

# NEURAL NETWORKS

# SPEED-UP AND COMPRESSION

**Lecture 1: Introduction**

# THE TEAM



**Julia Gusak**, PhD  
Senior research scientist, Skoltech  
(main collaborations: Prof. Oseledets, Prof. Cichocki, Prof. Phan)



**Stanislav Abukhovich**  
PhD student, Skoltech  
(supervisor: Prof.Cichocki)



**Dmitry Ermilov**  
PhD student, Skoltech  
(supervisor: Prof.Cichocki)



**Konstantin Sobolev**  
PhD student, Skoltech  
(supervisor: Prof.Phan)

# OUTLINE

- Motivation for neural networks speed-up and compression
- Estimation of neural networks effectiveness
- Overview of main compression techniques
- Course schedule, outcomes, and completion criteria

# MOTIVATION

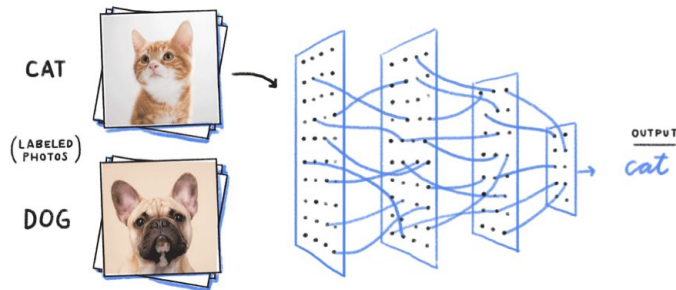
# WHAT IS A NEURAL NETWORK?

- Function that is defined by its parameters.
- Inference phase: given an input sample neural network provides the output.
- Training phase: neural networks' parameters are tuned using a set of input samples, loss function, and iterative gradient-based parameter updates.

# EXAMPLE: COMPUTER VISION

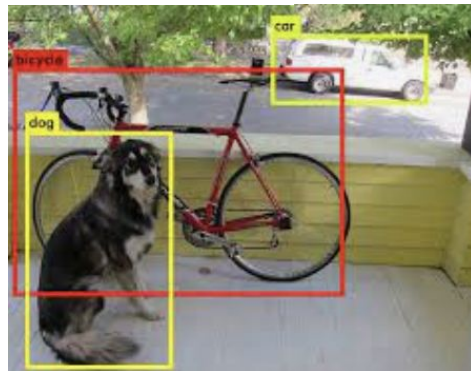
- Image recognition

- image  $\rightarrow$  NN  $\rightarrow$  class label



- Object detection

- image  $\rightarrow$  NN  $\rightarrow$  image with frames around objects



- Action recognition

- video frame  $\rightarrow$  NN  $\rightarrow$  action description

# EXAMPLE: SPEECH RECOGNITION

- Speaker recognition
  - audio frame  $\rightarrow$  NN  $\rightarrow$  speaker ID
- Speech recognition
  - audio frame  $\rightarrow$  NN  $\rightarrow$  text corresponding to the audio frame
- Speech generation
  - audio frame  $\rightarrow$  NN  $\rightarrow$  next audio frame

# EXAMPLE: NATURAL LANGUAGE PROCESSING

- Sentiment analysis
  - sentence → NN → positive/negative
- Question & Answer
  - sentence-question → NN → sentence-answer
- Translation
  - sentence → NN → sentence translated into foreign language
- Text generation
  - sentence → NN → next sentence



# WHAT'S THE PROBLEM?

- Most state of the art deep neural networks are overparameterized and exhibit a high computational cost
- The size of neural networks affects
  - power and memory consumption
  - running time
  - CO2 emission
  - money

# WHAT'S THE PROBLEM?

DL model limitations:

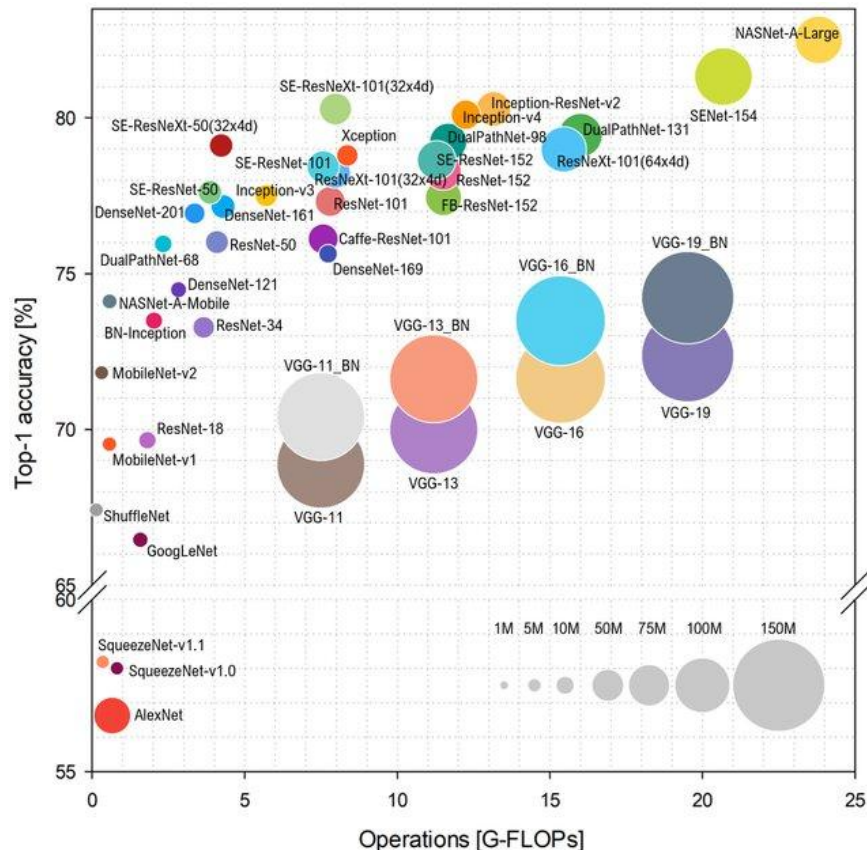
- High memory consumption
- Huge computational requirements
- Great power consumption



Difficult to deploy on portable devices  
(e.g. laptops and smartphones)



Efficient architecture design is required



# WHAT'S THE PROBLEM?

Energy: 190 MW\*hours



Air emissions : 85 tonnes of CO<sub>2</sub>



GPT-3 training



Heating for 126 houses  
in Denmark



Car drive to the  
Moon

# DIFFERENT MODALITIES, SIMILAR TECHNIQUES

- Image
- Video
- Audio
- Text
- Medical signals (EEG, ECG)
- Physical measurements

-> NN -> prediction

! Similar Deep Learning techniques are used to solve a vast variety of tasks

# WHAT WE WANT FROM NEURAL NETWORKS?

- High predictive quality
- Robustness to the shifts in input data
- Efficient training (fast convergence to optimal parameters)
- **Efficient inference** (fast execution, low memory usage)

# ESTIMATION OF EFFECTIVENESS

# REPRESENTATION OF NEURAL NETWORK PARAMETERS

- Bits & Bytes
  - bit: 0 or 1
  - 1 B(=byte) = 8 bits
  - 1 KB = 1024(=2<sup>10</sup>) bytes
  - 1 MB = 1024 KB = 1024\*1024 bytes
  - 1 GB = 1024 MB = 1024\*1024\*1024 bytes
- Types of representations
  - int (sign|absolute value)
    - int8: 8 bits = 1 byte
    - int32: 32 bits = 4 bytes
  - float (sign|exponent|mantissa)
    - float32: 32 bits = 4 bytes
    - float64: 64 bits = 8 bytes

# KEY FACTORS

- Speed
  - Theoretical: FLOP, MAC, FLOP/second, Bytes/second
  - Empirical: wall-clock inference time, throughput
- Memory
  - Theoretical: (# of model parameters) \* (# of bytes in one parameter)
  - Empirical: allocated memory



# FLOP & MAC

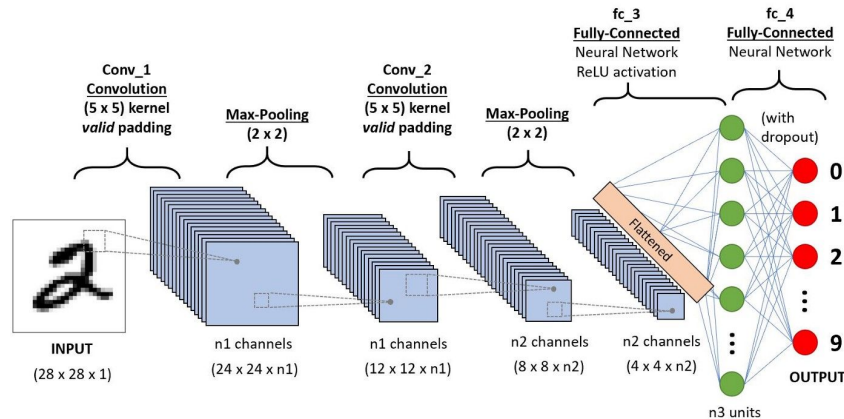
- FLOP - floating point operations
- MAC - multiply-accumulate operations

Example:  $c = c + (a * b)$ , where  $a, b, c$  - scalars

- 2 FLOP
- 1 MAC

# RECAP: NEURAL NETWORKS

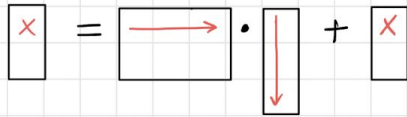
- Main operations
  - Linear
  - Non-linear
- Main neural networks types
  - Feed-forward
  - Residual
  - Attention-based



# LINEAR (FULLY-CONNECTED) LAYER

$X$  - input,  $N_{in} \times 1$   
 $Y$  - output,  $N_{out} \times 1$   
 $W$  - weight matrix,  $N_{in} \times N_{out}$   
 $b$  - bias vector,  $N_{out} \times 1$

$$Y = W^T X + b$$



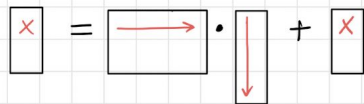
$$\text{Params}_{\text{Linear}} = N_{in} N_{out} + N_{out}$$

# LINEAR (FULLY-CONNECTED) LAYER: FLOP

Example

$$N_{in} = 3, N_{out} = 2 \Rightarrow \begin{array}{ll} X: 3 \times 1 & W: 3 \times 2 \\ Y: 2 \times 1 & b: 2 \times 1 \end{array}$$

$$Y = W^T X + b$$



$$\text{Params}_{\text{Linear}} = N_{in} N_{out} + N_{out} = 3 \cdot 2 + 2 = 8$$

FLOP<sub>Linear</sub> - ?

$$y_i = \underbrace{(w_{i1} \cdot x_1 + w_{i2} \cdot x_2 + w_{i3} \cdot x_3)}_{\substack{\text{mult: } 3 \\ \text{add: } 2}} + \underbrace{b_i}_{\text{add: } 1}$$

$$\text{FLOP}_{y_i} = N_{in} + (N_{in} - 1) + 1 = 2 N_{in}$$

$$\text{FLOP}_{\text{Linear}} = 2 N_{in} N_{out}$$

# LINEAR (FULLY-CONNECTED) LAYER: MAC

Example

$$N_{in} = 3, N_{out} = 2 \Rightarrow$$

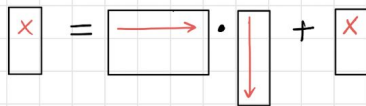
$$X: 3 \times 1$$

$$Y: 2 \times 1$$

$$W: 3 \times 2$$

$$b: 2 \times 1$$

$$Y = W^T X + b$$



MAC  
Linear - ?

$$y_i := b_i$$

$$y_i = y_i + w_{i1} \cdot x_1$$

$$y_i = y_i + w_{i2} \cdot x_2$$

$$y_i = y_i + w_{i3} \cdot x_3$$

$$MAC_{y_i} = N_{in}$$

$$MAC_{Linear} = N_{in} N_{out}$$

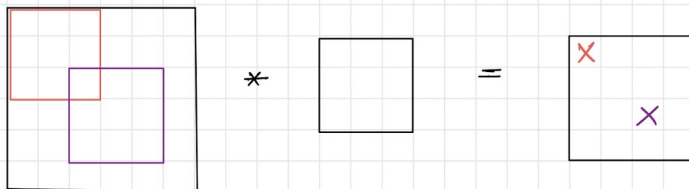
# CONVOLUTIONAL LAYER

$X$  - input,  $C_{in} \times H_{in} \times W_{in}$   
 $Y$  - output,  $C_{out} \times H_{out} \times W_{out}$   
 $W$  - weight (kernel),  $C_{out} \times C_{in} \times k_h \times k_w$   
 $b$  - bias,  $C_{out}$

$$Y = X * W + b$$

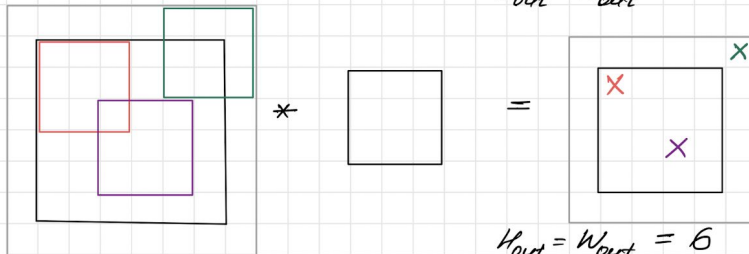
example

$C_{in} = C_{out} = 1$   
 $H_{in} = W_{in} = 6$   
 $k_h = k_w = 3$



$H_{out} = W_{out} = 4$

$padding_h = 1$   
 $padding_w = 1$

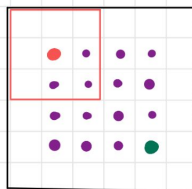


$H_{out} = W_{out} = 6$

# CONVOLUTIONAL LAYER

Example

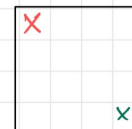
$$\begin{aligned} C_{in} &= C_{out} = 1 \\ H_{in} &= W_{in} = 6 \\ k_h &= k_w = 3 \end{aligned}$$



\*

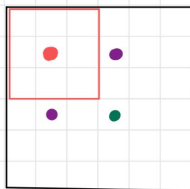


=



$$H_{out} = W_{out} = 4$$

$$\begin{aligned} \text{stride}_h &= 2 \\ \text{stride}_w &= 2 \end{aligned}$$



\*



=



$$H_{out} = W_{out} = 2$$

$$MAC_{y_{ij}} = k_h k_w$$

$$C_{in} = C_{out} = 1 \Rightarrow MAC_{conv} = H_{out} W_{out} k_h k_w$$

$$C_{in} = 1, C_{out} > 1 \Rightarrow MAC_{conv} = H_{out} W_{out} k_h k_w C_{out}$$

# CONVOLUTIONAL LAYER

## Standard

$$\text{Params} = k_h k_w C_{in} C_{out}$$

$$\text{MAC} = \text{Hart Wout} \cdot k_h k_w C_{in} \cdot C_{out}$$

## Depth-wise ( $C_{out} = C_{in} = C$ )

$$\text{Params} = k_h k_w 1 C$$

$$\text{MAC} = \text{Hart Wout} \cdot k_h k_w 1 \cdot C$$

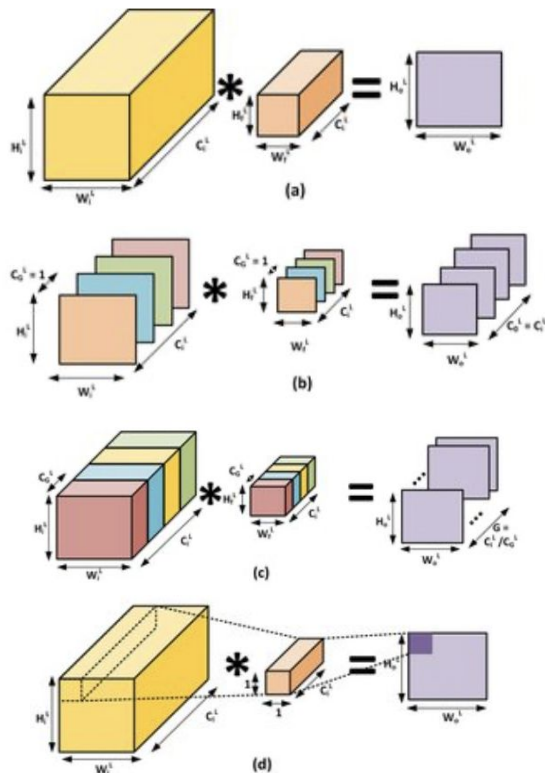
## Group ( $G$ - # of groups)

$$\text{Params} = k_h k_w \frac{C_{in}}{G} \cdot C_{out}$$

$$\text{MAC} = \text{Hart Wout} \cdot k_h k_w \frac{C_{in}}{G} \cdot C_{out}$$

## Point-wise

$$\text{Standard with } k_h = k_w = 1$$





# PRACTICAL EXERCISES

- Compute Parameters/ FLOP/ MAC for different Linear layers
  - theoretically
  - empirically, using Python packages
    - FlopCo
    - time

# PRACTICAL EXERCISES

- Non-linear layers
  - ReLU
  - GeLU
  - etc.
- Normalization layers
  - Layer Norm
  - Batch Norm
  - etc.

# NORMALIZATION LAYERS

$$\hat{x}_i = \frac{1}{\hat{\sigma}_i} (x_i - \mu_i), \quad \mu_i = \frac{1}{m} \sum_{k \in S_i} x_k$$

$$\hat{\sigma}_i = \sqrt{\frac{1}{m} \sum_{k \in S_i} (x_k - \mu_i)^2 + \epsilon}$$

$S_i$  - set of pixels in which mean & std are computed

$$i = (i_N, i_C, i_H, i_W)$$

$$y_i = \gamma \hat{x}_i + \beta, \quad \gamma, \beta - \text{trainable scale \& shift}$$

BN  $S_i = \{k \mid k_C = i_C\}$

Same statistics for:

all spatial coords  $(h, w)$  in one channel (across batch)

LN  $S_i = \{k \mid k_N = i_N\}$

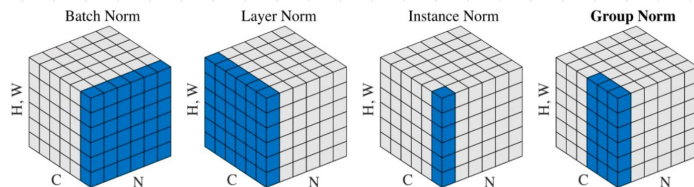
all channel-spatial coords  $(c, h, w)$

IN  $S_i = \{k \mid k_N = i_N, k_C = i_C\}$

all spatial coords  $(h, w)$  in one channel

GN  $S_i = \left\{k \mid k_N = i_N, \left\lfloor \frac{k_C}{C/G} \right\rfloor = \left\lfloor \frac{i_C}{C/G} \right\rfloor \right\}$

all spatial coords  $(h, w)$  in one group of channels



# INTENSITY OF OPERATIONS

- Real speed is limited by processor characteristics
  - time of math operation ( $T_{\text{math}}$ )
  - the time of memory access operation ( $T_{\text{mem}}$ )
- $p = T_{\text{math}}/T_{\text{mem}} = \text{BW}_{\text{math}}/\text{BW}_{\text{mem}}$  - intensity of operation,
  - $\text{BW}_{\text{math}}$  - speed of floating point operations (measured in FLOP/second)
  - $\text{BW}_{\text{mem}}$  - speed of memory access (measured in bytes/second)
  - Modern GPUs:  $\text{BW}_{\text{math}} \gg \text{BW}_{\text{mem}}$
- $p_{\text{algo}} = \text{\#FLOP}/\text{\#Mem}$ 
  - $p_{\text{algo}} = p$ , - memory and processor are fully utilized during computations
  - $p_{\text{algo}} < p$ , - algo is limited by memory access speed
  - $p_{\text{algo}} > p$ , - algo is limited by math operation speed

# INTENSITY OF OPERATIONS: EXAMPLE

16-bit arithmetics, GPU Nvidia A100

Operation	Arithmetics' intensity	Time is limited by ....
Linear (4096 x 1024, batch 512)	315 FLOP/B	arithmetics
Linear (4096 x 1024, batch 1)	1 FLOP/B	memory access
MaxPooling with kernel 3x3	2.25 FLOP/B	memory access
ReLU	0.25 FLOP/B	memory access
LayerNorm	< 10 FLOP/B	memory access

# COMPRESSION TECHNIQUES: BRIEF OVERVIEW

# COMPRESSION METHODS

- **Pruning methods** (structured / unstructured)
  - redundant weights / neurons are pruned, hence, the whole model is compressed.
- **Tensor-based methods**
  - use matrix or tensor decomposition to estimate the informative parameters of deep neural networks;
  - In most cases, a much lower total computational cost can be achieved by replacing a convolutional layer with a sequence of several smaller convolutional layers.
- **Quantization methods**
  - use low-bit representations for weights / activations;
  - can significantly accelerate networks, but they usually require special hardware to reach a theoretical speed-up in practice.
- **Knowledge distillation methods**
  - deal with a pre-trained network (*teacher* network), and an accelerated network (*student* network);
  - outputs (resulting and/or intermediate) of the teacher network are used to guide the student network training.

# COURSE INFO



# SCHEDULE

1. Introduction. Measures of neural networks effectiveness.
2. Pruning.
3. Tensor methods.
4. Quantization.
5. Inference on mobile devices.
6. 6.1. Knowledge distillation.  
6.2. Large neural networks training acceleration

# COMPLETION CRITERIA

- 5 Homework assignments (HAs) + 1 optional HA
- Acceptance rules
  - HA submitted before 11.59 am next day - max 100%
  - HA submitted before 11.59 am next next day - max 75%
  - HA submitted before 11.59 am last lecture day - max 50%
- To pass the course:
  - Get  $\geq 50\%$  for any 4 of required 5 HAs

# RESOURCES

- FlopCo <https://github.com/juliagusak/flopco-pytorch>
  - Python library that aims to make FLOPs and MACs counting simple and accessible for PyTorch neural networks.
  - FlopCo allows to collect other useful model statistics, such as number of parameters, shapes of layer inputs/outputs, etc.
- Papers & code links  
<https://github.com/juliagusak/model-compression-and-acceleration-progress>

# SEE YOU NEXT LECTURE!

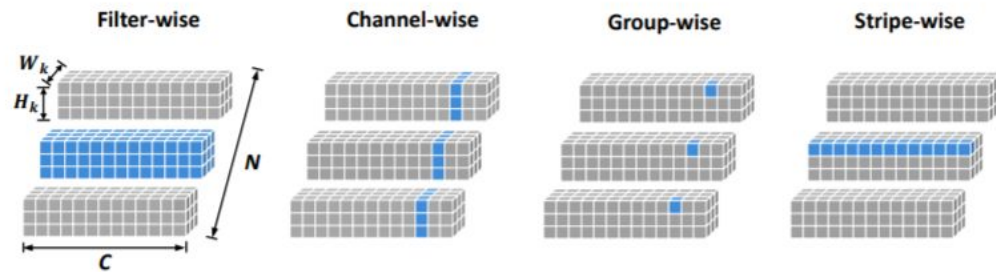
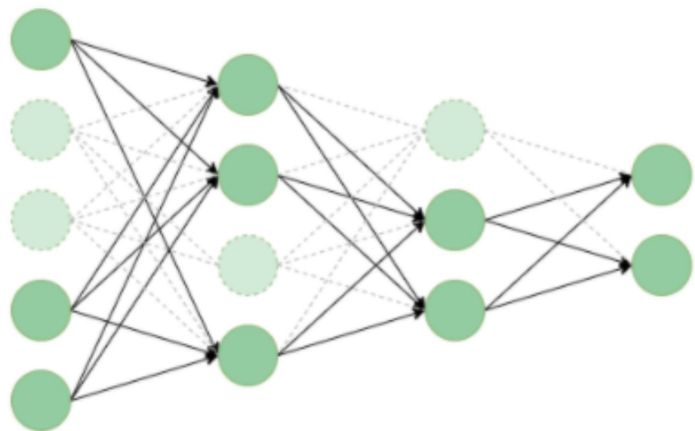


Figure 2: The visualization of different types of pruning.