

ML for natural and physical scientists 2023 9

Deep Learning 2 - Convolutional NNs

this slide deck:

https://slides.com/federicabianco/mlpns23_9

- Machine Learning basic concepts
 - interpretability
 - parameters vs hyperparameters
 - supervised/unsupervised
- CART methods
- Clustering methods
- Neural Networks

- Neural Networks
 - the brain connection
 - perceptron
 - learning
 - activation functions
 - shallow nets
 - deep nets architecture
 - back-propagation
 - preprocessing and whitening (minibatch)

legend
networks

perceptrons

Perceptrons are **linear classifiers**:

makes its predictions based on a linear predictor function

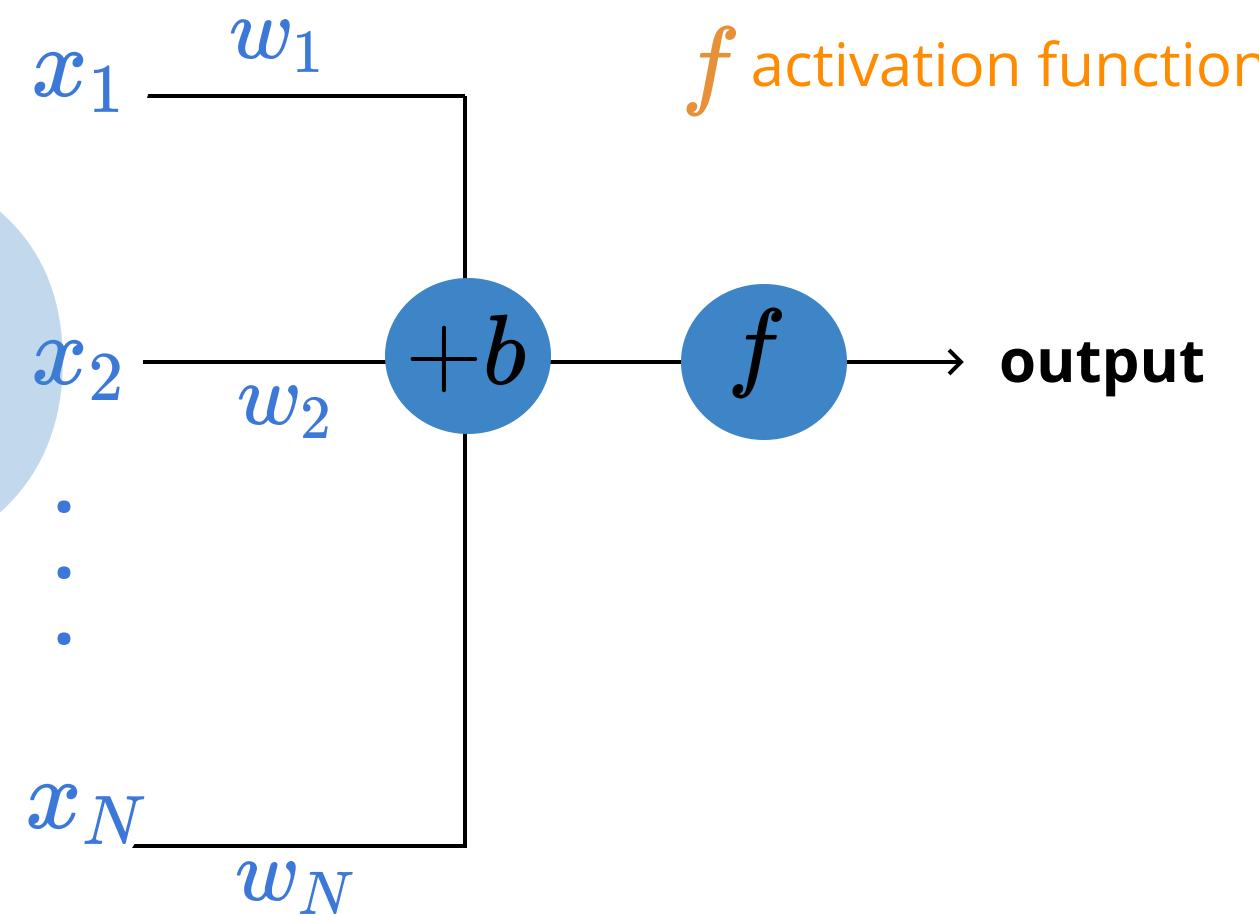
combining a set of weights

(=parameters) with the feature vector.

$$y = w \cdot x + b$$

$$y = \sum_i w_i x_i + b$$

$$y = f(\sum_i w_i x_i + b)$$



w_i weights

b bias

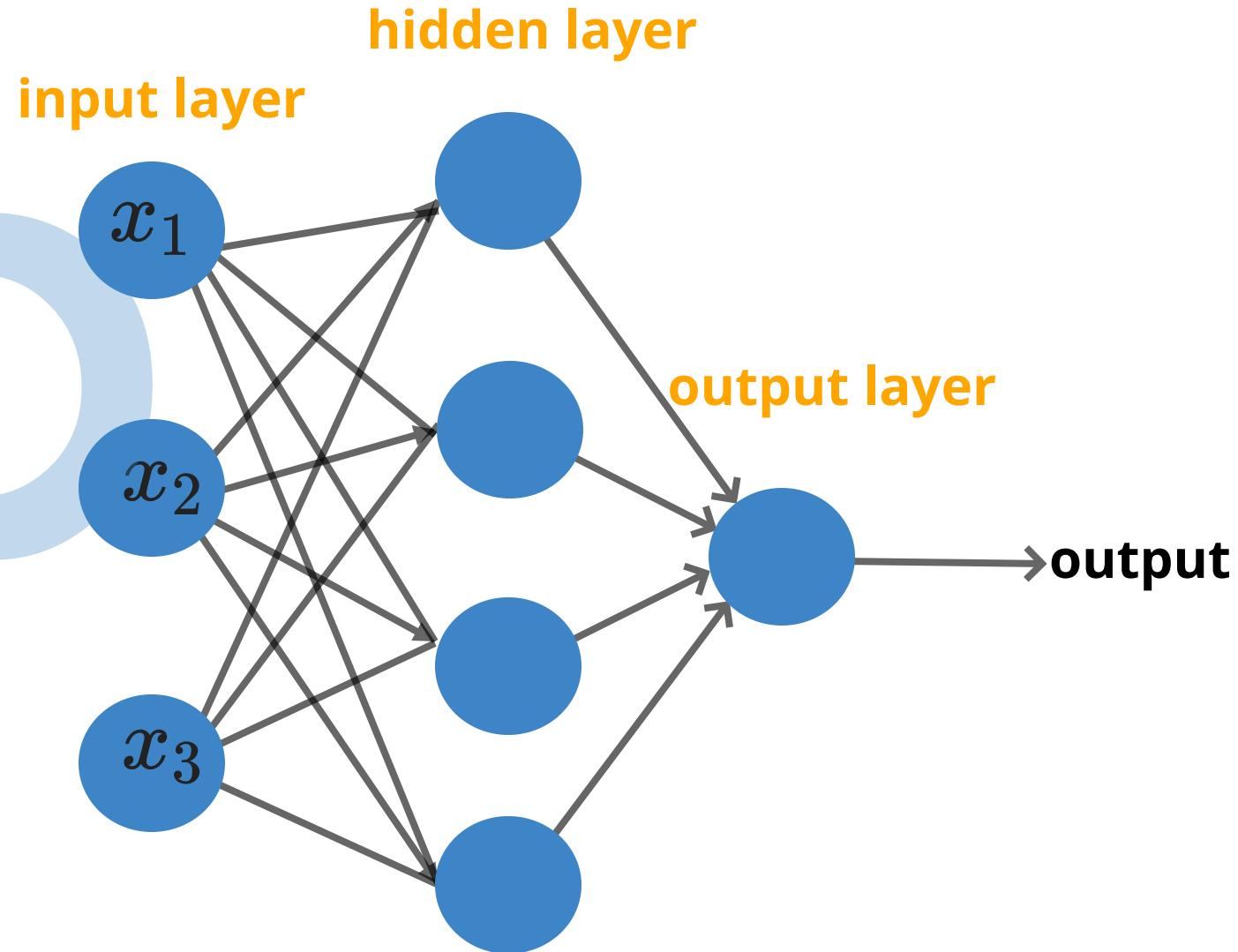
f activation function

multilayer perceptron

1970: multilayer
perceptron architecture

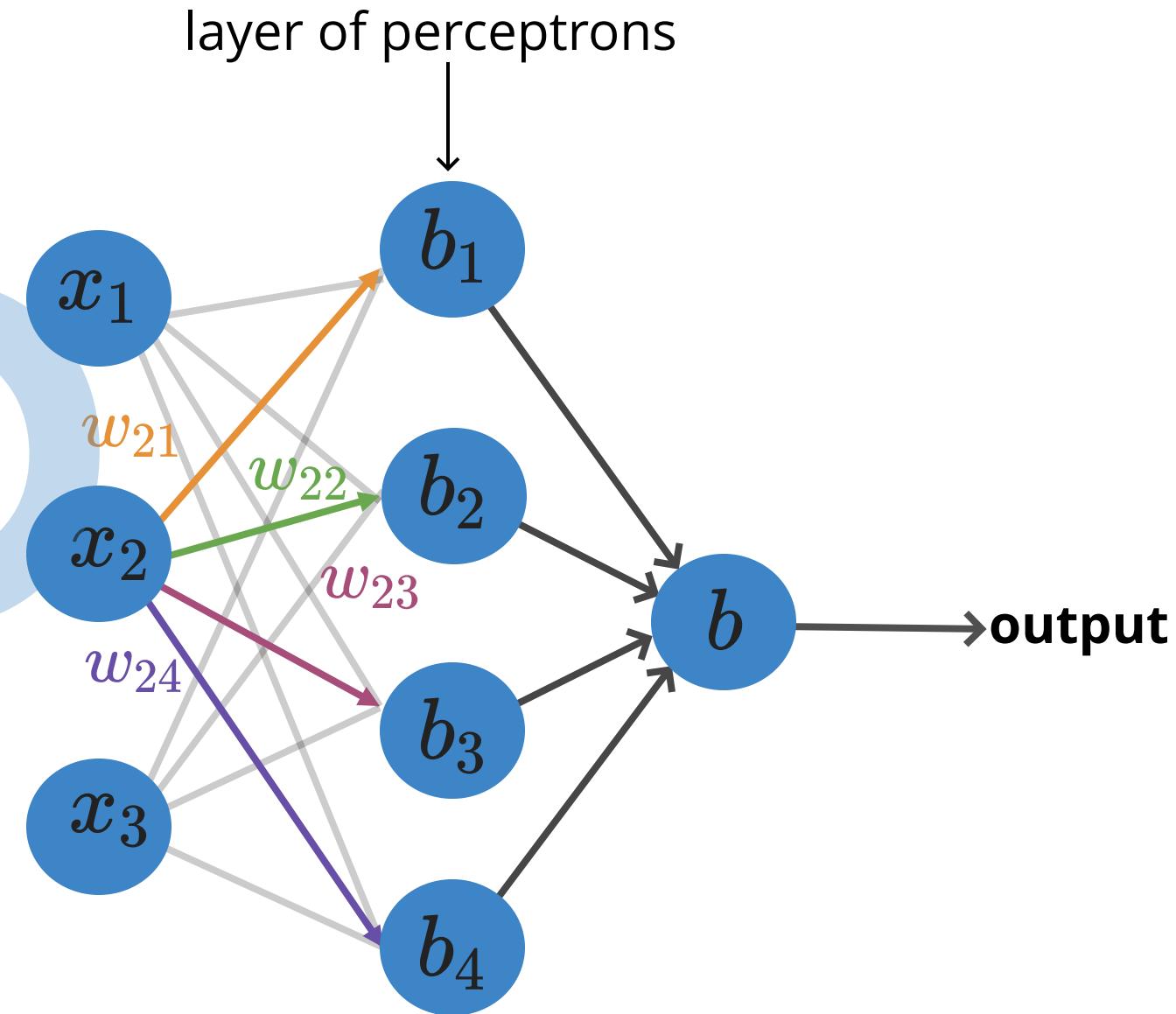
repo

Fully connected: all nodes go to
all nodes of the next layer.



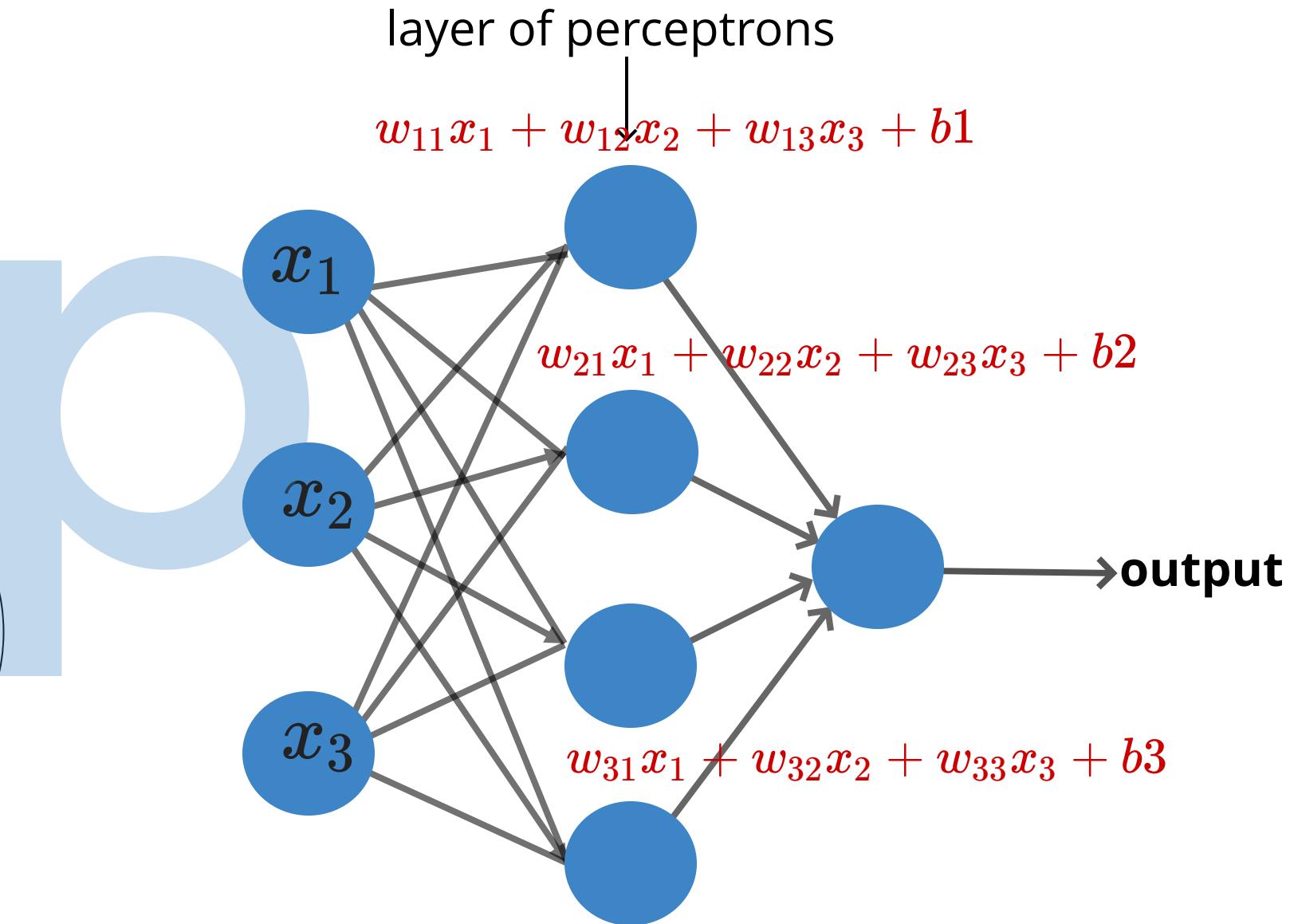
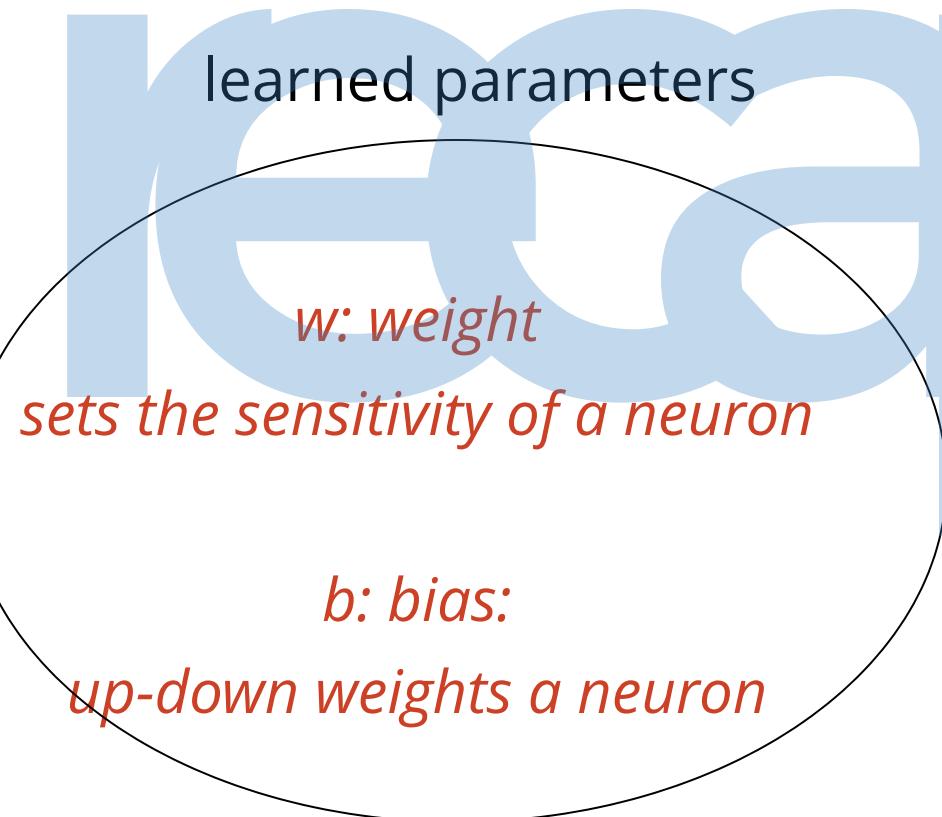
multilayer perceptron

recap



multilayer perceptron

Fully connected: all nodes go to all nodes of the next layer.



multilayer perceptron

what we are doing is exactly a series of matrix multiplications.

$$\begin{bmatrix} a_1 & a_2 & a_3 & \dots & a_n \\ b_1 & b_2 & b_3 & \dots & b_n \\ c_1 & c_2 & c_3 & \dots & c_n \\ \dots & \dots & \dots & \dots & \dots \\ m_1 & m_2 & m_3 & \dots & m_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{bmatrix} = \begin{bmatrix} (a_1 x_1) + (a_2 x_2) + (a_3 x_3) + \dots + (a_n x_n) \\ (b_1 x_1) + (b_2 x_2) + (b_3 x_3) + \dots + (b_n x_n) \\ (c_1 x_1) + (c_2 x_2) + (c_3 x_3) + \dots + (c_n x_n) \\ \dots \\ (m_1 x_1) + (m_2 x_2) + (m_3 x_3) + \dots + (m_n x_n) \end{bmatrix}$$

multilayer perceptron

Fully connected: all nodes go to all nodes of the next layer.

f: activation function:

turns neurons on-off

w: weight

sets the sensitivity of a neuron

b: bias:

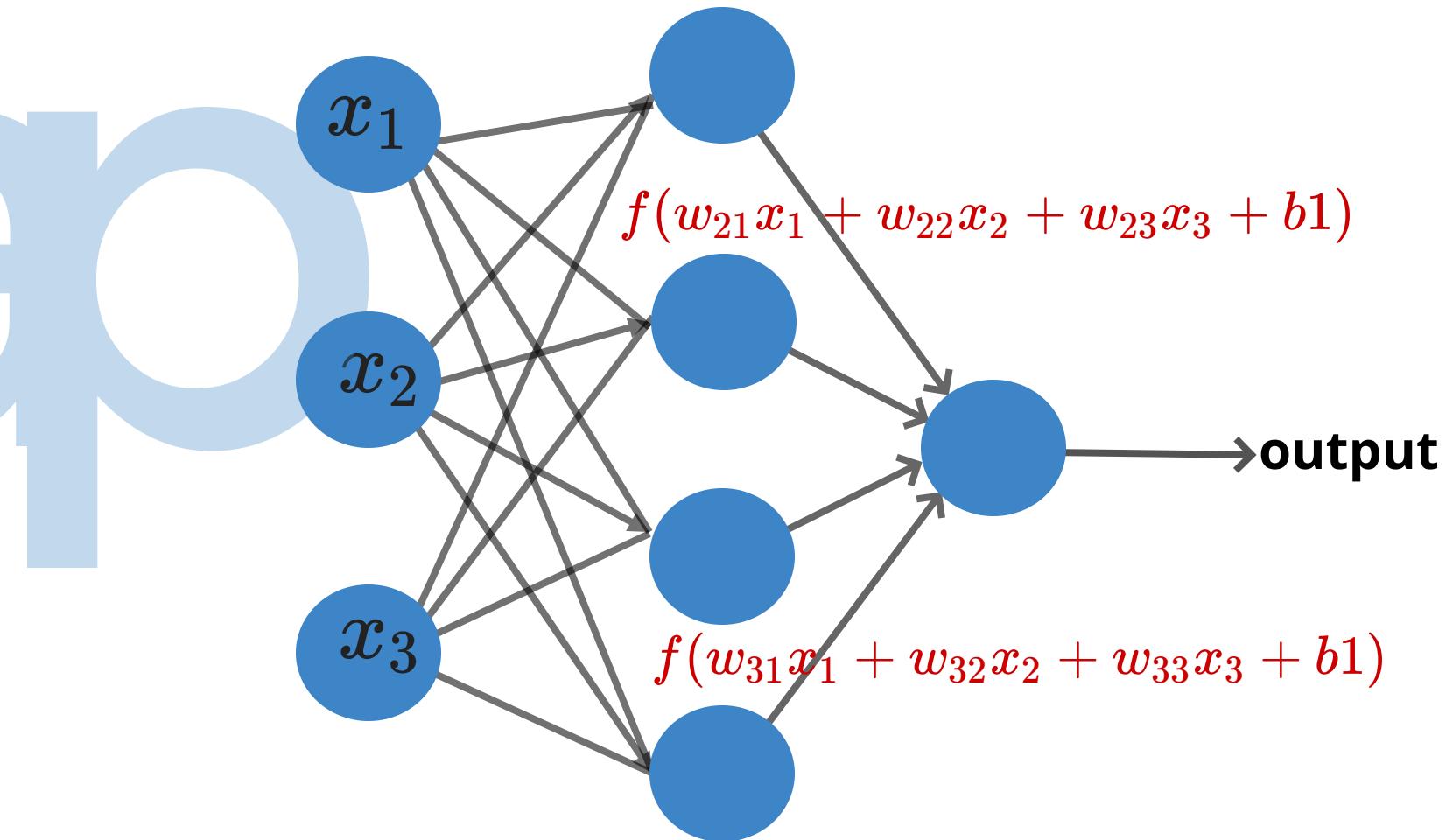
up-down weights a neuron

layer of perceptrons

$$f(w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + b_1)$$

$$f(w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + b_1)$$

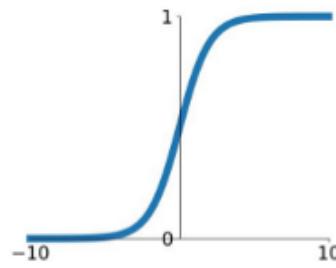
$$f(w_{31}x_1 + w_{32}x_2 + w_{33}x_3 + b_1)$$



activation functions

Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

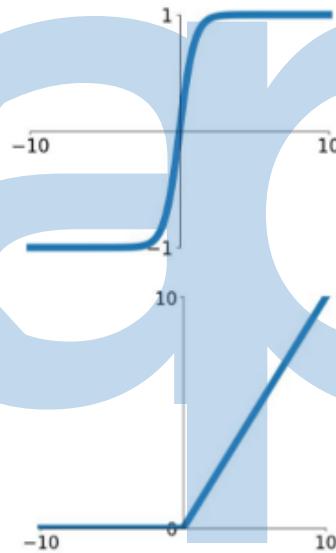


tanh

$$\tanh(x)$$

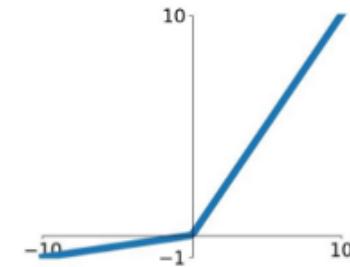
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

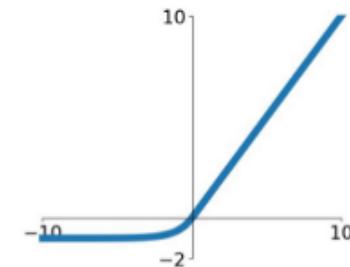


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



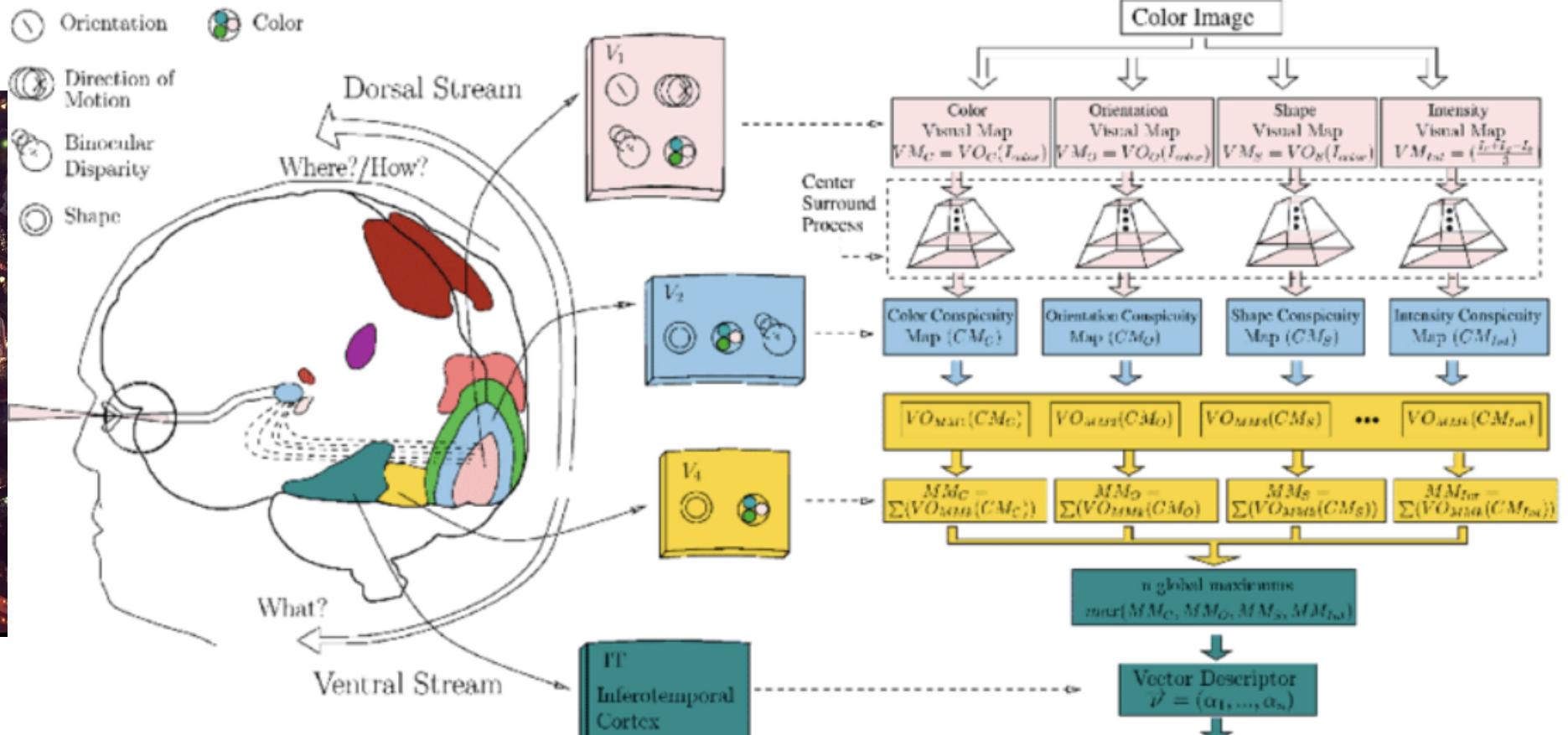
CNN

1

Convolutional Neural Nets

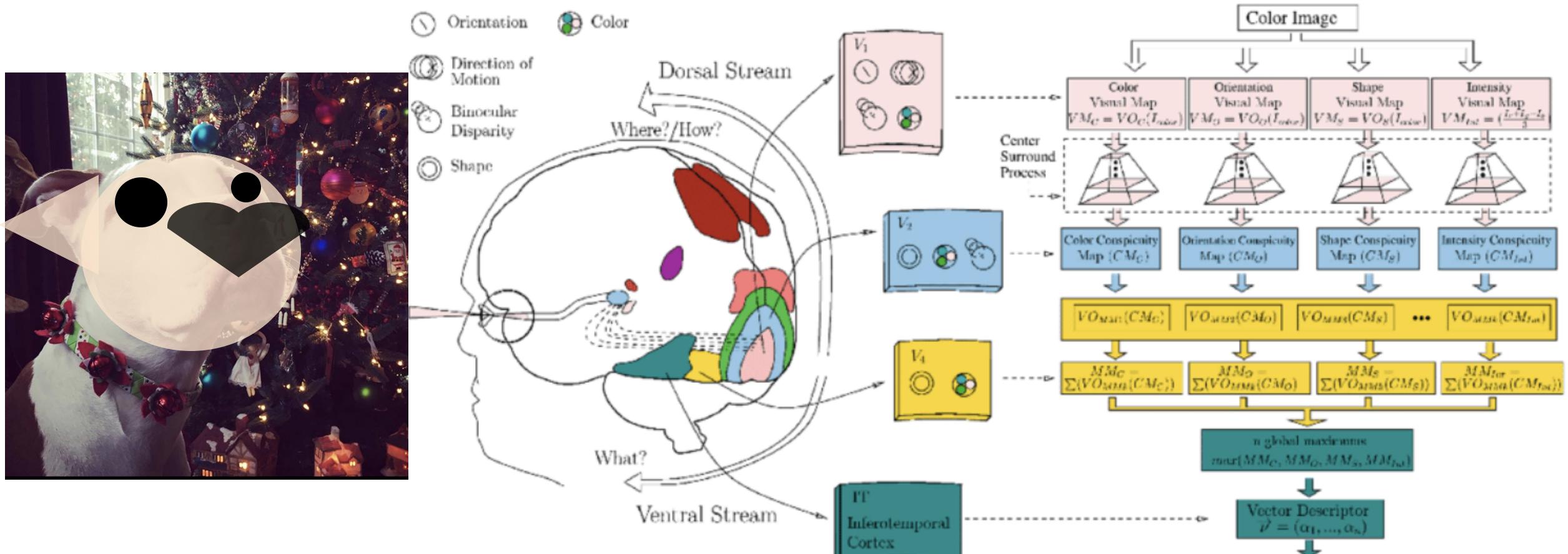


@akumadog



Brain Programming and the Random Search in Object Categorization Olague et al 2017

The visual cortex learns hierarchically: first detects simple features, then more complex features and ensembles of features



CNN

1a

Convolution

Convolution

convolution is a mathematical operator on two functions

f and g

that produces a third function

$f \otimes g$

expressing *how the shape of one is modified by the other.*

Convolution Theorem

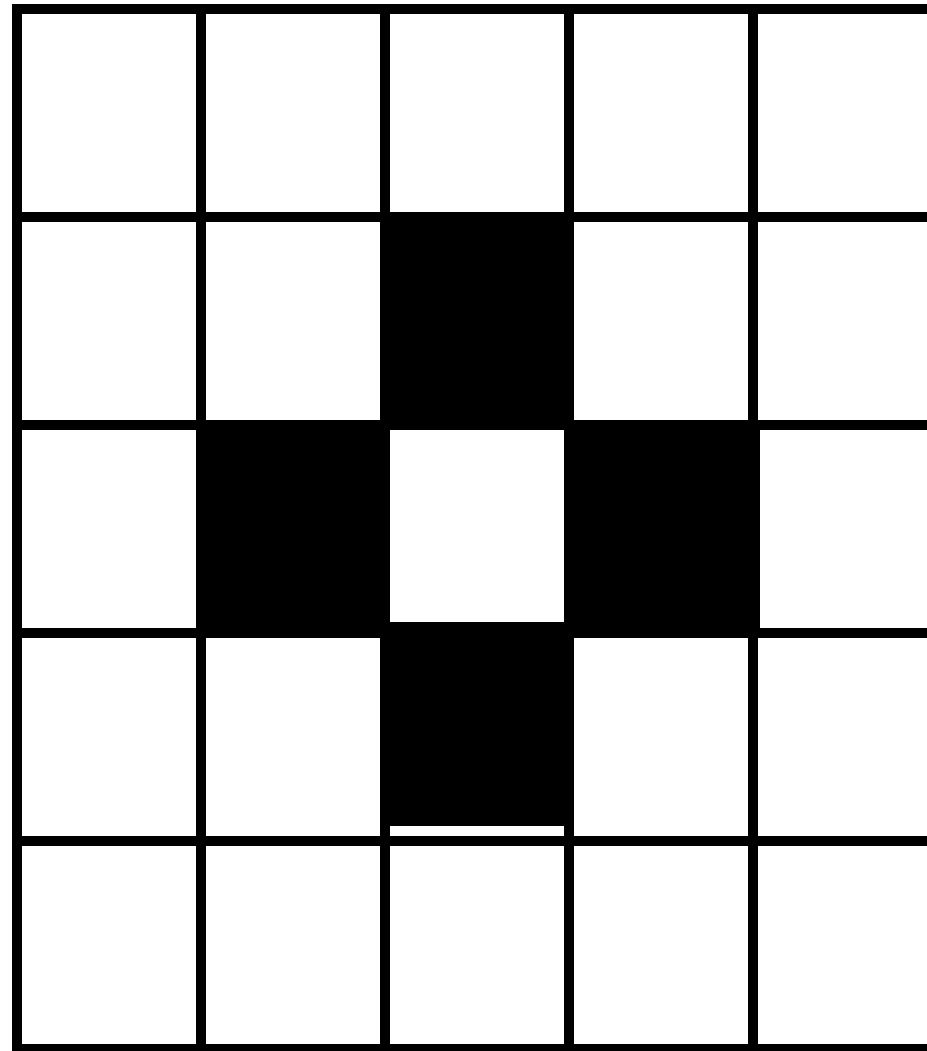
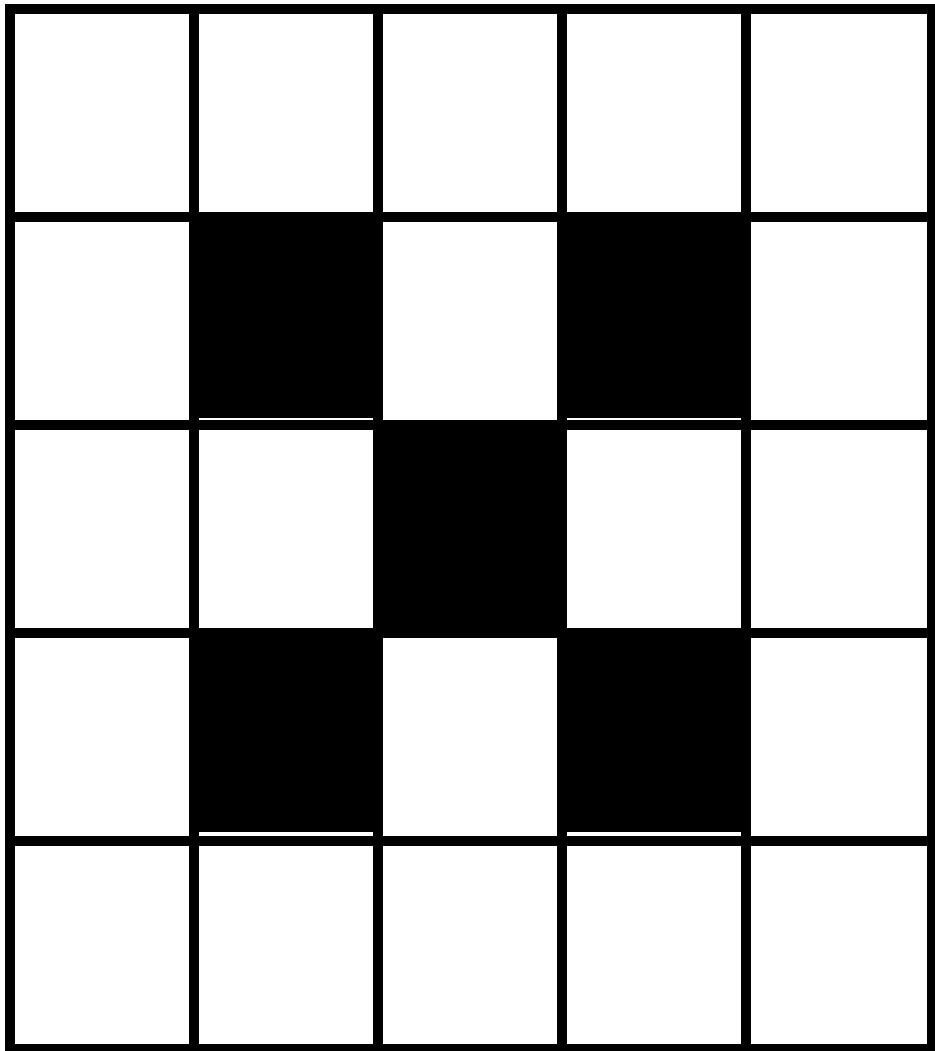
$$F(\nu) = \int_{\mathbb{R}^n} f(x) e^{-2\pi i x \cdot \nu} dx,$$

$$G(\nu) = \int_{\mathbb{R}^n} g(x) e^{-2\pi i x \cdot \nu} dx,$$

$$f * g = \mathcal{F}^{-1} \{ \mathcal{F}\{f\} \cdot \mathcal{F}\{g\} \}$$

\mathcal{F} fourier transform

two images.



-1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	-1

-1	-1	-1	-1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	-1	-1	-1	-1

-1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	-1

feature maps

1	-1	-1
-1	1	-1
-1	-1	1

-1	-1	1
-1	1	-1
1	-1	-1

-1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	-1

⊗
convolution

1	-1	-1
-1	1	-1
-1	-1	1

$$\begin{aligned}
 & (-1 * 1) + (-1 * -1) + (-1 * -1) + \\
 & (-1 * -1) + (1 * 1) + (-1 * -1) \\
 & (-1 * -1) + (-1 * -1) + (1 * 1) \\
 = 7
 \end{aligned}$$

1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

7		

$$\begin{aligned}
 & (-1 * 1) + (-1 * -1) + (-1 * -1) + \\
 & (-1 * 1) + (-1 * 1) + (-1 * 1) \\
 & (-1 * -1) + (-1 * 1) + (-1 * 1) \\
 & = -3
 \end{aligned}$$

-1	-1	-1	-1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

7	-3	

-1	-1			
-1	1			
-1	-1			
-1	1	-1	1	-1
-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

7	-3	3

-1	-1	-1	-1	-1
1	-1	-1	1	-1
-1	1	-1	-1	-1
-1	-1	1	1	-1
-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

7	-3	3
?		

-1	-1	-1	-1	-1
-1	1	-1	-1	-1
-1	-1	1	-1	-1
-1	-1	-1	1	-1
-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

7	-3	3
?	?	

-1	-1	-1	-1	-1
-1	1	-1	-1	-1
-1	-1	-1	-1	-1
-1	1	-1	-1	-1
-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

7	-3	3
?	?	

-1	-1	-1	-1	-1
-1	1	-1	1	-1
1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

7	-3	3
?	?	

-1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

7	-3	3
?	?	

-1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	1



1	-1	-1
-1	1	-1
-1	-1	1

=

7	-3	3
?	?	

input layer

-1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	-1

convolution layer

1	-1	-1
-1	1	-1
-1	-1	1



feature map

7	-3	3
-3		

=

input layer

-1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	-1

convolution layer



1	-1	-1
-1	1	-1
-1	-1	1

feature map

=

7	-3	3
-3	5	-3
3	-3	7

the feature map is "richer": we went from binary to R

input layer

-1	-1	-1	-1	-1
-1	1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	1	-1
-1	-1	-1	-1	-1

convolution layer



1	-1	-1
-1	1	-1
-1	-1	1

feature map

7	-3	3
-3	5	-3
3	-3	7

the feature map is "richer": we went from
binary to R
and it is reminiscent of the original layer

Convolve with different feature: each neuron is 1 feature

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1



1	-1	1
-1	1	-1
1	-1	1

=

0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.11	0.33	-0.77	1.00	-0.77	0.33	-0.11
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1



-1	-1	1
-1	1	-1
1	-1	-1

=

0.33	-0.11	0.55	0.33	0.11	-0.11	0.77
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.77	-0.11	0.11	0.33	0.55	-0.11	0.33

CNN

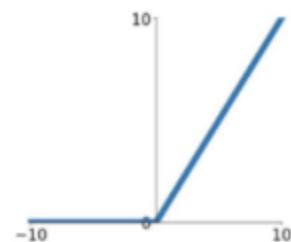
tb

ReLU

ReLU: normalization that replaces negative values with 0's

7	-3	3
-5	5	-3
-6	-1	7

ReLU
 $\max(0, x)$



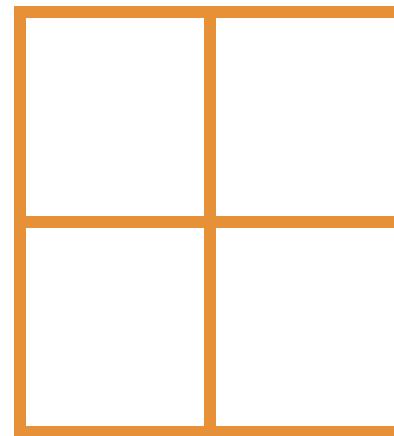
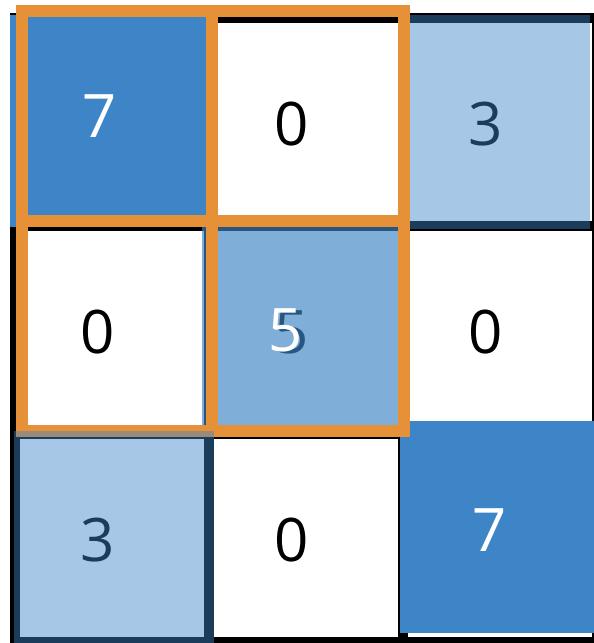
7	0	3
0	5	0
3	0	7

CNN

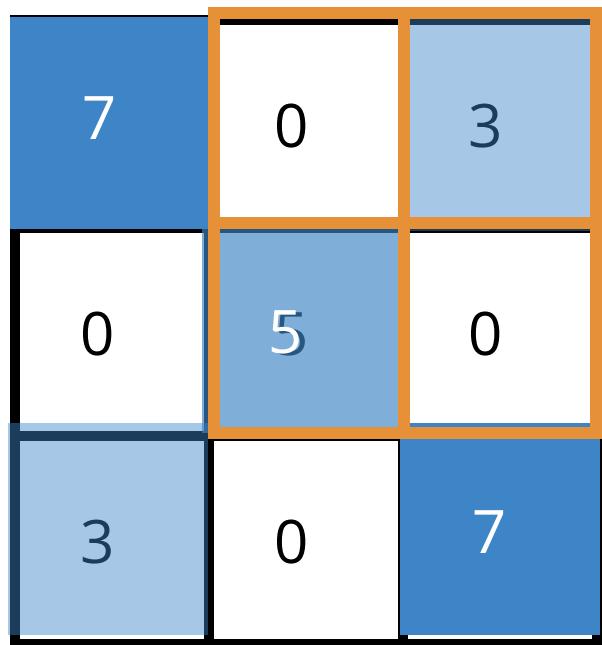
tc

Max-Pool

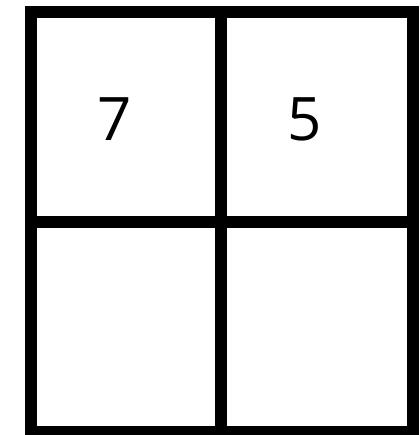
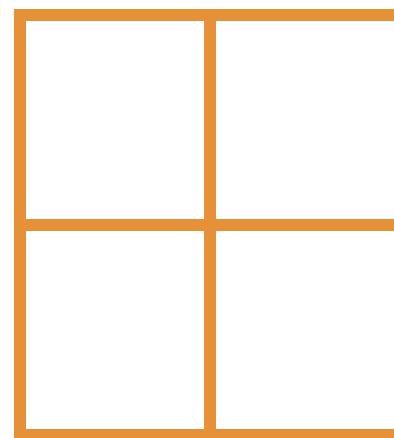
MaxPooling: reduce image size, generalizes result



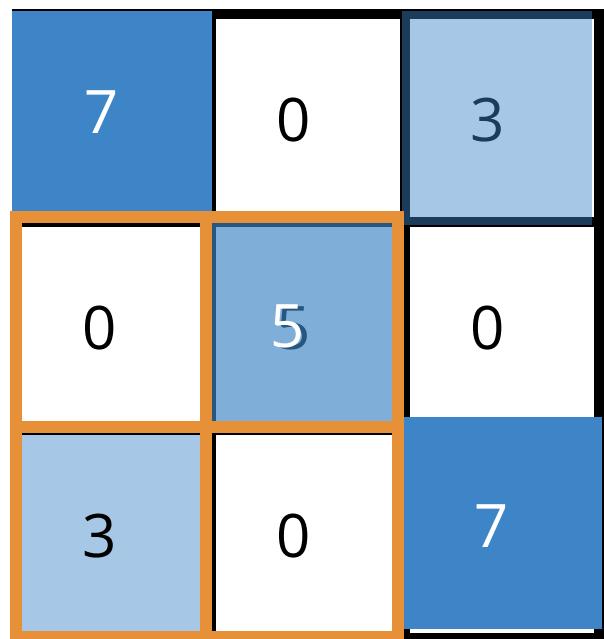
MaxPooling: reduce image size, generalizes result



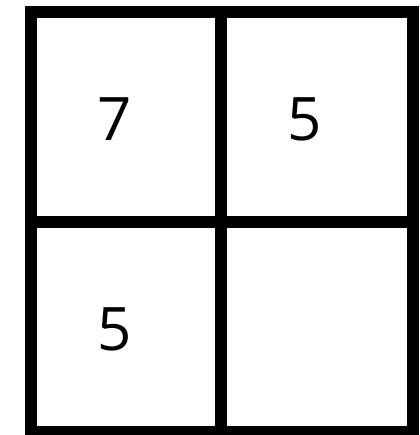
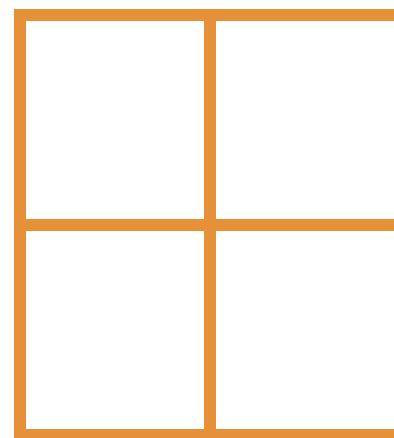
2x2 Max Poll



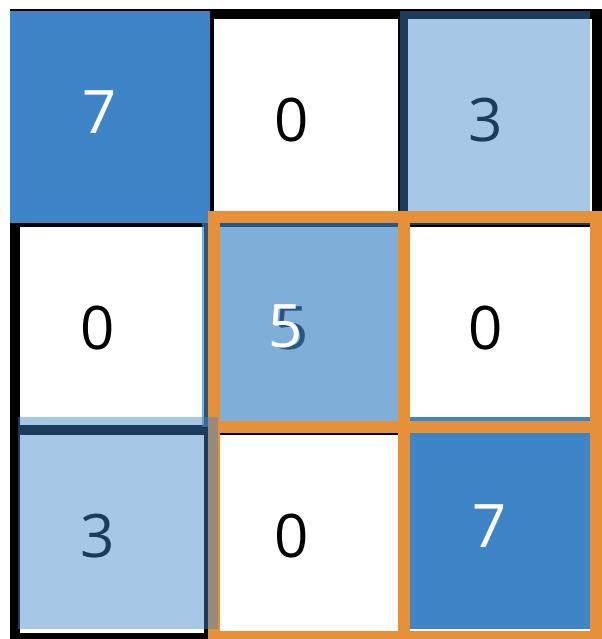
MaxPooling: reduce image size, generalizes result



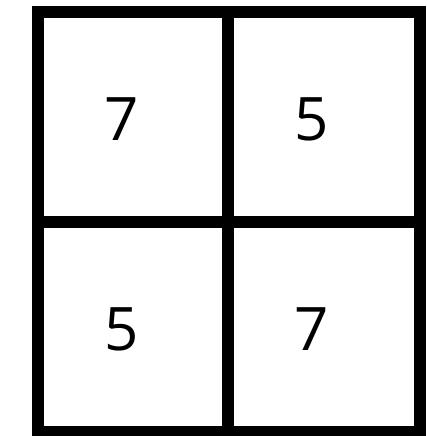
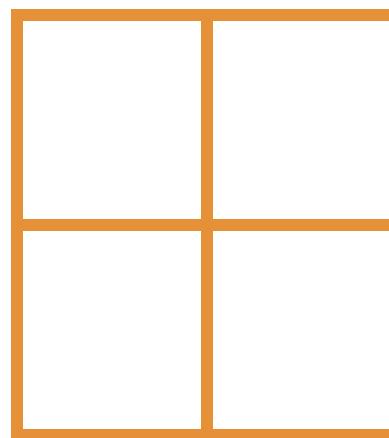
2x2 Max Poll



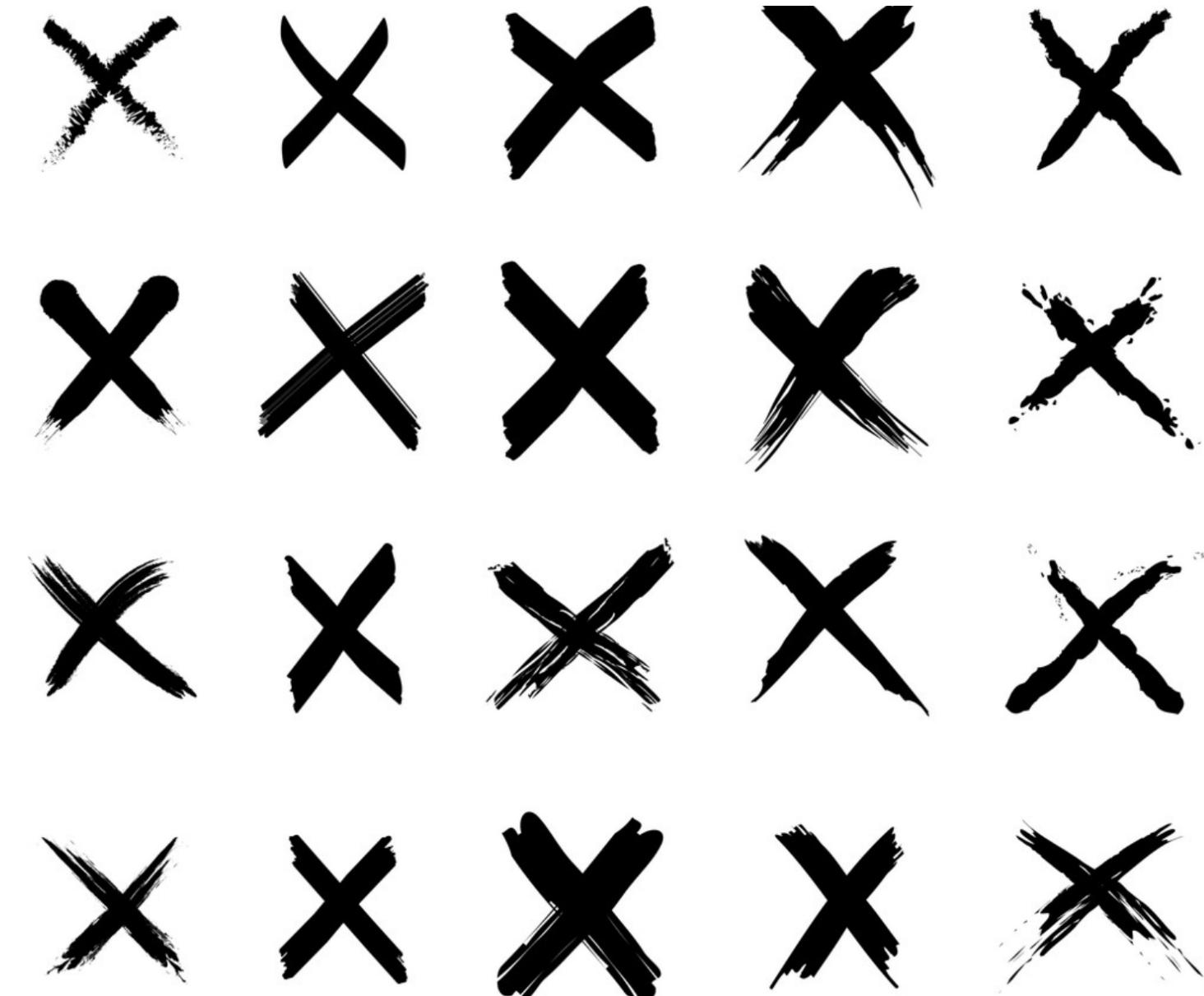
MaxPooling: reduce image size, generalizes result



2x2 Max Poll



MaxPooling:
reduce image
size & *generalizes*
result

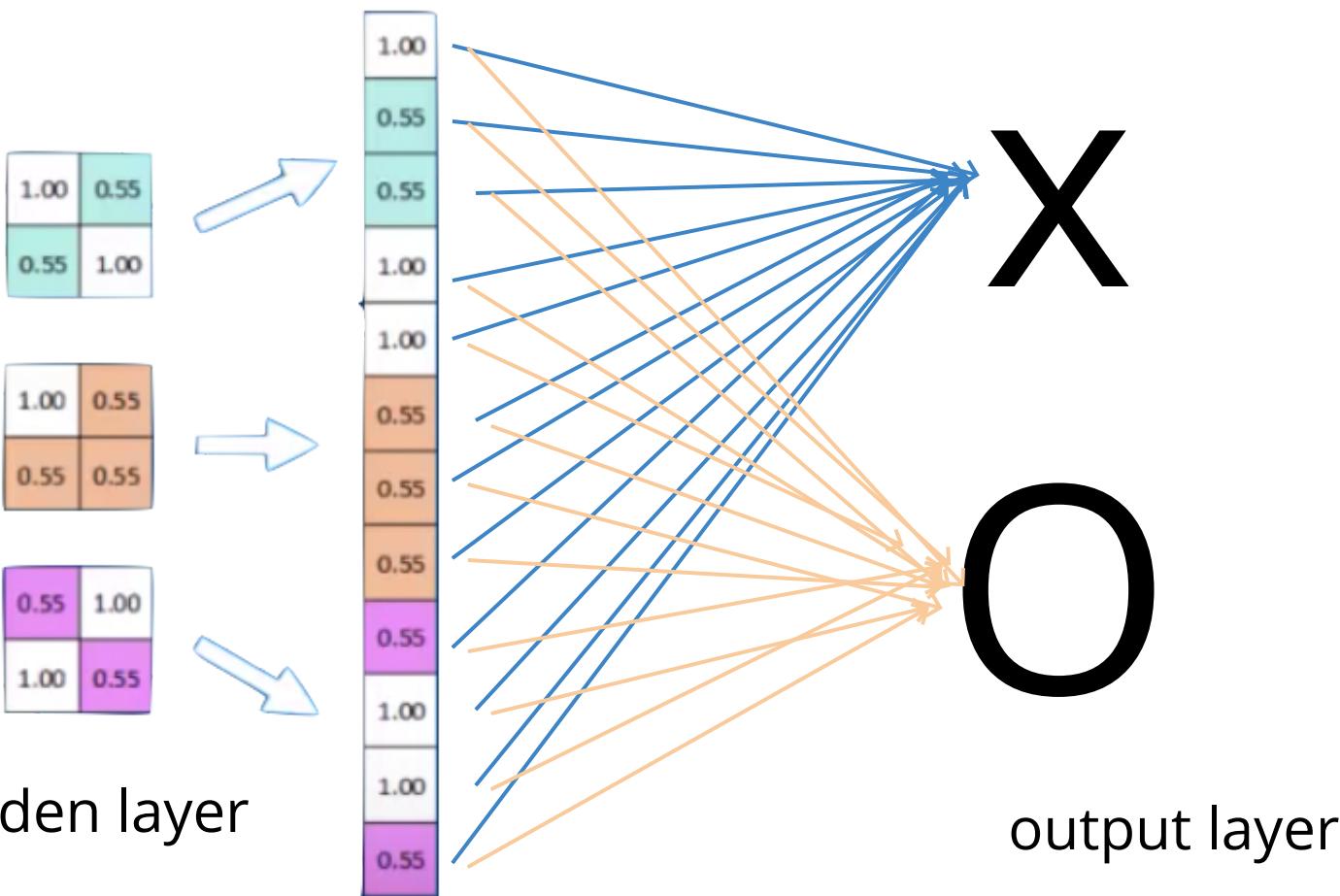


By reducing the size and picking the maximum of a sub-region we make the network less sensitive to specific details

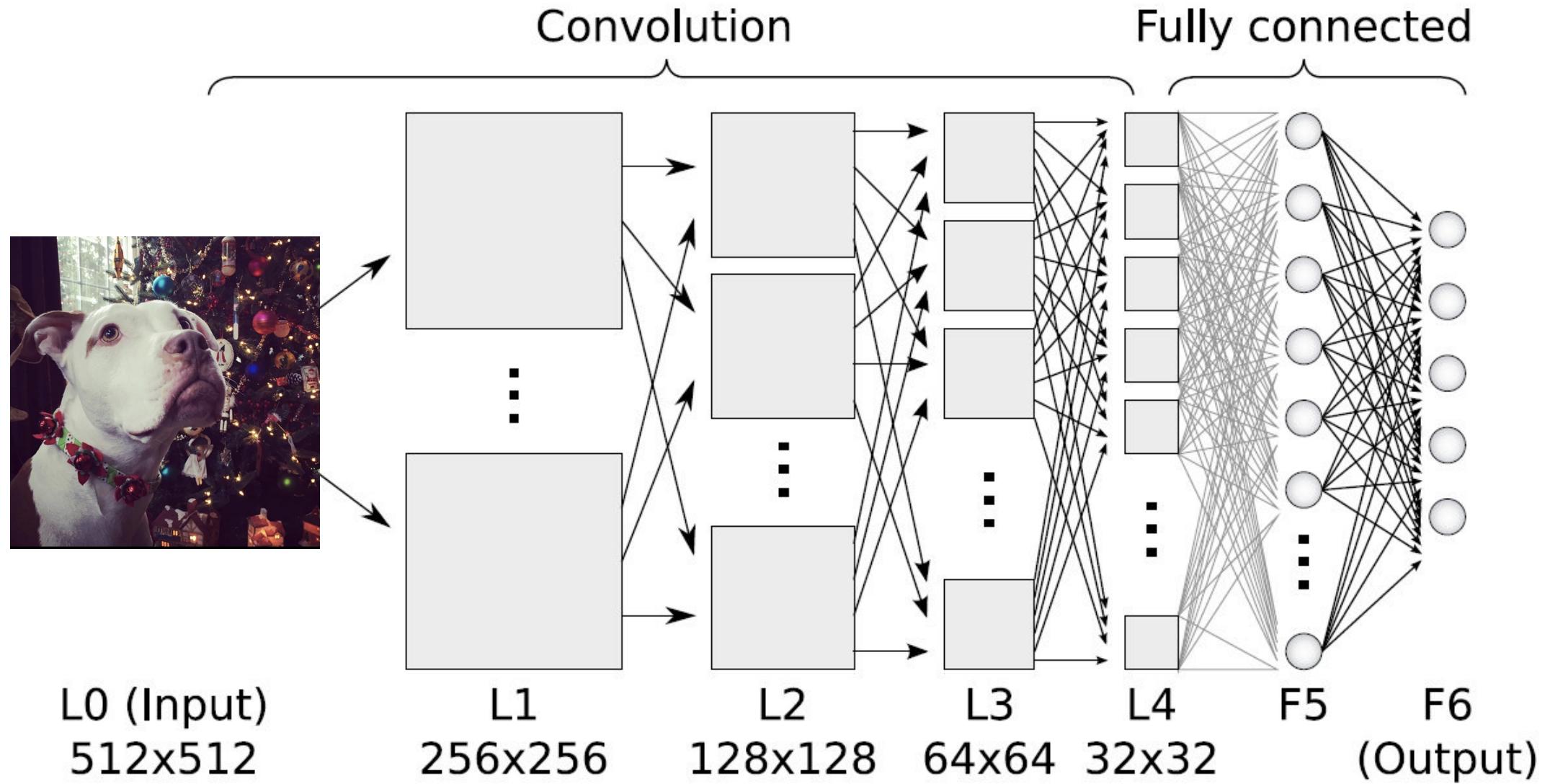
CNN



final layer: the final layer is fully connected



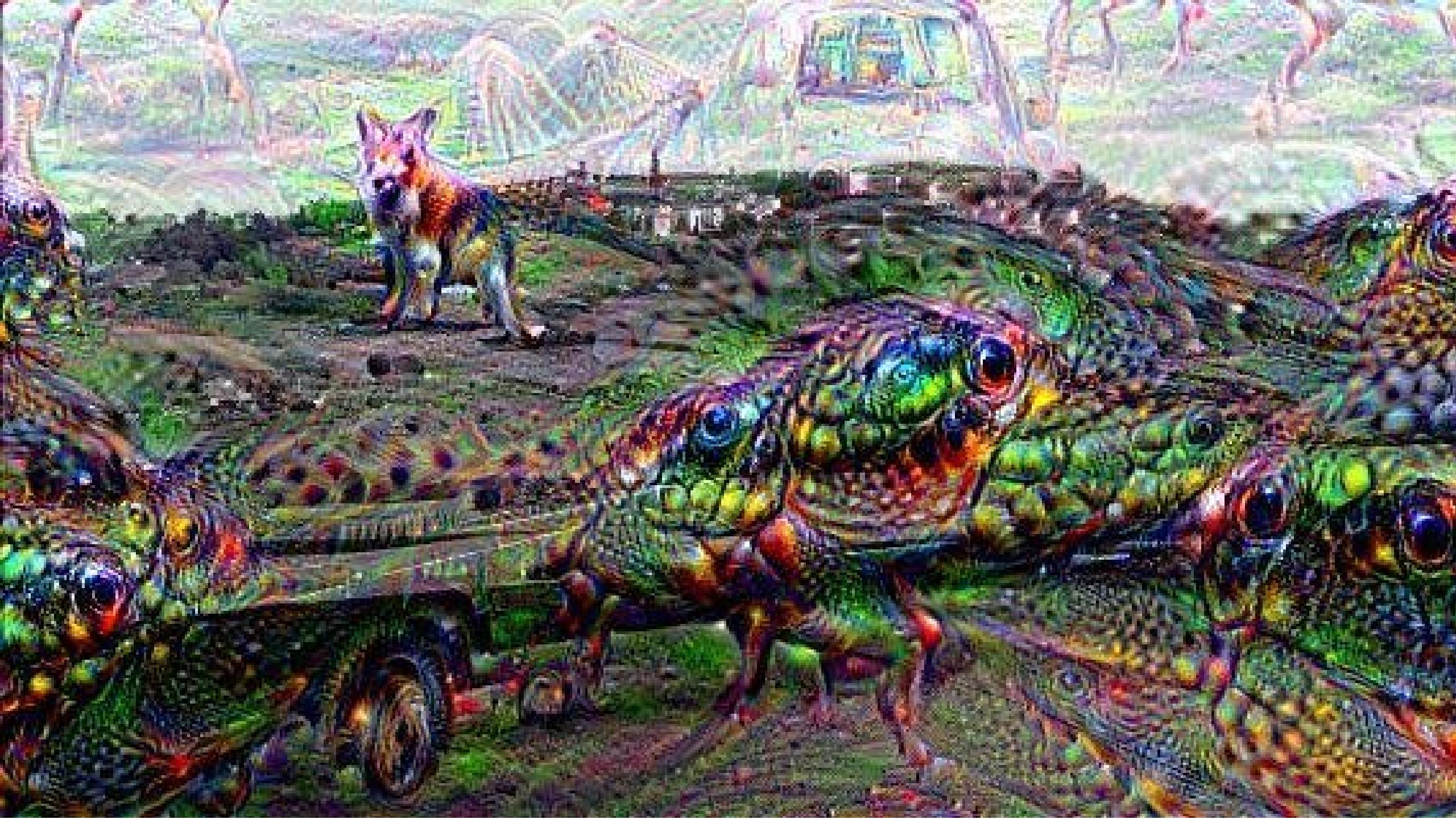
Stack multiple convolution layers



A deep dream image generated by a neural network. The central focus is a dog's head, which has been transformed into a vibrant, multi-colored pattern of blues, reds, yellows, and greens. The dog's eyes are particularly prominent, showing a mix of blue and orange hues. In the background, a small dog is standing on a grassy hillside, looking towards the viewer. The sky above is filled with a soft, hazy light.

deep dreams





what is happening in DeepDream?

Deep Dream (DD) is a google software, a pre-trained NN (originally created on the Cafe architecture, now imported on many other platforms including tensorflow).

The high level idea relies on training a convolutional NN to recognize common objects, e.g. dogs, cats, cars, in images. As the network learns to recognize those objects it develops its convolutional layers to pick out "features" of the NN, like lines at a certain orientations, circles, etc.

Each neuron, is a filters: e.g. edge finders.

The DeepDream software runs this NN on an image you give it, and it loops on some hidden layers, thus "manifesting" the things it knows how to recognize in the image. The output of an inner layer (input of the next inner layer) is called a "feature map". We are taking a peek into the feature maps of a deep neural network trained to recognize common objects.

<https://colah.github.io/posts/2015-08-Backprop/>

excellent blog post on BP: <http://colah.github.io/posts/2015-08-Backprop/>

Training a DNN

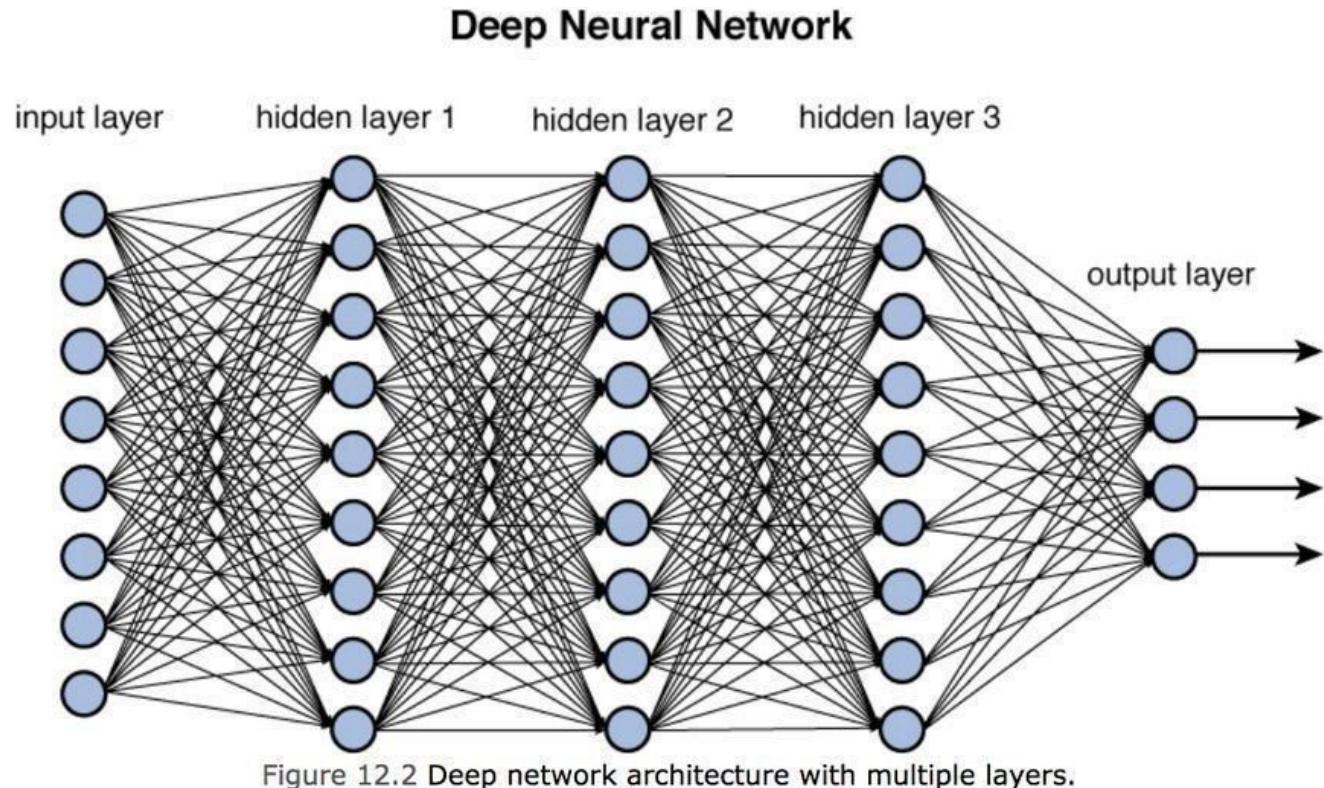
Training models with this many parameters requires a lot of care:

- . defining the metric
- . optimization schemes
- . training/validation/testing sets

But just like our simple linear regression case, the fact that small changes in the parameters leads to small changes in the output for the right activation functions.

define a cost function, e.g.

$$C = \frac{1}{2} |y - a^L|^2 = \frac{1}{2} \sum_j (y_j - a_j^L)^2$$



feed data forward through network and calculate cost metric

for each layer, calculate effect of small changes on next layer

$$\vec{y} = f_N(\dots(f_1(\vec{x}W_i + b_1\dots W_N + b_N)))$$

overfitting

2

Minibatch

&

Dropout

overfitting

Split your training set into many smaller subsets and train on each small set separately

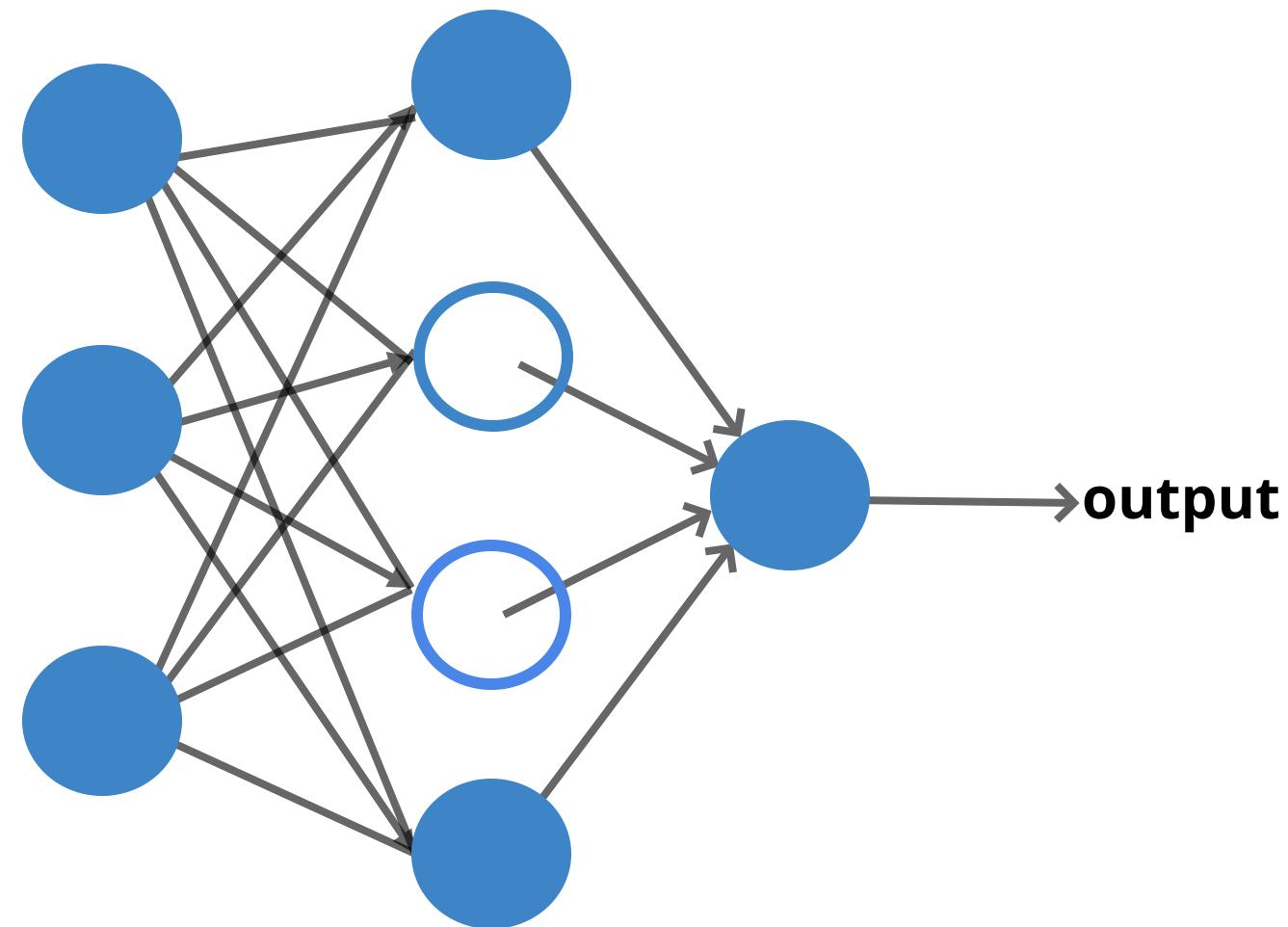
- If one updates model parameters after processing the whole training data (i.e., epoch), it would take too long to get a model update in training, and the entire training data probably won't fit in the memory.
- If one updates model parameters after processing every instance (i.e., stochastic gradient descent), model updates would be too noisy, and the process is not computationally efficient.
- Therefore, minibatch gradient descent is introduced as a trade-off between fast model updates (memory efficiency) and accurate model updates (computational efficiency).

<https://www.quora.com/What-is-a-minibatch-in-a-neural-network>

overfitting

Artificially remove some neurons for different minibatches to avoid overfitting

Dropout



Architecture components: neurons, activation function

- basically each neuron is a multivariate regression with an activation function that turns the output into a probability
- changing the weights and biases in the linear regression gives different results

Single layer NN: perceptrons

- perceptrons were developed in the 50s but a long time passed since then till people figured out how to build complex layered architectures and especially how to train them

Deep NN:

- DNN are multi-layer architectures of neurons. They can be fully connected (each neuron goes to each neuron of the next layer) or not (a neuron goes only to some neurons in the next layer)
- DNN have a lot of parameters (thousands!) which makes the interpretability and feature extraction of NN difficult.

Convolutional NN

- convolutional NN are DNN with three types of layers:
 - convolutional layers: run filters through an image to detect features like edges or colors
 - maxpool layers: decrease the size of the previous layer outputs and removes some details
 - ReLU : rectified linear units: normalizes the output of conv layers so that it is all positive (sets negatives to 0)
- CNNs are great for the stud of structure in large datasets (images are large datasets)

Training an NN:

- most MI methods are trained by gradient descent: change weights and biases based

Lots of parameters and lots of hyperparameters! What to choose? cheatsheet

1. **architecture - wide networks tend to overfit, deep networks are hard to train**
2. number of epochs - the sweet spot is when learning slows down, but before you start overfitting... it may take DAYS! jumps may indicate bad initial choices
3. loss function - needs to be appropriate to the task, e.g. classification vs regression
4. activation functions - needs to be consistent with the loss function
5. optimization scheme - needs to be appropriate to the task and data
6. learning rate in optimization - balance speed and accuracy
7. batch size - smaller batch size is faster but leads to overtraining

An article that compares various DNNs

AN ANALYSIS OF DEEP NEURAL NETWORK MODELS FOR PRACTICAL APPLICATIONS

Alfredo Canziani & Eugenio Culurciello

Weldon School of Biomedical Engineering
Purdue University
`{canziani,euge}@purdue.edu`

Adam Paszke

Faculty of Mathematics, Informatics and Mechanics
University of Warsaw
`a.paszke@students.mimuw.edu.pl`

ABSTRACT

Since the emergence of Deep Neural Networks (DNNs) as a prominent technique in the field of computer vision, the ImageNet classification challenge has played a major role in advancing the state-of-the-art. While accuracy figures have steadily increased, the resource utilisation of winning models has not been properly taken into account. In this work, we present a comprehensive analysis of important metrics in practical applications: accuracy, memory footprint, parameters, operations count, inference time and power consumption. Key findings are: (1) power consumption is independent of batch size and architecture; (2) accuracy and inference time are in a hyperbolic relationship; (3) energy constraint is an upper bound on the maximum achievable accuracy and model complexity; (4) the number of operations is a reliable estimate of the inference time. We believe our analysis provides a compelling set of information that helps design and engineer efficient DNNs.

<https://arxiv.org/pdf/1605.07678.pdf>

An article that compares various DNNs

AN ANALYSIS OF DEEP NEURAL NETWORK MODELS FOR PRACTICAL APPLICATIONS

Alfredo Canziani & Eugenio Culurciello
Weldon School of Biomedical Engineering
Purdue University
`{canziani,euge}@purdue.edu`

Adam Paszke
Faculty of Mathematics, Informatics and Mechanics
University of Warsaw
`a.paszke@students.mimuw.edu.pl`

ABSTRACT

Since the emergence of Deep Neural Networks (DNNs) as a prominent technique in the field of computer vision, the ImageNet classification challenge has played a major role in advancing the state-of-the-art. While accuracy figures have steadily increased, the resource utilisation of winning models has not been properly taken into account. In this work, we present a comprehensive analysis of important metrics in practical applications: accuracy, memory footprint, parameters, operations count, inference time and power consumption. Key findings are: (1) power consumption is independent of batch size and architecture; (2) accuracy and inference time are in a hyperbolic relationship; (3) energy constraint is an upper bound on the maximum achievable accuracy and model complexity; (4) the number of operations is a reliable estimate of the inference time. We believe our analysis provides a compelling set of information that helps design and engineer efficient DNNs.

<https://arxiv.org/pdf/1605.07678.pdf>

accuracy comparison

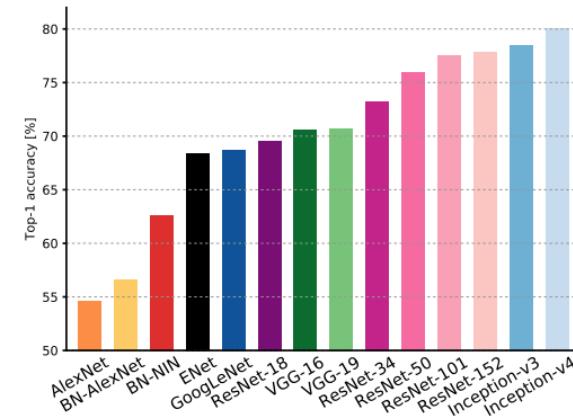


Figure 1: **Top1 vs. network.** Single-crop top-1 validation accuracies for top scoring single-model architectures. We introduce with this chart our choice of

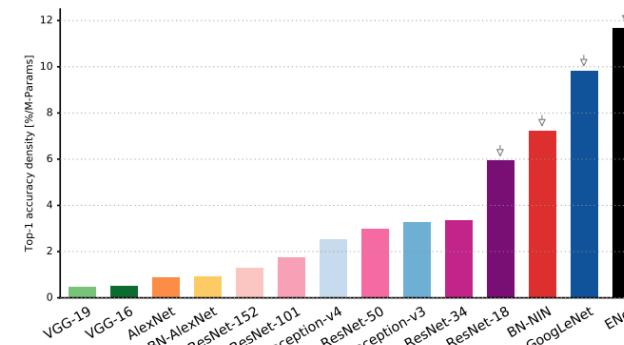


Figure 10: **Accuracy per parameter vs. network.** Information density (accuracy per parameters) is an efficiency metric that highlight that capacity of a specific architecture to better utilise its parametric space. Models like VGG and AlexNet are clearly oversized, and do not take fully advantage of their potential learning ability. On the far right, ResNet-18, BN-NIN, GoogLeNet and ENet (marked by grey arrows) do a better job at “squeezing” all their neurons to learn the given task, and are the winners of this section.

An article that compares various DNNs

AN ANALYSIS OF DEEP NEURAL NETWORK MODELS FOR PRACTICAL APPLICATIONS

Alfredo Canziani & Eugenio Culurciello
Weldon School of Biomedical Engineering
Purdue University
{canziani,euge}@purdue.edu

Adam Paszke
Faculty of Mathematics, Informatics and Mechanics
University of Warsaw
a.paszke@students.mimuw.edu.pl

ABSTRACT

Since the emergence of Deep Neural Networks (DNNs) as a prominent technique in the field of computer vision, the ImageNet classification challenge has played a major role in advancing the state-of-the-art. While accuracy figures have steadily increased, the resource utilisation of winning models has not been properly taken into account. In this work, we present a comprehensive analysis of important metrics in practical applications: accuracy, memory footprint, parameters, operations count, inference time and power consumption. Key findings are: (1) power consumption is independent of batch size and architecture; (2) accuracy and inference time are in a hyperbolic relationship; (3) energy constraint is an upper bound on the maximum achievable accuracy and model complexity; (4) the number of operations is a reliable estimate of the inference time. We believe our analysis provides a compelling set of information that helps design and engineer efficient DNNs.

<https://arxiv.org/pdf/1605.07678.pdf>

accuracy comparison

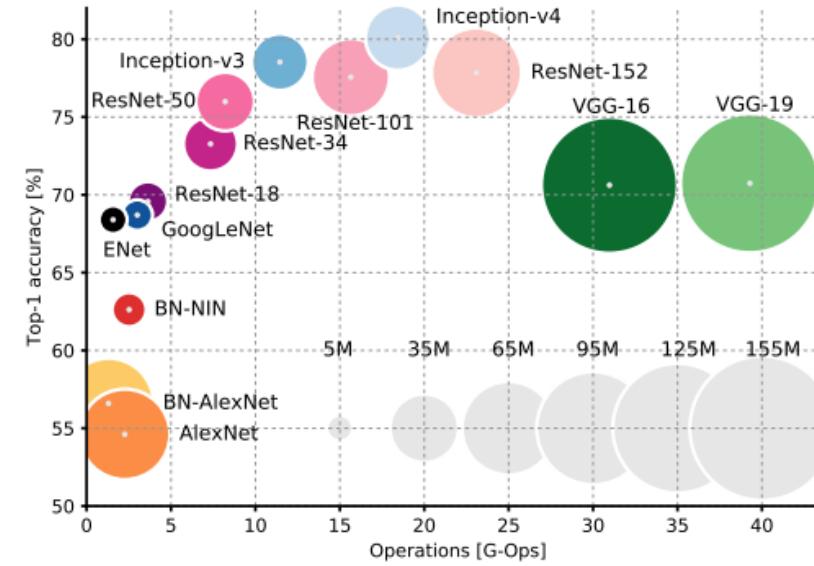


Figure 2: Top1 vs. operations, size \propto parameters.
Top-1 one-crop accuracy versus amount of operations required for a single forward pass. The size of the blobs is proportional to the number of network parameters; a legend is reported in the bottom right corner, spanning from 5×10^6 to 155×10^6 params. Both these figures share the same y-axis, and the grey dots highlight the centre of the blobs.

An article that compares various DNNs

AN ANALYSIS OF DEEP NEURAL NETWORK MODELS FOR PRACTICAL APPLICATIONS

Alfredo Canziani & Eugenio Culurciello
Weldon School of Biomedical Engineering
Purdue University
{canziani,euge}@purdue.edu

Adam Paszke
Faculty of Mathematics, Informatics and Mechanics
University of Warsaw
a.paszke@students.mimuw.edu.pl

ABSTRACT

Since the emergence of Deep Neural Networks (DNNs) as a prominent technique in the field of computer vision, the ImageNet classification challenge has played a major role in advancing the state-of-the-art. While accuracy figures have steadily increased, the resource utilisation of winning models has not been properly taken into account. In this work, we present a comprehensive analysis of important metrics in practical applications: accuracy, memory footprint, parameters, operations count, inference time and power consumption. Key findings are: (1) power consumption is independent of batch size and architecture; (2) accuracy and inference time are in a hyperbolic relationship; (3) energy constraint is an upper bound on the maximum achievable accuracy and model complexity; (4) the number of operations is a reliable estimate of the inference time. We believe our analysis provides a compelling set of information that helps design and engineer efficient DNNs.

<https://arxiv.org/pdf/1605.07678.pdf>

batch size

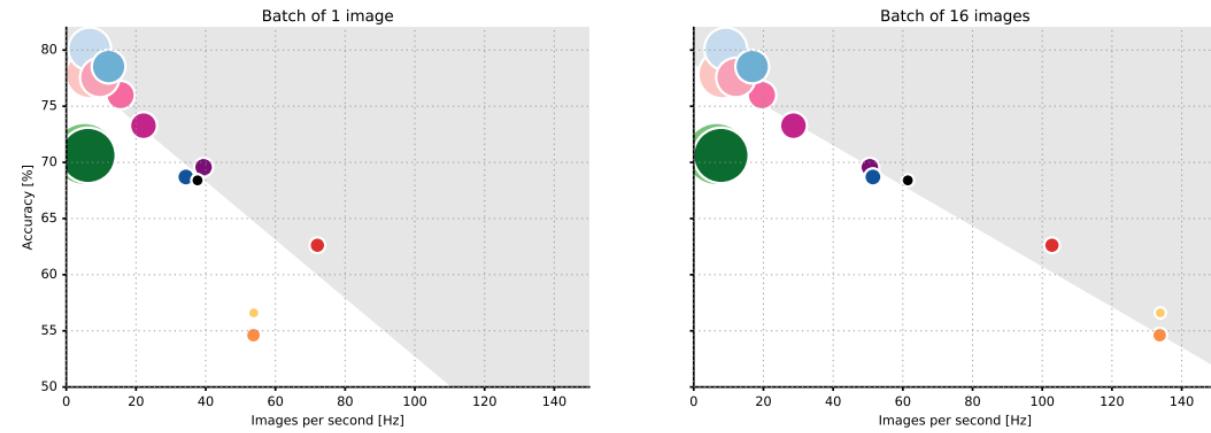


Figure 9: Accuracy vs. inferences per second, size \propto operations. Non trivial linear upper bound is shown in these scatter plots, illustrating the relationship between prediction accuracy and throughput of all examined architectures. These are the first charts in which the area of the blobs is proportional to the amount of operations, instead of the parameters count. We can notice that larger blobs are concentrated on the left side of the charts, in correspondence of low throughput, *i.e.* longer inference times. Most of the architectures lay on the linear interface between the grey and white areas. If a network falls in the shaded area, it means it achieves exceptional accuracy or inference speed. The white area indicates a suboptimal region. *E.g.* both AlexNet architectures improve processing speed as larger batches are adopted, gaining 80 Hz.

Lots of parameters and lots of hyperparameters! What to choose? cheatsheet

1. architecture - wide networks tend to overfit, deep networks are hard to train
2. **number of epochs** - the sweet spot is when learning slows down, but before you start overfitting... it may take DAYS! jumps may indicate bad initial choices
3. **loss function** - needs to be appropriate to the task, e.g. classification vs regression
4. activation functions - needs to be consistent with the loss function
5. optimization scheme - needs to be appropriate to the task and data
6. learning rate in optimization - balance speed and accuracy
7. batch size - smaller batch size is faster but leads to overtraining

Lots of parameters and lots of hyperparameters! What to choose?



https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

3.0.1 regression

- loss='mean_squared_error' L2: default loss to use for regression problems. => linear activation funct in output layer, one node out

alternatives: loss='mean_squared_logarithmic_error', 'mean_absolute_error' (which is L1 instead of L2)

3.0.2 binary classification

- loss='binary_crossentropy' => sigmoid activation function in output layer, one node out

alternatives: 'hinge'

3.0.3 multiclass classification

categorical encoded as numerical

- loss='categorical_crossentropy' => softmax n nodes out

onehot encoded categorical

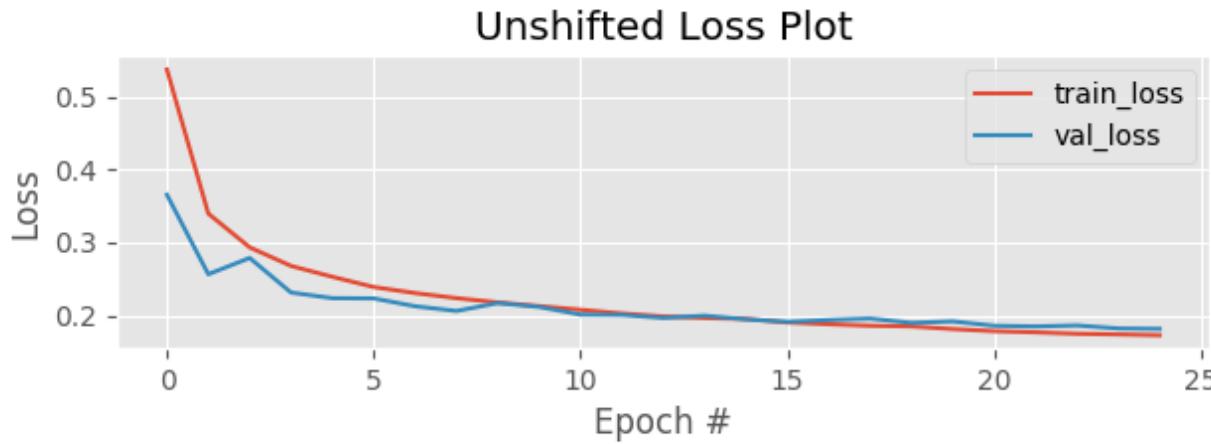
- 'parse_categorical_crossentropy' => softmax n nodes out
- 'kullback Leibler Divergence Loss' => probabilistic categorical classification; $\log(P/Q)$

What should I choose for the loss function and how does that relate to the activation function and optimization?

*Lots of parameters and lots of
hyperparameters! What to choose?*



cheatsheet

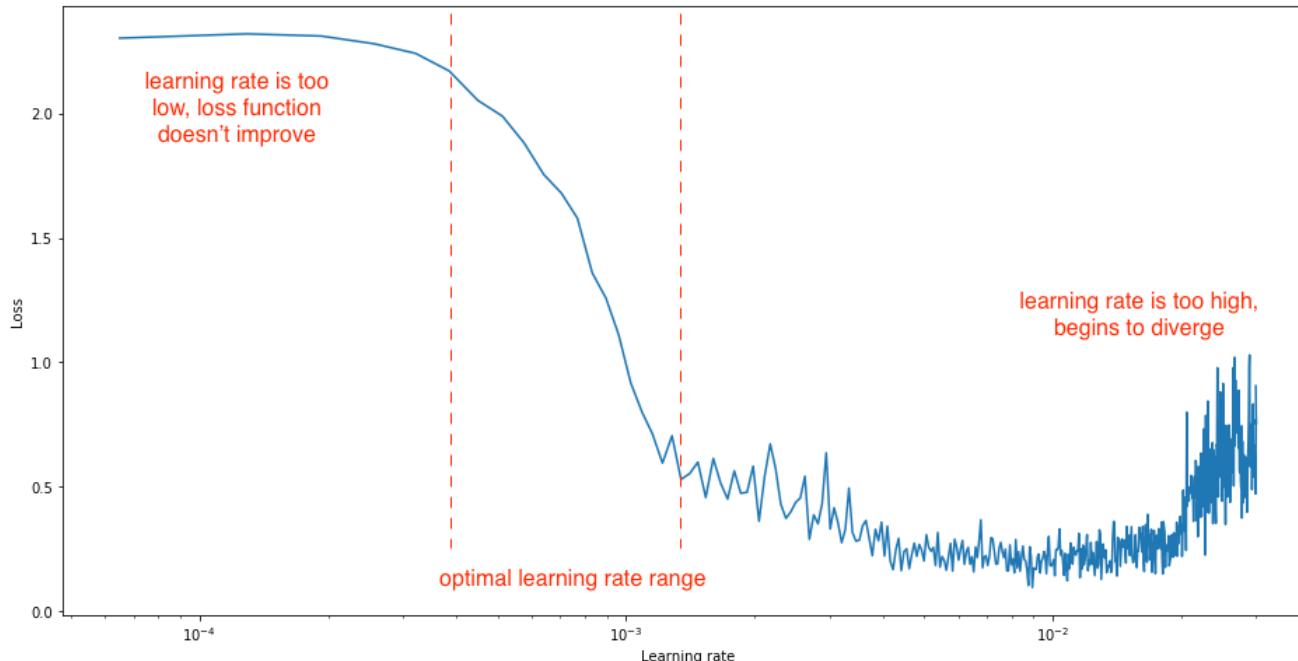


always check your loss function! it should go down smoothly and flatten out at the end of the training.

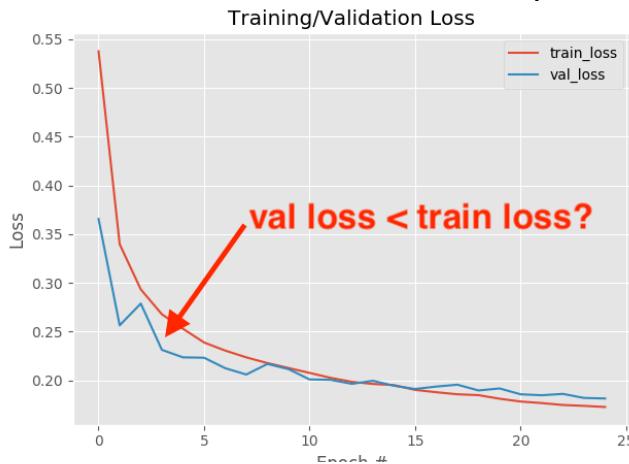
not flat? you are still learning!
too flat? you are overfitting...

loss (gallery of horrors)

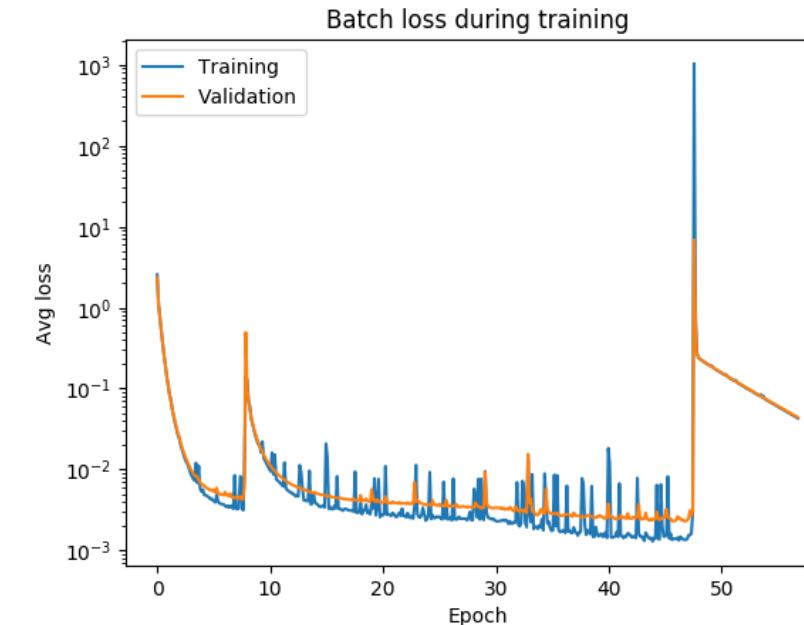
https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb



loss and learning rate (not that the appropriate learning rate depends on the chosen optimization scheme!)



when you use validation you are introducing regularizations (e.g. dropout) so the loss can be smaller than for the training set



jumps are not unlikely (and not necessarily a problem) if your activations are discontinuous (e.g. relu)

Building a DNN with keras and tensorflow autoencoder for image reconstruction



https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

What should I choose for the loss function and how does that relate to the activation function and optimization?

loss	good for	activation last layer	size last layer
mean_squared_error	regression	linear	one node
mean_absolute_error	regression	linear	one node
mean_squared_logarithmit_error	regression	linear	one node
binary_crossentropy	binary classification	sigmoid	one node
categorical_crossentropy	multiclass classification	sigmoid	N nodes
Kullback_Divergence	multiclass classification, probabilistic interpretation	sigmoid	N nodes

On the interpretability of DNNs

Windows (4b:237)
excite the car detector
at the top and inhibit
at the bottom.

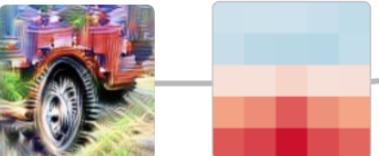


● positive (excitation)
● negative (inhibition)

Car Body (4b:491)
excites the car
detector, especially at
the bottom.



Wheels (4b:373) excite
the car detector at the
bottom and inhibit at
the top.



A car detector (4c:447)
is assembled from
earlier units.

<https://distill.pub/2020/circuits/zoom-in/>

... [[Jerusalem Report]] ... [http://www.jrep.com] Left-of-center Eng
* "[hToausal maogurt]" "(http://www.bsinioom/ -iat af tenter (ng
[' [Cassmene] Beaonds s a [ad : xne. waaaoca. s &ato- nfhlsum-ouc
' s mFurnls iaeltsa' : . i ' cdw- 2tpiisoeg. er / . a] (oseswr- ciddrs [mt
: : AqDenebiutn | Cipreel . b1emr. 9:ahb- npumughnmp) Teiretu: eoseodsald
T&Tf Siwrpe] aluveleurus : -mprts < moa2deyshilrjc. Augl. 1p. larc : fae

glish [[weekly newspaper]] " [[YNet News]] " [http://www.ynetnews. c
lish c [Caakly cawspaper]] " [[hTAA at]] " (http://www.bacahets. co
iaci- lhSoipli sec lenpls . ' ' [Co * wessl s a [ad : xne. waea. awatoa
eena. pCci et nedlo x] gicill s' [sAmFeSahon] t' : . imomw- 2 #piisoessis. /er
syz . spenn alruellrra . '# : oDuFreieu p . : b1edr. < : ahb- nptwt. xigh
adpeamArbdeorpitee] dts - | T [[BaAvTp oSwao, . . oacstp, tcoa2drulwoclens

om/] English-language website of Israel's largest newspaper ' [[Yed
m/ -xglis hlinguagesairsite of tsraelis singlaawsaperso' [[Tel
. s &ntiaca- sardeelh oantbianfanreif ' aatdir scoe ena. i ThAoai
. n. c] (deen epesaaiki ieledh,irthraonse. coseus. setlgors. satCare
/ ma) Tvdryzil couedsu: tha- oo tu, stu, tveper y- tuaevrtid, tBAmSusy
r] p. llvaod, eytc- n dm- oibuvb] bb imsulta lybna, d, iiuiticp.] (IsvHytu

loth Ahronoth] " . Hebrew-language periodicals: ' ' " [[Globes] ,
t i (feanemti) : . ' [errewsleenguage: arosodical : . ' ' [Taaba] : red
nnh Srmuw] ey s [' ineia'si wddh' s olrif: stl ' [hAeovelt s
eg' aC lrisz] i e' : . # : TAAaat Baseeilo' ianfvl tt' ' & &mCoerone':
ut] Asaoigs] . . . : sMBolous: Toua- n: d woapnu a'n: , C: & # : afDrusu] ,
suiedNoegan o . . : { CCui bohe Cybksis: r- epcntsnk i < : & 11s T Guitrsi .

The highlighted neuron here gets very excited when the RNN is inside the [[]] markdown environment and turns off outside of it. Interestingly, the neuron can't turn on right after it

sees the character "[", it must wait for the second "[" and then activate. This task of counting whether the model has seen one or two "[" is likely done with a different neuron.

... [[Jerusalem Report]] ... [http://www.jrep.com] Left-of-center Eng
* "[hToausal maogurt]" "(http://www.bsinioom/ -iat af tenter (ng
[' [Cassmene] Beaonds s a [ad : xne. waaaoca. s &ato- nfhlsum-ouc
' s mFurnls iaeltsa' : . i ' cdw- 2tpiisoeg. er / . a] (oseswr- ciddrs [mt
: : AqDenebiutn | Cipreel . b1emr. 9:ahb- npumughnmp) Teiretu: eoseodsald
T&Tf Siwrpe] aluveleurus : -mprts < moa2deyshilrjc. Augl. 1p. larc : fae

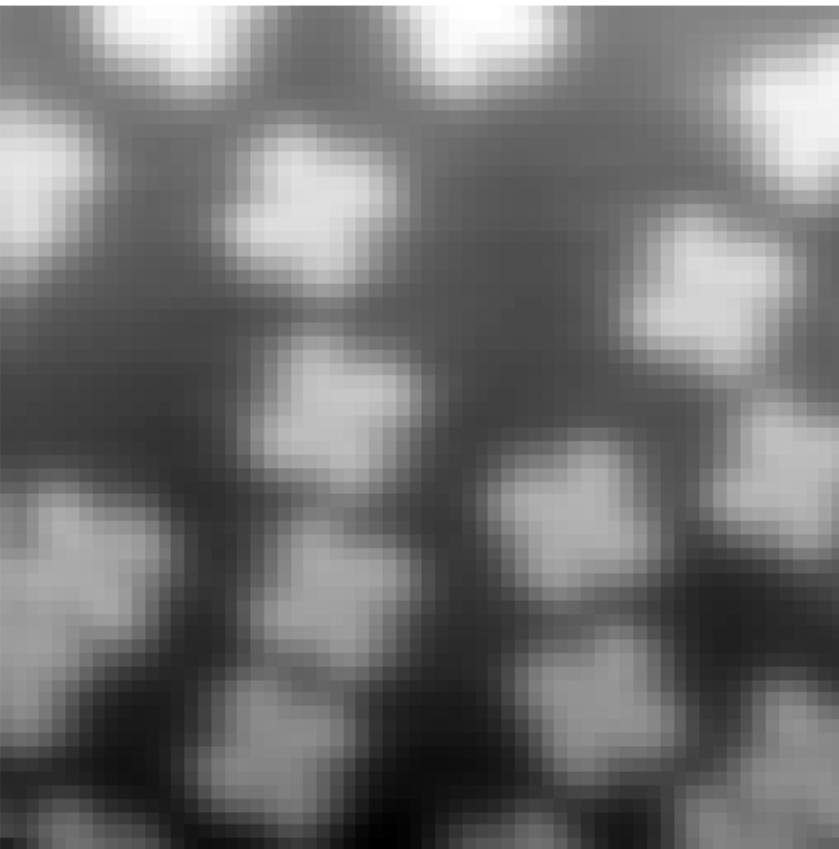
glish [[weekly newspaper]] " [[YNet News]] " [http://www.ynetnews. c
lish c [Caakly cawspaper]] " [[hTAA at]] " (http://www.bacahets. co
iaci- lhSoipli sec lenpls . ' ' [Co * wessl s a [ad : xne. waea. awatoa
eena. pCci et nedlo x] gicill s' [sAmFeSahon] t' : . imomw- 2 #piisoessis. /er
syz . spenn alruellrra . '# : oDuFreieu p . : b1edr. < : ahb- nptwt. xigh
adpeamArbdeorpitee] dts - | T [[BaAvTp oSwao, . . oacstp, tcoa2drulwoclens

Here we see a neuron that varies seemingly linearly across the [[]] environment. In other words its activation is giving the RNN a time-aligned coordinate system across the [[]] scope. The RNN can use this information to make different characters more or less likely depending on how early/late it is in the [[]] scope (perhaps?).

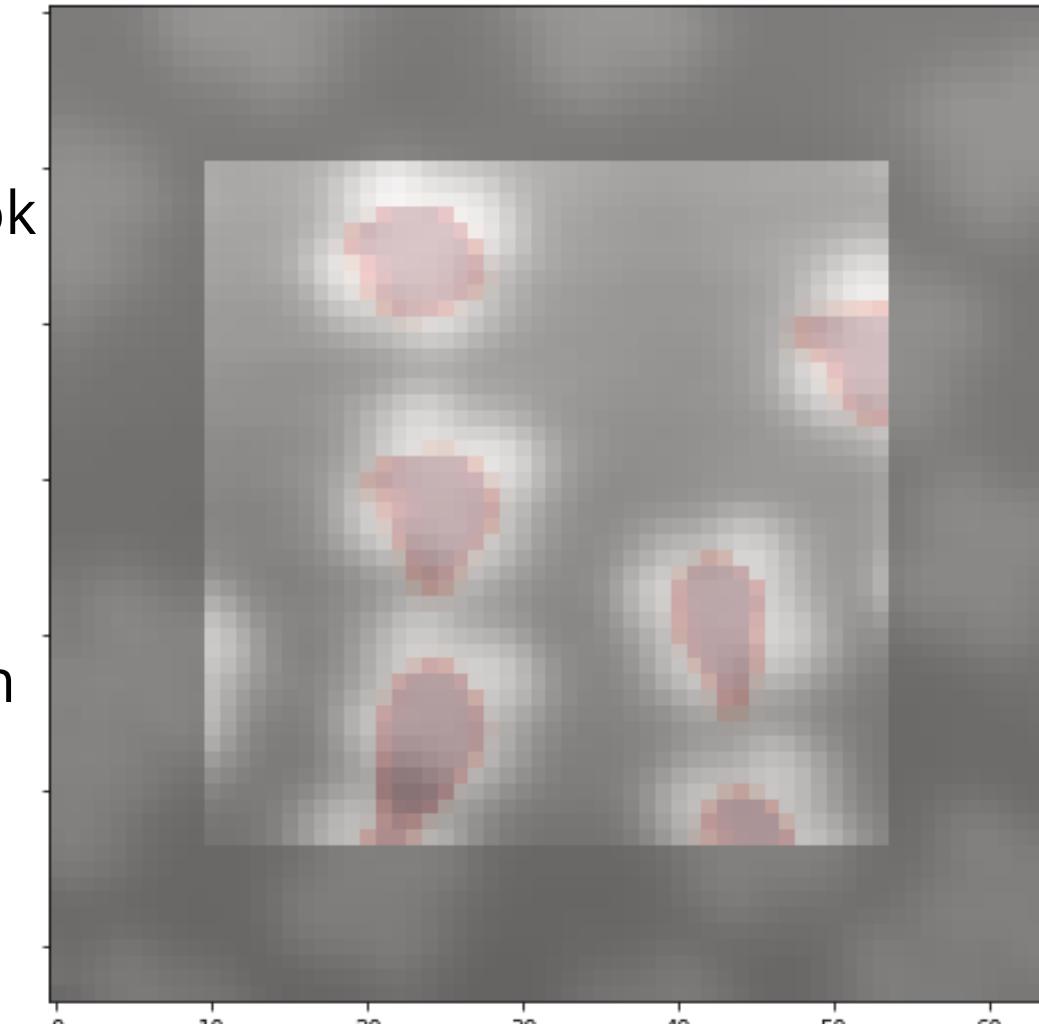
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

YOLO and R-CNN

Object detection

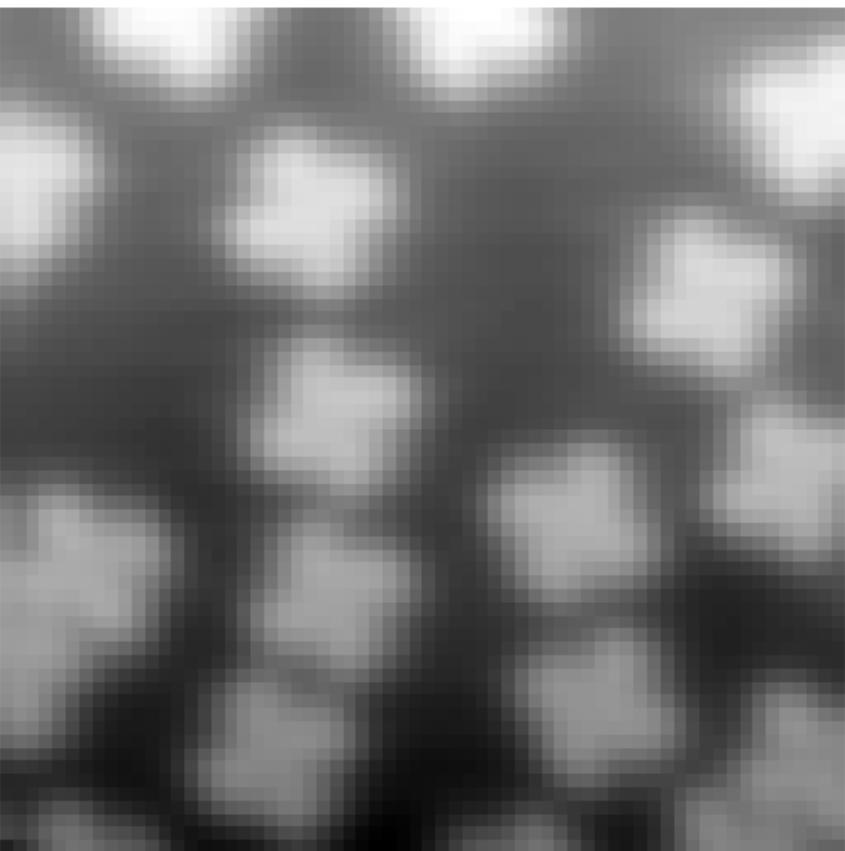


Naive model: we took different region of the image and measured the probability of presence of the object in that region

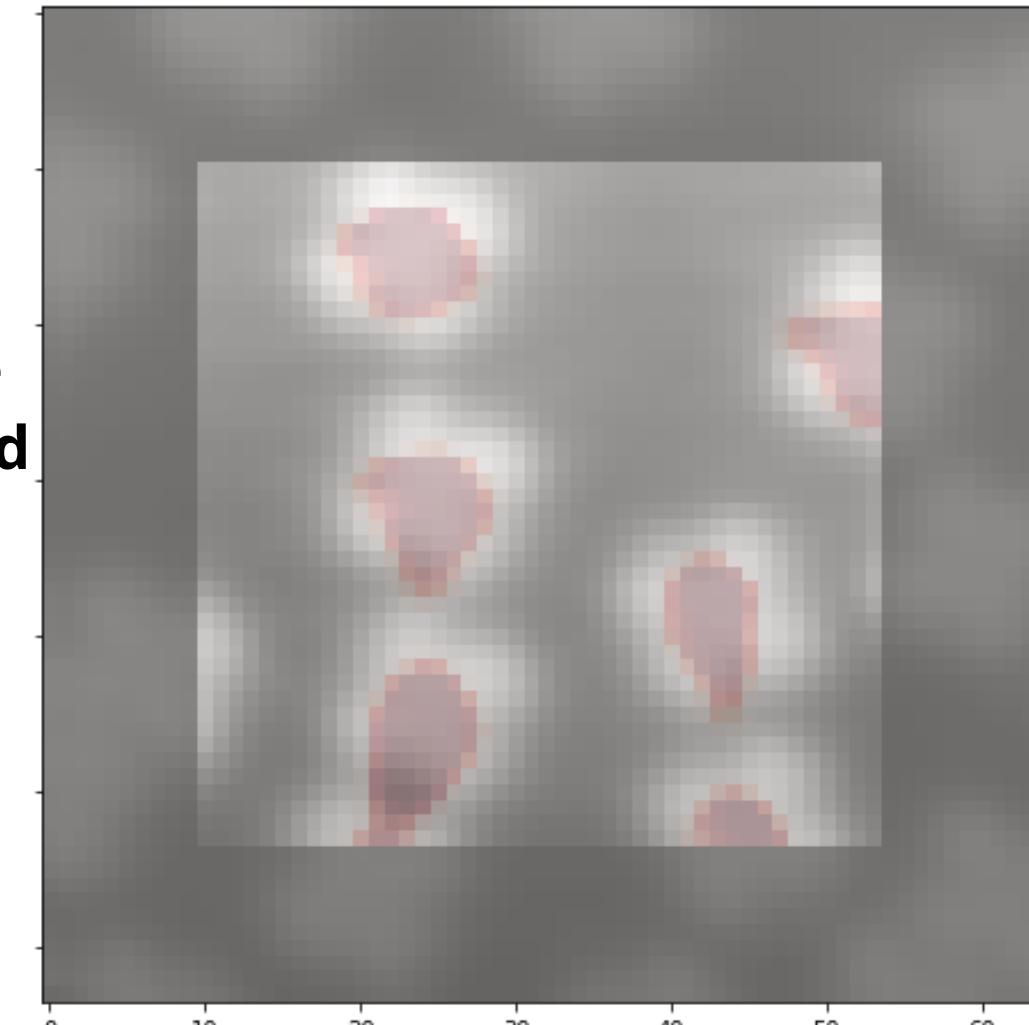


YOLO and R-CNN

Object detection

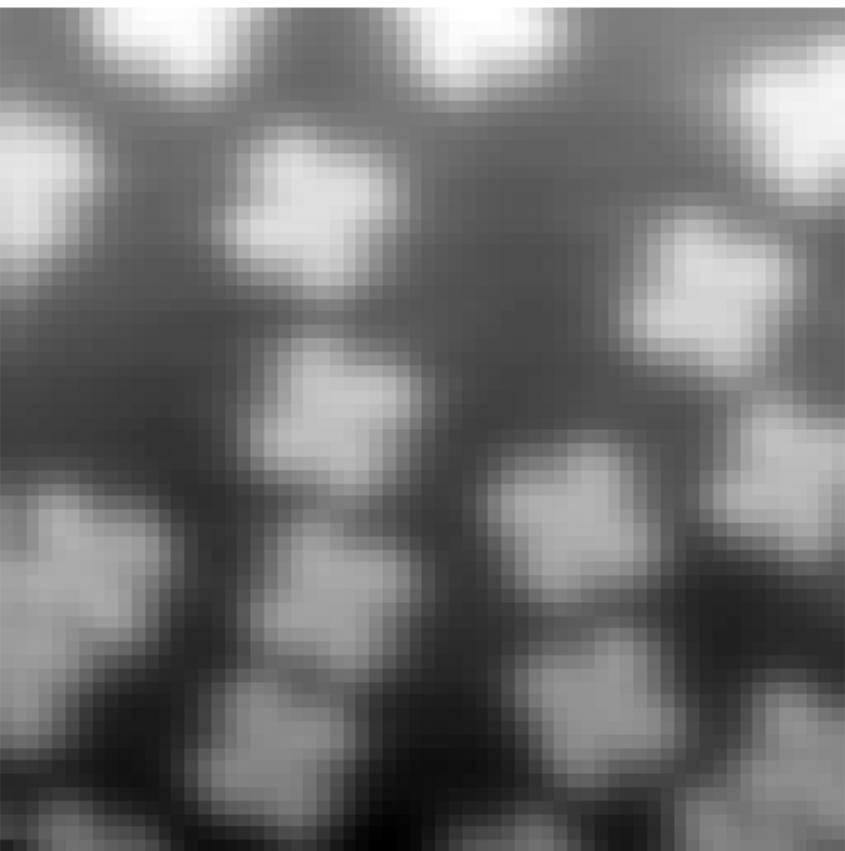


Problem: we had to search the whole image which is time consuming, **we could only find 1 kind of object at 1 scale**

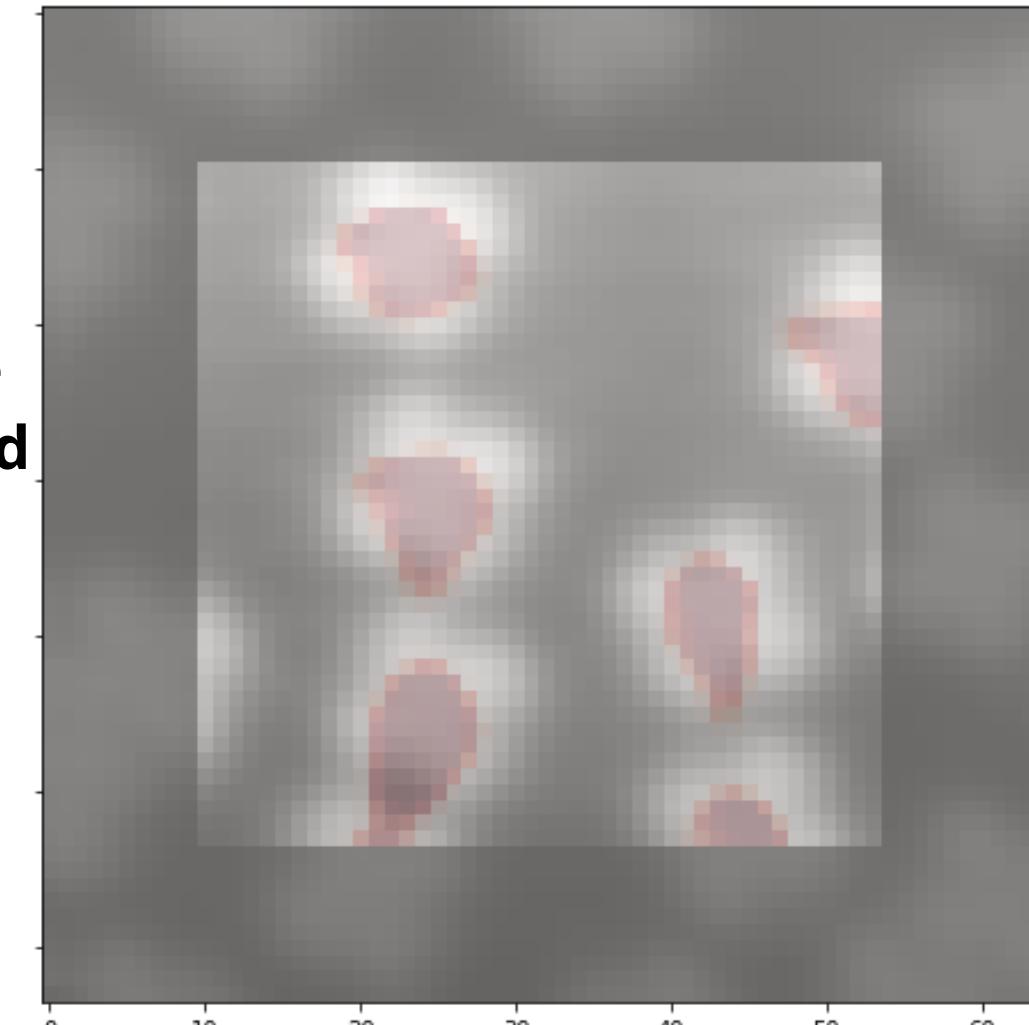


YOLO and R-CNN

Object detection



Problem: we had to search the whole image which is time consuming, **we could only find 1 kind of object at 1 scale**



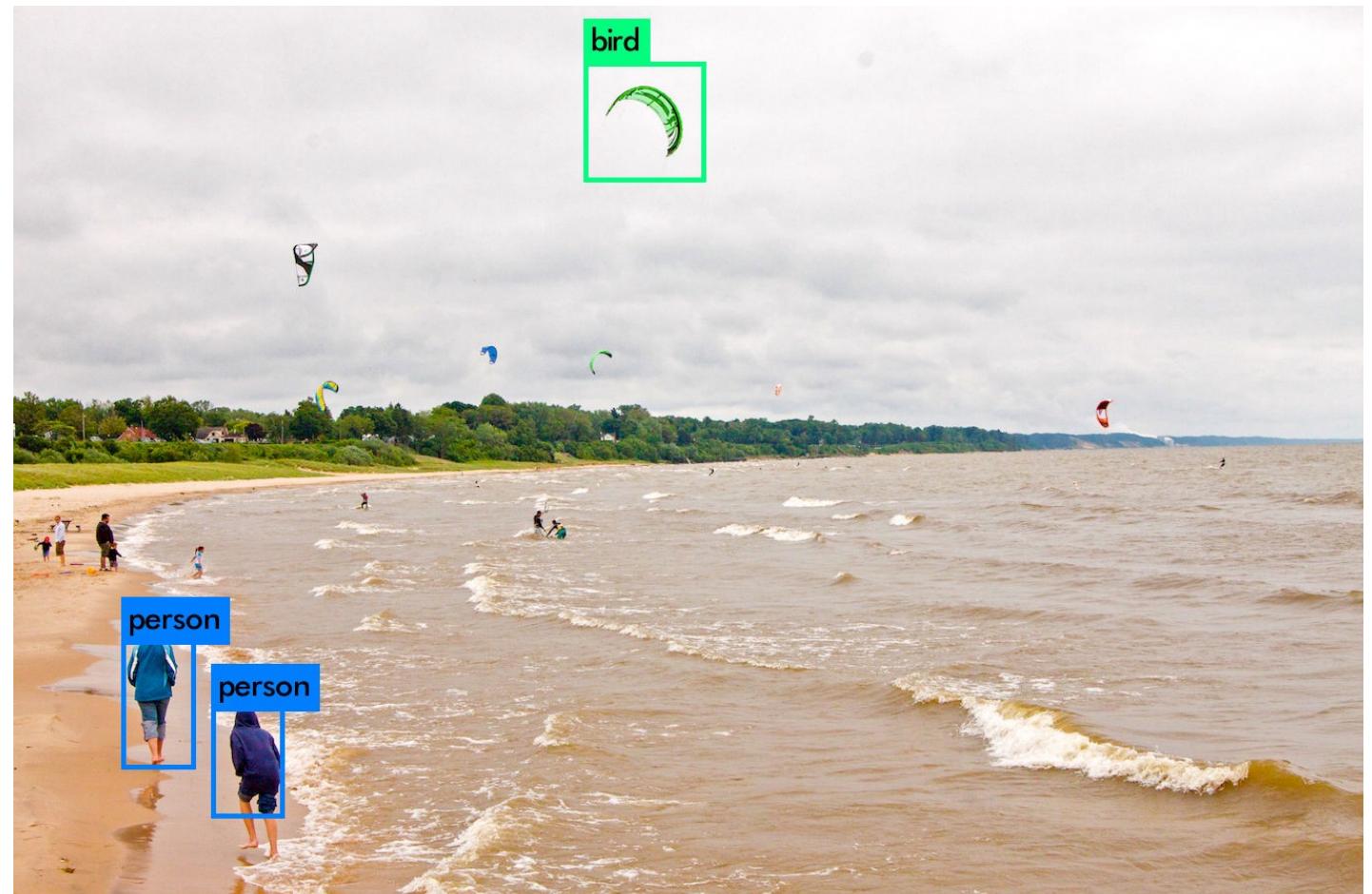
YOLO and R-CNN

Object detection

What if you do not know what is in the mage?

Final Dense layer has undefined size (one per kind of object in the region)

Objects can have different scale or axis ration: how many regions can you search before the problem blows up computationally??



R-CNN

Girshick et al. 2013

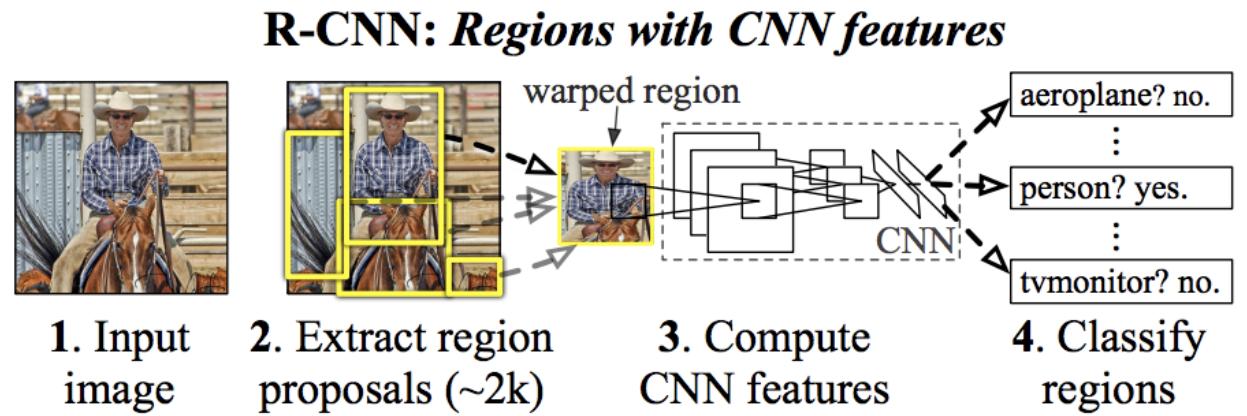
Extract 2000 regions from the image "region proposals."

Feature Extraction CNN produces a 4096-dimensional feature vector in an output dense layer

SVM classify the presence of the object within that candidate region proposal.

TOO SLOW (47 sec to test 1 image)

1. Generate initial sub-segmentation, we generate many candidate regions
2. Use greedy algorithm to recursively combine similar regions into larger ones
3. Use the generated regions to produce the final candidate region proposals



Fast R-CNN

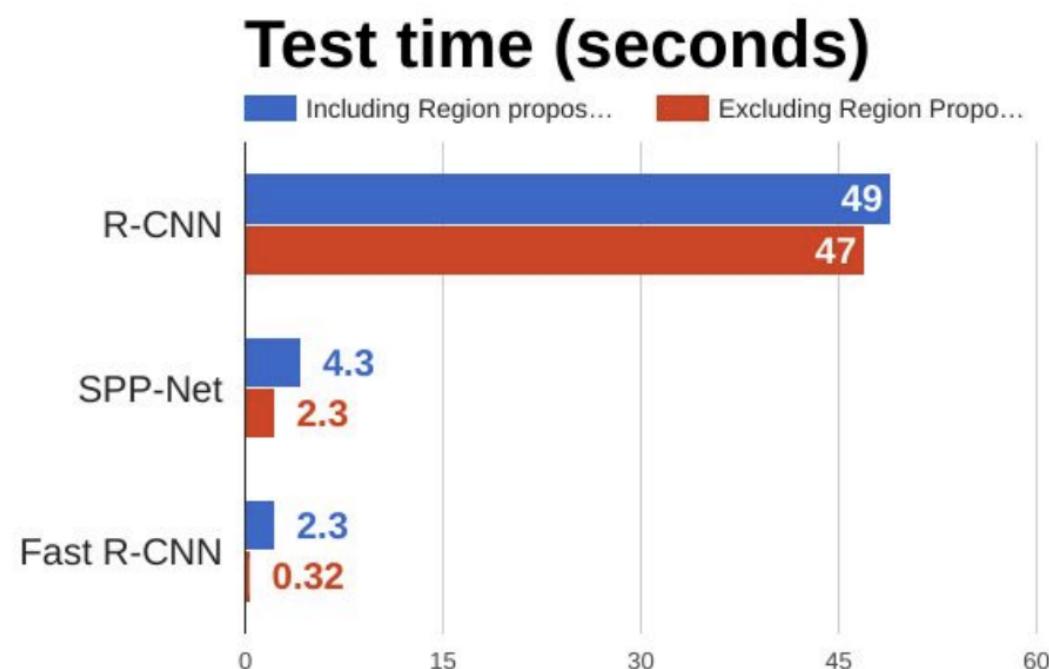
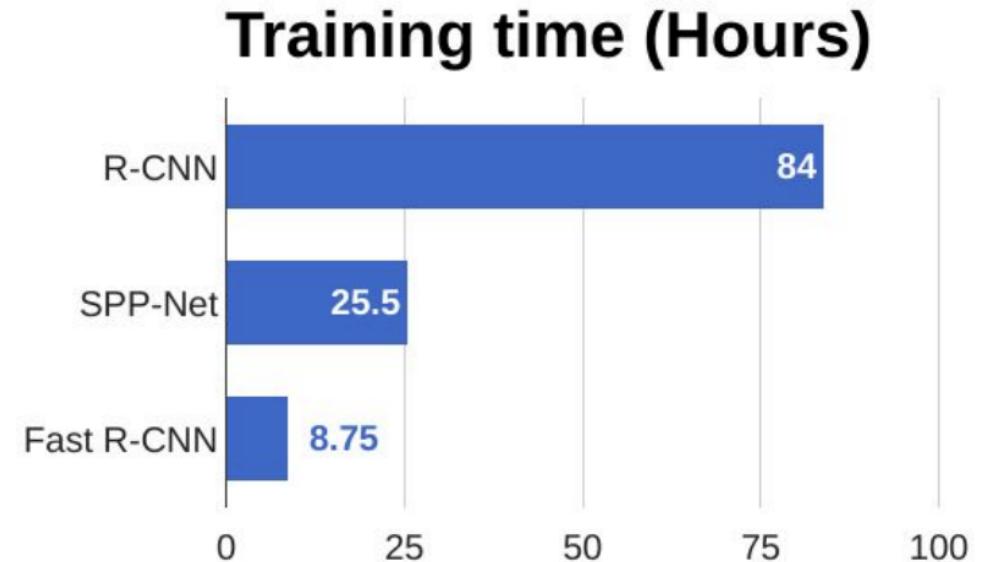
Girshick et al. 2015

Use a CNN to generate convolutional feature maps

Use Selective Search Algorithm to identify the RPs and warp them into squares

Using an ROI ***pooling layer*** to reshape them into a fixed size so that it can be fed into a fully connected layer - predict box offset

Softmax layer to predict the class of the proposed



Fast R-CNN

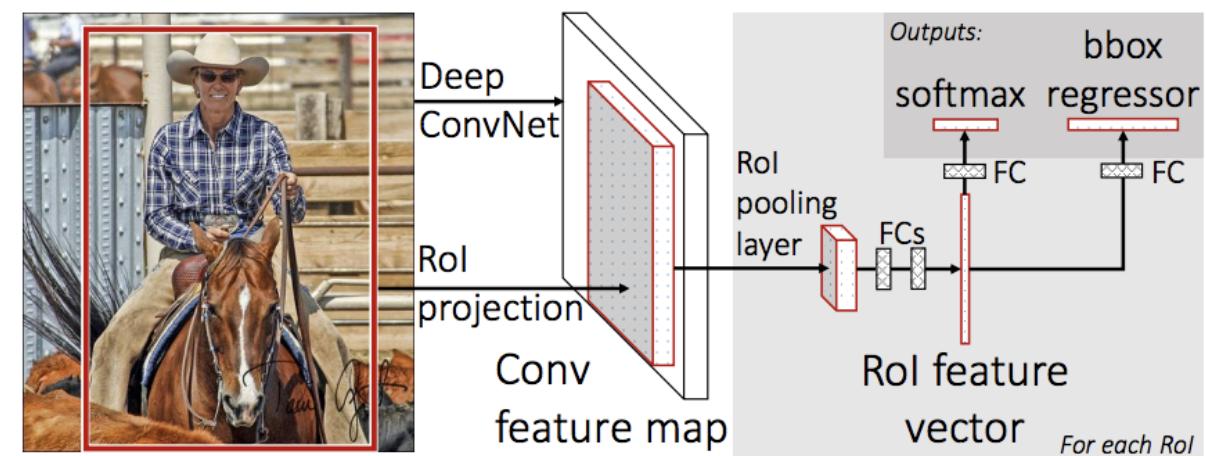
Girshick et al. 2015

Use a CNN to generate convolutional feature maps

Use Selective Search Algorithm to identify the RPs and warp them into squares

Using an Roi ***pooling layer*** to reshape them into a fixed size so that it can be fed into a fully connected layer - predict box offset

Softmax layer to predict the class of the proposed



Faster R-CNN

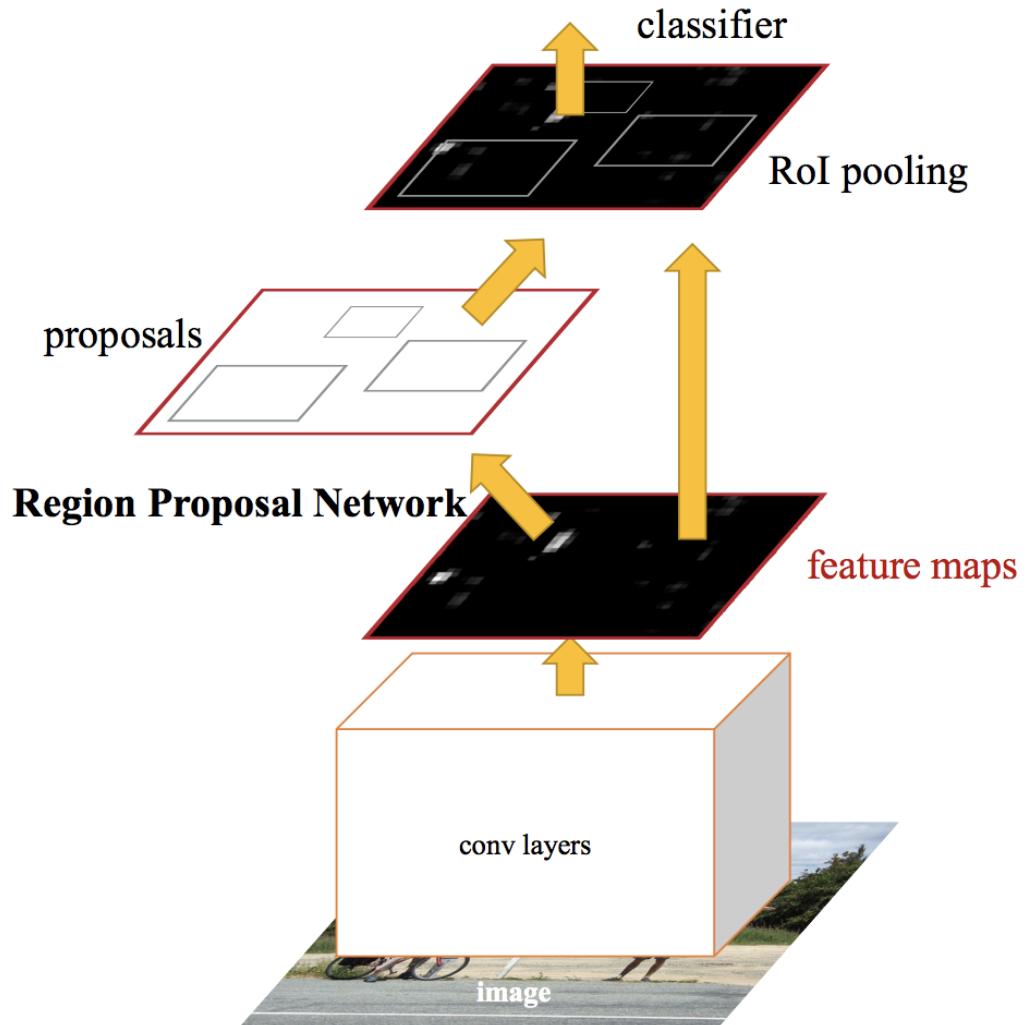
Ren et al. 2015

Use a CNN to generate convolutional feature maps

Use CNN to predict RPs and warp them into squares

Using an RoI ***pooling layer*** to reshape them into a fixed size so that it can be fed into a fully connected layer - predict box offset

Softmax layer to predict the class of the proposed



Faster R-CNN

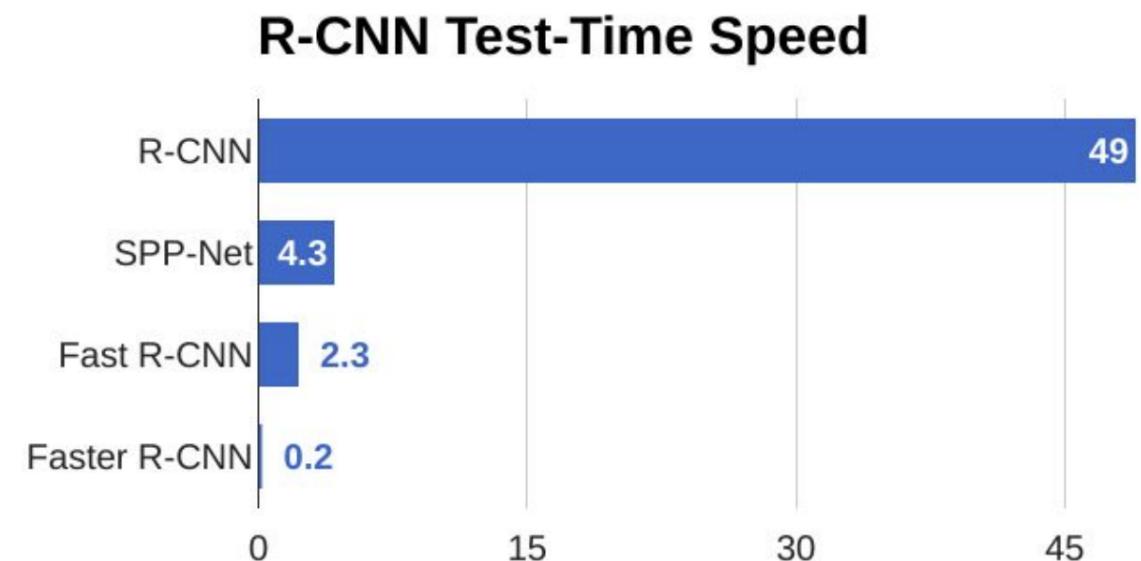
Ren et al. 2015

Use a CNN to generate convolutional feature maps

Use CNN to predict RPs and warp them into squares

Using an ROI ***pooling layer*** to reshape them into a fixed size so that it can be fed into a fully connected layer - predict box offset

Softmax layer to predict the class of the proposed



Yolo

Redmon et al 2016

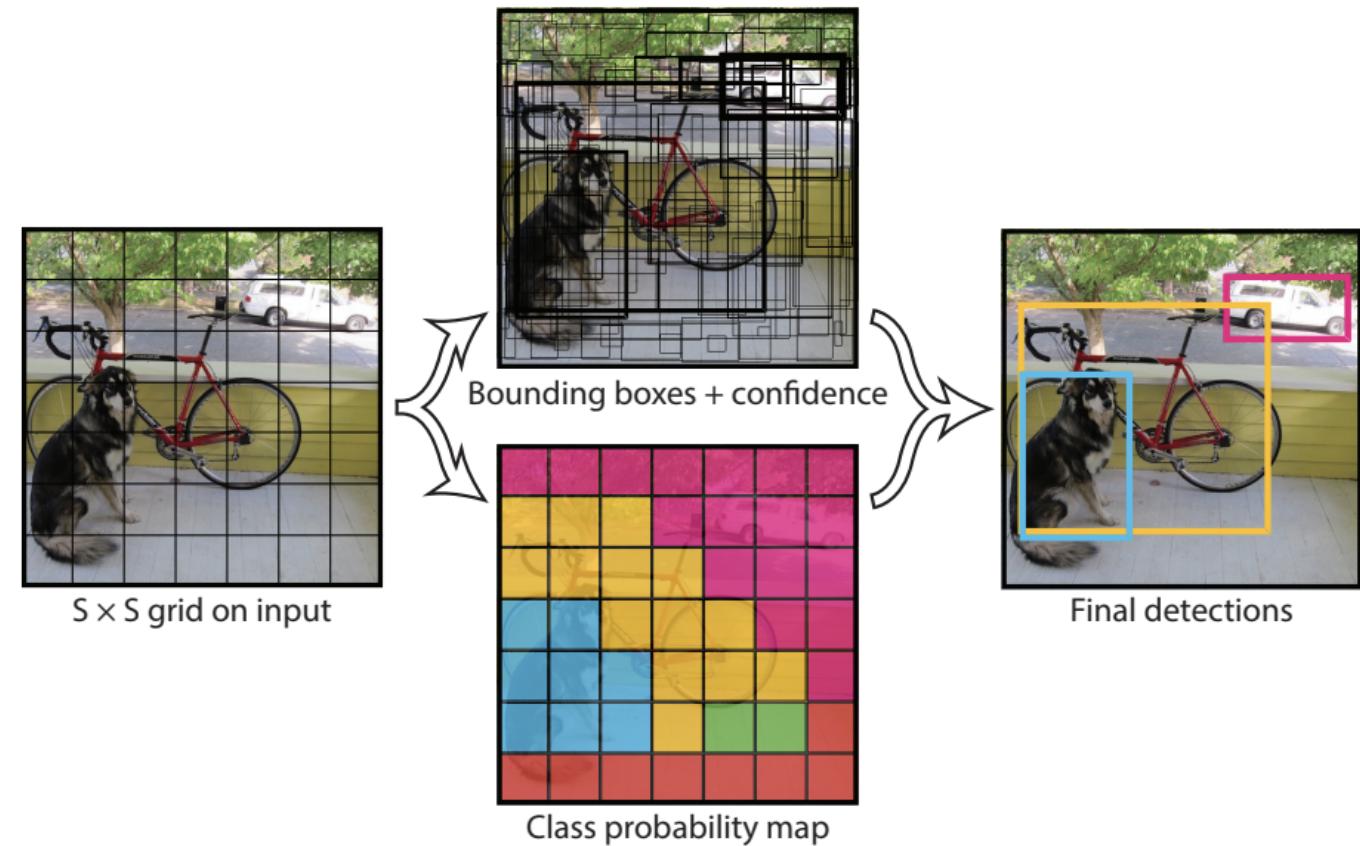
What if you looked at the whole image instead of Rols in the image??

Split an image into a $S \times S$ grid

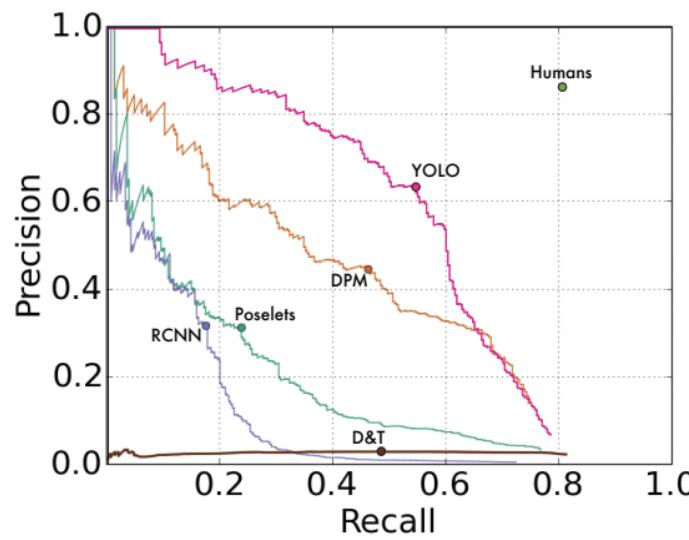
For each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities.

CNN outputs probability that BB has an object (+ offset)

High prob BBs are classified



Redmon et al 2016



(a) Picasso Dataset precision-recall curves.

	VOC 2007	Picasso		People-Art
	AP	AP	Best F_1	AP
YOLO	59.2	53.3	0.590	45
R-CNN	54.2	10.4	0.226	26
DPM	43.2	37.8	0.458	32
Poselets [2]	36.5	17.8	0.271	
D&T [4]	-	1.9	0.051	

(b) Quantitative results on the VOC 2007, Picasso, and People-Art Datasets. The Picasso Dataset evaluates on both AP and best F_1 score.

Figure 5: Generalization results on Picasso and People-Art datasets.

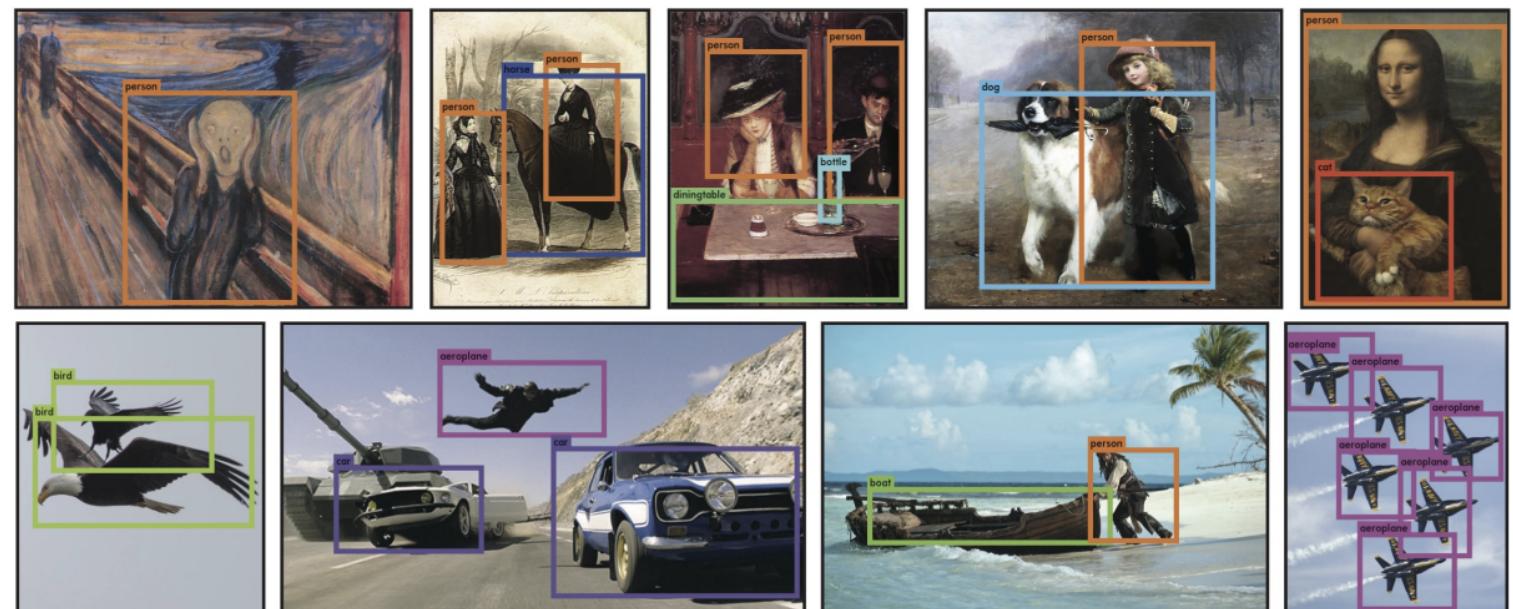
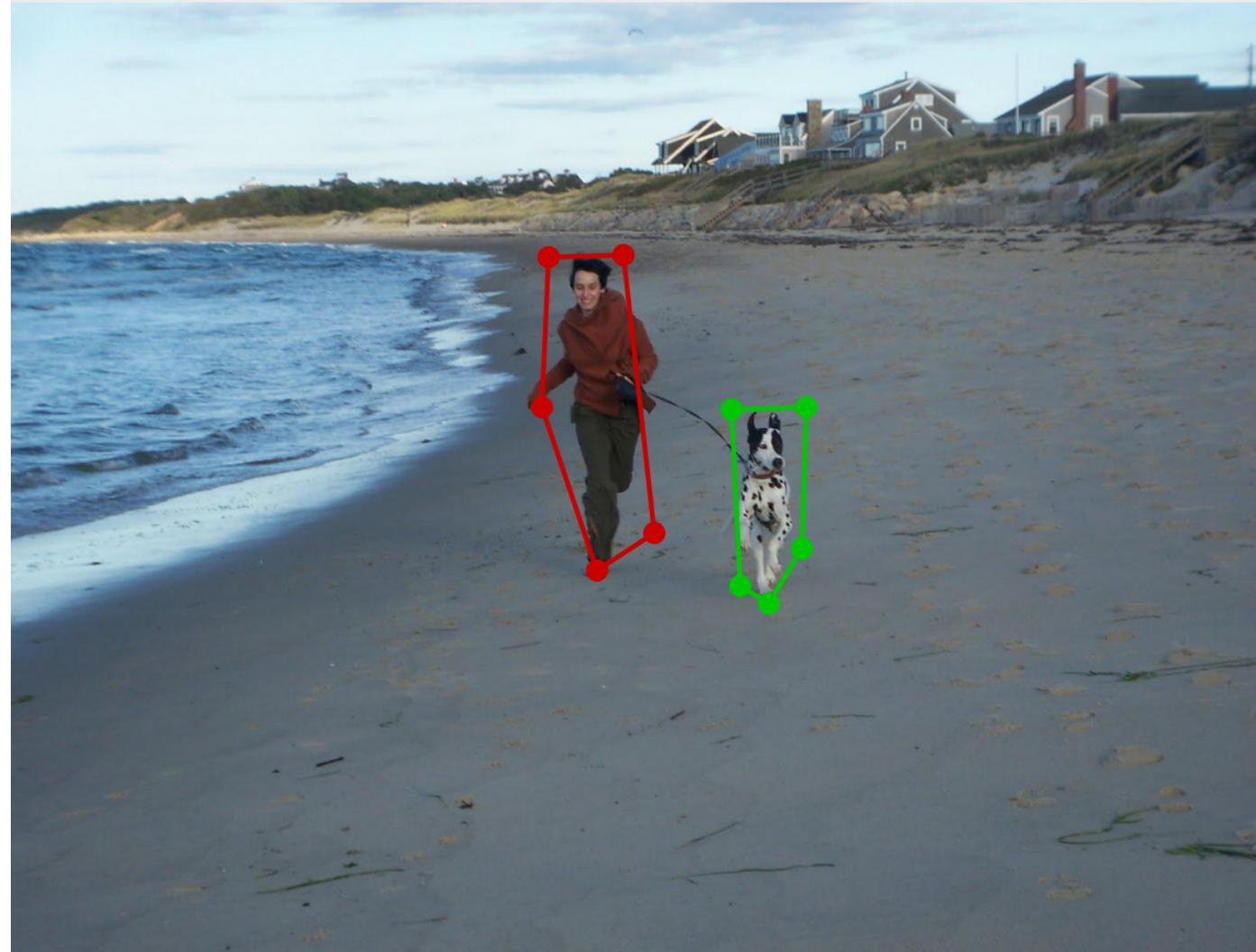


Figure 6: Qualitative Results. YOLO running on sample artwork and natural images from the internet. It is mostly accurate although it does think one person is an airplane.

Labling tools

<https://www.v7labs.com/blog/best-image-annotation-tools>



Neural Network and Deep Learning

an excellent and free book on NN and DL

<http://neuralnetworksanddeeplearning.com/index.html>

History of NN

<https://cs.stanford.edu/people/eroberts/courses/soco/projects/neural-networks/History/history2.html>

resources

Gradient Descent

https://ml-cheatsheet.readthedocs.io/en/latest/gradient_descent.html

Backpropagation

<http://colah.github.io/posts/2015-08-Backprop/>

resources



Physics Informed NN

Pin

-infinity - 1950's

Observations Descriptions			
2. S. Janis	Mar 11. '12	O **	
30. mons'		** O *	
2. Feb:		O *** *	
3. mons'		O * *	
3. Mar. 5.		* O *	
4. mons'		* O **	
6. mons'		** O *	
8. mar 11. '13.		*** O	
10. mons'		* * * O *	
11.		* * O *	
12. H. 4 negg:		* O *	
13. mons'		* * O *	
14. Mar. 6.		* * * O *	

Application regime:

theory driven: little data, mostly theory,
falsifiability and all that...

PinN

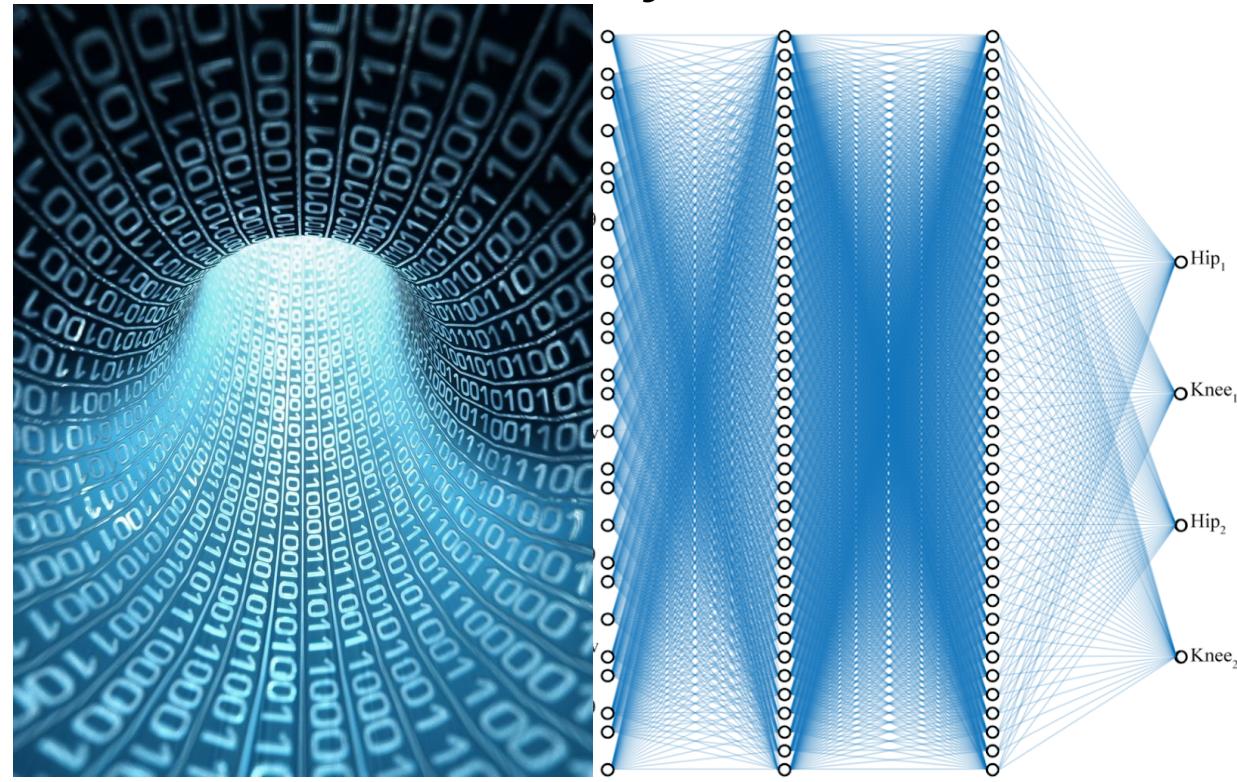
-infinity - 1950's

Observations Deviations			
2. S. Joris marl H. 12	O **		
30. none	** O *		
2. xon:	O *** *		
3. more	O * *		
3. Ho. s.	* O *		
4. more!	* O **		
6. manc	** O *		
8. manc H. 13.	*** * O		
10. manc.	* * * O *		
11.	* * O *		
12. H. 4 negg:	* O *		
13. manc'	* * * O *		
14. manc.	* * * O *		

theory driven: little data, mostly theory, falsifiability and all that...

Application regime:

-1980's - today



data driven: lots of data, drop theory and use associations, black-box models

PINN

-infinity - 1950's

	Observations Deviations	
2. S. Jan. 1951	O ***	
30. mard.	** O *	
2. Feb.	O *** *	
3. mard.	O * *	
3. Febr. 5.	* O *	
9. mard.	* O **	
6. mard.	** O *	
8. mard. H. 13.	*** * O	
10. mard.	* * * O *	
11.	* * O *	
12. H. 4 wegg.	* O *	
13. mard.	* ** O *	
14. Febr.	* *** O *	

theory driven: little data, mostly theory, falsifiability and all that...

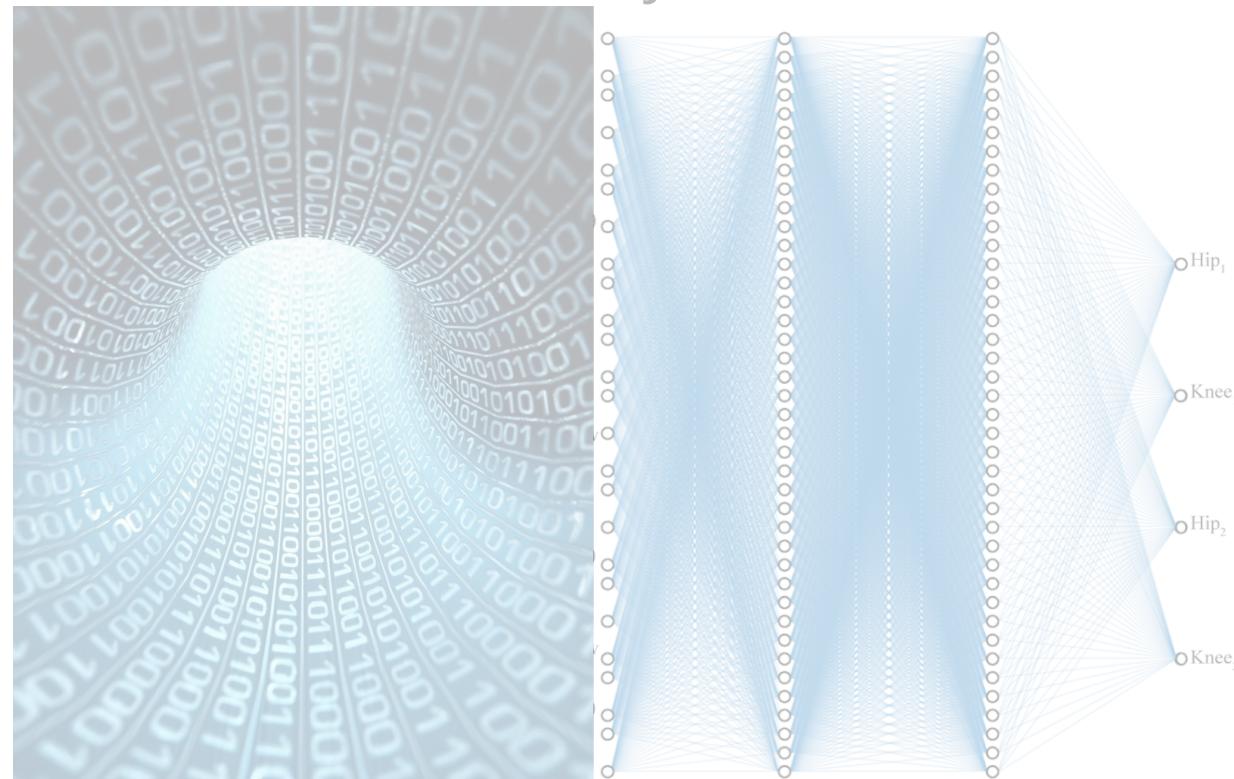
lots of data yet not enough for entirely automated decision making

complex theory that cannot be solved analytically

combine it with some theory

Application regime:

-1980's - today



data driven: lots of data, drop theory and use associations, black-box models



$\mathcal{N}[\cdot]$ is a nonlinear differential operator

$$\partial_t u(t, x) + \mathcal{N}[u](t, x) = 0$$

$$u(0, x) = u_0(x)$$

$$u : [0, T] \times D \Rightarrow \mathbb{R}$$

$$(t, x) \in (0, T] \times D$$

General conservation law

$$u(t) + \frac{\partial f(u)}{\partial x} = 0$$

e.g. flux function (linear)

$$f(u) = a \cdot u$$

Burgers eq (non-linear)

$$f(u) = \frac{1}{2}u^2$$



Non Linear PDEs are hard to solve!

Existence and uniqueness of solutions

A fundamental question for any PDE is the existence and uniqueness of a solution for given boundary conditions. open problem of [existence \(and smoothness\) of solutions](#) to the Navier–Stokes equations is one of the seven [Millennium Prize problems](#) in mathematics.

- Raissi et al. *Physics Informed Deep Learning (Part I): Data-driven Solutions of Nonlinear Partial Differential Equations.* arXiv [1711.10561](https://arxiv.org/abs/1711.10561)
- Raissi et al. *Physics Informed Deep Learning (Part II): Data-driven Discovery of Nonlinear Partial Differential Equations.* arXiv [1711.10566](https://arxiv.org/abs/1711.10566)
- Raissi et al. *Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations.* J. Comp. Phys. 378 pp. 686-707 DOI: [10.1016/j.jcp.2018.10.045](https://doi.org/10.1016/j.jcp.2018.10.045)

https://en.wikipedia.org/wiki/Nonlinear_partial_differential_equation



Non Linear PDEs are hard to solve!

Linear approximation

The solutions in a neighborhood of a known solution can sometimes be studied by linearizing the PDE around the solution. This corresponds to studying the tangent space of a point of the moduli space of all solutions.

- Raissi et al. *Physics Informed Deep Learning (Part I): Data-driven Solutions of Nonlinear Partial Differential Equations*. arXiv 1711.10561
- Raissi et al. *Physics Informed Deep Learning (Part II): Data-driven Discovery of Nonlinear Partial Differential Equations*. arXiv 1711.10566
- Raissi et al. *Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations*. J. Comp. Phys. 378 pp. 686-707 DOI: 10.1016/j.jcp.2018.10.045

https://en.wikipedia.org/wiki/Nonlinear_partial_differential_equation



Non Linear PDEs are hard to solve!

Exact solutions

It is often possible to write down some special solutions explicitly in terms of elementary functions (though it is rarely possible to describe all solutions like this). One way of finding such explicit solutions is to reduce the equations to equations of lower dimension, preferably ordinary differential equations, which can often be solved exactly.

- Raissi et al. *Physics Informed Deep Learning (Part I): Data-driven Solutions of Nonlinear Partial Differential Equations*. arXiv 1711.10561
- Raissi et al. *Physics Informed Deep Learning (Part II): Data-driven Discovery of Nonlinear Partial Differential Equations*. arXiv 1711.10566
- Raissi et al. *Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations*. J. Comp. Phys. 378 pp. 686-707 DOI: 10.1016/j.jcp.2018.10.045

https://en.wikipedia.org/wiki/Nonlinear_partial_differential_equation



Non Linear PDEs are hard to solve!

Numerical solutions

Numerical solution on a computer is almost the only method that can be used for getting information about arbitrary systems of PDEs. There has been a lot of work done, but a lot of work still remains on solving certain systems numerically, especially for the Navier–Stokes and other equations related to weather prediction.

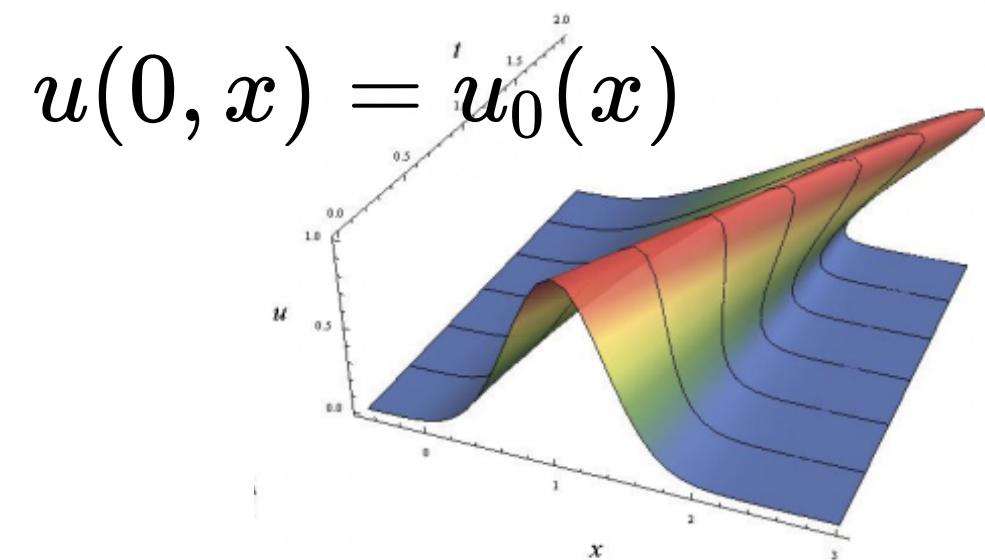
- Raissi et al. *Physics Informed Deep Learning (Part I): Data-driven Solutions of Nonlinear Partial Differential Equations.* arXiv 1711.10561
- Raissi et al. *Physics Informed Deep Learning (Part II): Data-driven Discovery of Nonlinear Partial Differential Equations.* arXiv 1711.10566
- Raissi et al. *Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations.* J. Comp. Phys. 378 pp. 686-707 DOI: 10.1016/j.jcp.2018.10.045

https://en.wikipedia.org/wiki/Nonlinear_partial_differential_equation



**Burgers equation:
second order non-linear PDE**

$$\partial_t u + u \partial_x u - \nu \partial_{xx} u = 0,$$



**Applications of Burgers eq:
shock weave formation, turbulence, the
weather problem, traffic flow and
acoustic transmission**

$$u : [0, T] \times D \Rightarrow \mathbb{R}$$

$$(t, x) \in (0, T] \times D$$

x : spatial coordinate

t : temporal coordinate

$u(x, t)$: speed of fluid at x, t

ν : viscosity



**Burgers equation:
second order non-linear PDE**

$$\partial_t u + u \partial_x u - (0.01/\pi) \partial_{xx} u = 0,$$

How to solve analytically

<https://www.youtube.com/watch?v=5ZrwxQr6aV4>

Domain

$$(t, x) \in (0, 1] \times (-1, 1),$$
$$x \in [-1, 1],$$
$$t \in (0, 1]$$

Boundary Conditions

$$u(0, x) = -\sin(\pi x),$$
$$u(t, -1) = u(t, 1) = 0.$$



**Burgers equation:
second order non-linear PDE**

$$\partial_t u + u \partial_x u - (0.01/\pi) \partial_{xx} u = 0,$$

How to solve analytically

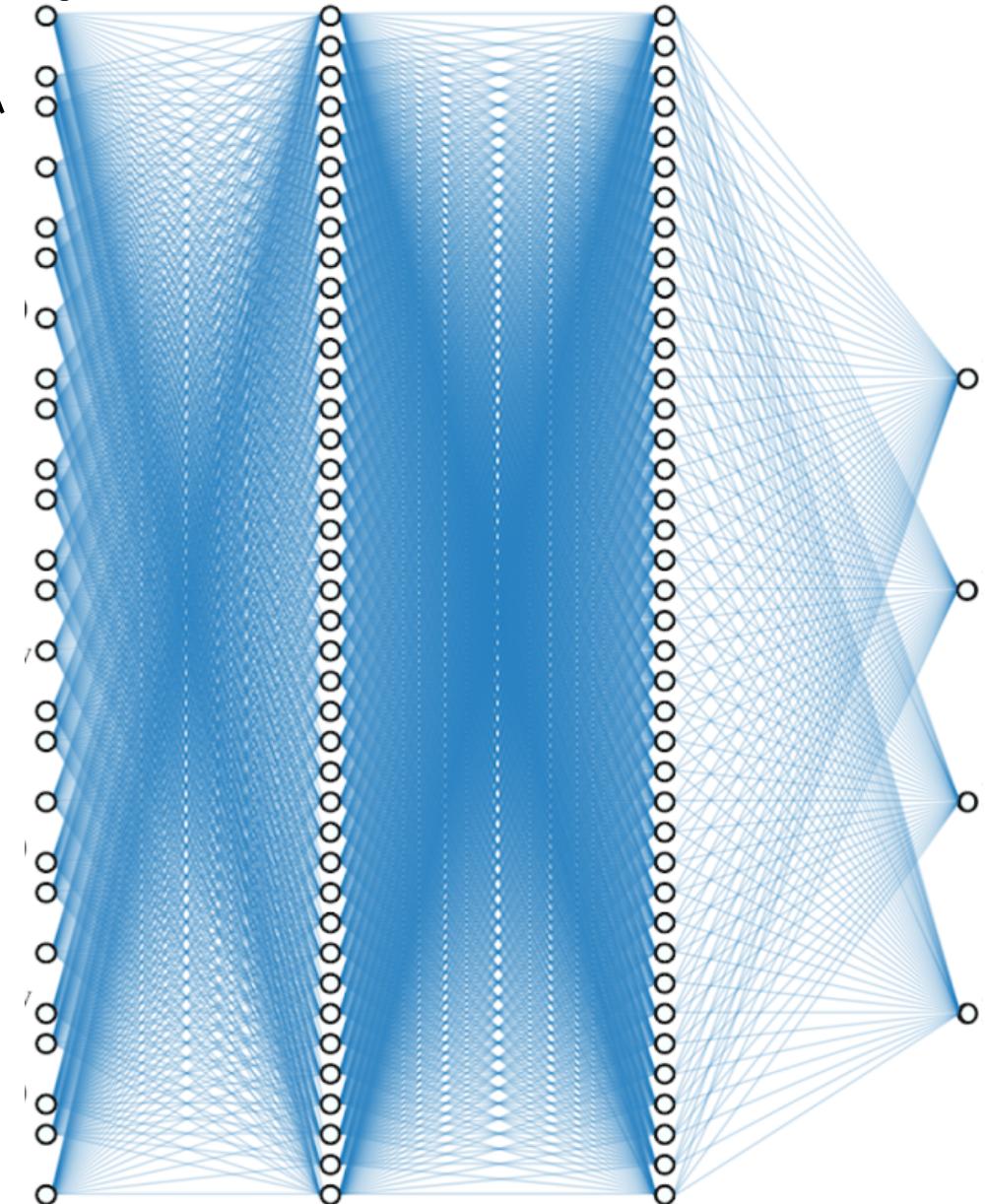
[https://www.youtube.com/watch?
v=5ZrwxQr6aV4](https://www.youtube.com/watch?v=5ZrwxQr6aV4)

[https://www.youtube.com/embed/nFBvAd81lds?
enablejsapi=1](https://www.youtube.com/embed/nFBvAd81lds?enablejsapi=1)

PinN

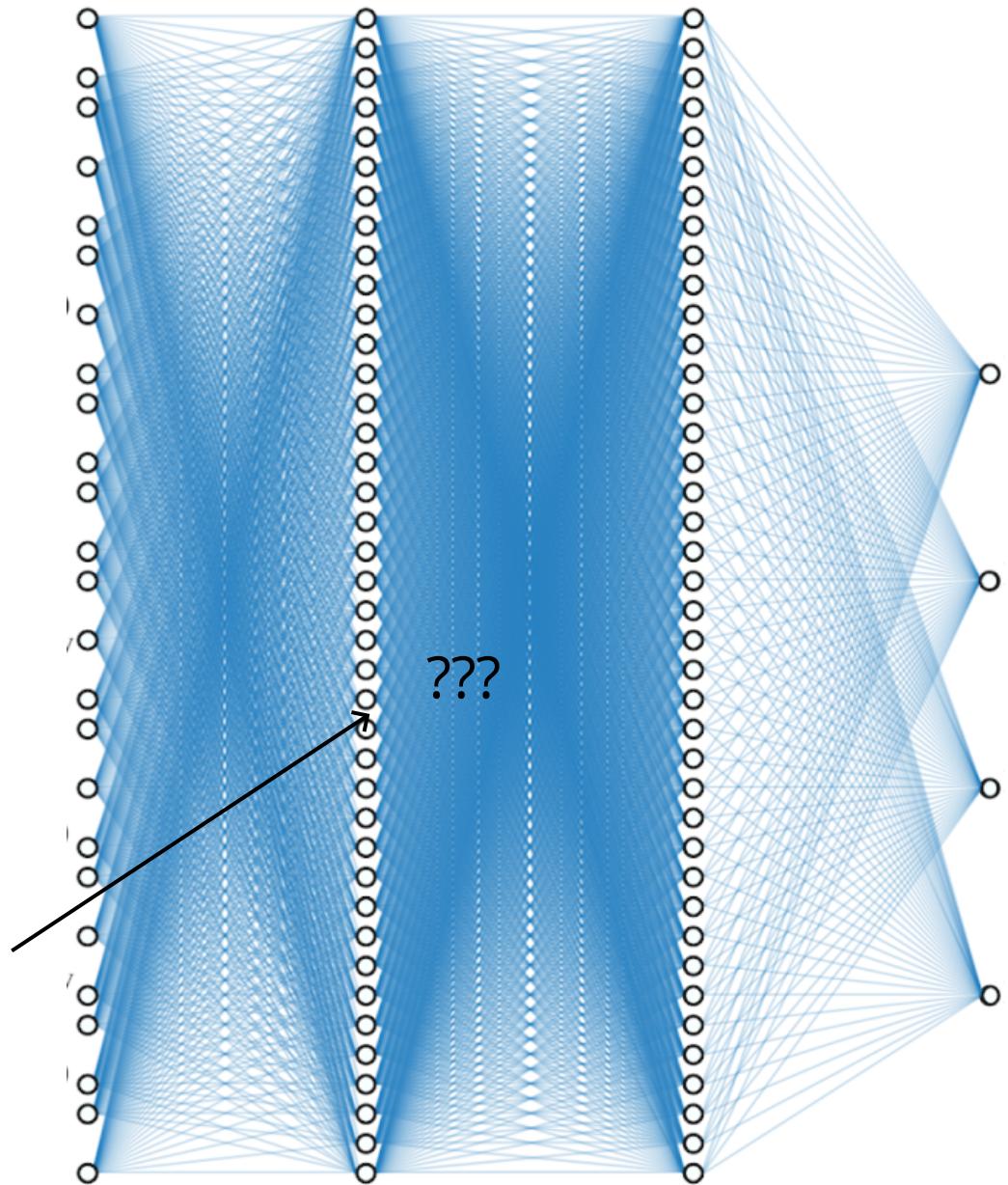
input layer

- Provide training points at the boundary with calculated solution (trivial cause we have boundary conditions)



PINN

- Provide training points at the boundary with calculated solution (trivial cause we have boundary conditions)
- Provide the physical constraint: make sure the solution satisfies the PDE

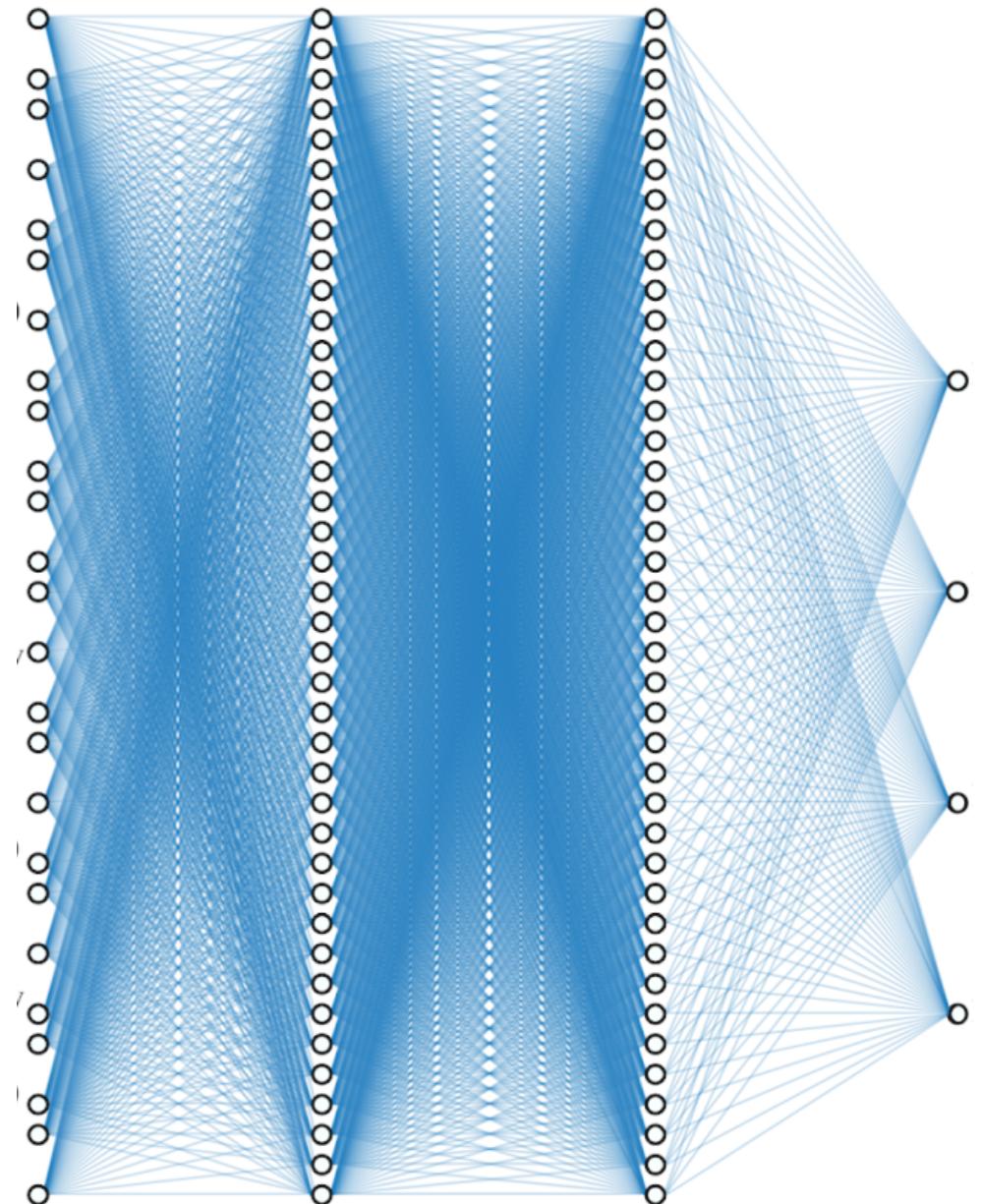


PinN

- Provide training points at the boundary with calculated solution (trivial cause we have boundary conditions)

- Provide the physical constraint: make sure the solution satisfies the PDE

via a modified loss function that includes
residuals of the prediction and residual of
the PDE



PinN

- Provide training points at the boundary with calculated solution (trivial cause we have boundary conditions)

- Provide the physical constraint: make sure the solution satisfies the PDE

via a modified loss function that includes
residuals of the prediction and residual of
the PDE

$$\begin{aligned} \text{loss} = L2 + PDE = \\ \sum(u_\theta - u)^2 + \\ (\partial_t u_\theta + u_\theta \partial_x u_\theta - (0.01/\pi) \partial_{xx} u_\theta)^2 \end{aligned}$$



<https://maziarraissi.github.io/PINNs/>

PINN



[https://github.com/fedhere/MLPNS_FBianco
/blob/main/PINN/PINN_Burgers.ipynb](https://github.com/fedhere/MLPNS_FBianco/blob/main/PINN/PINN_Burgers.ipynb)

- Provide training points at the boundary with calculated solution (trivial cause we have boundary conditions)

$$\begin{aligned} \text{loss} = & L2 + PDE = \\ & \sum(u_\theta - u)^2 + \\ & \partial_t u_\theta + u_\theta \partial_x u_\theta - (0.01/\pi) \partial_{xx} u_\theta \end{aligned}$$

- Provide the physical constraint: make sure the solution satisfies the PDE

via a modified loss function that includes
residuals of the prediction and residual of
the PDE