# Deep Learning for Computer Vision

Andrii Liubonko

Grammarly*

# Logistics

4 units

2 types of homework
- paper review
- mini-project

01  december, 23:59, approve for a paper
09 december, 23:59, paper review                    (SOFT, PENALTY 30%)
*30 december, 23:59, deadline*                       *(HARD)*

https://github.com/lyubonko/ucu2020cv

# Overview of the course

**Unit I**                                                  (26 Nov, 11.30-14.30)
    [T] Intro to Convolution Neural Networks (CNNs)
    [P] pytorch

**Unit II**                                                 (28 Nov, 15.00-18.00)
    [T] CNNs in depth
    [P] classification

**Unit III**                                                (10 Dec, 14.00-17.00)
    [T] Attention in CV
    [P] Transformers

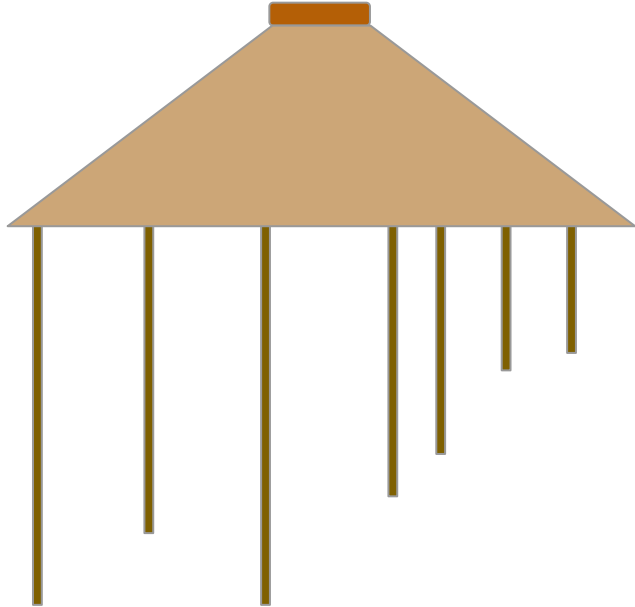**Unit IV**                                                 (11 Dec, 10.00-13.00)
    [T] Object Detection
    [P] project structure, detection

# Overview of the course

**Unit I**                                        (26 Nov, 11.30-14.30)
 [T] Intro to Convolution Neural Networks (CNNs)
 [P] pytorch

**Unit II**                                       (28 Nov, 15.00-18.00)
 [T] CNNs in depth
 [P] simple net, classification

**Unit III**                                      (10 Dec, 14.00-17.00)
 [T] Attention in CV
 [P] Transformers

**Unit IV**                                       (11 Dec, 10.00-13.00)
 [T] Object Detection
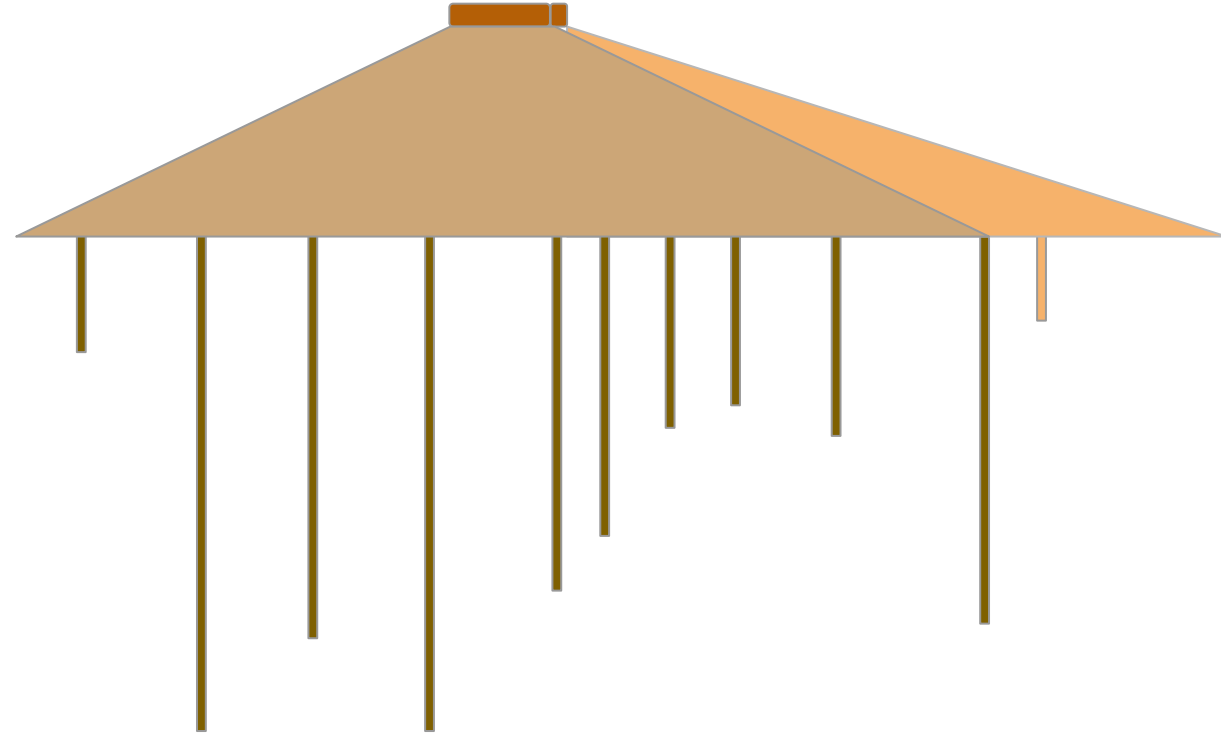 [P] project structure, detection

# Goals of the Course

- working knowledge of essential elements/blocks of Convolutional Neural Networks (CNNs)

- modern CNNs architectures

- attention in Computer Vision (current trend)

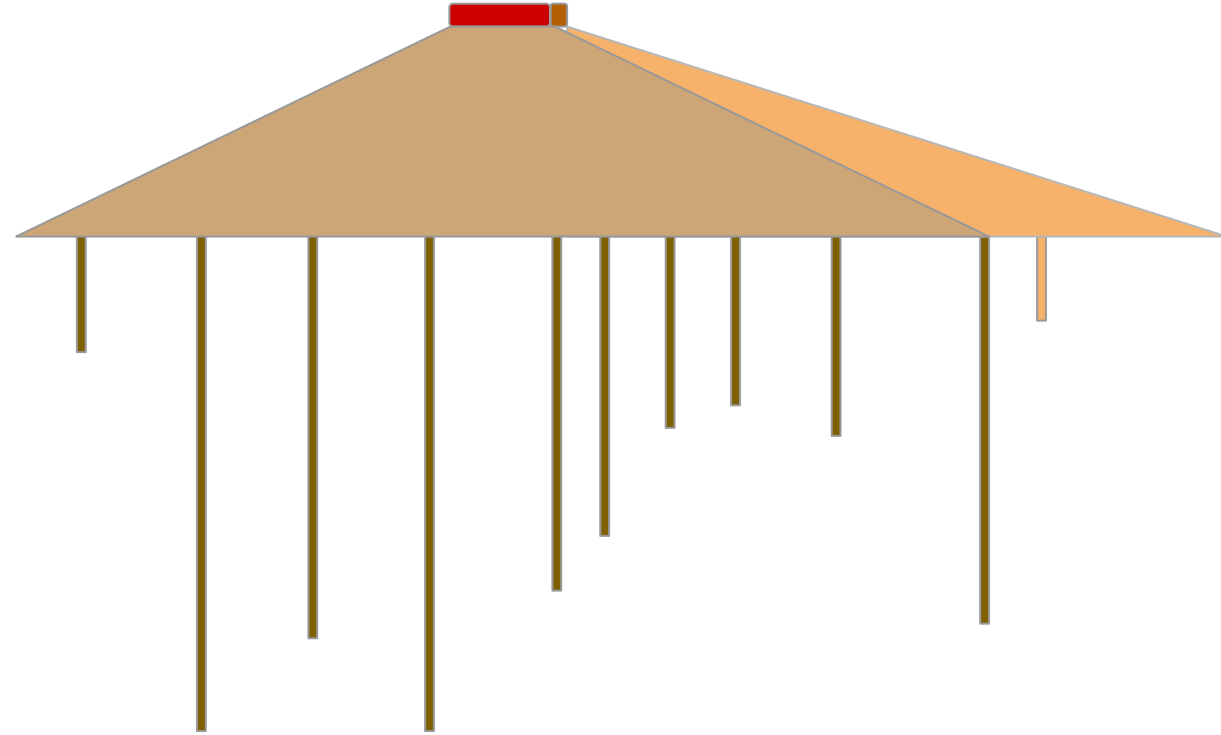- get deeper with one particular problem (Object Detection)

# Goals of the Course
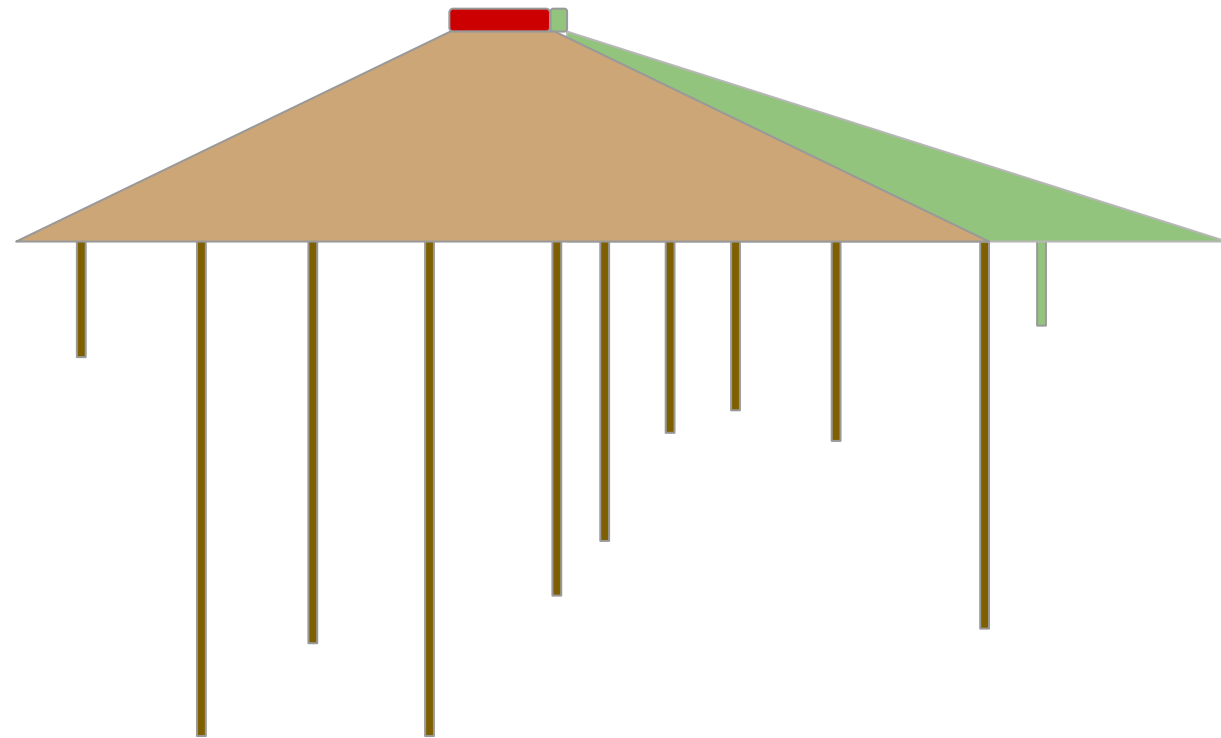
# Goals of the Course

# Goals of the Course

- essential CNNs elements/blocks
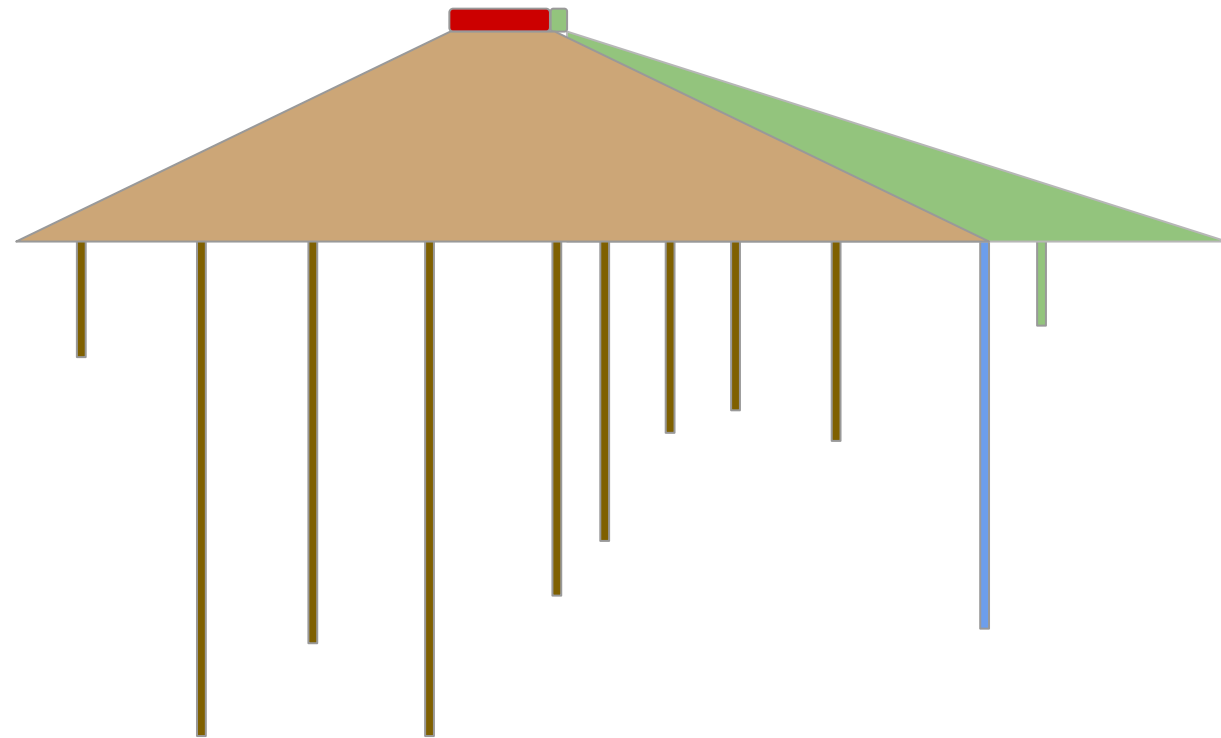
# Goals of the Course

- essential CNNs elements/blocks

- modern CNNs architectures

# Goals of the Course



- essential CNNs elements/blocks

- modern CNNs architectures
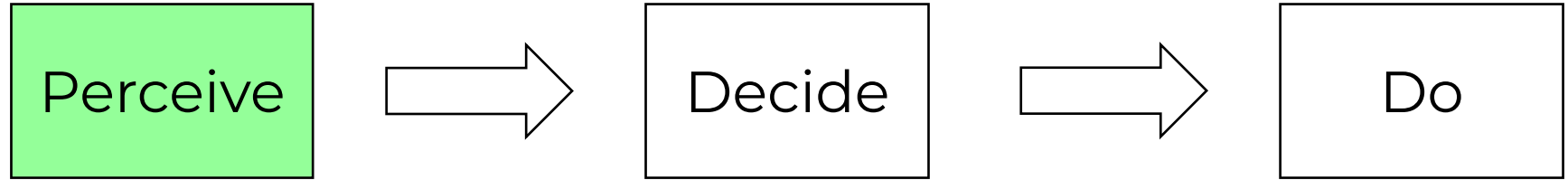
- attention in CV

# Goals of the Course



- essential CNNs elements/blocks

- modern CNNs architectures
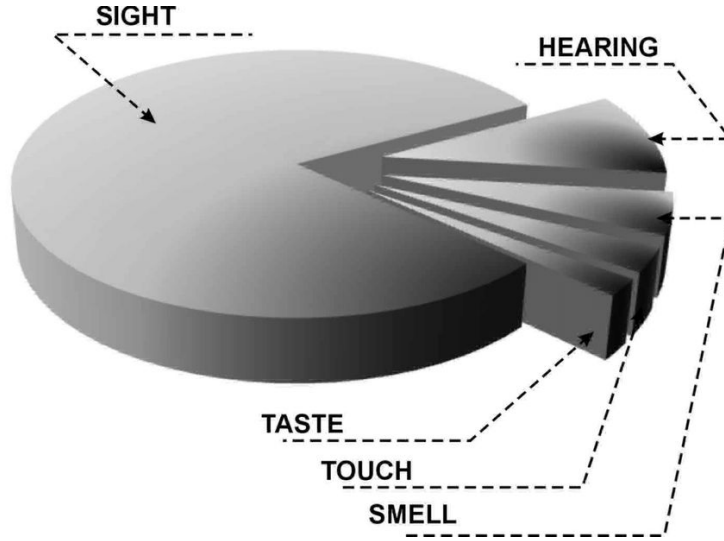
- attention in CV

- Object Detection

# Content of today lecture

- **Intro**

- DL review
  - neural networks
  - training & testing
  - supervised, semi-supervised, self-training

- Main components of CNNs (motivation and details)
  - convolutional layer
  - pooling layer

- Datasets
  - ImageNet

# Intro

Perceive ⟹ Decide ⟹ Do
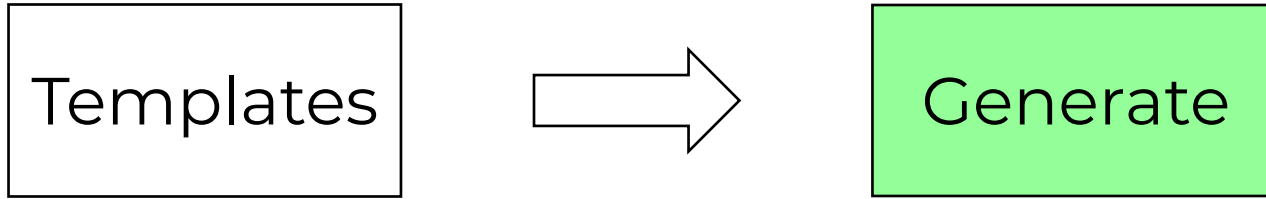
# Intro





*The goal of **computer vision** is to extract useful information from visual input (images, 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

Templates $\Rightarrow$ Generate

# Intro

## Generative adversarial network (GAN)

# Intro



what humans see



what computers see

# Intro

```
Image  →  Human Crafted Features  →  ML methods  →  Result
```

Image → Human Crafted Features → ML methods → Result

Descriptors: SIFT, SURF, etc.

Linear, SVM, Trees

**Intro**



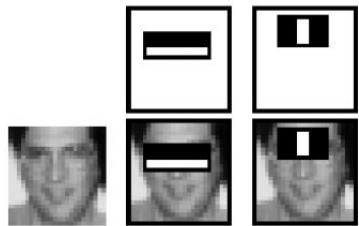Human Crafted Features → ML methods (crossed out)

Image → Neural Networks → Result
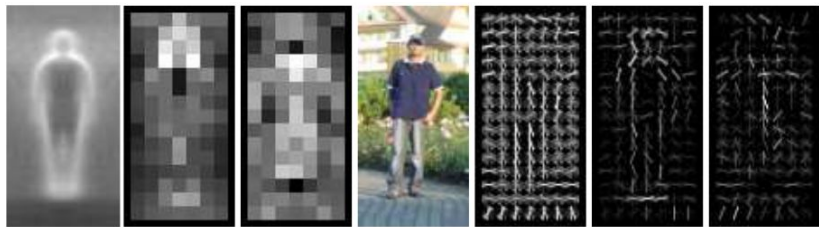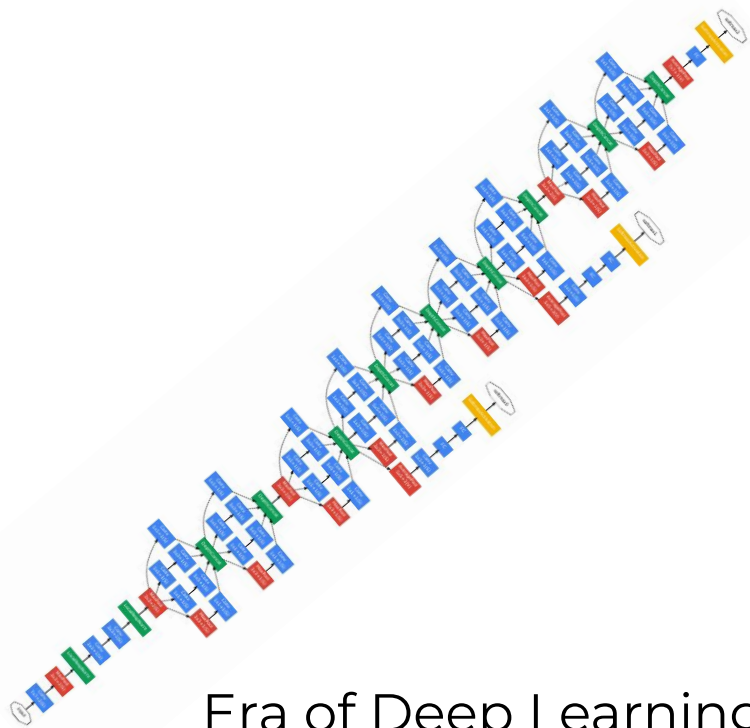
# Intro



Era of Human-Crafter Features

Era of Deep Learning

1986
BackProp

1998
LeNet

2012
AlexNet

# Intro

Models/
Algorithms

Big
Dataset

Computational
resources

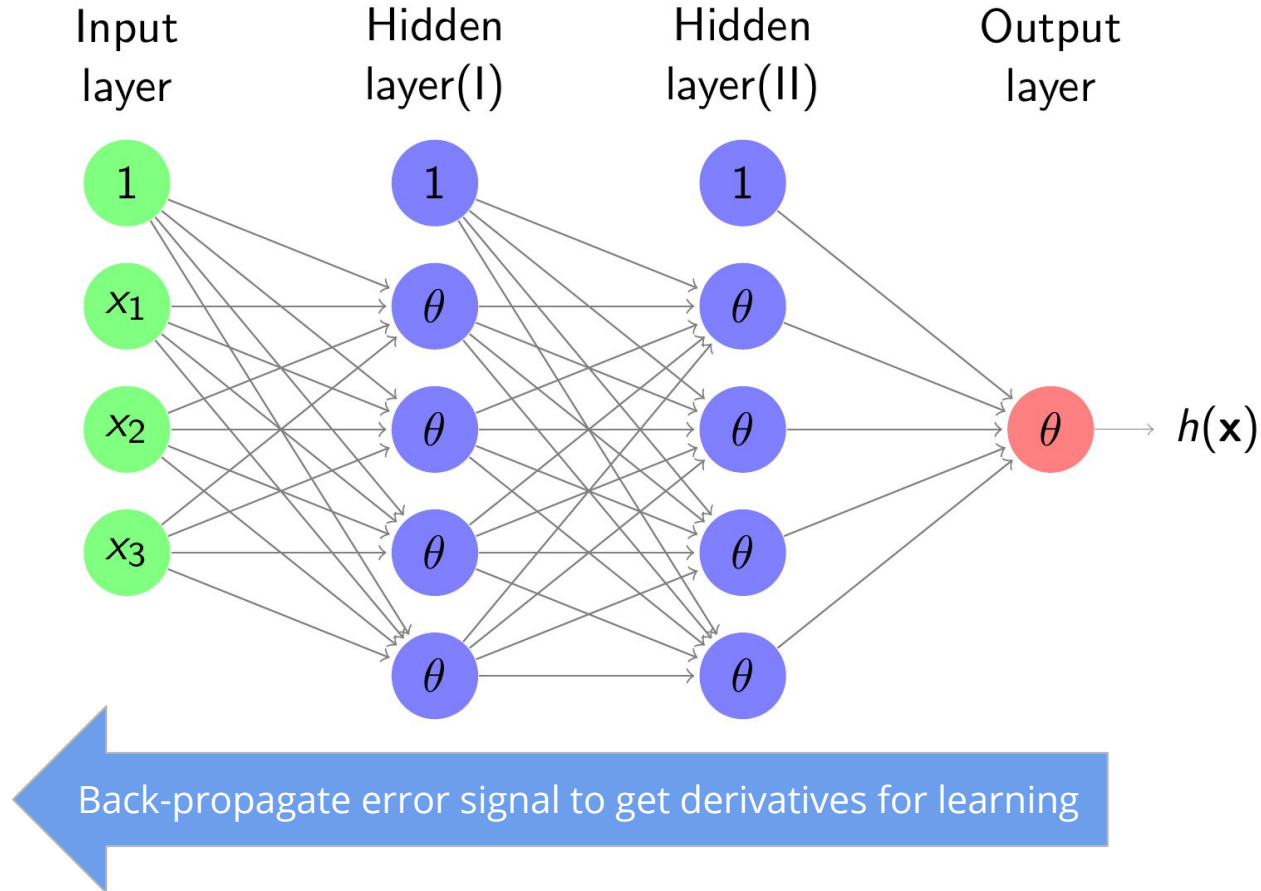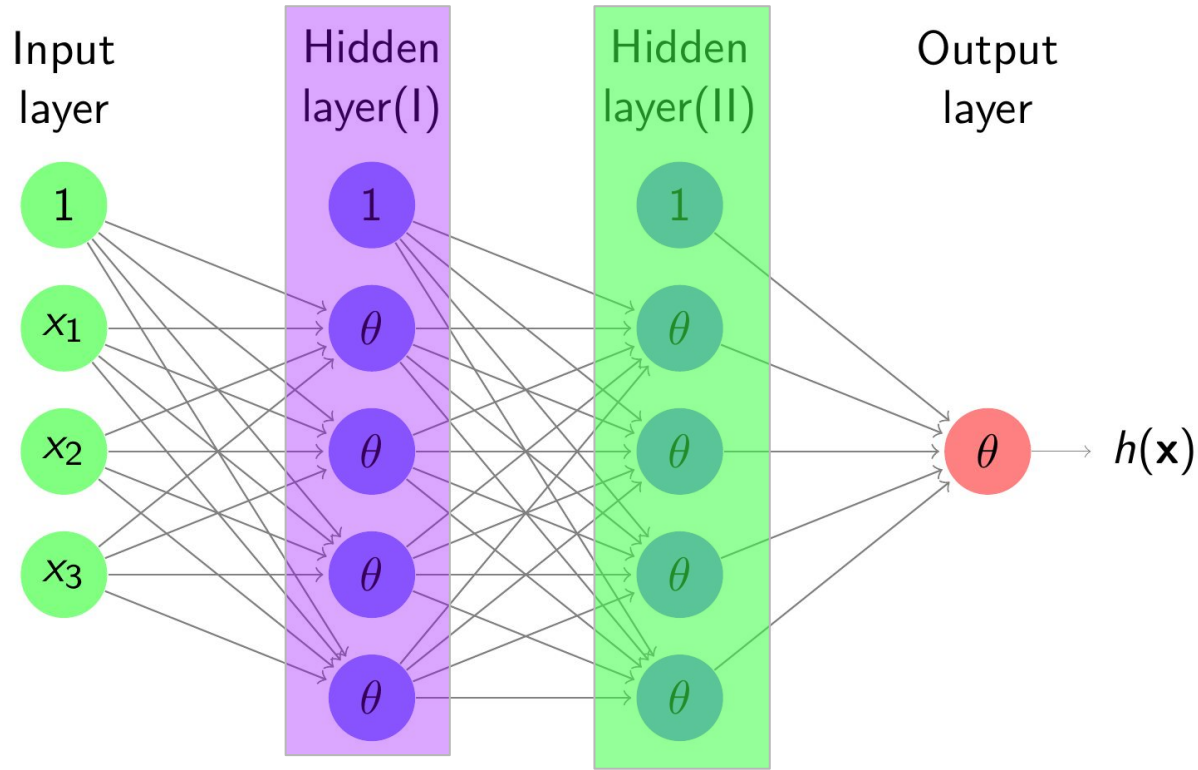# Content

- Intro

- **DL review**
  - neural networks
  - training & testing
  - loss functions
  - supervised, semi-supervised

- Main components of CNNs (motivation and details)
  - convolutional layer
  - pooling layer

- Datasets & Metrics
  - ImageNet

# Neural Networks

# Neural Networks



$$\left\{ f(\boldsymbol{x}; \boldsymbol{\theta}) = \mathbf{W}_L \boldsymbol{\sigma}_L (\mathbf{W}_{L-1} \cdots \boldsymbol{\sigma}_2 (\mathbf{W}_2 \boldsymbol{\sigma}_1 (\mathbf{W}_1 \boldsymbol{x})) ) \mid \boldsymbol{\theta} = \{ \mathbf{W}_1, \ldots, \mathbf{W}_L \} \right\}$$
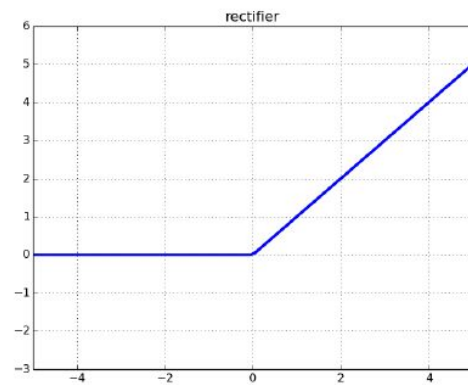
# Neural Networks

parameters in NN:

$$W_l^{ij} = \begin{cases} 1 \le l \le L & \text{layers} \\ 0 \le i \le d^{(l-1)} & \text{inputs} \\ 1 \le j \le d^{(l)} & \text{outputs} \end{cases}$$

activation:



rectifier

$$x_j^{(l)} = \sigma(s_j^{(l)}) = \sigma\left( \sum_{i=0}^{d^{(l-1)}} W_l^{ij} x_i^{(l-1)} \right)$$

$$\sigma(s) = RELU(s) = max(0, x)$$

# Neural Networks

Define **Loss (or Cost) function**:

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

**Gradient Descent (GD)** minimizes:

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

by iterative steps along $-\nabla L_{train}$:

$$\Delta \omega = -\eta \nabla L_{train}(\omega)$$
$$\omega_{prev} = \omega_{next} + \Delta \omega$$

# Neural Networks

# Neural Networks

# Neural Networks (supervised way)



{ [plane image] , plane }

{ [cat image] , cat }

...

$F(\mathbf{x_n})$

$e(F(\mathbf{x_n}), y_n)$

$e(F(\mathbf{x_n}), y_n)$

...

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

# Neural Networks (self-supervised way)

SimCLR



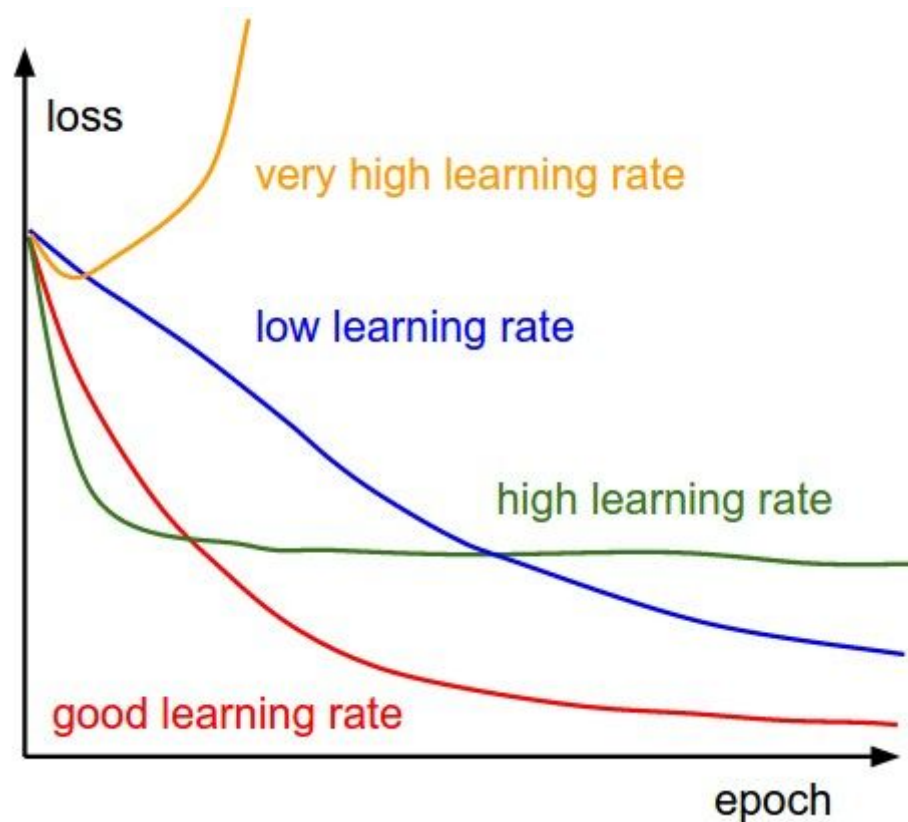$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\boldsymbol{z}_i, \boldsymbol{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\boldsymbol{z}_i, \boldsymbol{z}_k)/\tau)}$$

$$\boldsymbol{z}_i = g(\boldsymbol{h}_i) = W^{(2)}\sigma(W^{(1)}\boldsymbol{h}_i)$$

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

2002.05709

# Self-supervised

SimCLR

# Self-supervised

SimCLR



Representation

Original Image → Data Augmentation (T) → Transformed Images ($x_i$, $x_j$) → Base Encoder f(.) (Encoder) → $h_i$, $h_j$ → Projection Head g(.) (Dense | Relu | Dense) → $z_i$, $z_j$ → Maximize similarity

Finetuning

classification, detection, ...

# Neural Networks (semi-supervised)



{ [plane image], plane }

{ [cat image], cat }

Train teacher model with labeled data → Infer pseudo-labels on unlabeled data

Data augmentation
Dropout
Stochastic depth
→ Train **equal-or-larger** student model with combined data and noise injected → Make the student a new teacher

1911.04252

# Content

- Intro

- DL review
  - neural networks
  - training & testing
  - supervised, semi-supervised

- **Main components of CNNs (motivation and details)**
  - convolutional layer
  - pooling layer

- Datasets & Metrics
  - ImageNet

# Intro to CNN



Input layer
$3 \times 256 \times 256 = 200K$

Weights ($\omega_0$):
$200K \times$(size of next layer)

$=>$ millions of parameters
(for just one layer)

Need for alternative architecture

# Intro to CNN

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



INPUT 32x32

C1: feature maps 6@28x28

S2: f. maps 6@14x14

C3: f. maps 16@10x10

S4: f. maps 16@5x5

C5: layer 120

F6: layer 84

OUTPUT 10

Convolutions

Subsampling

Convolutions

Subsampling

Full connection

Full connection

Gaussian connections

# 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.
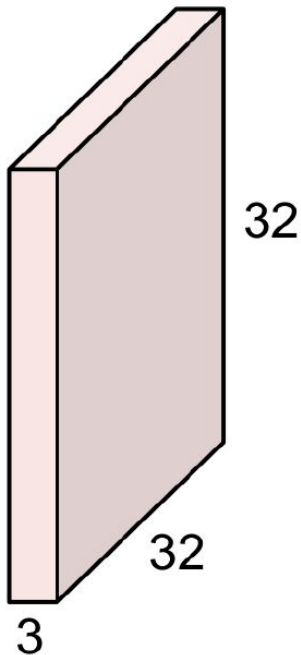
# Intro to CNN



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

32x32x3 image



32

32

3

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 $w$

32

32

3

1 number:
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

# Intro to CNN



32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

activation maps

28

28

1

# Intro to CNN



32 32 3 → CONV, ReLU e.g. 6 5x5x3 filters → 28 28 6 → CONV, ReLU e.g. 10 5x5x**6** filters → 24 24 10 → CONV, ReLU → ....

# Intro to CNN



32

32

3

*images from http://cs231n.stanford.edu/*

# Intro to CNN

| | |
|---|---|
| $1_{\times 1}$ $1_{\times 0}$ $1_{\times 1}$ 0 0 | 4 |

Image

Convolved Feature

# Intro to CNN



Image

Convolved Feature

# Intro to CNN



Image

Convolved Feature

# Intro to CNN



Image

Convolved Feature

# Intro to CNN

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



[ 1.  1.  1.]
[ 1.  1.  1.]
[ 1.  1.  1.]



[ 1.  2.  1.]
[ 0.  0.  0.]
[-1. -2. -1.]



[ 0. -1.  0.]
[-1.  5. -1.]
[ 0. -1.  0.]

# Intro to CNN

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

# Intro to CNN

- ► Accepts a volume of size $W1 \times H1 \times D1$
- ► 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 $W2 \times H2 \times D2$ where:
    - ► $W2 = (W1 - F + 2P)/S + 1$,
    - ► $H2 = (H1 - F + 2P)/S + 1$
    - ► $D2 = K$
- ► With parameter sharing, it introduces $F \times F \times D1$ weights per filter, for a total of $(F \times F \times D1) \times K$ weights and $K$ biases.
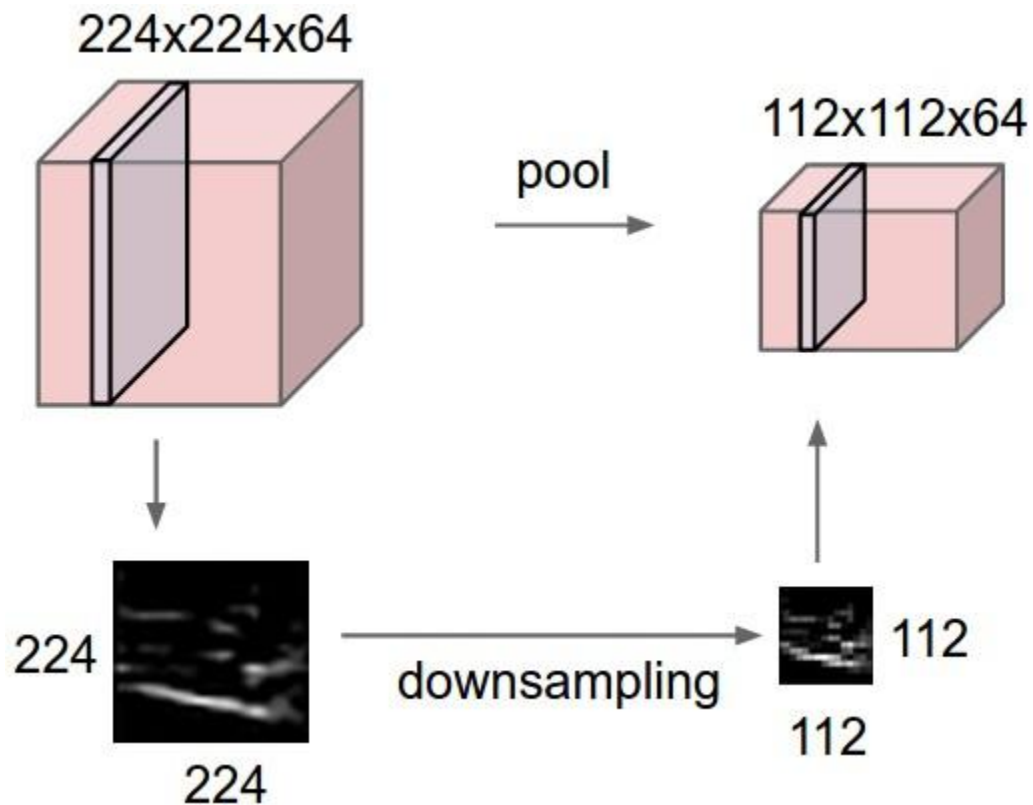
# Intro to CNN

► Convolution leverages four ideas that can help ML systems:

- Sparse interactions
- Parameter sharing
- Equivariant representations $f(g(\mathbf{x})) = g(f(\mathbf{x}))$
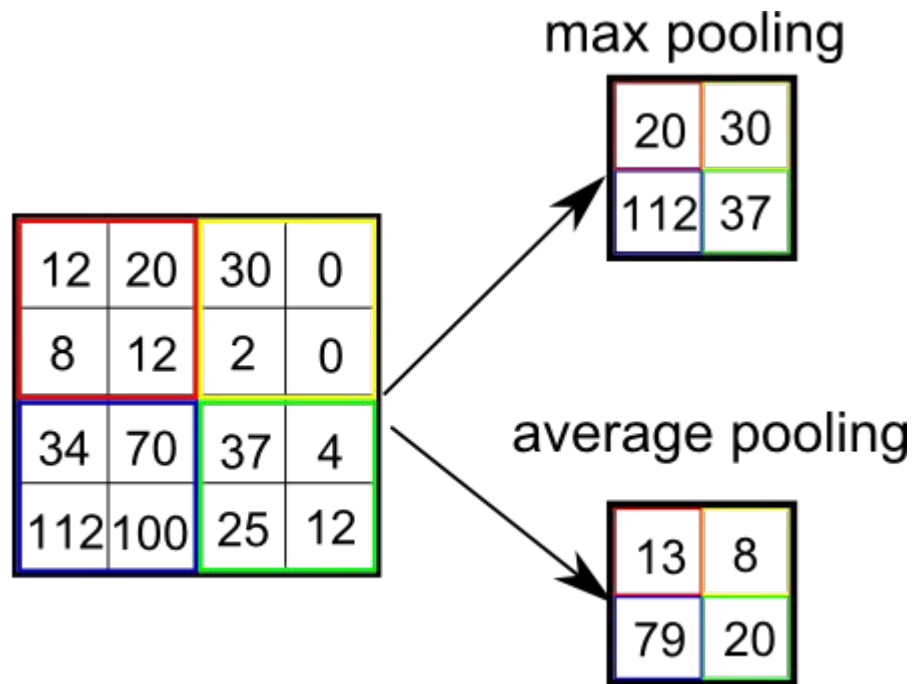- Ability to work with inputs of variable size

# Intro to CNN

▶ 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



224x224x64

pool →

112x112x64

224

224

downsampling →

112

112

*images from http://cs231n.stanford.edu/*

# Intro to CNN



max pooling

| 20 | 30 |
|----|----|
| 112 | 37 |

average pooling

| 13 | 8 |
|----|---|
| 79 | 20 |

Input:

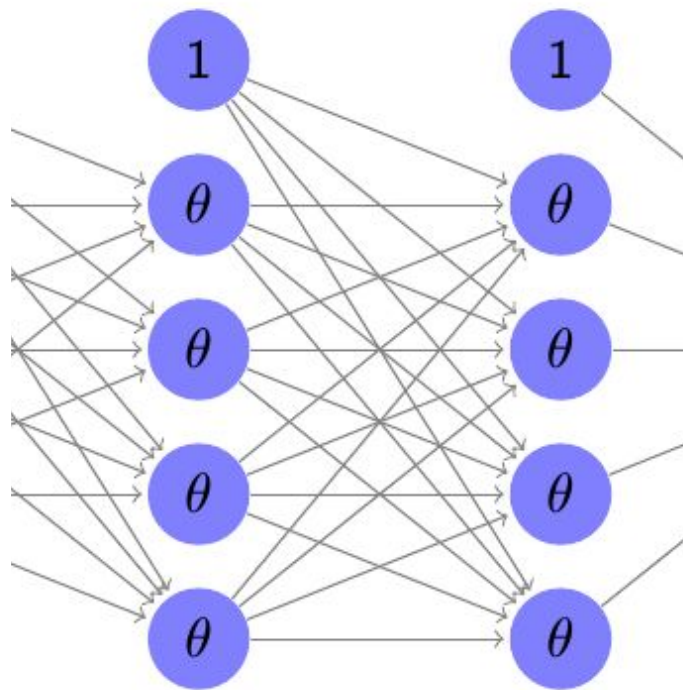| 12 | 20 | 30 | 0 |
|----|----|----|---|
| 8 | 12 | 2 | 0 |
| 34 | 70 | 37 | 4 |
| 112 | 100 | 25 | 12 |

# Intro to CNN

- Accepts a volume of size $W1 \times H1 \times D1$
- Requires three hyperparameters:
    - their spatial extent $F$,
    - the stride $S$,
- Produces a volume of size $W2 \times H2 \times D2$ where:
    - $W2 = (W1 - F)/S + 1$
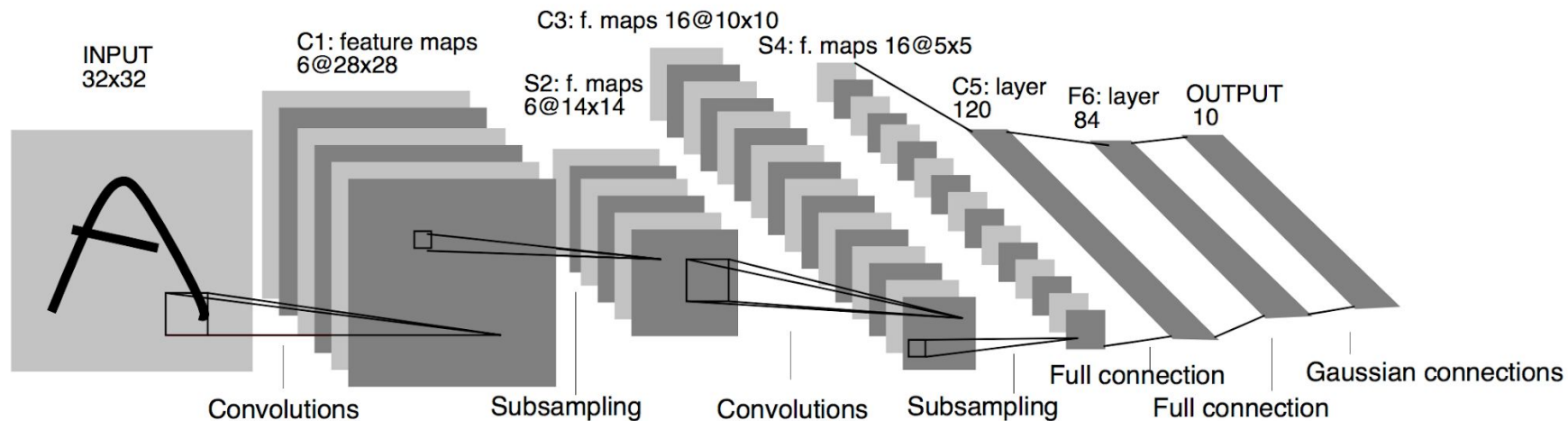    - $H2 = (H1 - F)/S + 1$
    - $D2 = D1$

# 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

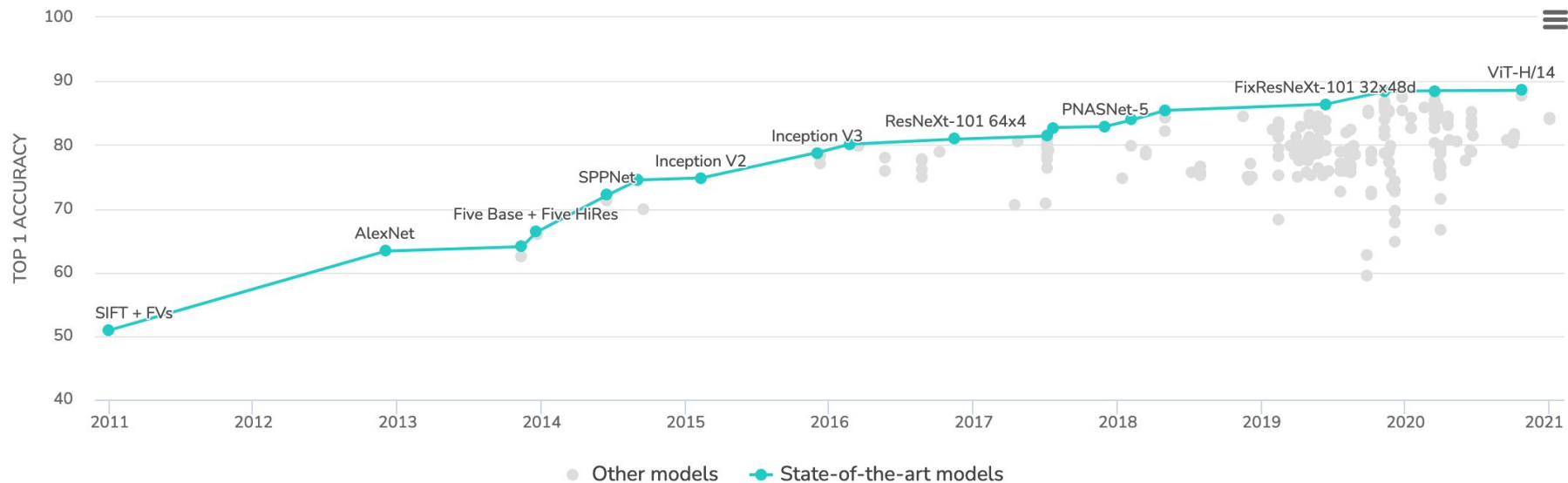LeNet-5 [1998, paper by LeCun et al.]

# Content

- Intro

- DL review
  - neural networks
  - training & testing
  - loss functions
  - supervised, semi-supervised

- Main components of CNNs (motivation and details)
  - convolutional layer
  - pooling layer

- **Dataset**
  - ImageNet

- Classes: 1000
- Problem: classification
- Set:
    - 1 200 000 examples

# Intro to CNN



TOP 1 ACCURACY

- SIFT + FVs
- AlexNet
- Five Base + Five HiRes
- SPPNet
- Inception V2
- Inception V3
- ResNeXt-101 64x4
- PNASNet-5
- FixResNeXt-101 32x48d
- ViT-H/14

Other models ● State-of-the-art models

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

# Intro to CNN