

Deep Learning for Computer Vision

Andrii Liubonko

Grammarly*



* The opinions expressed in this presentation and on the following slides are solely those of the presenter and not necessarily those of Grammarly

Logistics

4 units

2 types of homework:

- notebooks [30 % of FINAL SCORE]
- paper review [20 % of FINAL SCORE]
- mini-project [50 % of FINAL SCORE]

important dates (preliminary):

15 January, 23:59 : SOFT

29 January, 23:59 : HARD, PENALTY 30%

course repo:

<https://github.com/lyubonko/ucu2022cv>

Overview of the course

Unit I

[T] Intro

[P] pytorch

Unit II

[T] CNNs in depth, Object Detection

[P] simple nets

Unit III

[T] Attention in CV, Transformers

[P] classification

Unit IV

[T] Generative models, Diffusion Models

[P] project structure, stable diffusion

Overview of the course

Unit I

[T] Intro

[P] pytorch

Unit II

[T] CNNs in depth, Object Detection

[P] simple nets

Unit III

[T] Attention in CV, Transformers

[P] classification

Unit IV

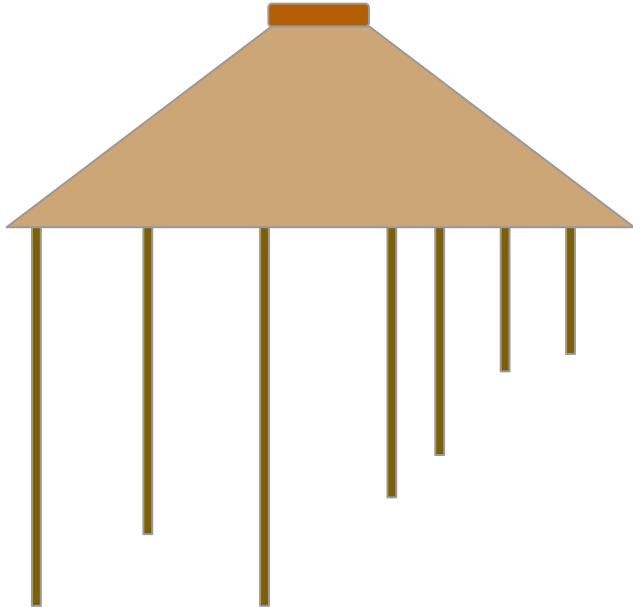
[T] Generative models, Diffusion Models

[P] project structure, stable diffusion

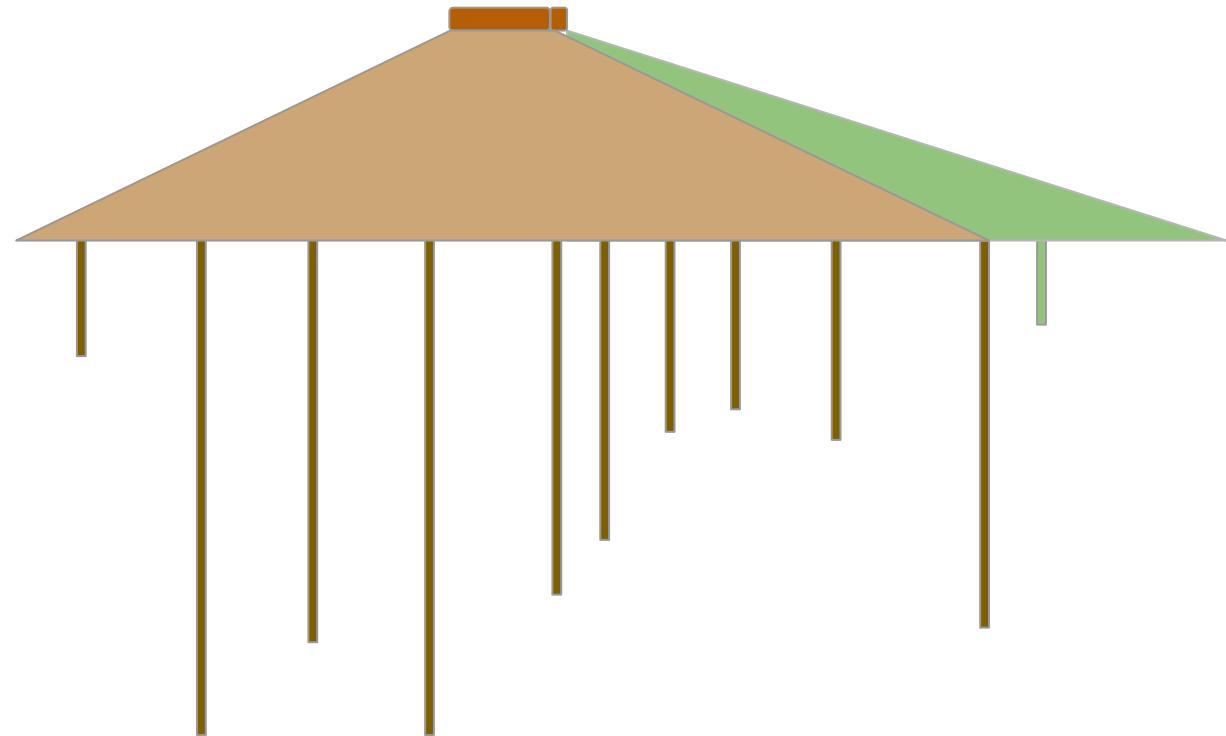
Goals of the Course

- working knowledge of essential elements/blocks of **Convolutional Neural Networks [CNNs]**
- modern **CNNs architectures**
- get deeper with one particular problem (**Object Detection**)
- attention in Computer Vision,
Transformers in Computer Vision
- get flavor of **Generative Modeling**, vis **Diffusion Models**

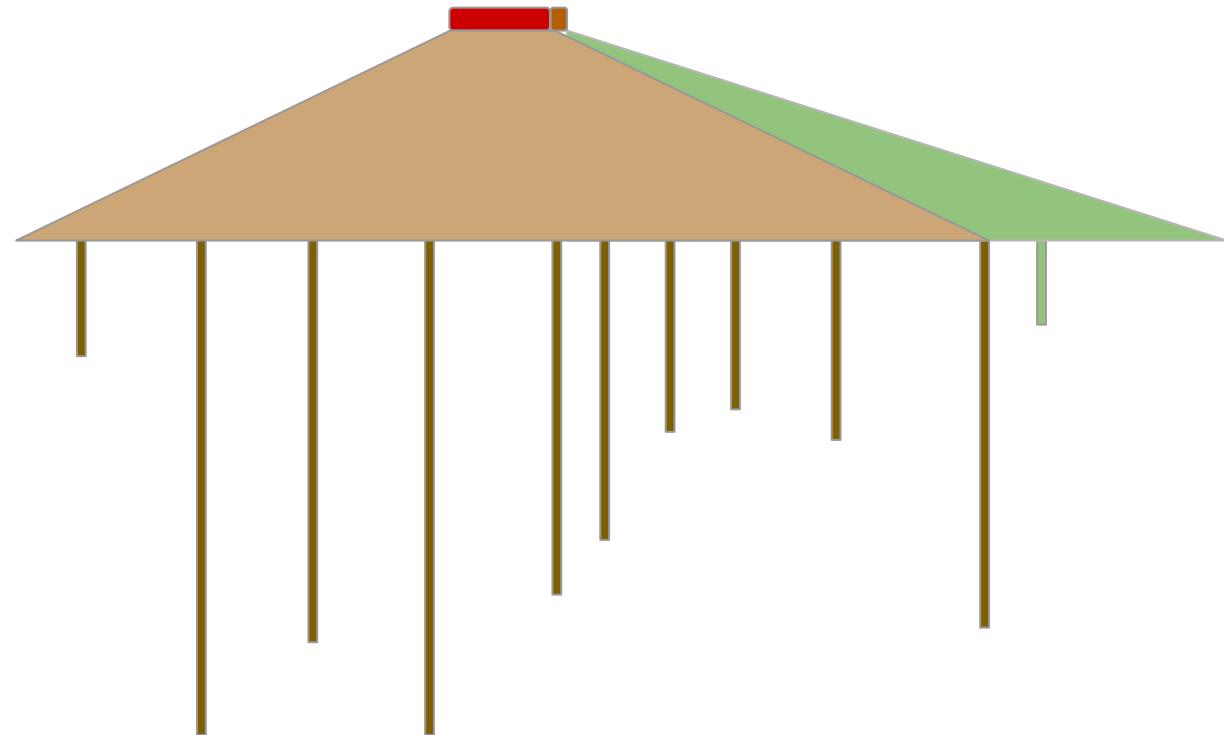
Goals of the Course



Goals of the Course

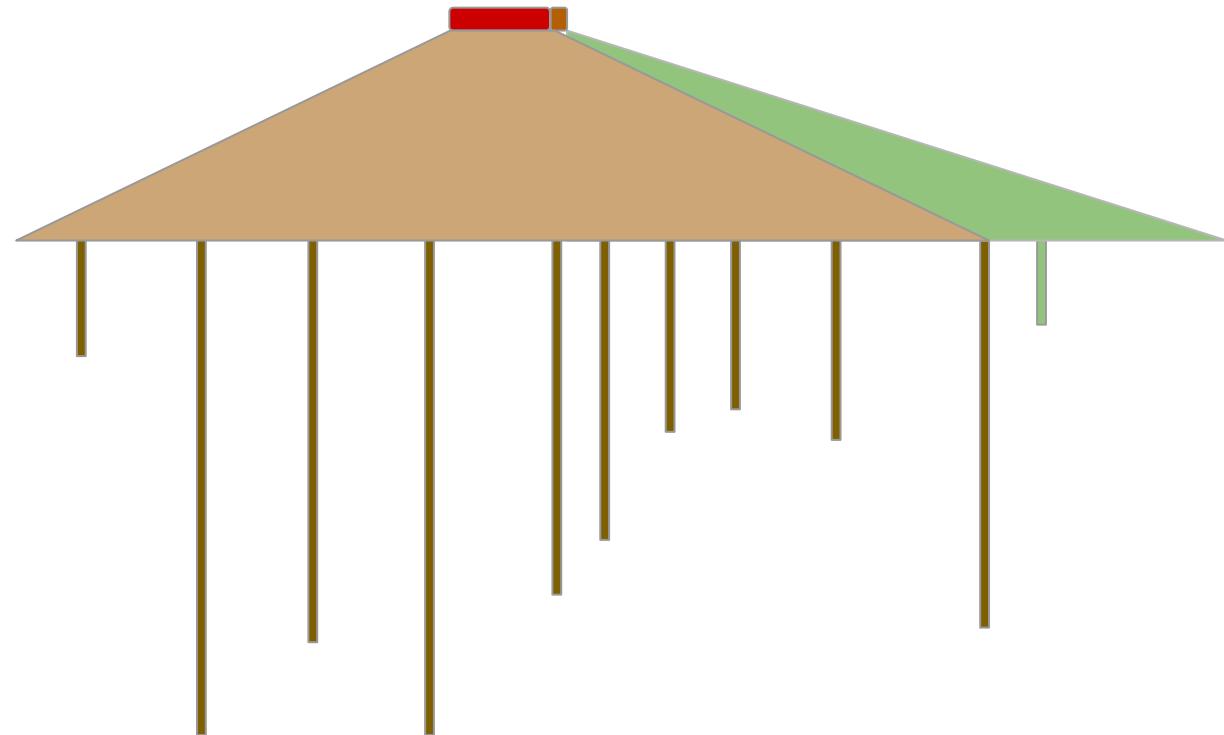


Goals of the Course



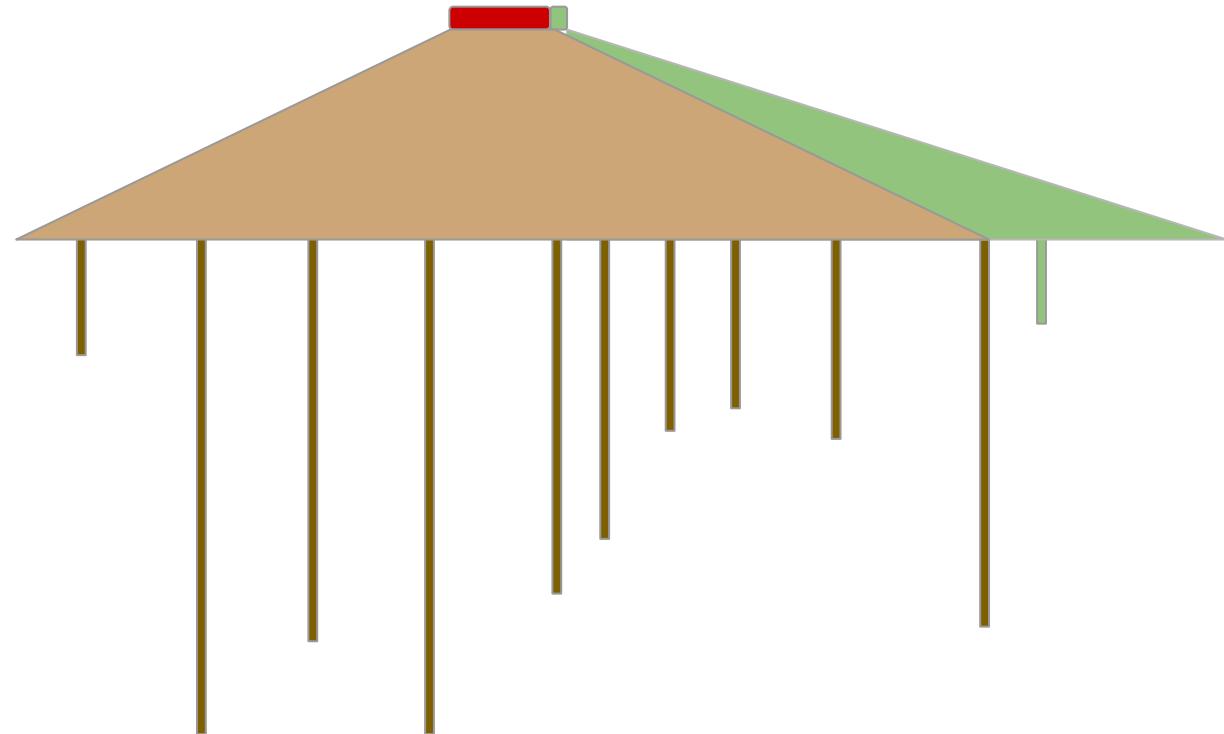
- essential elements/blocks

Goals of the Course



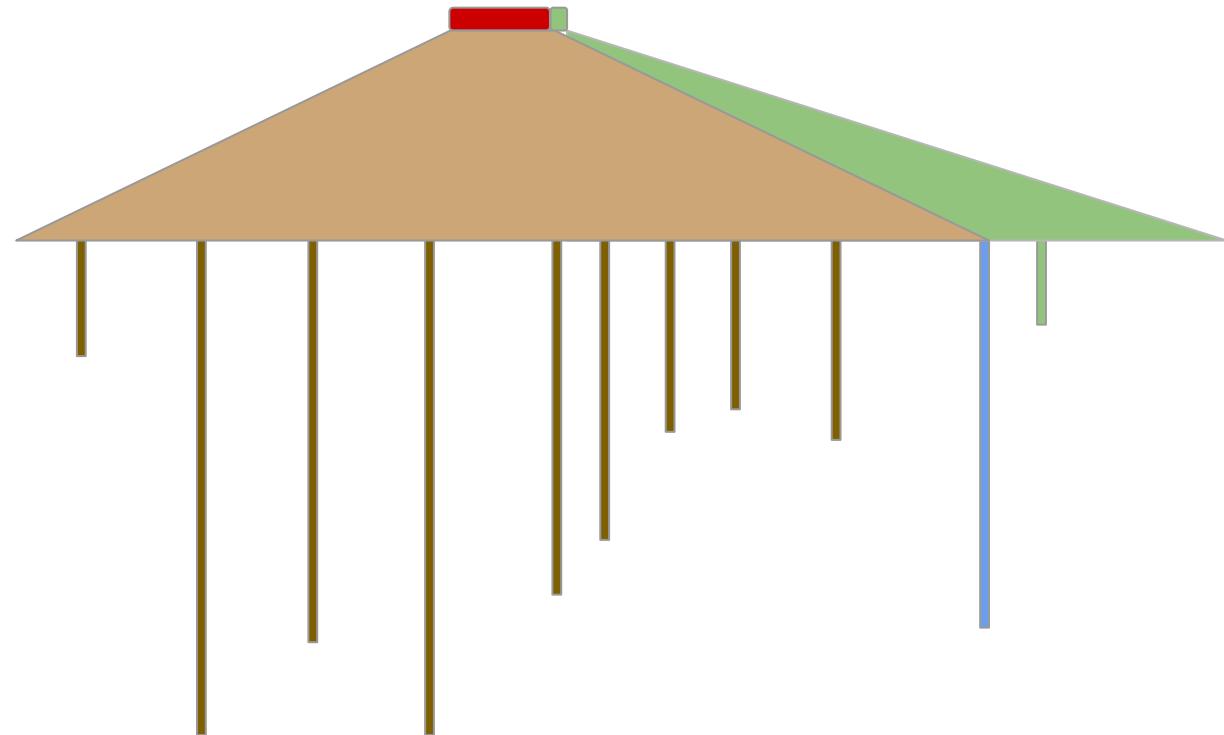
- essential elements/blocks
- modern architectures

Goals of the Course



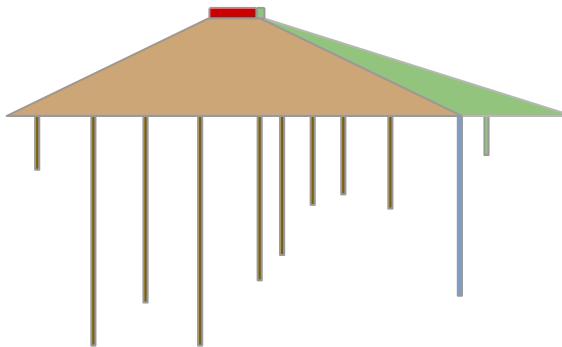
- essential elements(blocks)
- modern architectures
- attention in CV

Goals of the Course



- essential elements/blocks
- modern architectures
- attention in CV
- Object Detection

Goals of the Course



- working knowledge of **essential elements/blocks** of **Convolutional Neural Networks [CNNs]**
- modern **modern CNNs architectures**
- **attention** in Computer Vision, **Transformers** in Computer Vision
get deeper with one particular problem (**Object Detection**)
- get flavor of **Generative Modeling**, vis **Diffusion Models**

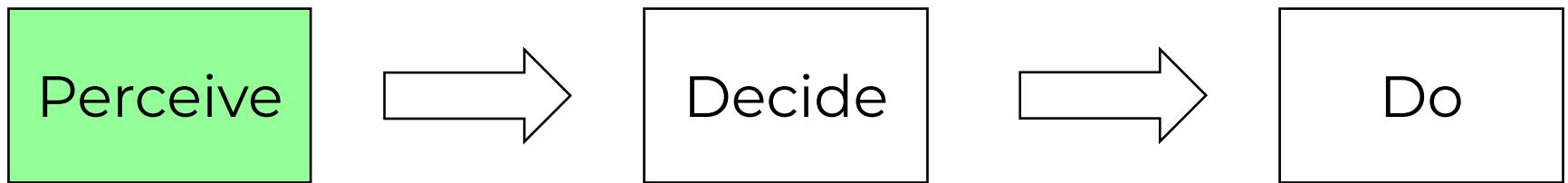
Content of today lecture

- **Intro, Context**
- ML/DL review
 - training & testing
 - neural networks
- CV specific setups
 - supervised, semi-supervised,
 - self-supervised, auto-encoders
- Main components of CNNs (motivation and details)
 - convolutional layer
 - pooling layer

Intro



Intro



Intro

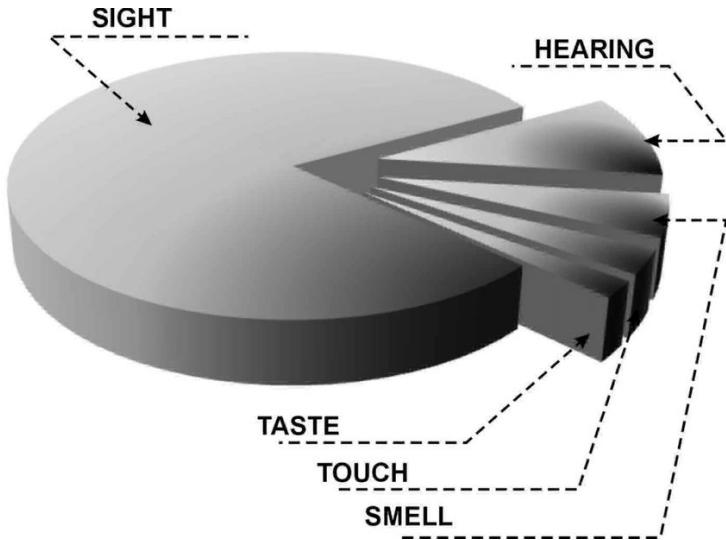
Generate

Intro

Generate



Intro



*The goal of **computer vision** is to perceive [extract useful information] from visual input, like images and video*

Intro



- indoor/outdoor? [image classification]
- Where are the objects? [object detection]
- How far is the object ? [depth estimation]
- What people are doing? [activity recognition]
- Is the state of the environment normal? [anomaly detection]
- ...

Intro



- scale variations
- occlusions
- illuminations
- deformations
-
-

Intro

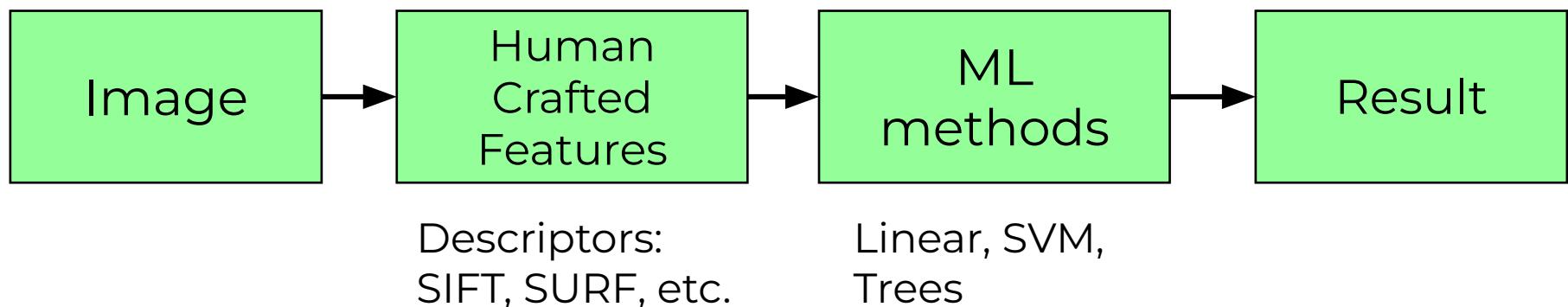


what humans see

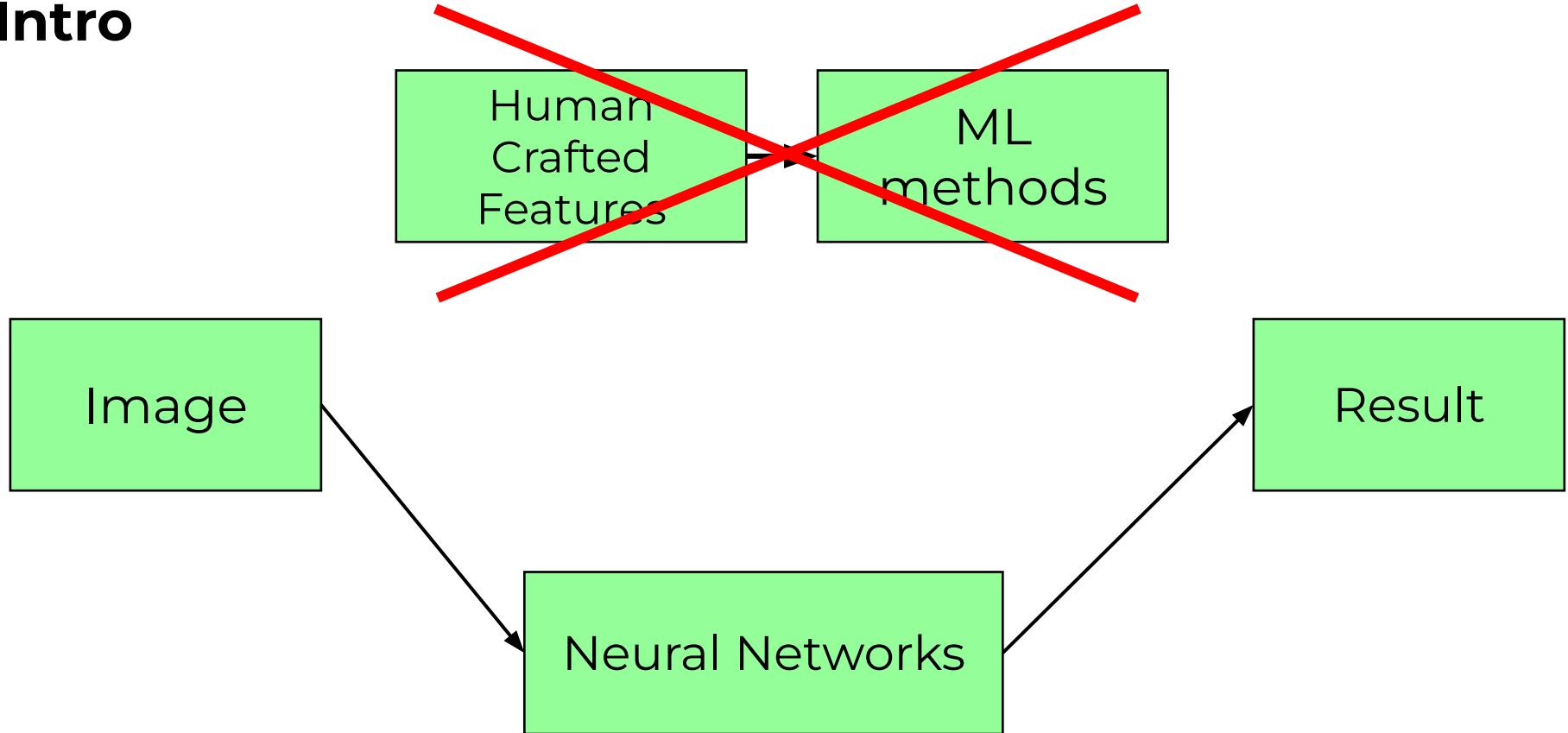
0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0

what computers see

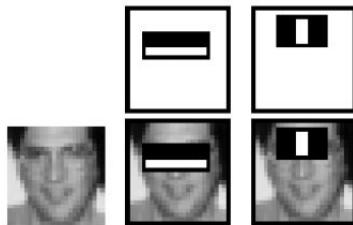
Intro



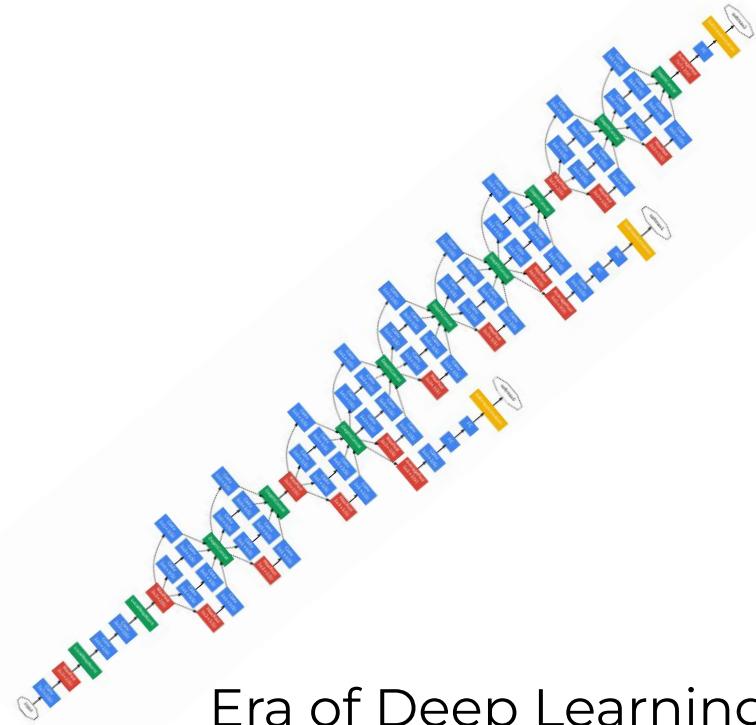
Intro



Intro



Era of Human-Crafter Features



Era of Deep Learning

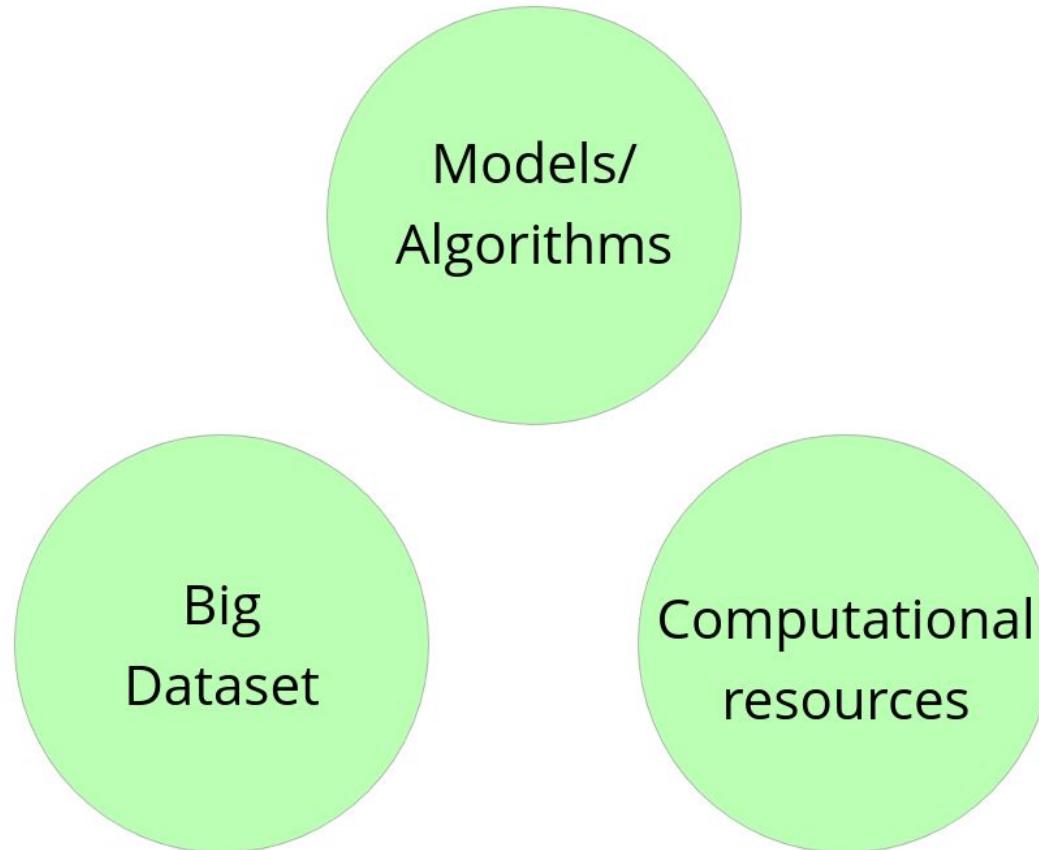


1986
BackProp

1998
LeNet

2012
AlexNet

Intro



Intro



Vladimir Vapnik Larry Jackel

1. Jackel bets (one fancy dinner) that by March 14, 2000, people will understand quantitatively why big neural nets working on large databases are not so bad. (Understanding means that there will be clear conditions and bounds)

Vapnik bets (one fancy dinner) that Jackel is wrong.

But .. If Vapnik figures out the bounds and conditions, Vapnik still wins the bet.

2. Vapnik bets (one fancy dinner) that by March 14, 2005, no one in his right mind will use neural nets that are essentially like those used in 1995.

Jackel bets (one fancy dinner) that Vapnik is wrong

A handwritten signature of V. Vapnik.

3/14/95

V. Vapnik

A handwritten signature of L. Jackel.

3/14/95

L. Jackel

A handwritten signature of Y. LeCun.

3/14/95

Witnessed by Y. LeCun

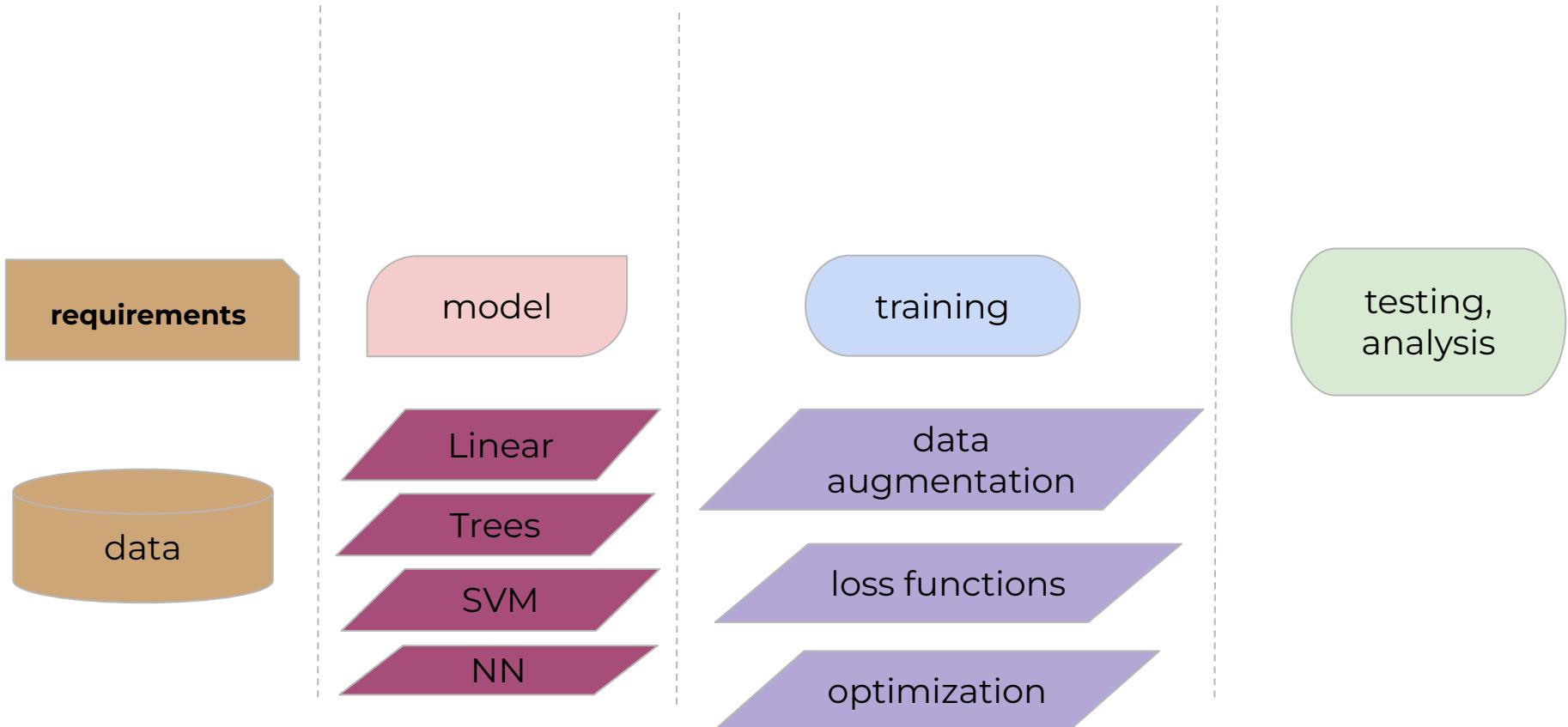
Content

- Intro, Context
- **ML/DL review**
 - training & testing
 - neural networks
- CV specific setups
 - supervised, semi-supervised,
 - self-supervised, auto-encoders
- Main components of CNNs (motivation and details)
 - convolutional layer
 - pooling layer

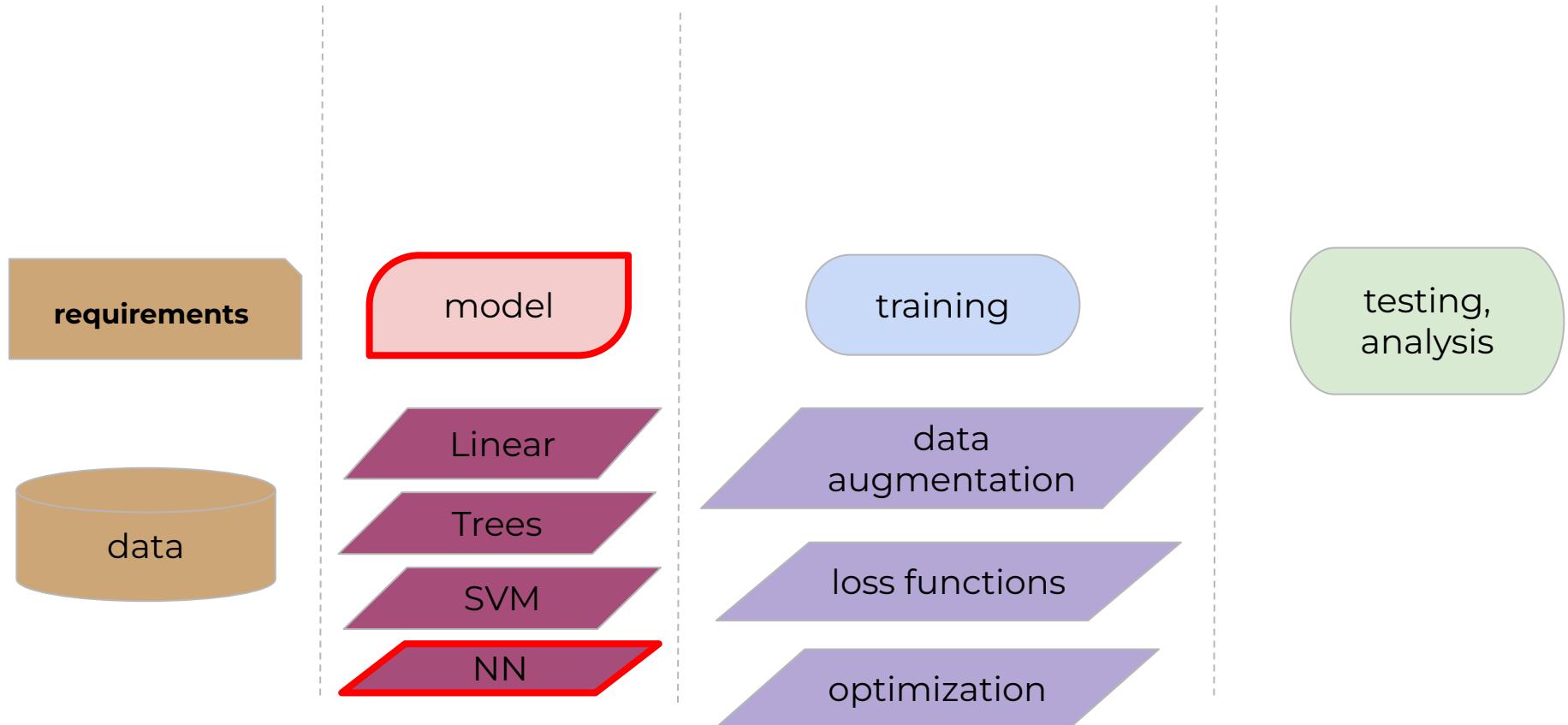
Machine Learning - types of learning

- Supervised
- Unsupervised
- Semi-supervised
- RL

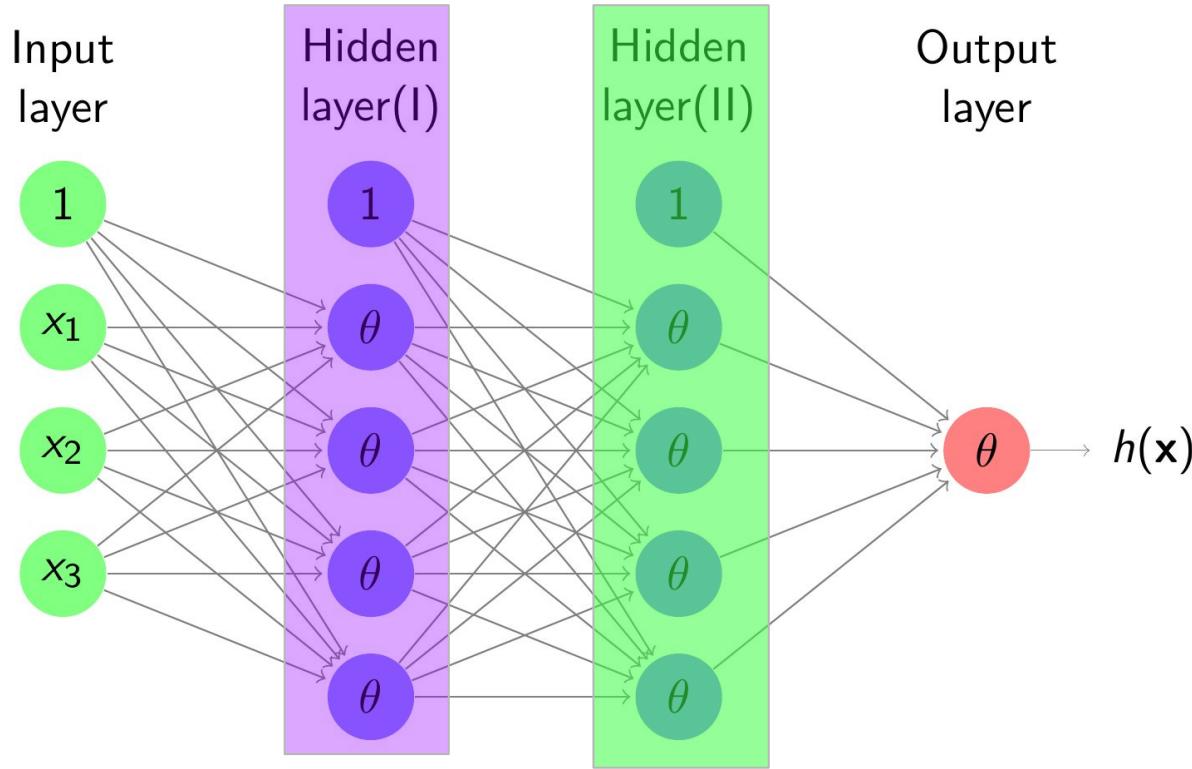
Machine Learning essential blocks



Machine Learning essential blocks

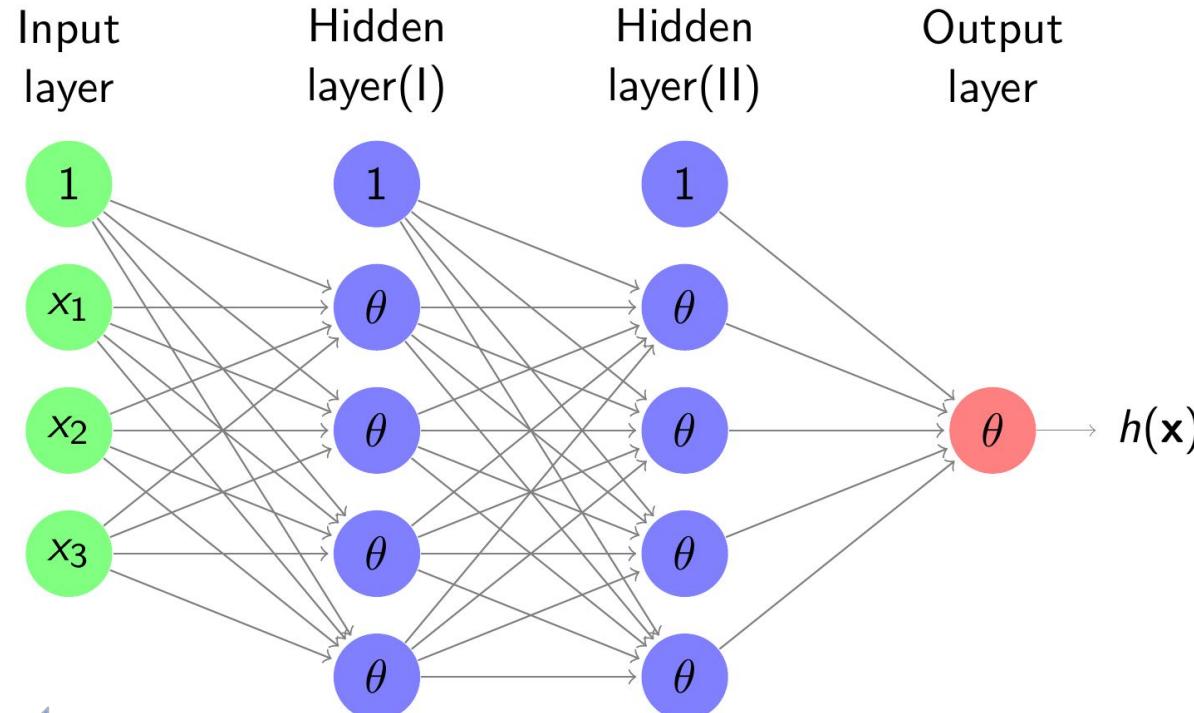


Model [Neural Networks]



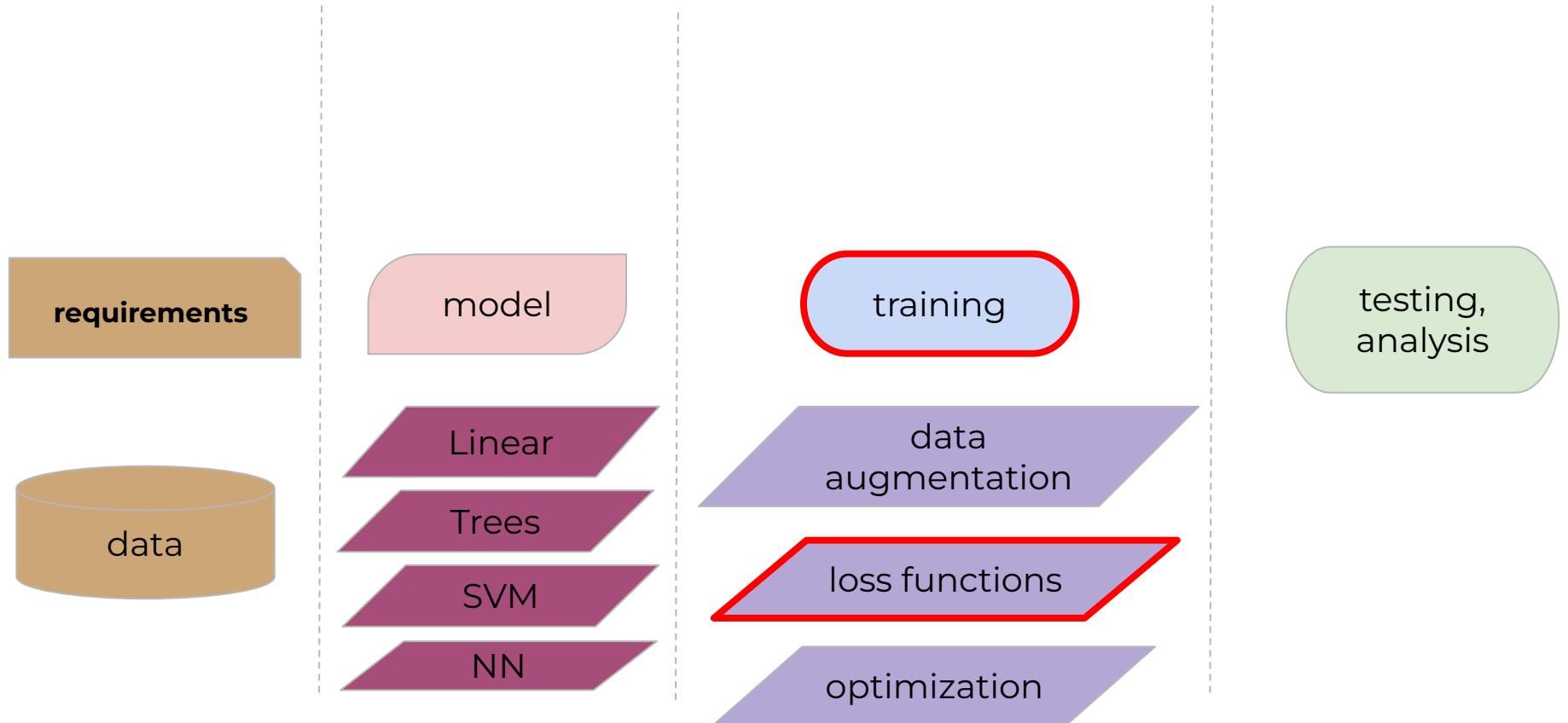
$$\{f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_L \sigma_L(\mathbf{W}_{L-1} \cdots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}))) \mid \boldsymbol{\theta} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}\}$$

Model [Neural Networks]



Back-propagate error signal to get derivatives for learning

Machine Learning essential blocks



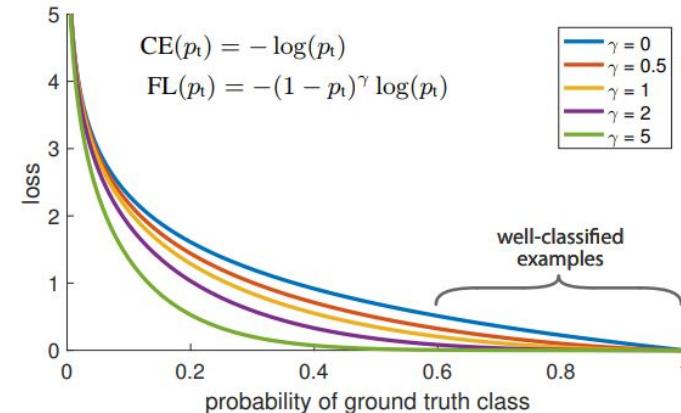
Loss function

Losses for classification, regression:

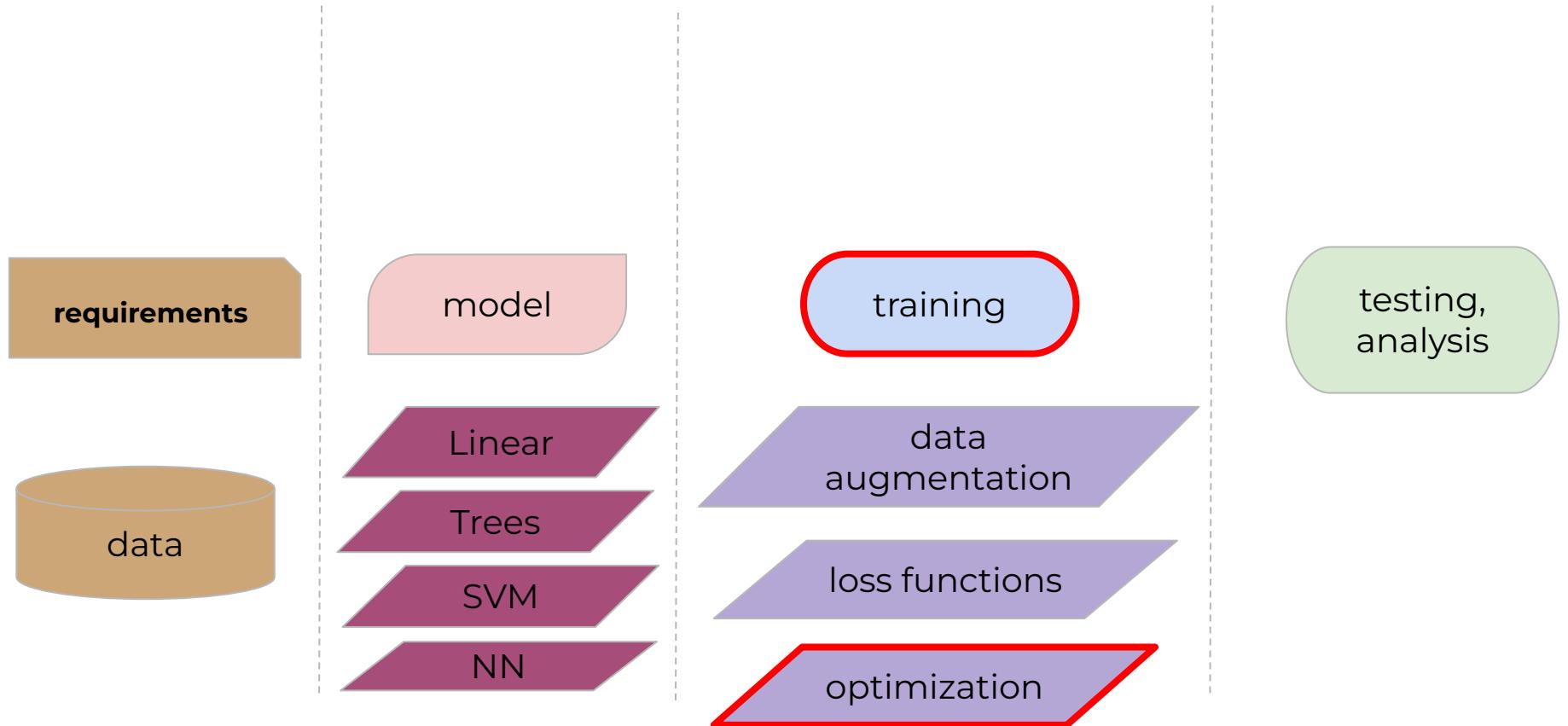
$$L_{logloss} = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^M y_{nk} \cdot \log(p_{nk})$$

$$L_{rmse} = \frac{1}{N} \sum_{n=1}^N \| F(\mathbf{x}_n) - y_i \|_2$$

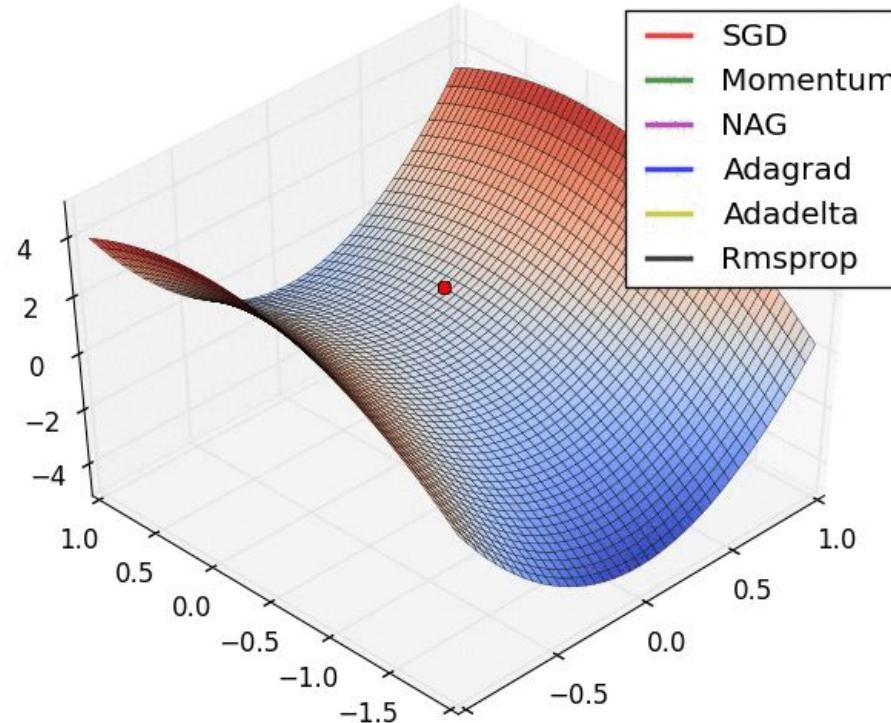
Focal Loss:



Machine Learning essential blocks



Optimization



Optimization

Input: Training set $\mathcal{S} \triangleq \cup_{i=1}^n \{(\mathbf{x}_i, \mathbf{y}_i)\}$, Loss function $l : \mathcal{W} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, Batch size b , Step size $\eta > 0$, Neighborhood size $\rho > 0$.

Output: Model trained with SAM

Initialize weights \mathbf{w}_0 , $t = 0$;

while not converged **do**

| Sample batch $\mathcal{B} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_b, \mathbf{y}_b)\}$;
| Compute gradient $\nabla_w L_{\mathcal{B}}(\mathbf{w})$ of the batch's training loss;
| Compute $\hat{\epsilon}(\mathbf{w})$ per equation 2;
| Compute gradient approximation for the SAM objective
(equation 3): $\mathbf{g} = \nabla_w L_{\mathcal{B}}(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})}$;
| Update weights: $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}$;
| $t = t + 1$;

end

return \mathbf{w}_t

Algorithm 1: SAM algorithm

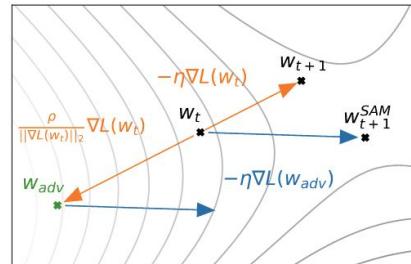
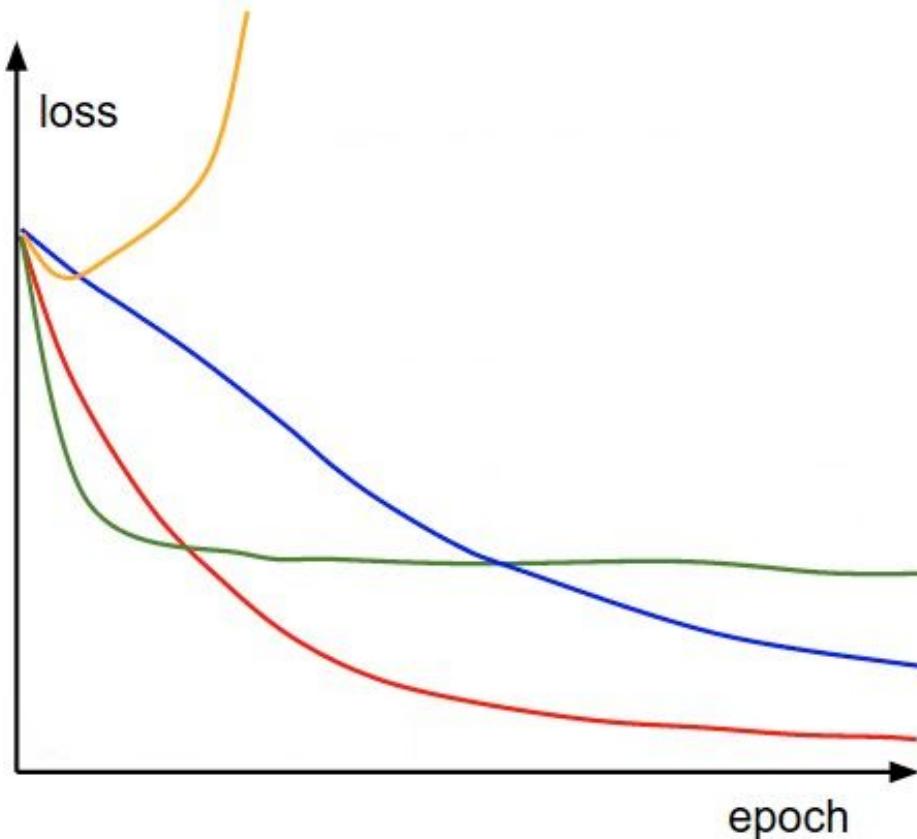


Figure 2: Schematic of the SAM parameter update.

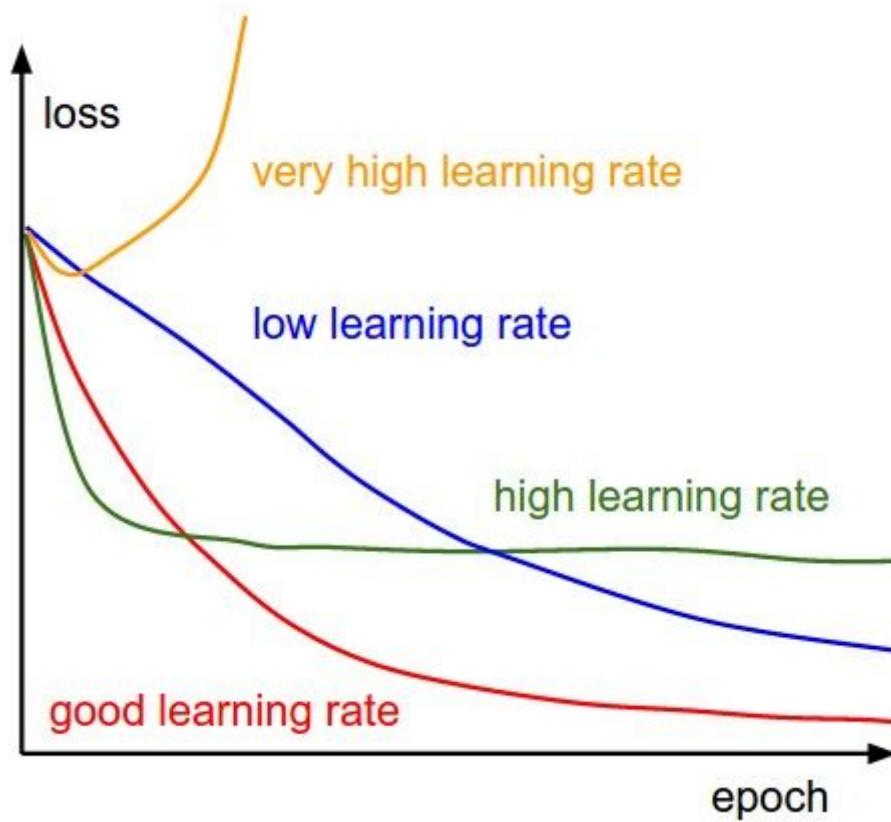


[2010.01412](#),
[JAX\(official\)](#), [pytorch](#)

Neural Networks



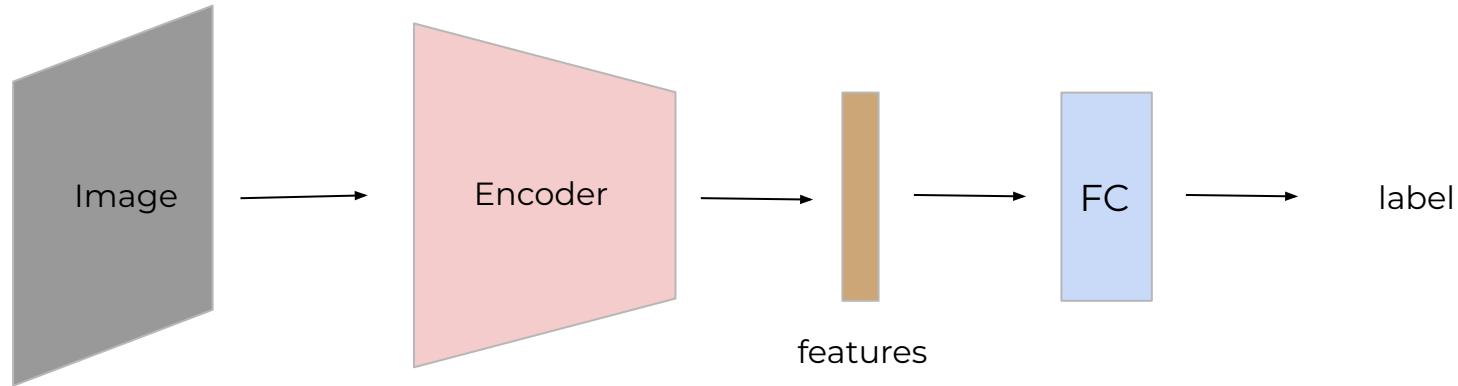
Neural Networks



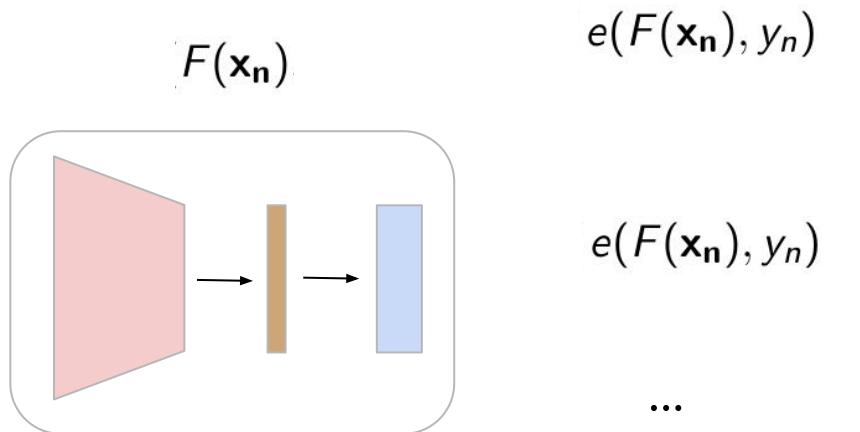
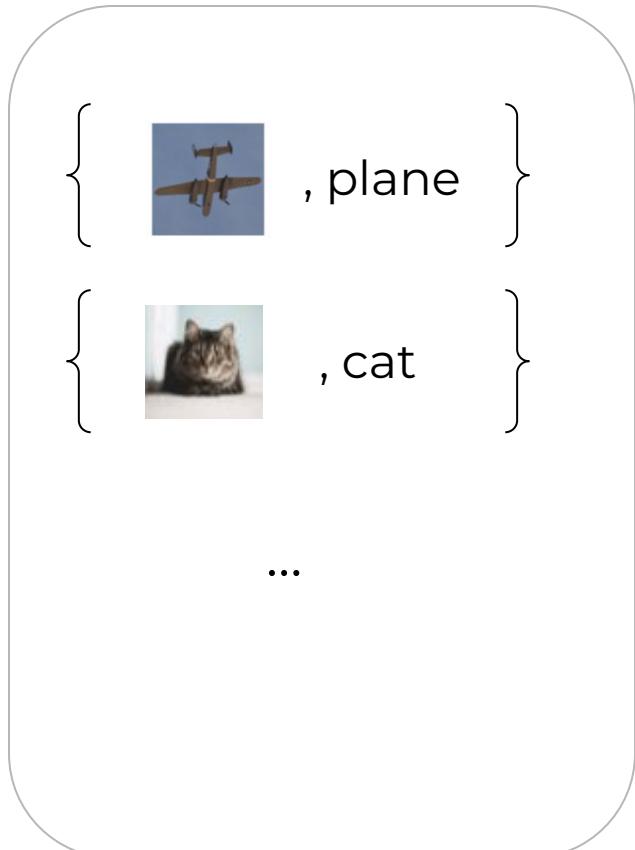
Content

- **Intro, Context**
- ML/DL review
 - training & testing
 - neural networks
- **CV specific setups**
 - supervised, semi-supervised,
 - self-supervised, auto-encoders
- Main components of CNNs (motivation and details)
 - convolutional layer
 - pooling layer

Neural Networks (supervised way)



Neural Networks (supervised way)



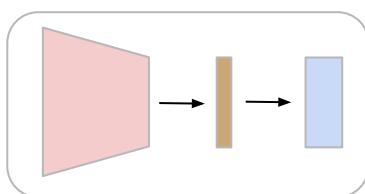
$$L_{train}(\omega) = \frac{1}{N} \sum_{n=1}^N e(F(\mathbf{x}_n), y_n)$$

Neural Networks (semi-supervised)

$\left\{ \begin{array}{l} \text{, plane} \\ \text{, cat} \end{array} \right\}$



...



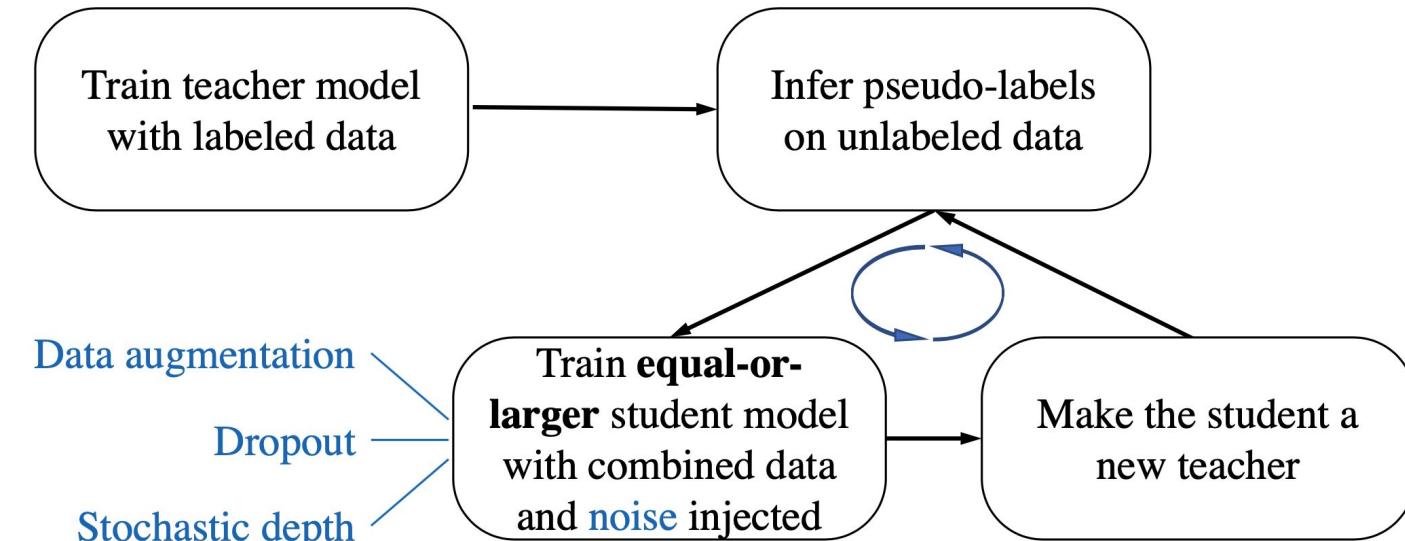
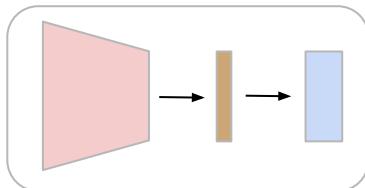
Train teacher model
with labeled data

Neural Networks (semi-supervised)

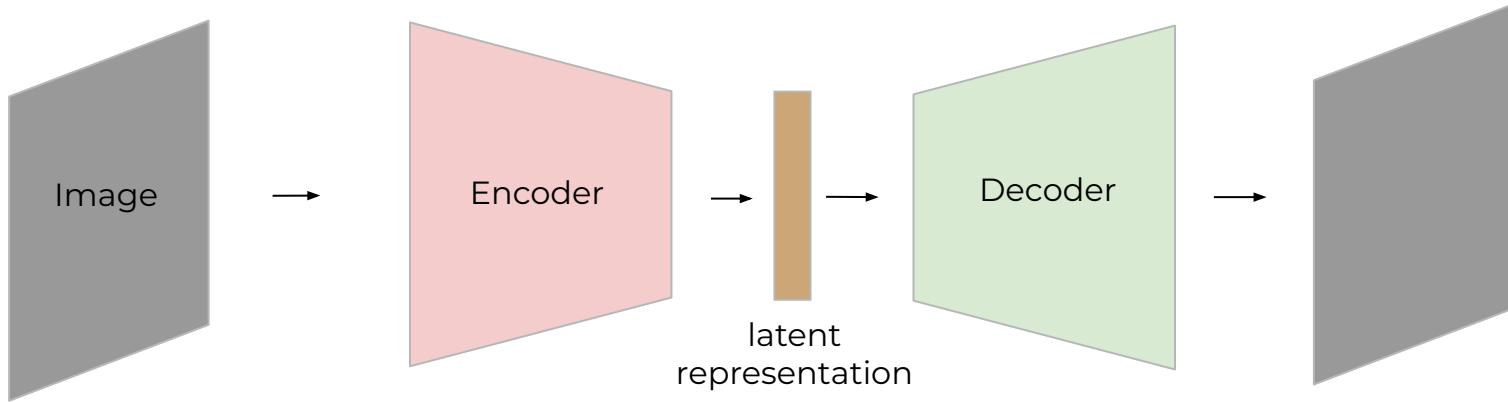
 , plane
 , cat



...

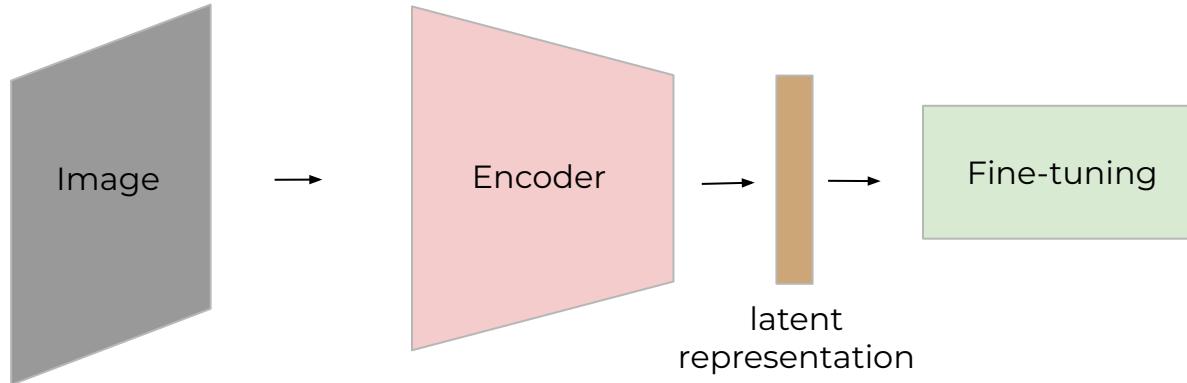


Neural Networks (auto-encoder way)



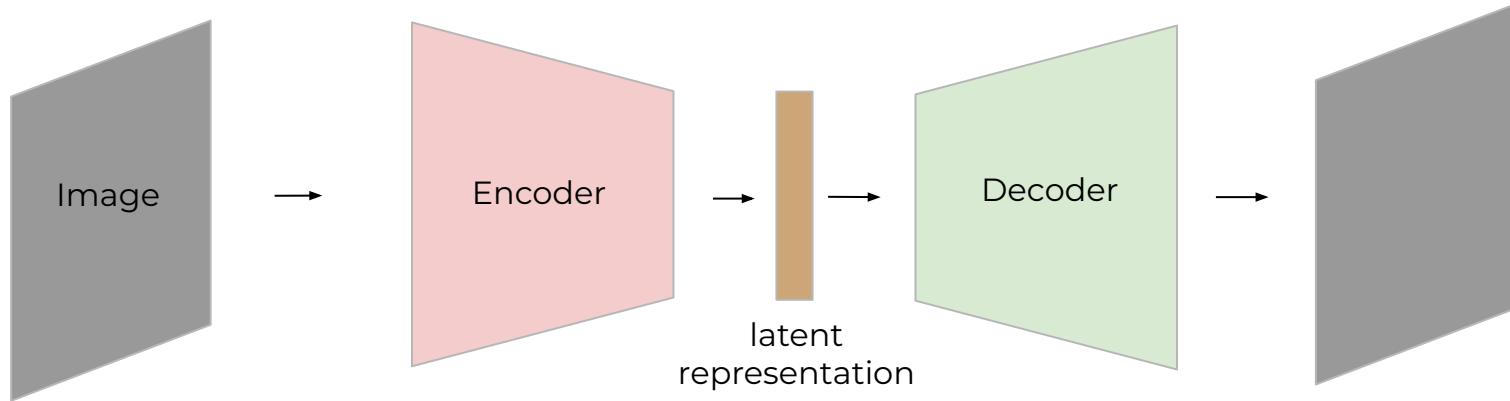
Neural Networks (discrimination)

Fine-tuning

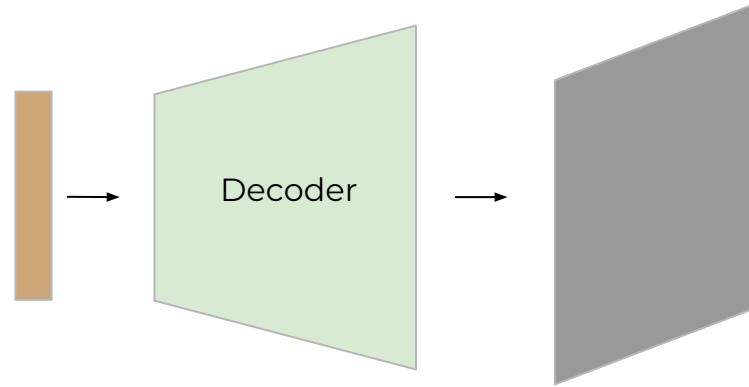


- Downstream tasks:
- classification
 - detection
 -

Neural Networks (generation)

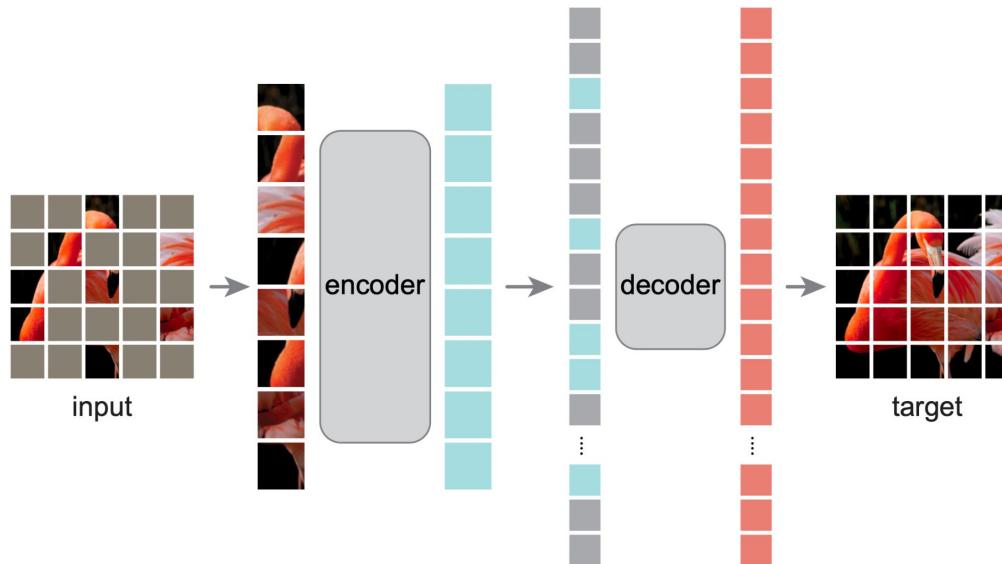


Neural Networks (generation)



Neural Networks (auto-encoder way)

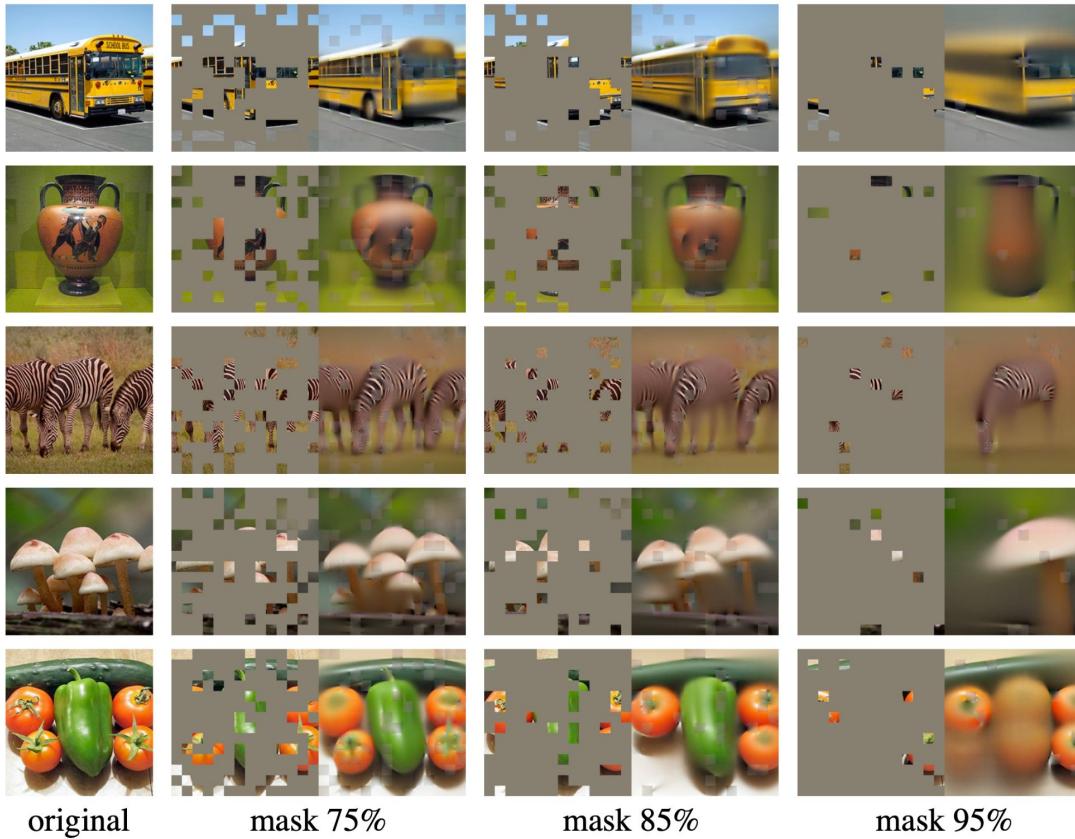
MAE



2111.06377

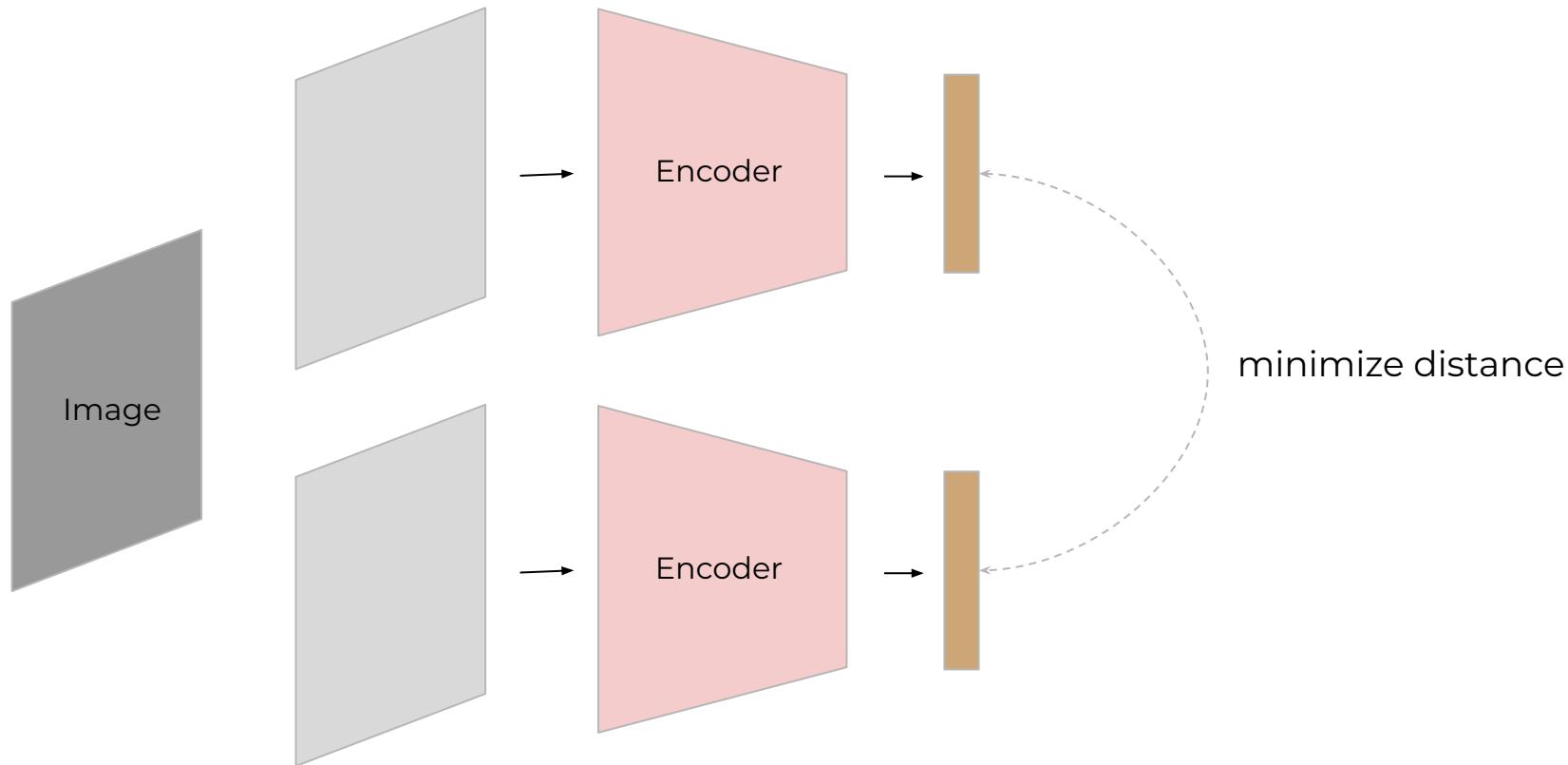
Neural Networks (auto-encoder way)

MAE



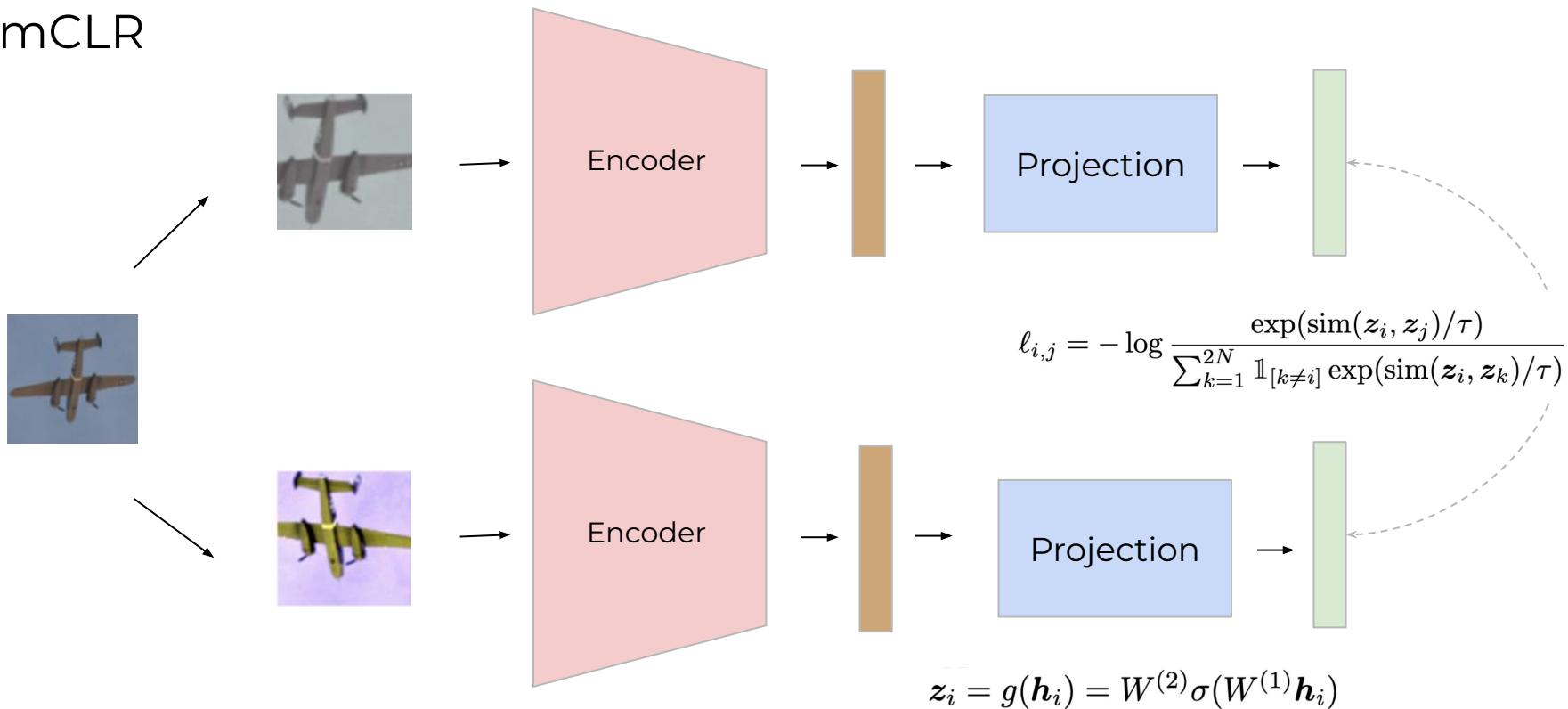
2111.06377

Neural Networks (self-supervised way)



Neural Networks (self-supervised way)

SimCLR



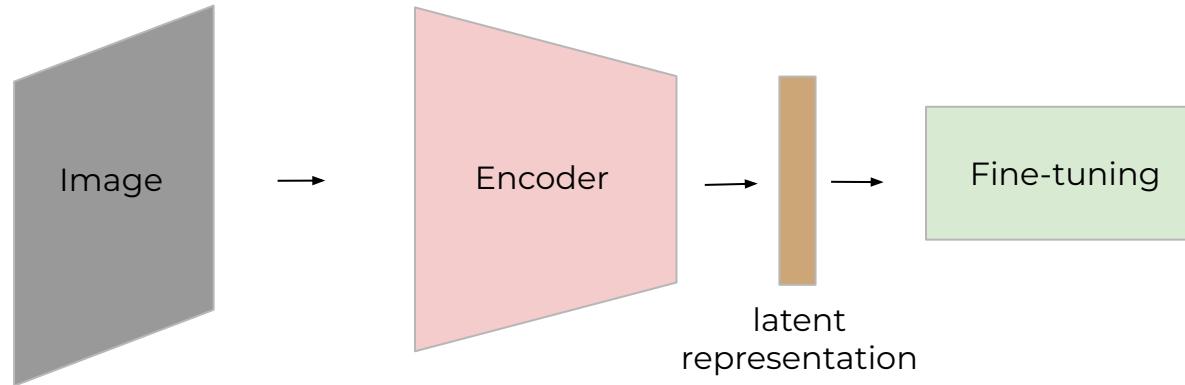
$$\mathbf{h}_i = f(\tilde{\mathbf{x}}_i) = \text{ResNet}(\tilde{\mathbf{x}}_i)$$

2002.05709

Neural Networks (self-supervised way)

SimCLR

Fine-tuning



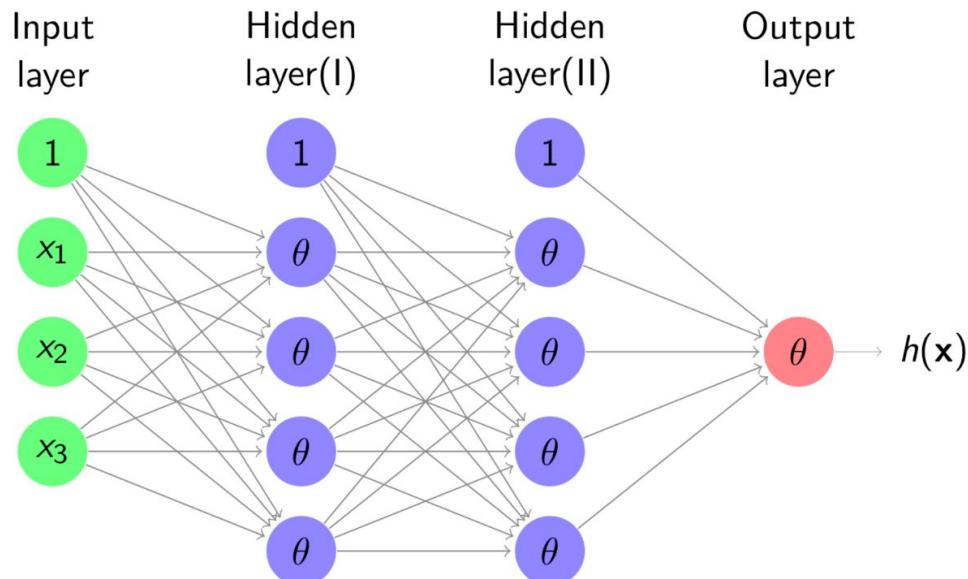
Downstream tasks:

- classification
- detection
-

Content

- Intro, Context
- ML/DL review
 - training & testing
 - neural networks
- CV specific setups
 - supervised, semi-supervised,
 - self-supervised, auto-encoders
- **Main components of CNNs (motivation and details)**
 - convolutional layer
 - pooling layer

Intro to CNN



Input layer

$$64 \times 64 = 4096$$

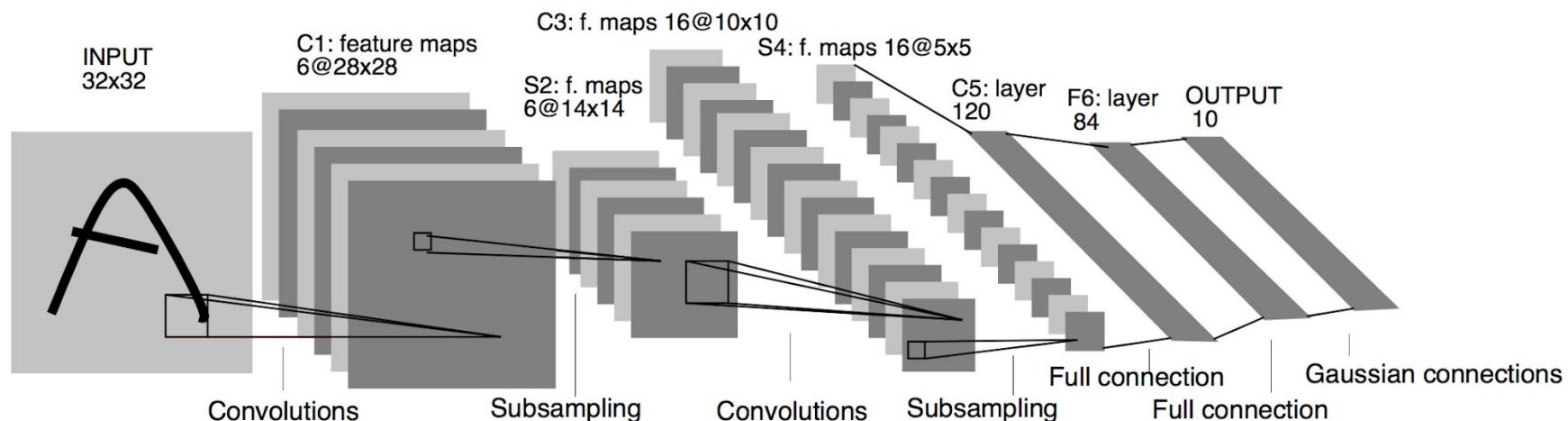
Weights: $4096 \times$
(size of second
layer)

=> millions of
parameters (for
just one layer)

Need for
alternative
architecture

Intro to CNN

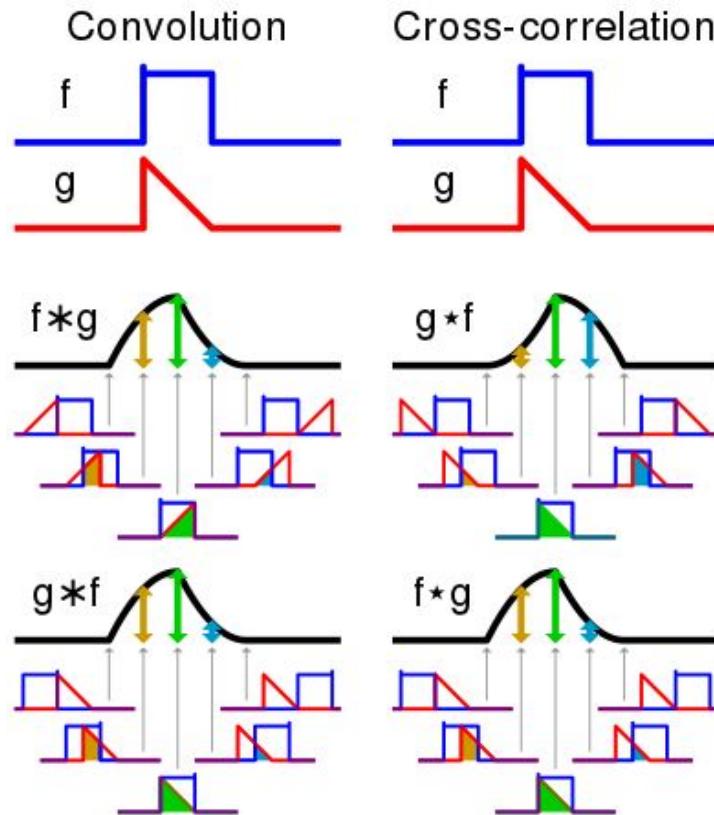
LeNet-5 [1998, paper by LeCun et al.]



Intro to CNN

- ▶ INPUT holds the raw pixel values of the image.
- ▶ CONV layer computes the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and the region they are connected to in the input volume.
- ▶ POOL layer performs a downsampling operation along the spatial dimensions (width, height).
- ▶ FC (i.e. fully-connected) layer computes the class scores. As with ordinary Neural Networks and as the name implies, each neuron in this layer is connected to all the numbers in the previous volume.

Intro to CNN



a function derived from two given functions by integration that expresses how the shape of one is modified by the other

$$(f * g)(t) \triangleq \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau.$$

Intro to CNN

1	0	1
0	1	0
1	0	1

Kernel

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

Intro to CNN

1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved
Feature

Intro to CNN

1	1 <small>$\times 1$</small>	1 <small>$\times 0$</small>	0 <small>$\times 1$</small>	0
0	1 <small>$\times 0$</small>	1 <small>$\times 1$</small>	1 <small>$\times 0$</small>	0
0	0 <small>$\times 1$</small>	1 <small>$\times 0$</small>	1 <small>$\times 1$</small>	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	

Convolved
Feature

Intro to CNN

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	4

Convolved
Feature

Intro to CNN

1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

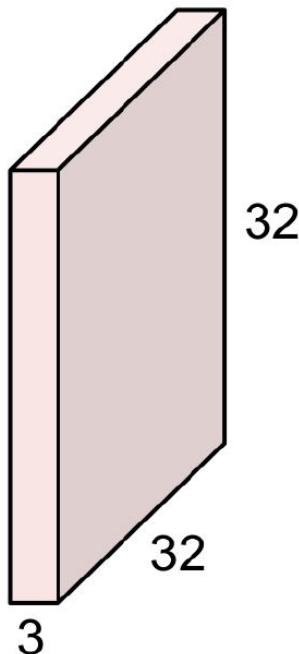
Image

4		

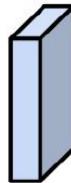
Convolved
Feature

Intro to CNN

32x32x3 image



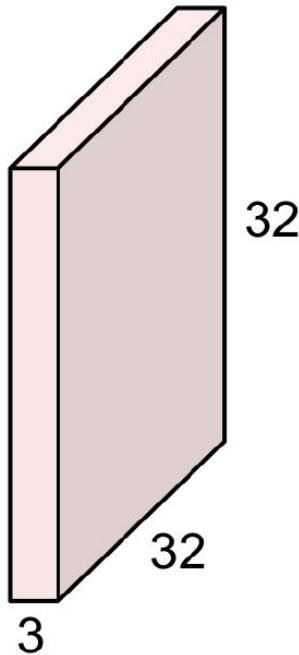
5x5x3 filter



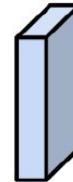
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Intro to CNN

32x32x3 image



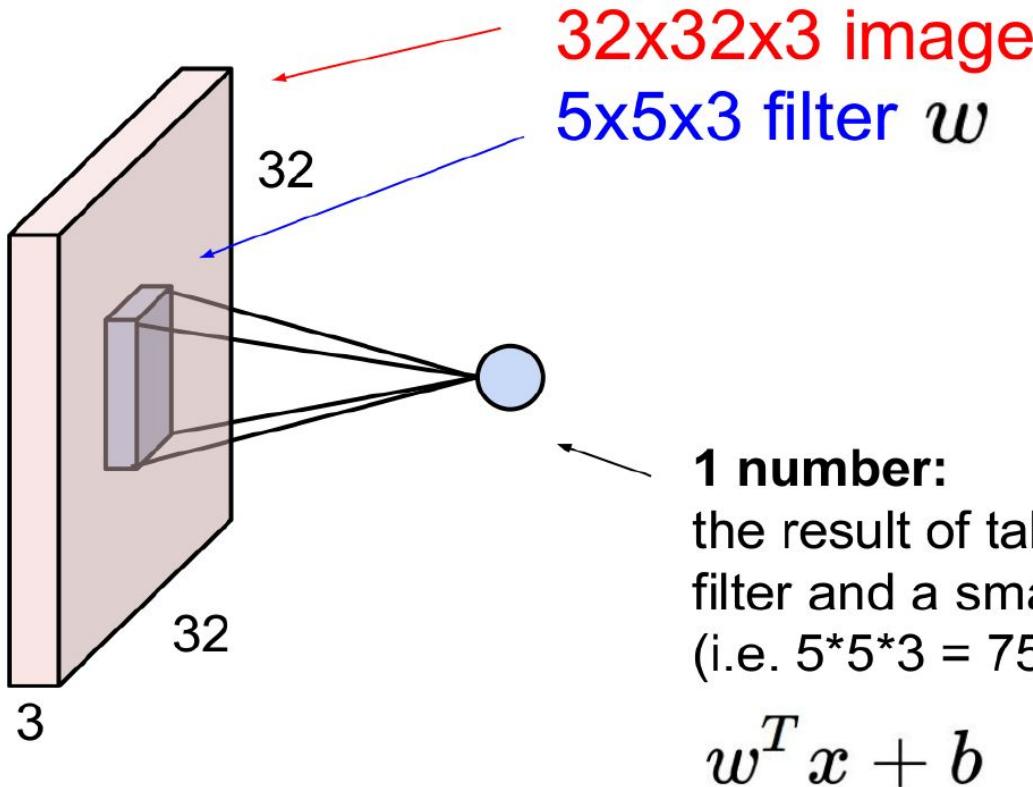
5x5x3 filter



1	0	1	0	1
0	1	0	1	0
1	0	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1

Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Intro to CNN

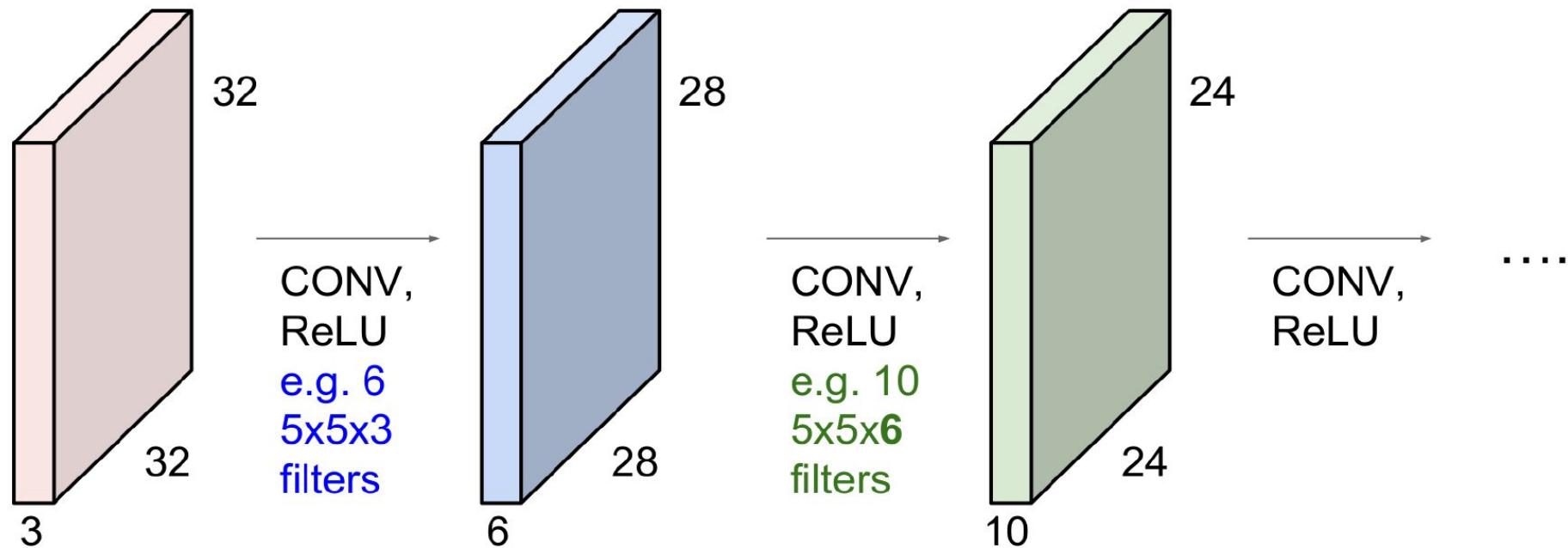


1	0	1	0	1
0	1	0	1	0
1	0	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1

Intro to CNN

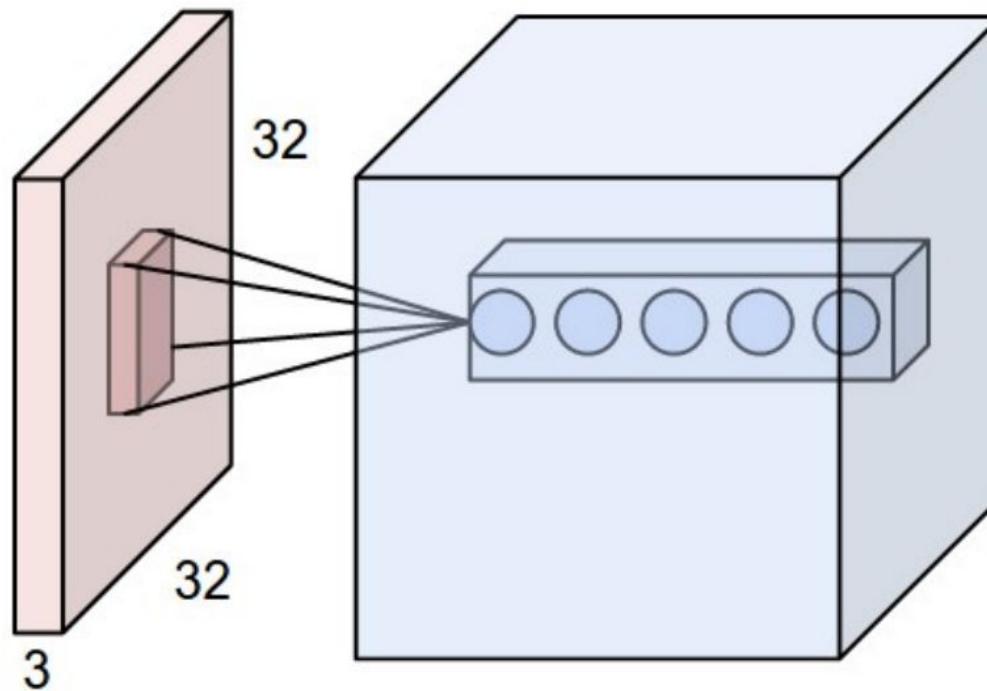


Intro to CNN



images from <http://cs231n.stanford.edu/>

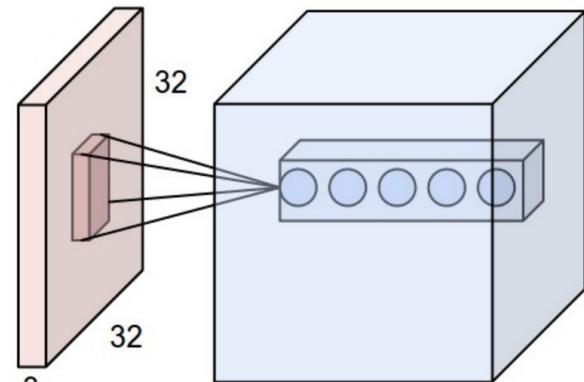
Intro to CNN



images from <http://cs231n.stanford.edu/>

Intro to CNN

- ▶ Accepts a volume of size $W_1 \times H_1 \times D_1$
- ▶ Requires four hyperparameters:
 - ▶ Number of filters K ,
 - ▶ their spatial extent F ,
 - ▶ the stride S ,
 - ▶ the amount of zero padding P .
- ▶ Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - ▶ $W_2 = (W_1 - F + 2P)/S + 1$,
 - ▶ $H_2 = (H_1 - F + 2P)/S + 1$
 - ▶ $D_2 = K$
- ▶ With parameter sharing, it introduces $F \times F \times D_1$ weights per filter, for a total of $(F \times F \times D_1) \times K$ weights and K biases.



Intro to CNN

- ▶ Convolution leverages four ideas that can help ML systems:
 - Sparse interactions
 - Parameter sharing
 - Equivariant representations $f(g(\mathbf{x})) = g(f(\mathbf{x}))$
 - Ability to work with inputs of variable size

Intro to CNN

We can use one single convolutional layer to modify a certain image

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$



$$\begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$



$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$



Intro to CNN

In training, we don't
specify kernels.
We learn kernels!

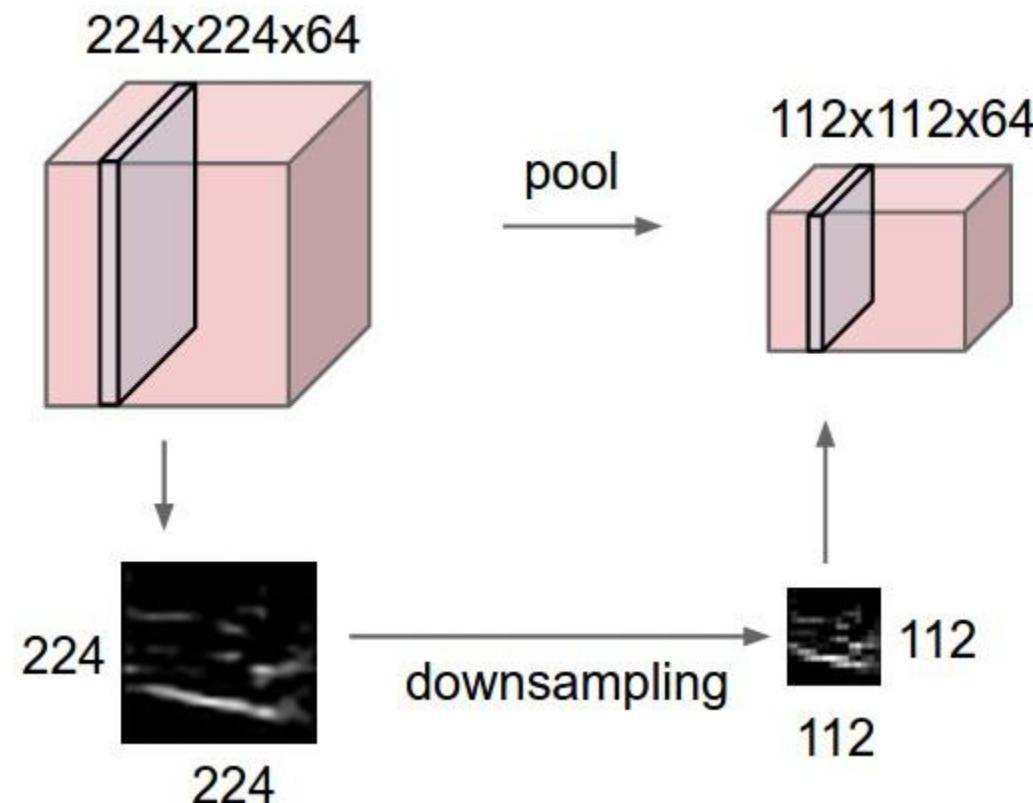


images from <http://cs231n.stanford.edu/>

Intro to CNN

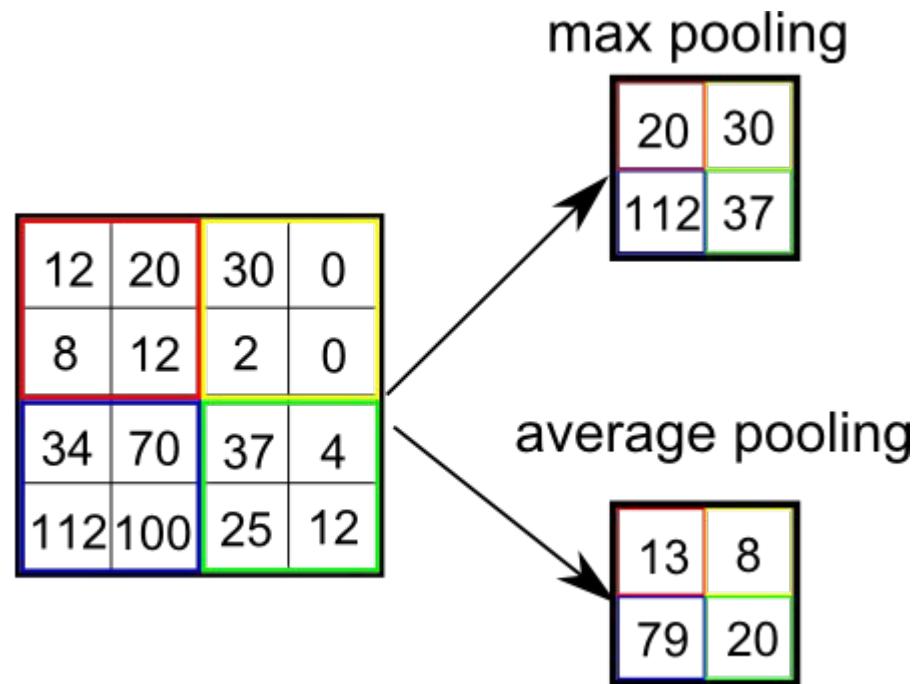
- ▶ INPUT holds the raw pixel values of the image.
- ▶ CONV layer computes the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and the region they are connected to in the input volume.
- ▶ POOL layer performs a downsampling operation along the spatial dimensions (width, height).
- ▶ FC (i.e. fully-connected) layer computes the class scores. As with ordinary Neural Networks and as the name implies, each neuron in this layer is connected to all the numbers in the previous volume.

Intro to CNN

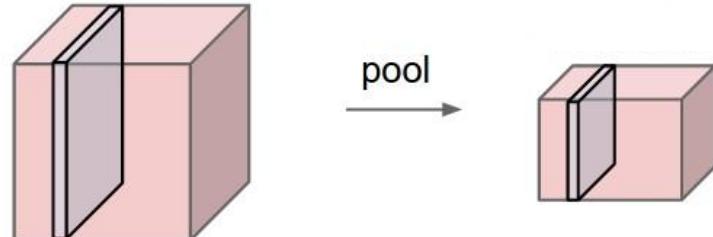


images from <http://cs231n.stanford.edu/>

Intro to CNN



Intro to CNN

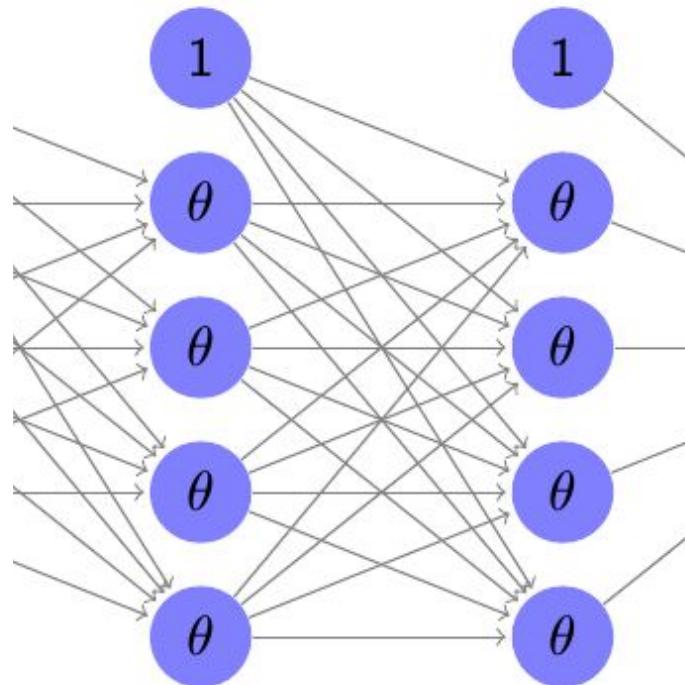


- ▶ Accepts a volume of size $W_1 \times H_1 \times D_1$
- ▶ Requires three hyperparameters:
 - ▶ their spatial extent F ,
 - ▶ the stride S ,
- ▶ Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - ▶ $W_2 = (W_1 - F)/S + 1$
 - ▶ $H_2 = (H_1 - F)/S + 1$
 - ▶ $D_2 = D_1$

Intro to CNN

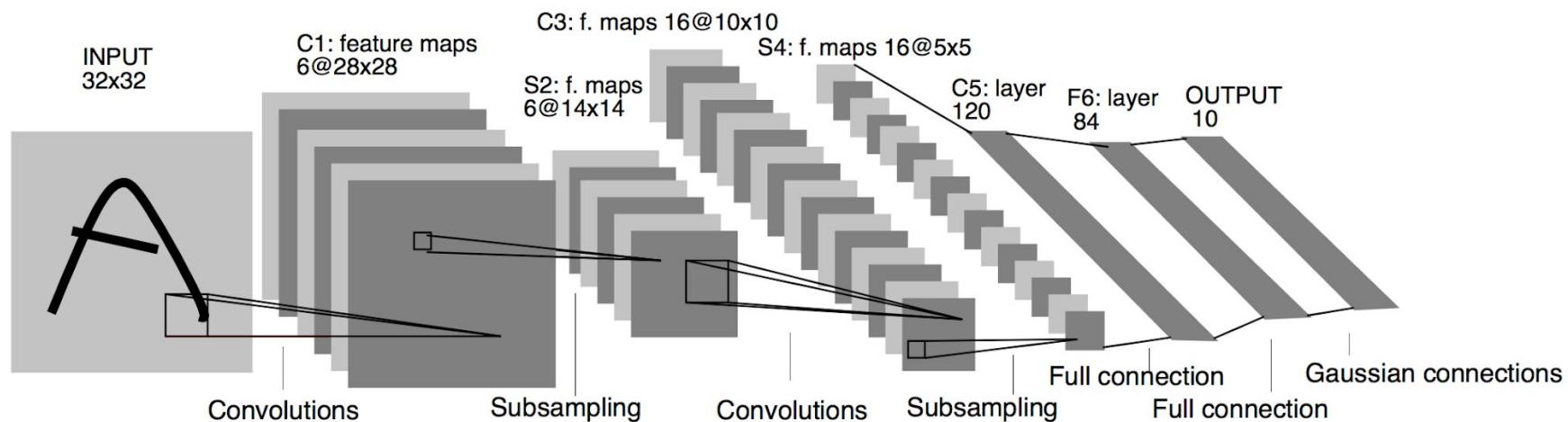
- ▶ INPUT holds the raw pixel values of the image.
- ▶ CONV layer computes the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and the region they are connected to in the input volume.
- ▶ POOL layer performs a downsampling operation along the spatial dimensions (width, height).
- ▶ FC (i.e. fully-connected) layer computes the class scores. As with ordinary Neural Networks and as the name implies, each neuron in this layer is connected to all the numbers in the previous volume.

Intro to CNN



Intro to CNN

LeNet-5 [1998, paper by LeCun et al.]



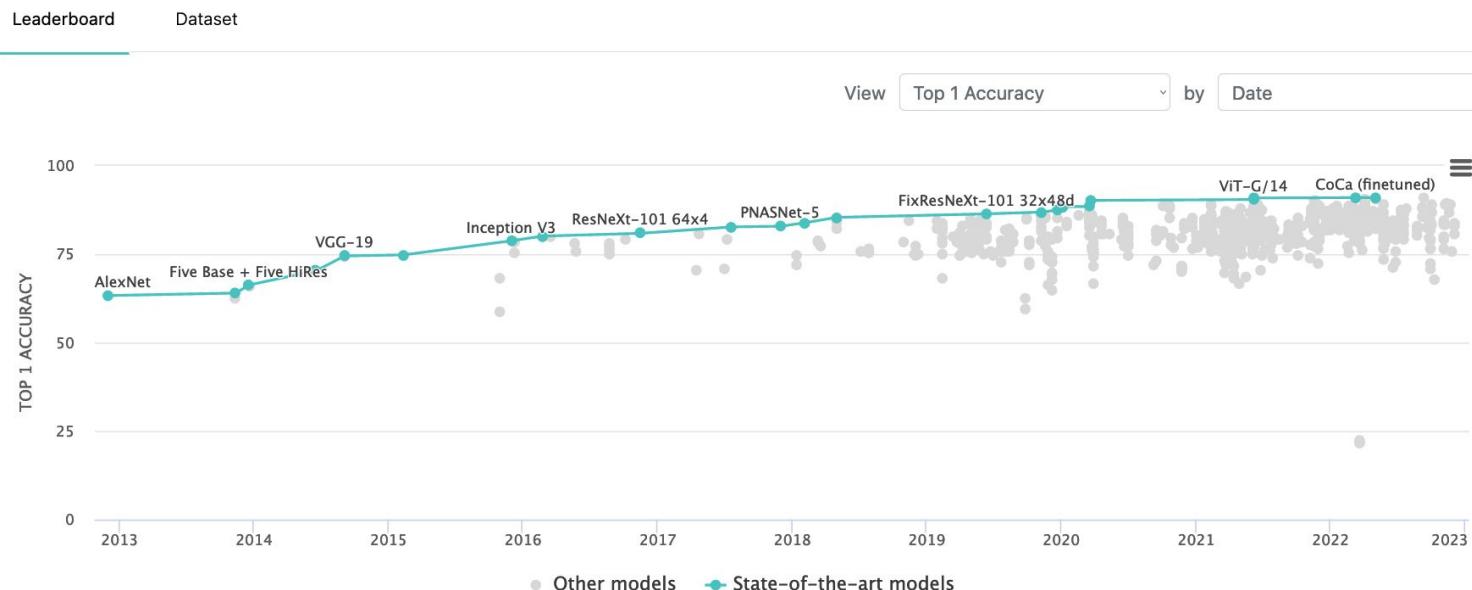
Intro to CNN

[CNN explainer](#)



Intro to CNN

Image Classification on ImageNet



<https://paperswithcode.com/sota/image-classification-on-imagenet>

Intro to CNN

