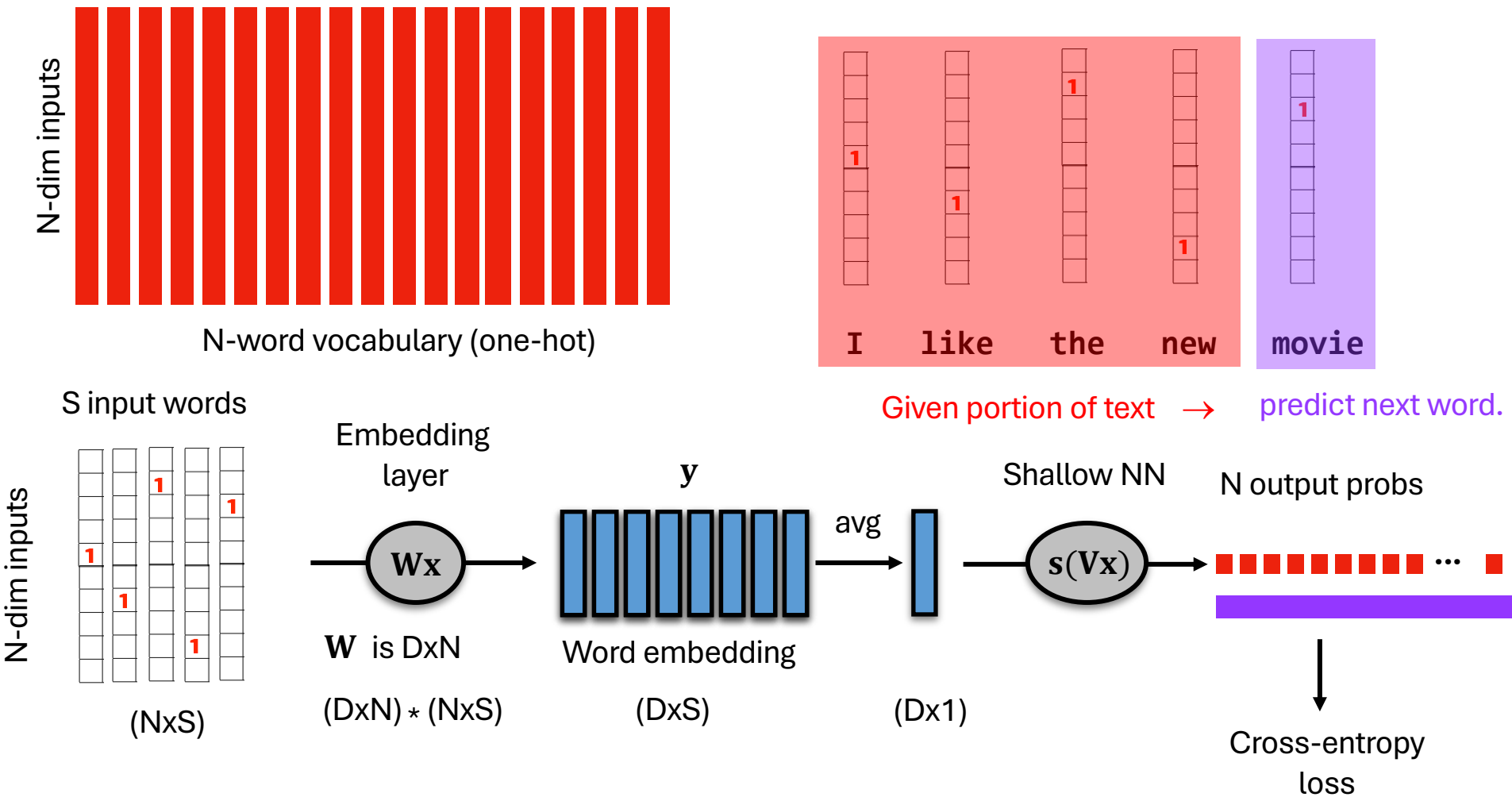Tan, Mingxing, and Q. Efficientnet Le. "Rethinking model scaling for convolutional neural networks." ICML 2019

# EfficientNet [2019]



Tan, Mingxing, and Q. Efficientnet Le. "Rethinking model scaling for convolutional neural networks." ICML 2019

# Word embedding

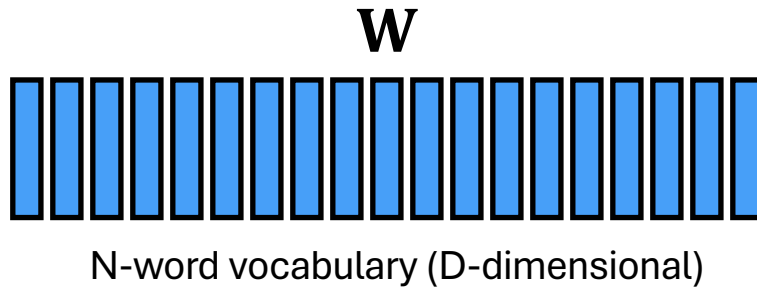- Word2vec represents words in low-dimensional continuous space



N-dim inputs

N-word vocabulary (one-hot)

I   like   the   new   movie

# Word embedding

- Word2vec represents words in low-dimensional continuous space



N-dim inputs

N-word vocabulary (one-hot)

I    like    the    new    movie

Given portion of text  →  predict next word.

S input words

N-dim inputs

(NxS)

Embedding layer

**Wx**

**W** is DxN

(DxN) ∗ (NxS)

**y**

Word embedding

(DxS)

avg

(Dx1)

Shallow NN

**s(Vx)**

N output probs

... 

Cross-entropy loss

# Word embedding

- Word2vec represents words in low-dimensional continuous space

**W**

N-word vocabulary (D-dimensional)

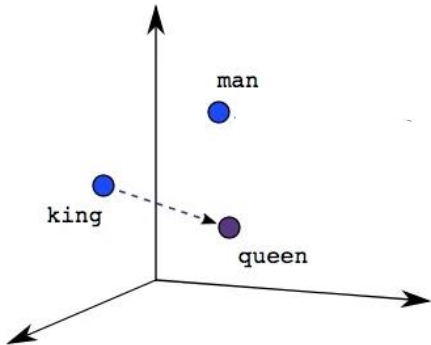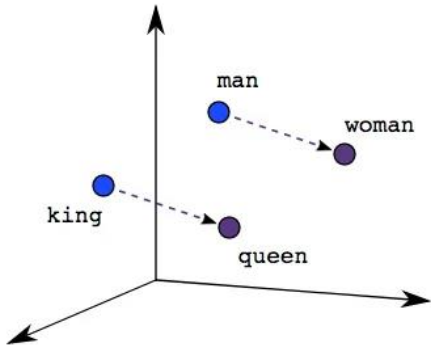| I | like | the | new | movie |
|------|------|------|------|------|
| 0.35 | 0.05 | 0.39 | 0.15 | 0.3 |
| 0.3 | 0.01 | 0.58 | 0.2 | 0.98 |
| 0.1 | 0.56 | 0.01 | 0.56 | 0.66 |
| 0.25 | 0.22 | 0.36 | 0.99 | 0.15 |

# Word embedding

- Word2vec represents words in low-dimensional continuous space

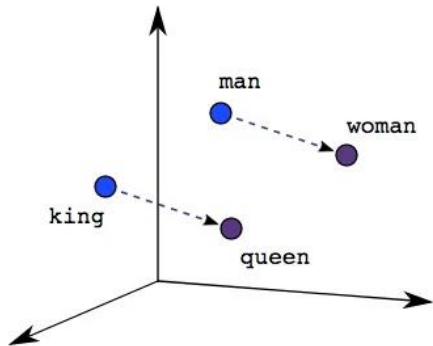- Word2vec represents words in low-dimensional continuous space



man

king

queen

Male-Female

# **Word embedding**
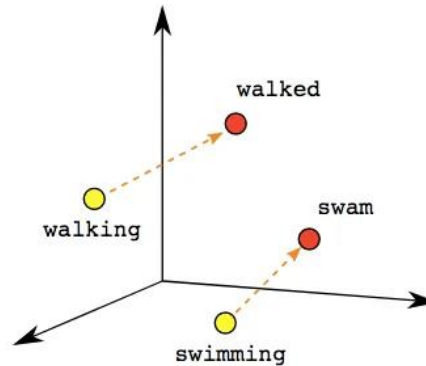
- Word2vec represents words in low-dimensional continuous space



Male-Female

"Word algebra": `king - man + woman = queen`

# Word embedding
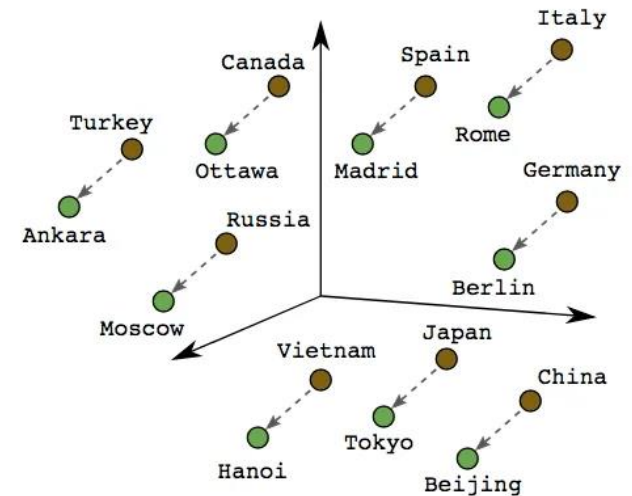
- Word2vec represents words in low-dimensional continuous space



Male-Female          Verb Tense          Country-Capital
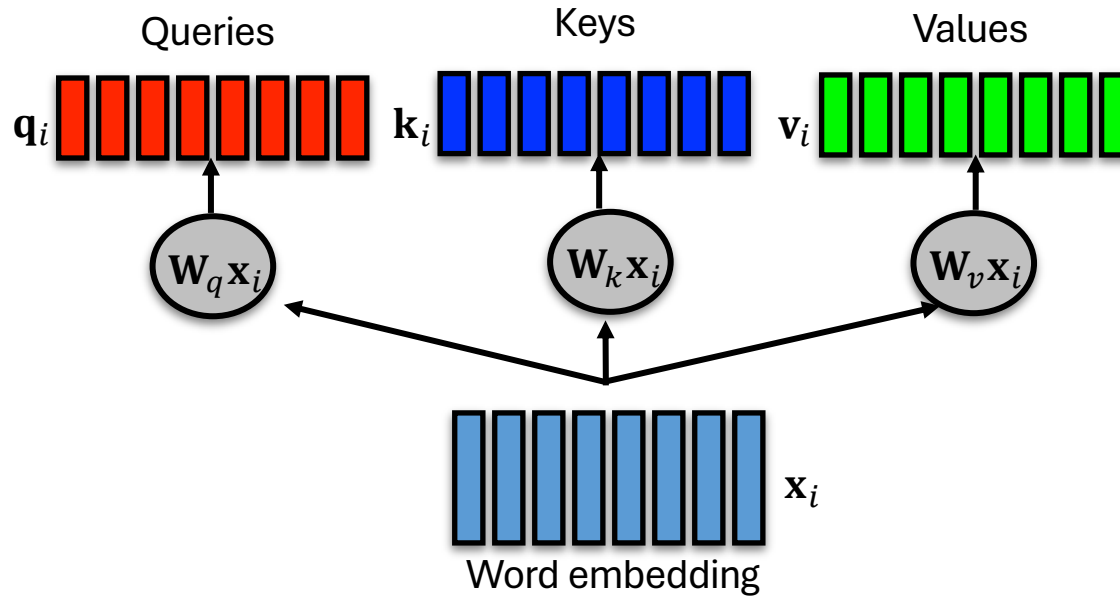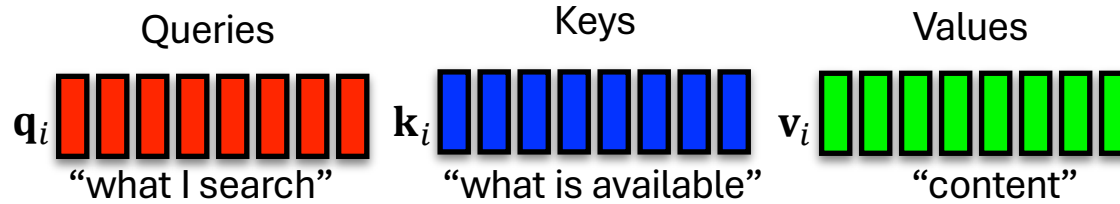
Queries  Keys  Values

$\mathbf{q}_i$ ▮▮▮▮▮▮▮▮▮  $\mathbf{k}_i$ ▮▮▮▮▮▮▮▮  $\mathbf{v}_i$ ▮▮▮▮▮▮▮▮

$\mathbf{W}_q\mathbf{x}_i$  $\mathbf{W}_k\mathbf{x}_i$  $\mathbf{W}_v\mathbf{x}_i$

▮▮▮▮▮▮▮▮ $\mathbf{x}_i$

Word embedding

# Self-attention

Queries | Keys | Values

$q_i$ "what I search" | $k_i$ "what is available" | $v_i$ "content"

Karel is teacher and Mario is plumber.

Which words contribute
to meaning of Karel?

$q_1$ $k_1$ $k_2$ $k_3$ $k_4$ $k_5$ $k_6$ $k_7$

$$\blacksquare = \blacksquare \times \blacksquare$$

$q_1^\top$ $k_1$

$q_1$
$k_3$

Scalar product measures similarity between vectors.

$k_5$ $k_7$

# Self-attention

Queries

$\mathbf{q}_i$ ▮▮▮▮▮▮▮

"what I search"

Keys

$\mathbf{k}_i$ ▮▮▮▮▮▮▮

"what is available"

Values

$\mathbf{v}_i$ ▮▮▮▮▮▮▮

"content"

**Attention**

$$\text{softmax} \left( \ \underset{\mathbf{q}_i^\top}{▬} \ \times \ \underset{\substack{\mathbf{k}_1 \quad \cdots \quad \mathbf{k}_n}}{▮▮▮▮▮▮▮} \ \right) = \underset{\substack{\mathbf{a}_1 \quad \cdots \quad \mathbf{a}_n}}{■■■■■■■}$$

**Attention-weighted sum of Values**

$$\underset{\substack{\mathbf{a}_1 \quad \cdots \quad \mathbf{a}_n}}{■■■■■■■} \ \underset{\substack{\mathbf{v}_1 \\ \mathbf{v}_n}}{▭} = ■▮+■▮+■▮+■▮+■▮+■▮+■▮+■▮ = \underset{\mathbf{y}_1}{▮}$$

**Output**

$\underset{\mathbf{y}_1}{▮}$

# Self-attention

Queries     Keys     Values

$\mathbf{q}_i$ ▮▮▮▮▮▮▮▮    $\mathbf{k}_i$ ▮▮▮▮▮▮▮▮    $\mathbf{v}_i$ ▮▮▮▮▮▮▮▮

"what I search"    "what is available"    "content"

**Attention**

$$\text{softmax} \left( \mathbf{Q}^\top \times \mathbf{K} \right) =$$

**Attention-weighted sum of values**

$$\mathbf{A} \times \mathbf{V}^\top = \text{softmax}(\mathbf{Q}^\top \mathbf{K})\mathbf{V}^\top =$$

**Output**

$$\mathbf{Y}$$

# Self-attention

# Self-attention

| | The |
|---|---|
| The | The |
| Doctor | Doctor |
| asked | asked |
| the | the |
| Nurse | Nurse |
| a | a |
| question | question |
| . | . |
| **She** | She |
| said | said |

Model assumes "she=nurse"

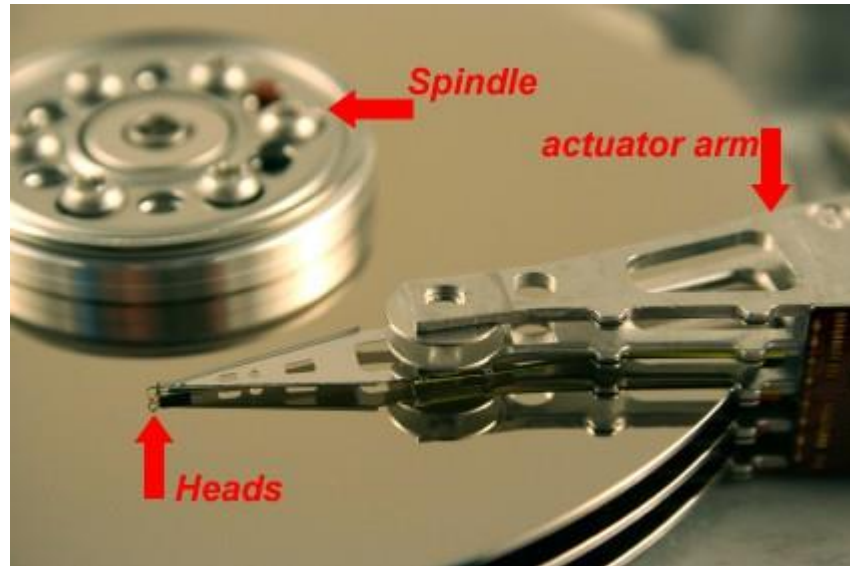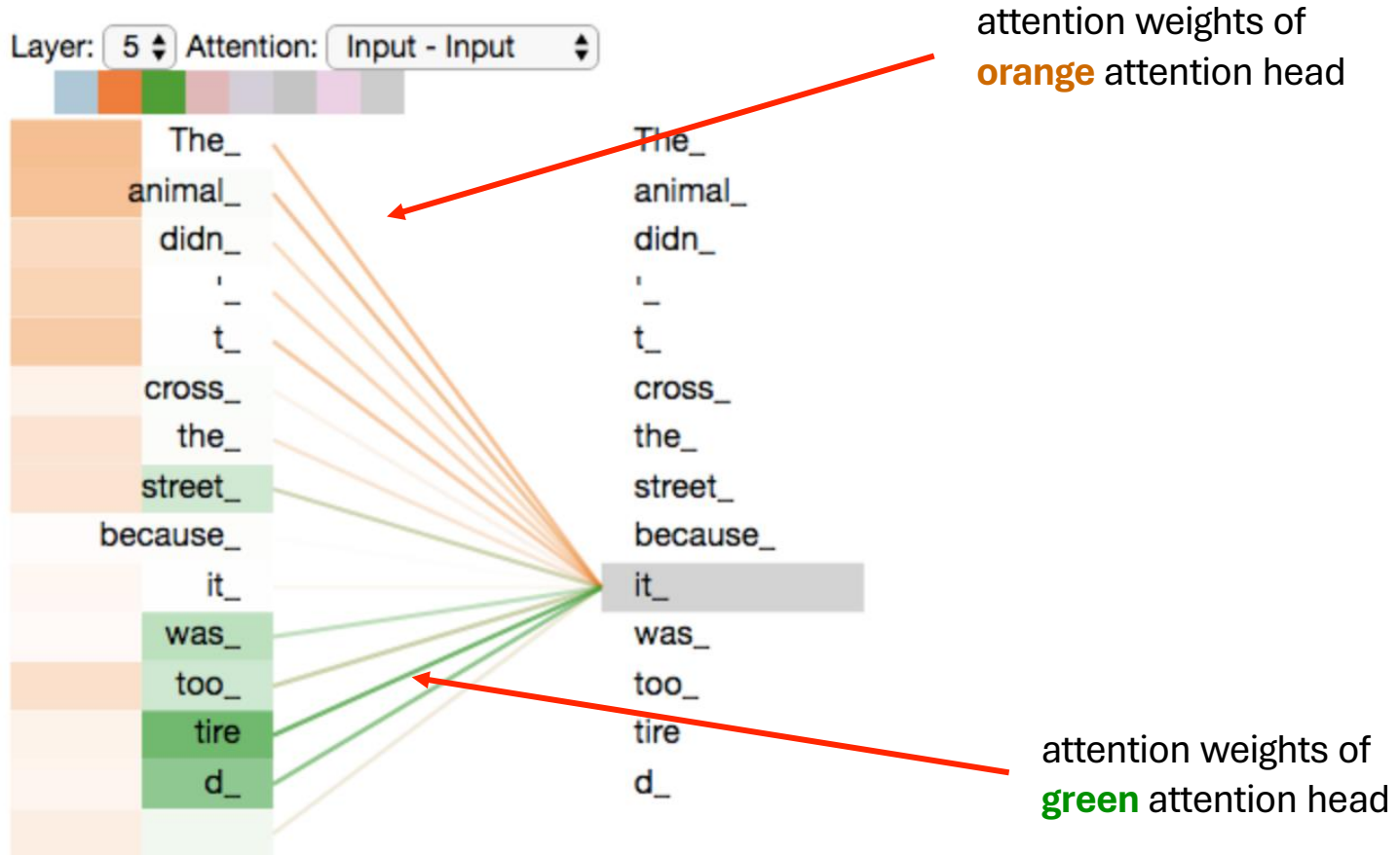| | The |
|---|---|
| The | The |
| Doctor | Doctor |
| asked | asked |
| the | the |
| Nurse | Nurse |
| a | a |
| question | question |
| . | . |
| **He** | He |
| asked | asked |

Model assumes "he=doctor"

# Multi-head self-attention (MHSA)

- Self-attention can be applied multiple times on the same sequence in parallel using more **heads**
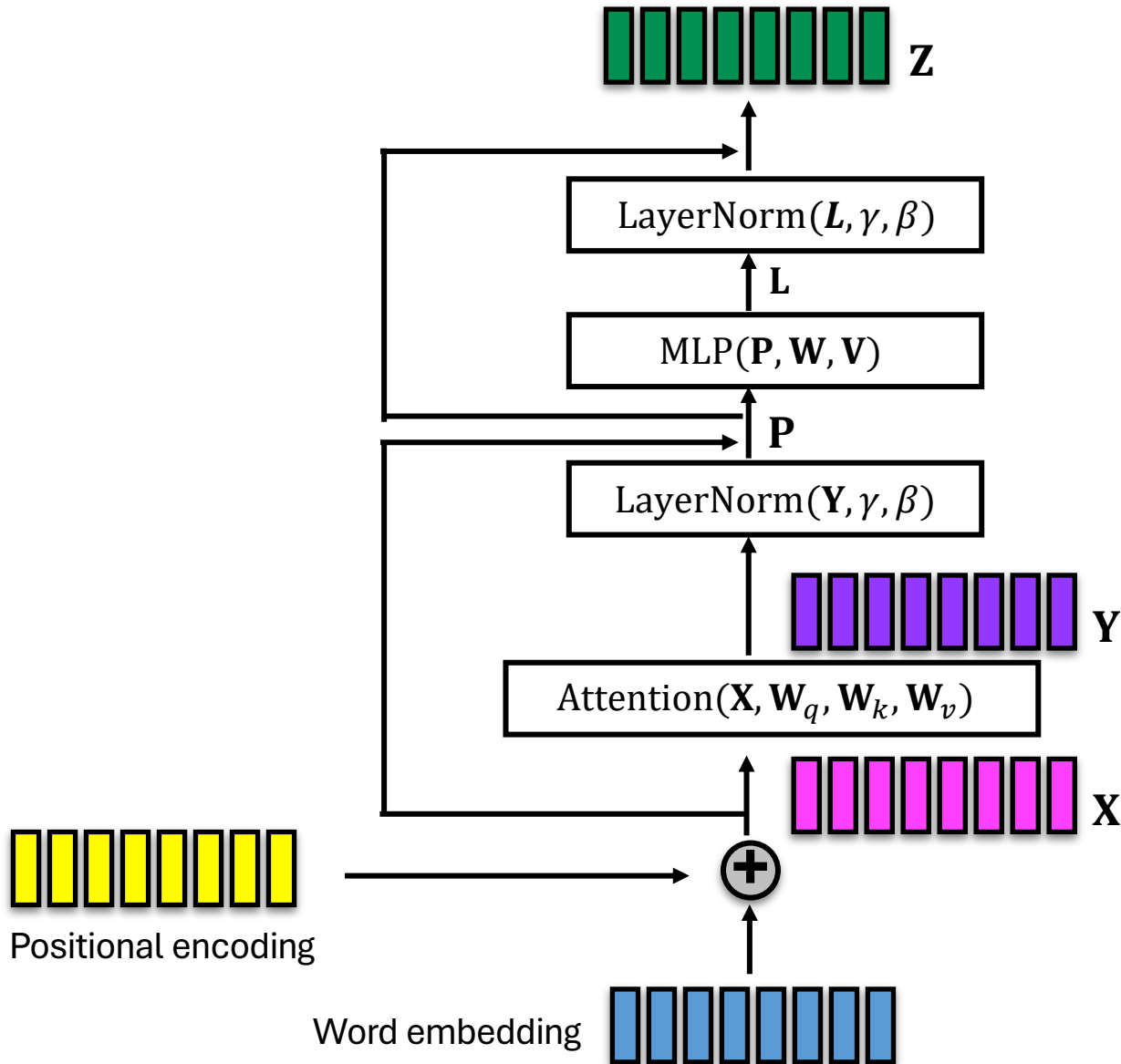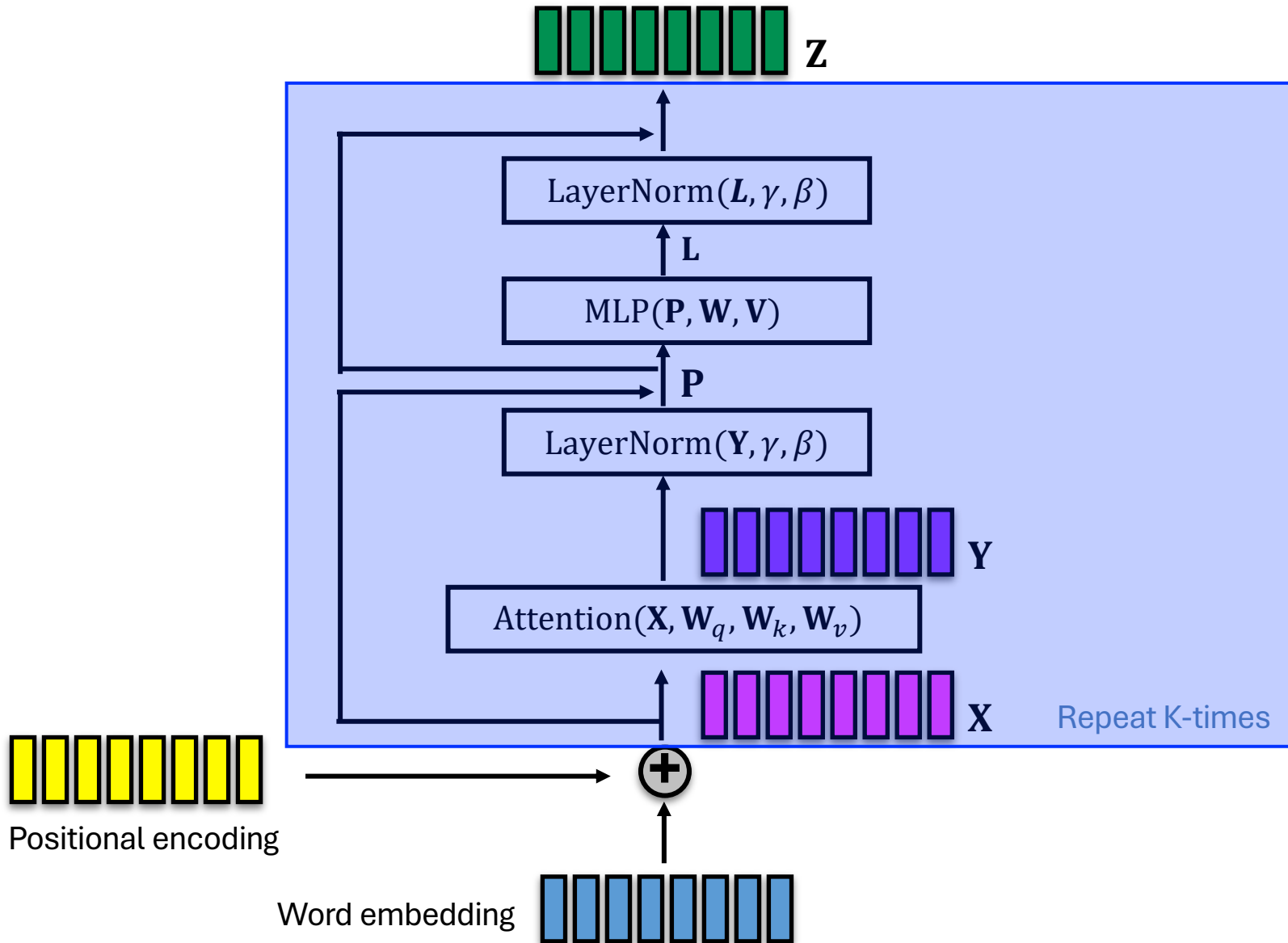
# Multi-head self-attention (MHSA)

The animal didn't cross the street because **it** was too tired.
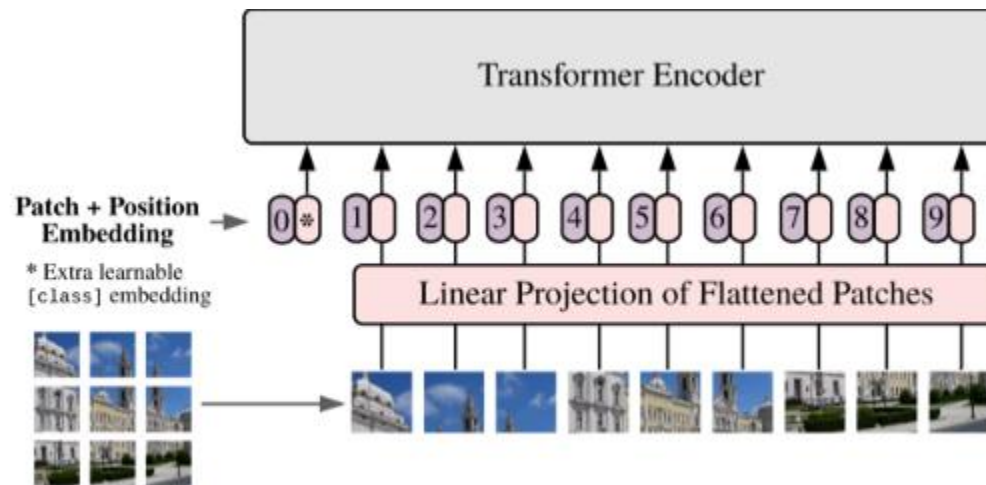
**Is "it" = animal vs "it" = street ?**



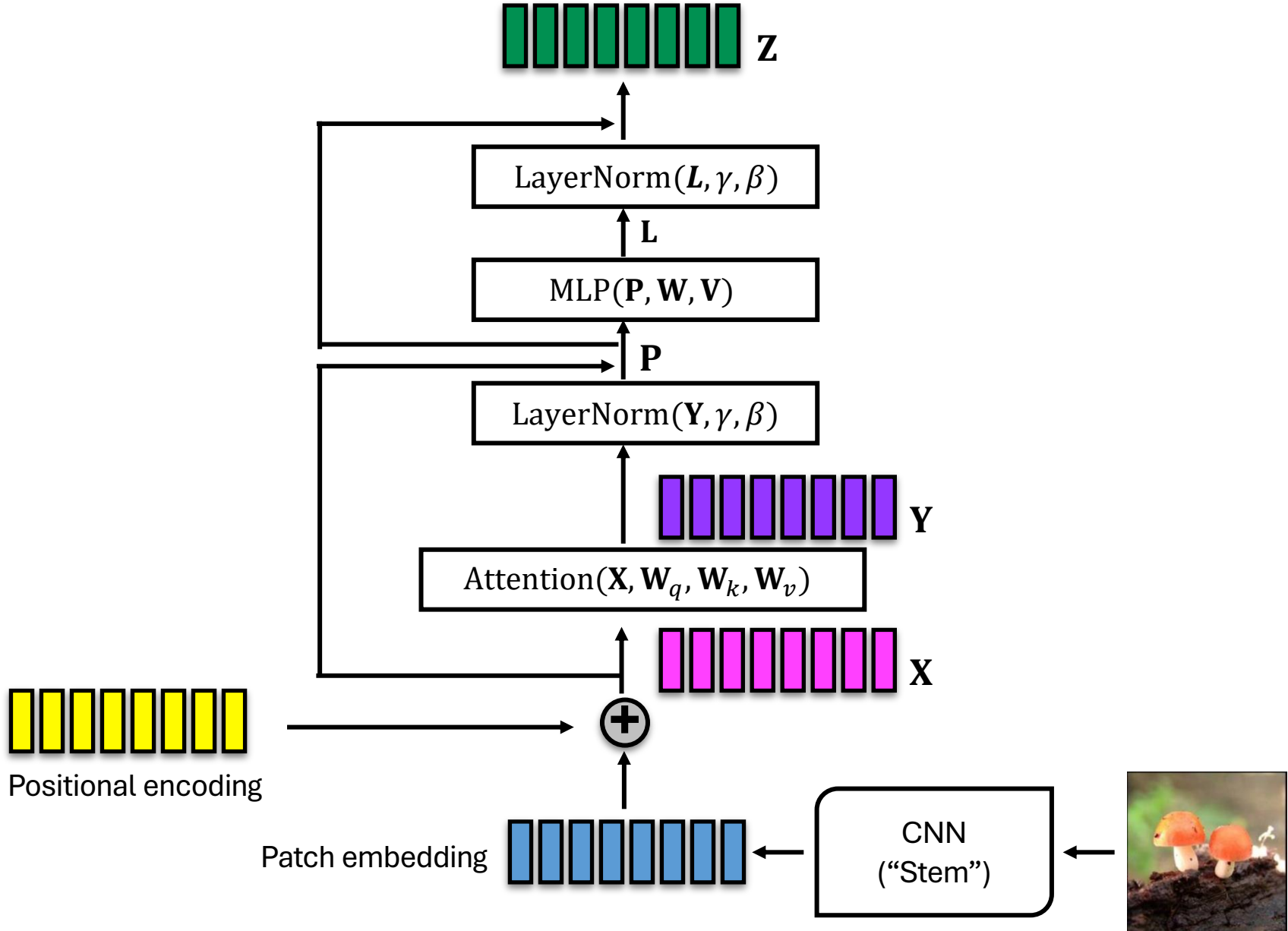attention weights of **orange** attention head

attention weights of **green** attention head

# Transformer

# Transformer



Z

LayerNorm($L, \gamma, \beta$)

L

MLP($\mathbf{P}, \mathbf{W}, \mathbf{V}$)

P

LayerNorm($\mathbf{Y}, \gamma, \beta$)

Y

Attention($\mathbf{X}, \mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v$)

X

Repeat K-times

Positional encoding

Word embedding

# Vision Transformer (ViT)

- Split image into 16x16 patches and treat each patch as a "word"



Dosovitskiy, Alexey et al. "An image is worth 16x16 words: Transformers for image recognition at scale.", arXiv 2020

# Vision Transformer (ViT)

# Vision Transformer (ViT)

# Competencies gained for the test

- Vanishing gradient problem, ResNet
- Self-attention, Transformers
- Vision Transformer