



LUND
UNIVERSITY

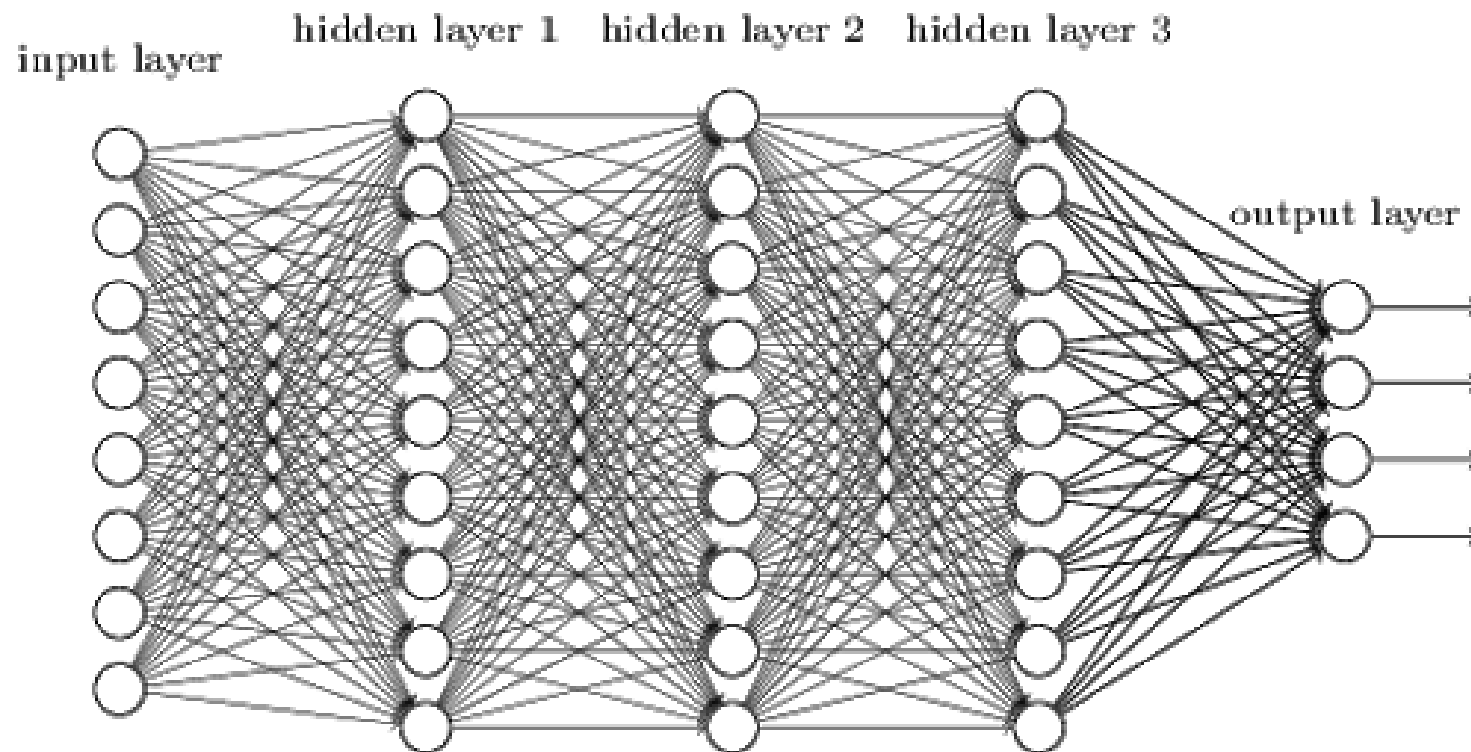
EITP25 2020

Lecture 4 – Machine Learning Models



Limitations of "regular" Neural Nets

Fully connected network

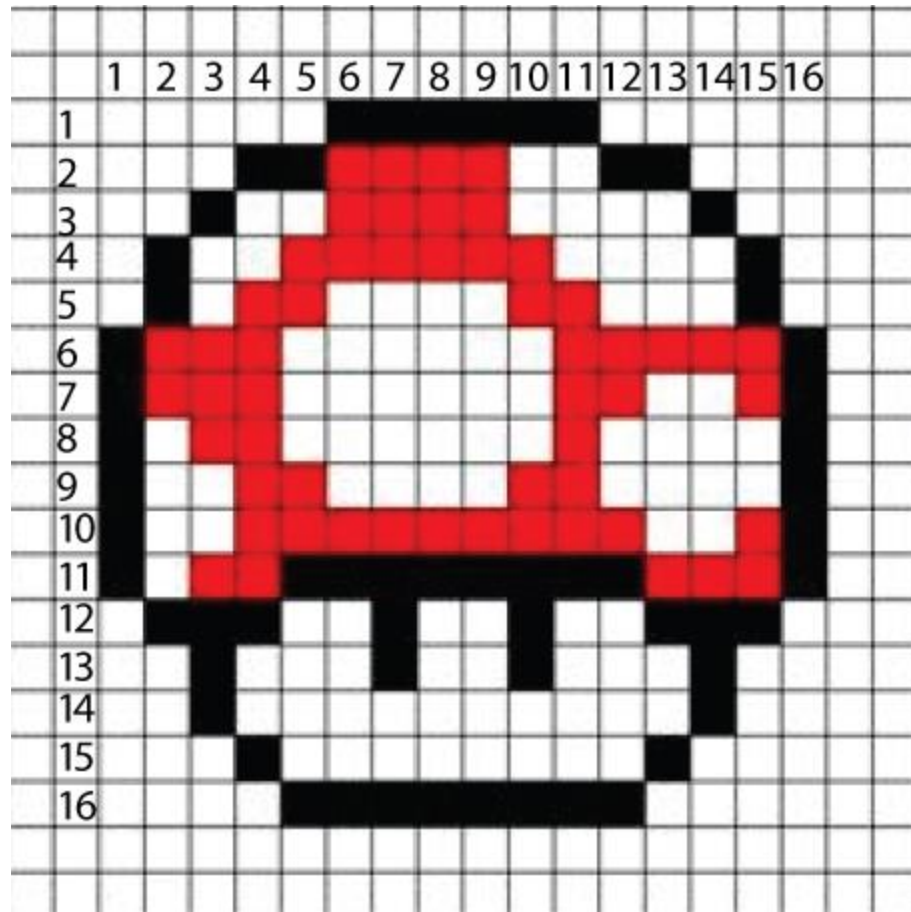


Outline – Lecture 4

Other types of neural network models

- Convolutional Neural Networks
- Recurrent Neural Networks (RNNs)
- Locally Connected Networks (LCNs)

Spatially structured data



An image is more than the sum of its pixels

**i.e. the relative position of data points
contain information, not just the pixel colors**

**How can we build an ANN that
takes advantage of this?**

The convolution operation

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

image

0	0	1
0	1	0
1	0	0

“Filter”

0	0	1
0	1	0
1	0	0

*

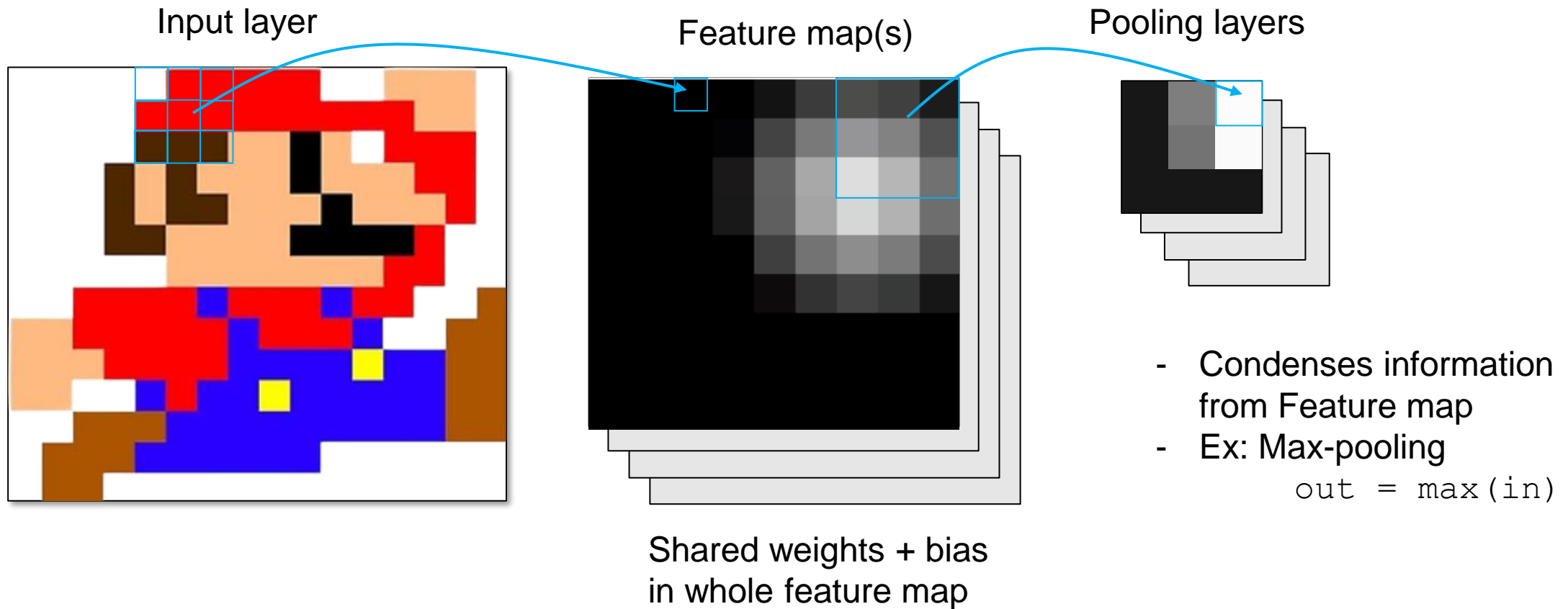
$$= 0*0 + 0*0 + 1*1 + 0*0 + 1*1 + 0*0 + 0*1 + 1*0 + 0*0 = 2$$

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

2	2	0	1
1	1	1	1
1	1	1	1
1	0	2	2

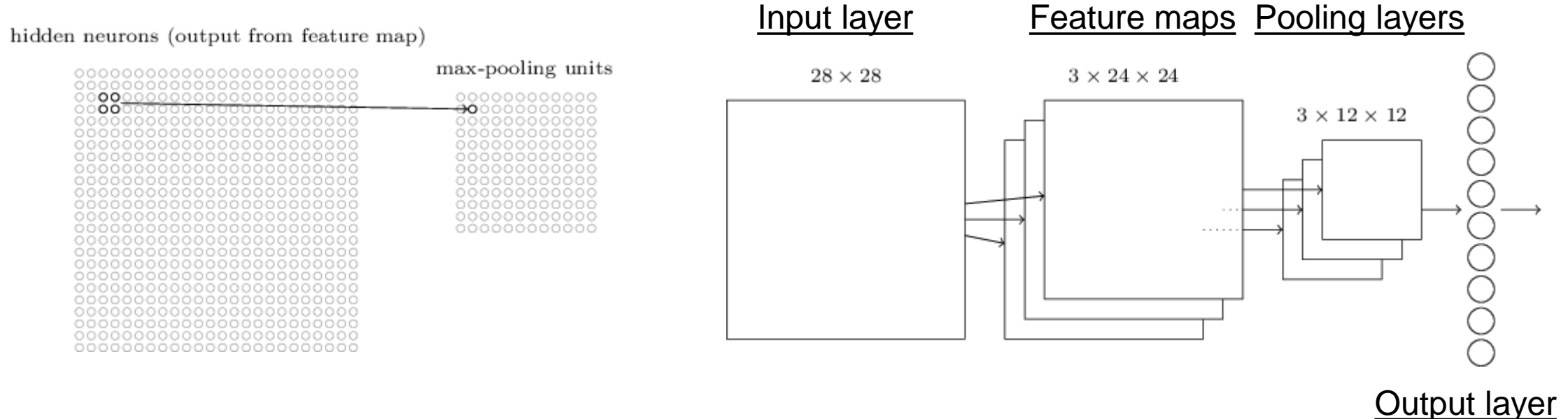
Feature map

Convolution in an ANN



Convolutional Neural Network (CNN)

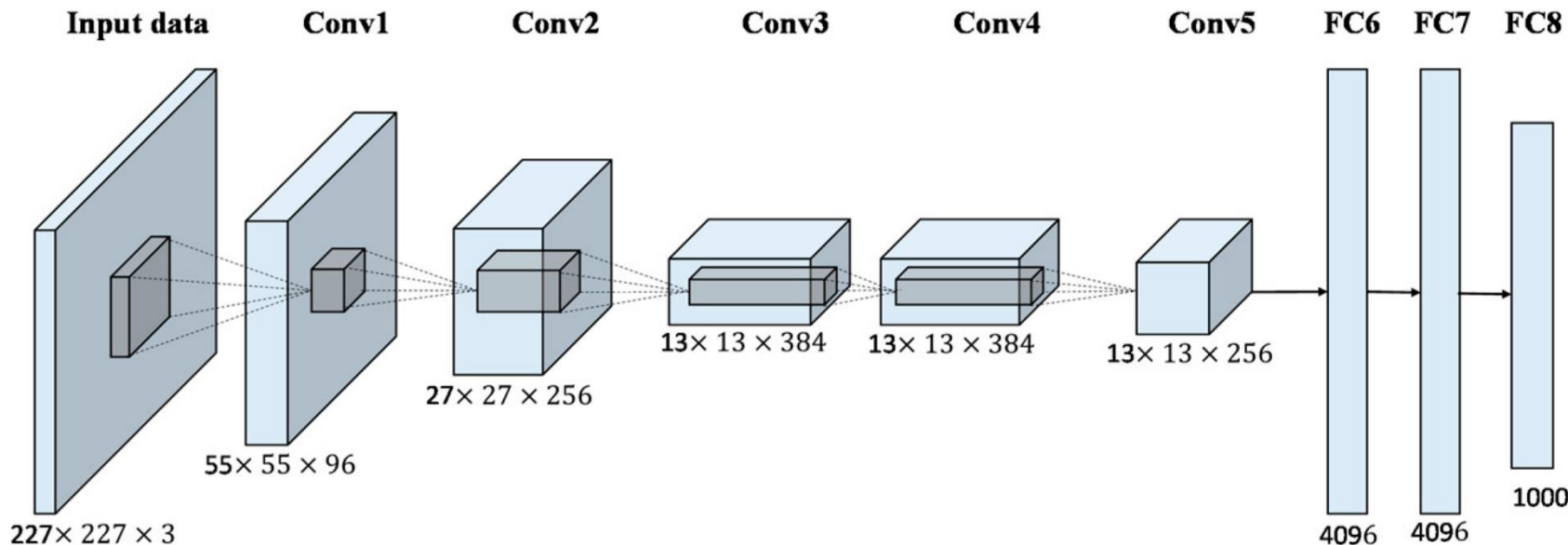
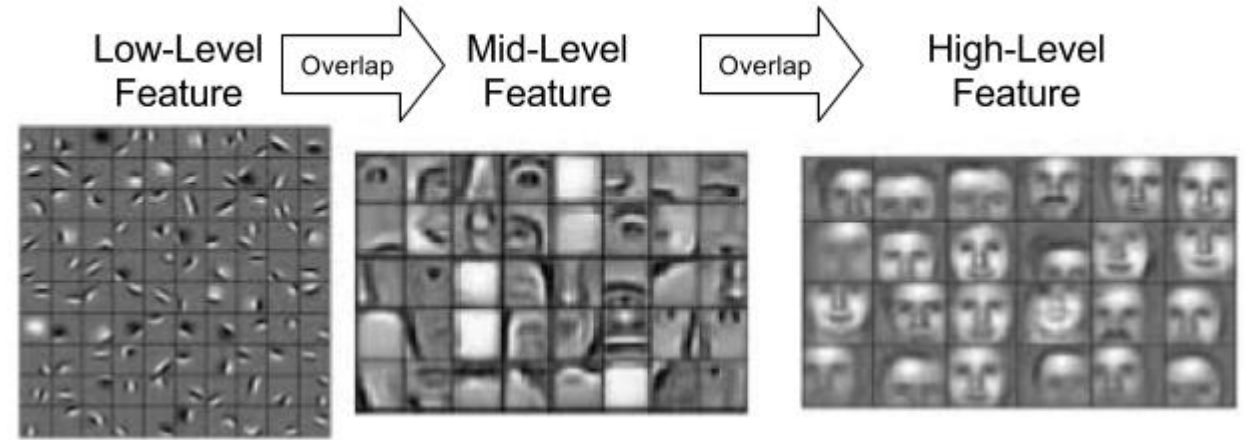
- A multitude of feature maps
- One pooling layer per feature map
 - L2 pooling: root of the sum of squares
- The Output layer is a normal completely connected layer to all neurons in the pooling layers



A full CNN

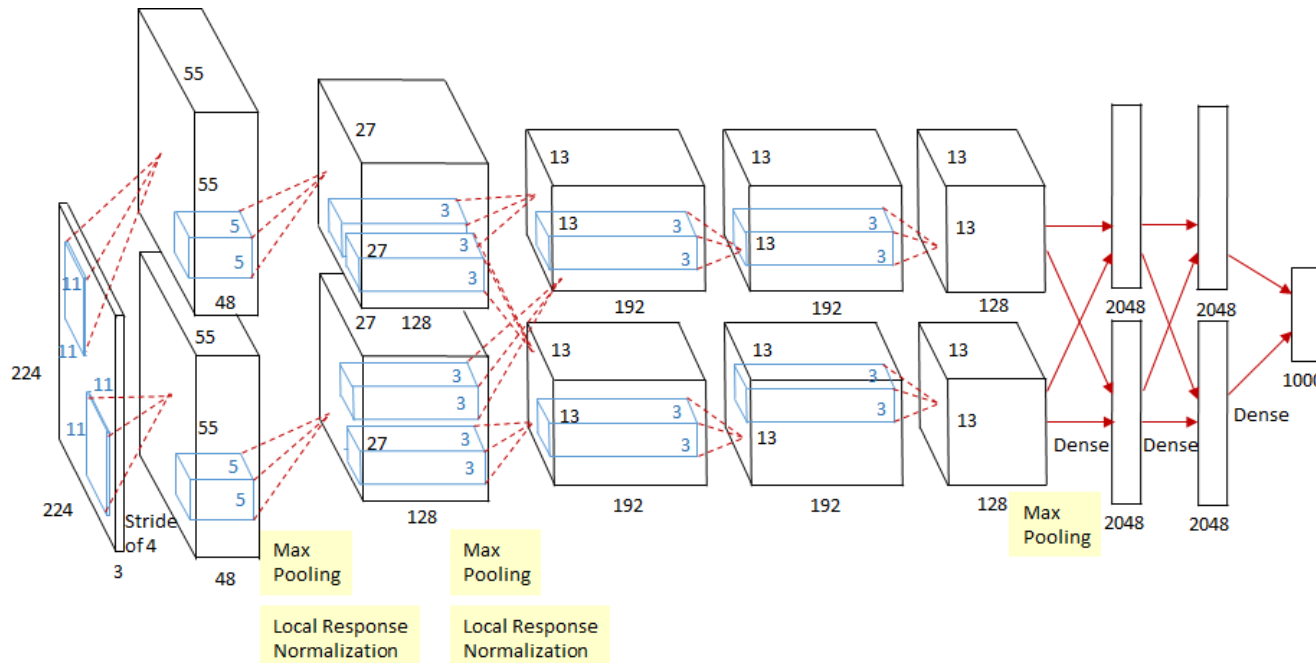
- Multiple convolutional layers for various levels of features.
- Data size reduced between layers
- For classification: Finish off with fully-connected layers

Feature Map in Convolutional Neural Networks (CNN)



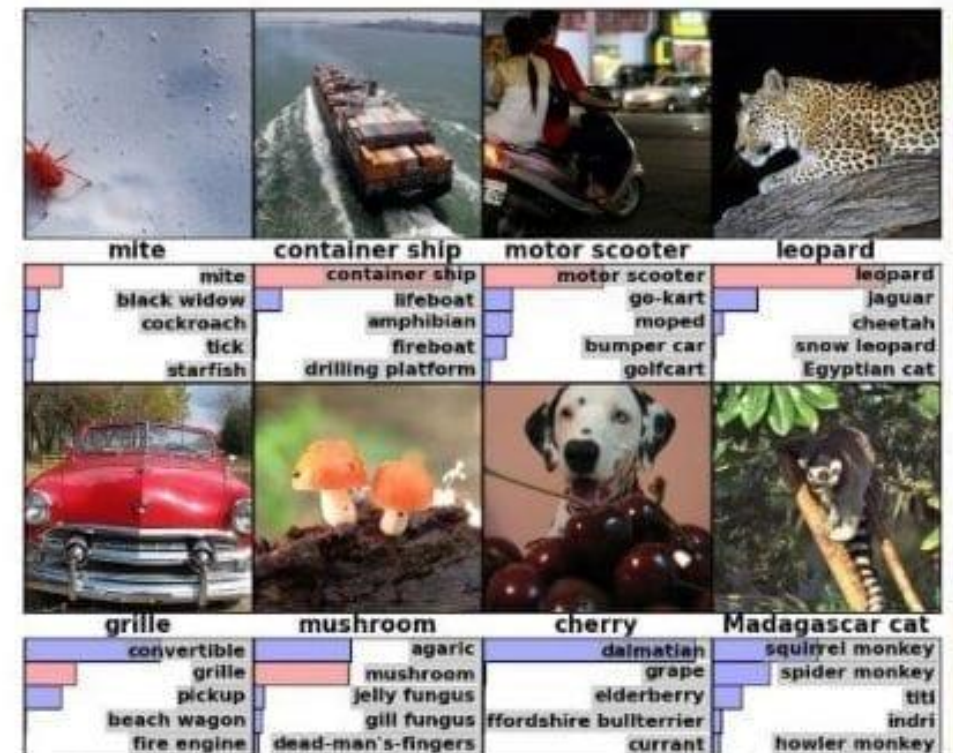
Classification Example

- AlexNet: Pioneer in CNNs, now standard for image classification
 - First to use ReLU on CNNs
 - 5 conv. Layers + 2 FCLs → 62.3 million parameters (!)
 - 6 days to train
- Won the 2012 ImageNet competition by a large margin
- Top 1 accuracy of 62.5% (Top 5: 83%)



ImageNet

15 million high-res images
in ~22000 categories

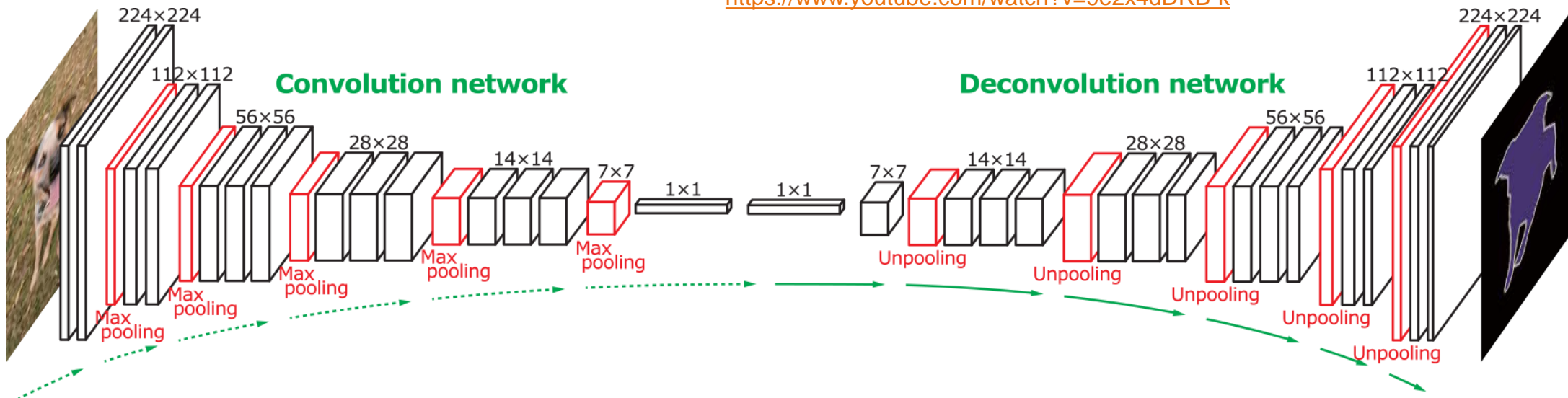


Other example

- Image segmentation
- Upscale features by “deconvolution”



<https://www.youtube.com/watch?v=9e2x4dDRB-k>



<https://medium.com/@wilburdes/semantic-segmentation-using-fully-convolutional-neural-networks-86e45336f99b>

Sequences

- CNNs see spatial correlations
- But what about if the order of data matters?
- **Q: Examples?**

How can one predict what comes next,
given an existing sequence of data?

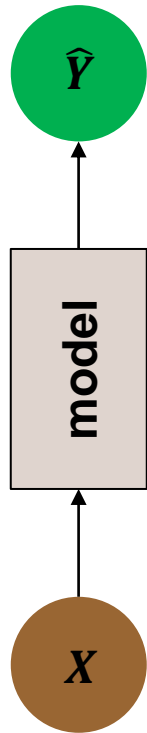
The **pony** ate the **hat** of my **father**

The **hat** of my **father** ate the **pony**

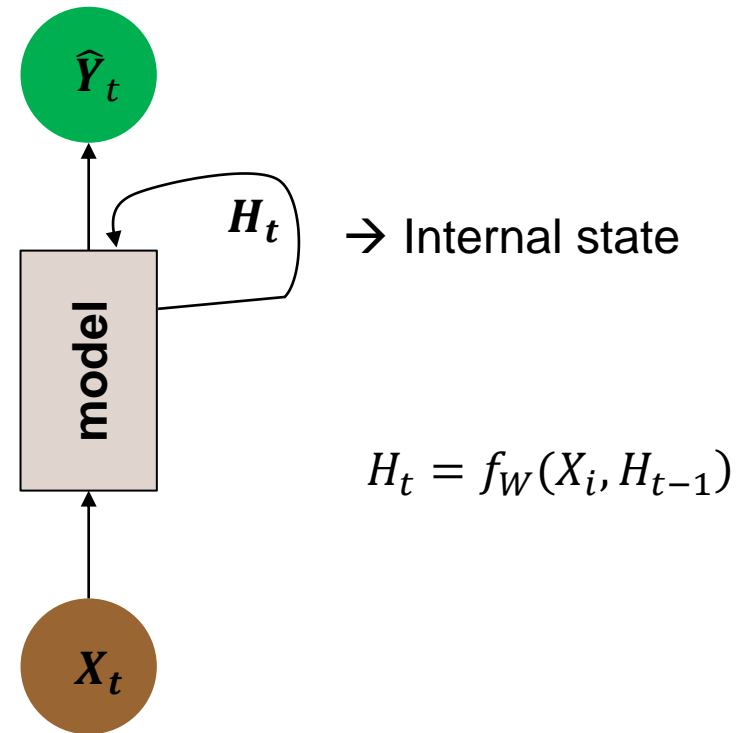
The of **hat** **pony** ate the **father** my

Recurrent Neural Networks

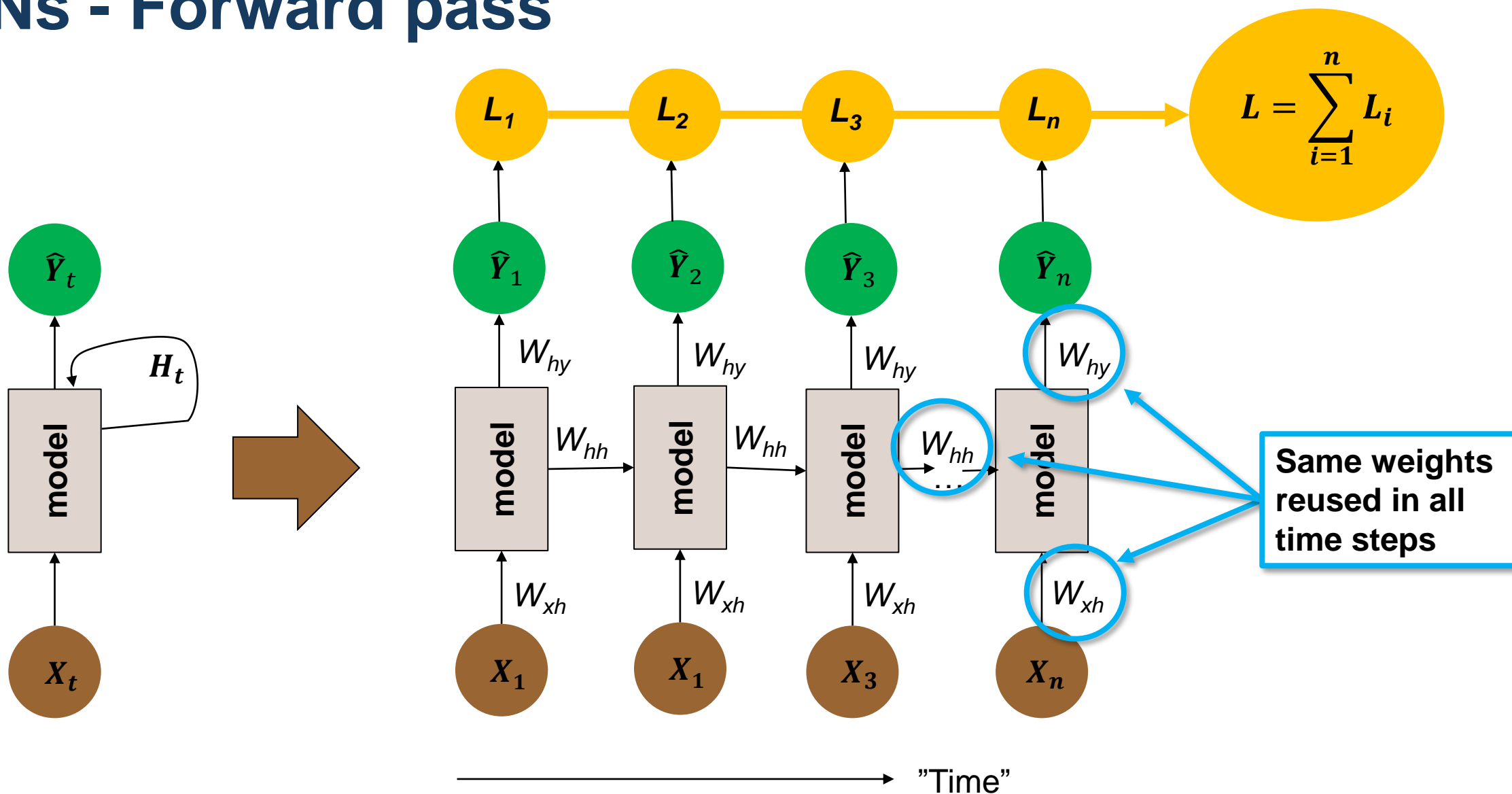
Standard "feed-forward"
network



Recurrent network

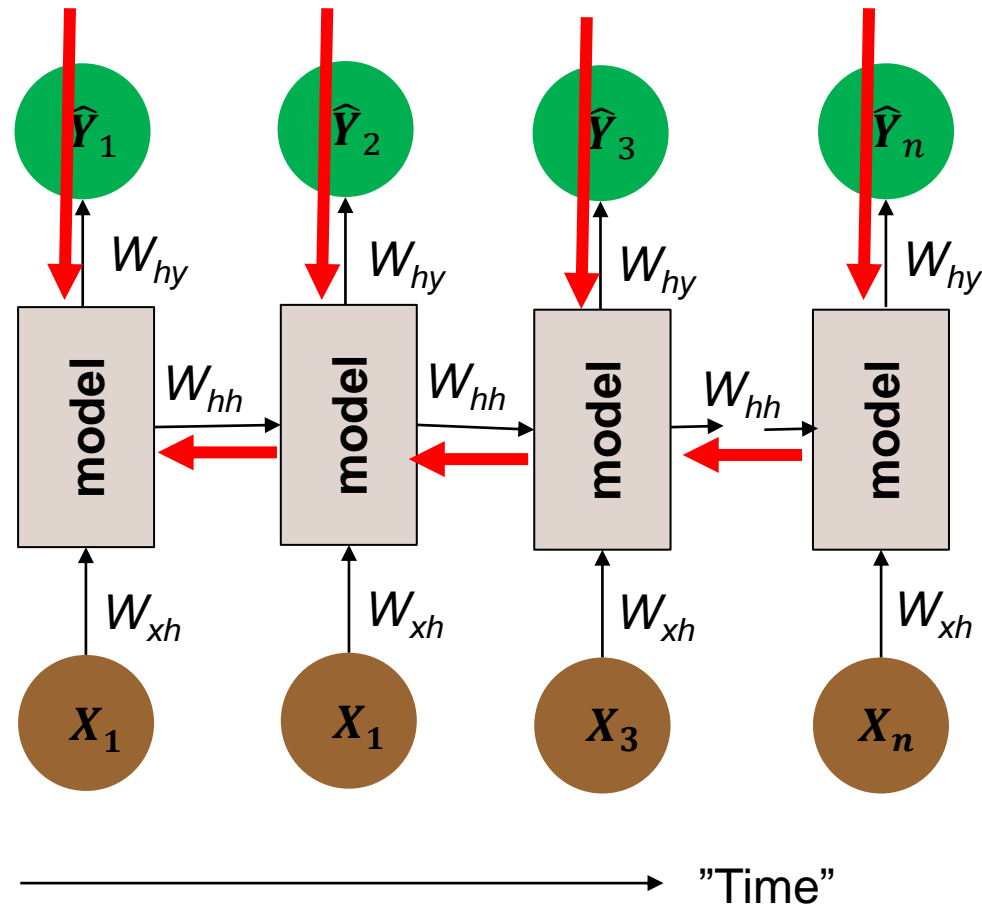


RNNs - Forward pass



RNNs - Backpropagation through time

- Training by backpropagation involves propagating gradients backwards through each time step.
- Each step consists of a matrix multiplication with W_{hh} and gradient of activation function.
- **Q: Potential issues?**



Vanishing gradient problem

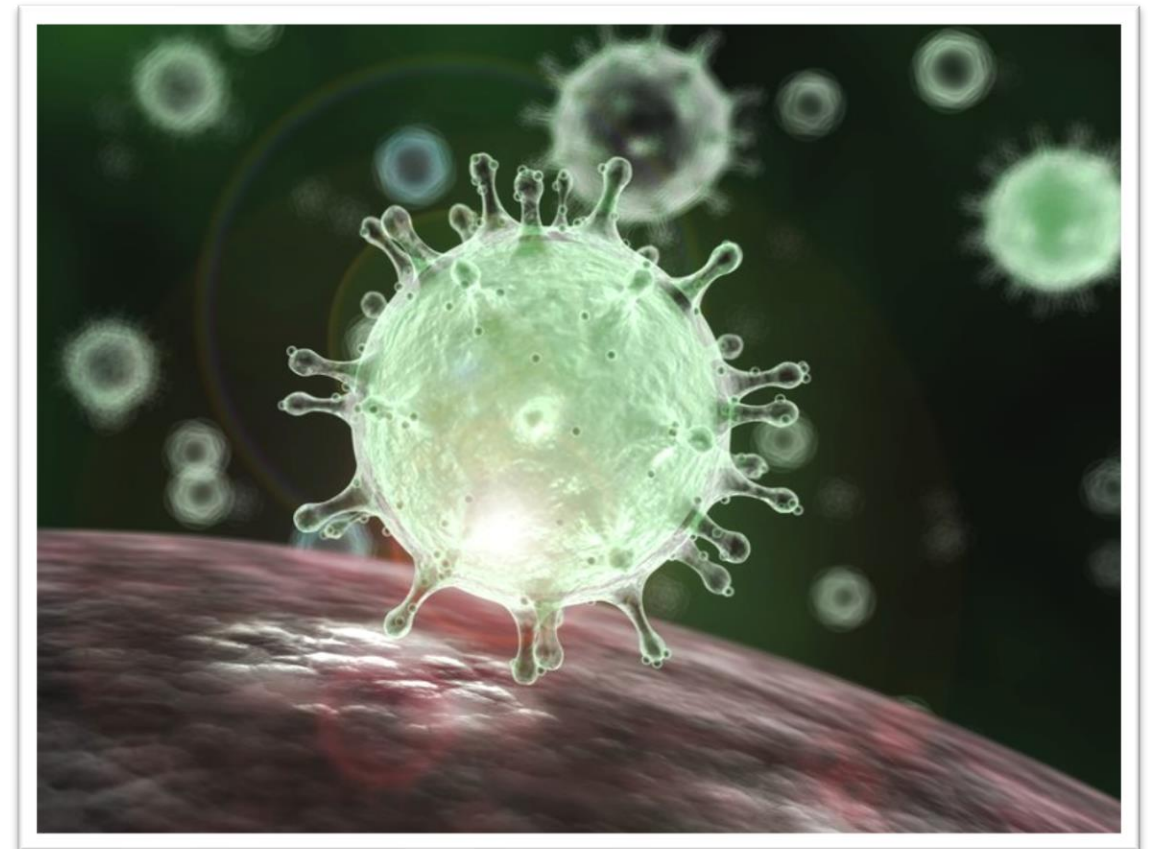
- Vanishing gradient → biasing to capture recent (short-term) dependencies

Short-term

The **virus** gave me the ____.

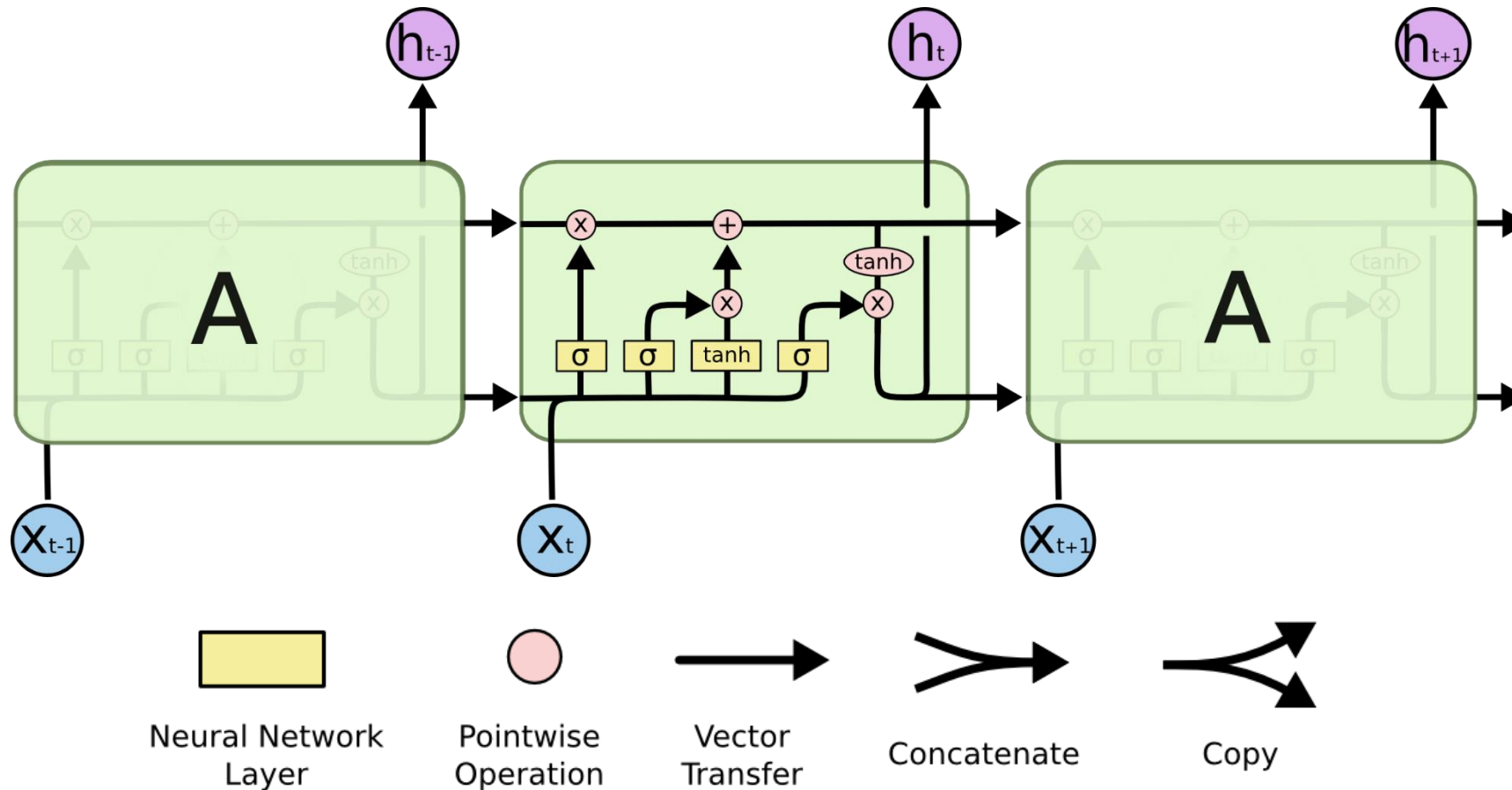
Long-term

A **virus** is an infective agent that typically consists of a nucleic acid molecule in a protein coat and can give a person the ____.



Long Short Term Memory Networks

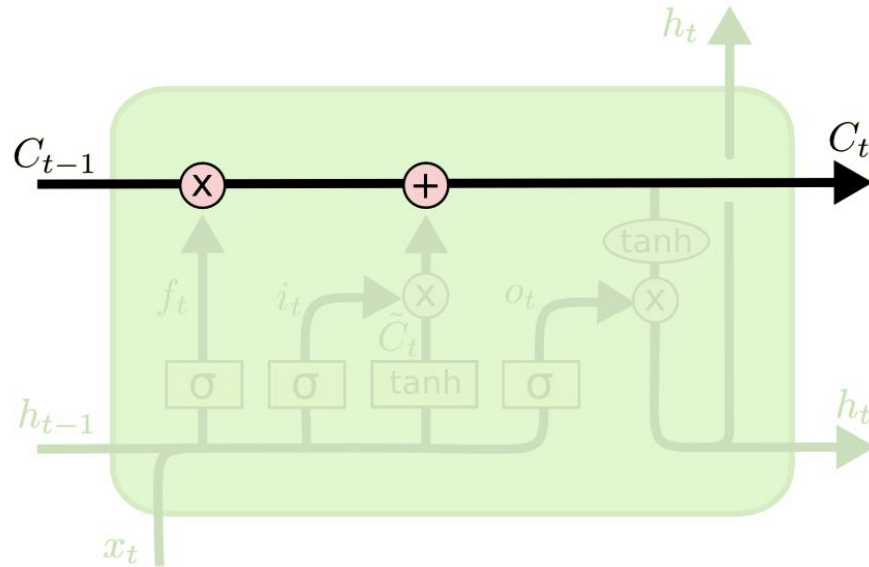
- LSTM: Recurrent units that are good for tracking long-term dependencies



Important features of LSTMs

Cell state

with easy information transfer through time



- Only elementwise multiplication and addition
- Avoids vanishing gradient problem!

Example:

"Mike likes his boat, but ..."

Subject

Gender

Singular/plural

Verb

Tense

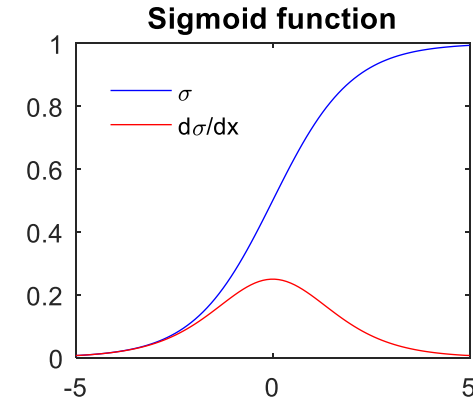
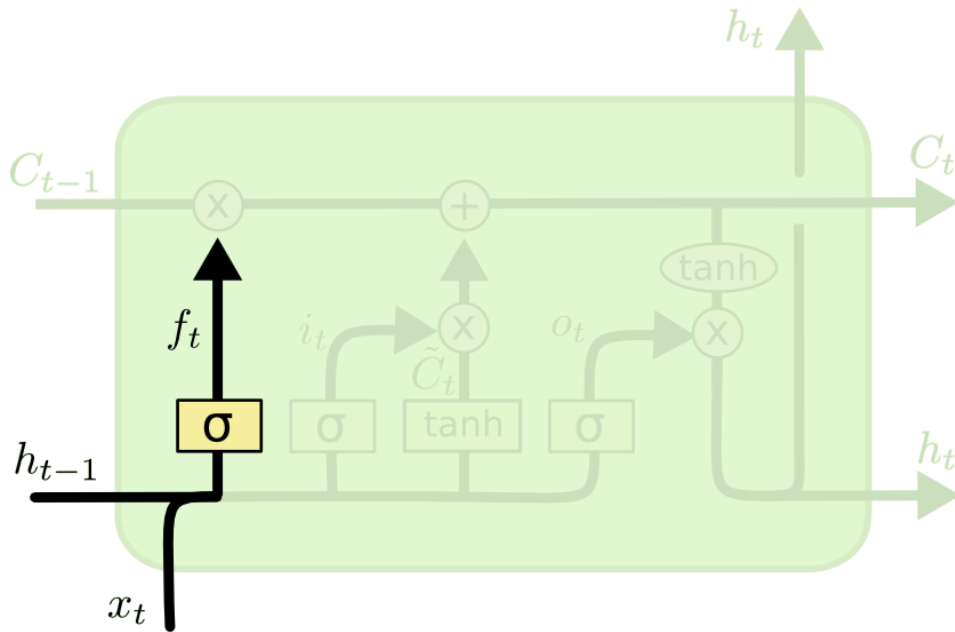
Inflection

Object

Type

Singular/plural

1. What to forget



$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

"Mike likes his boat, but Maria..."

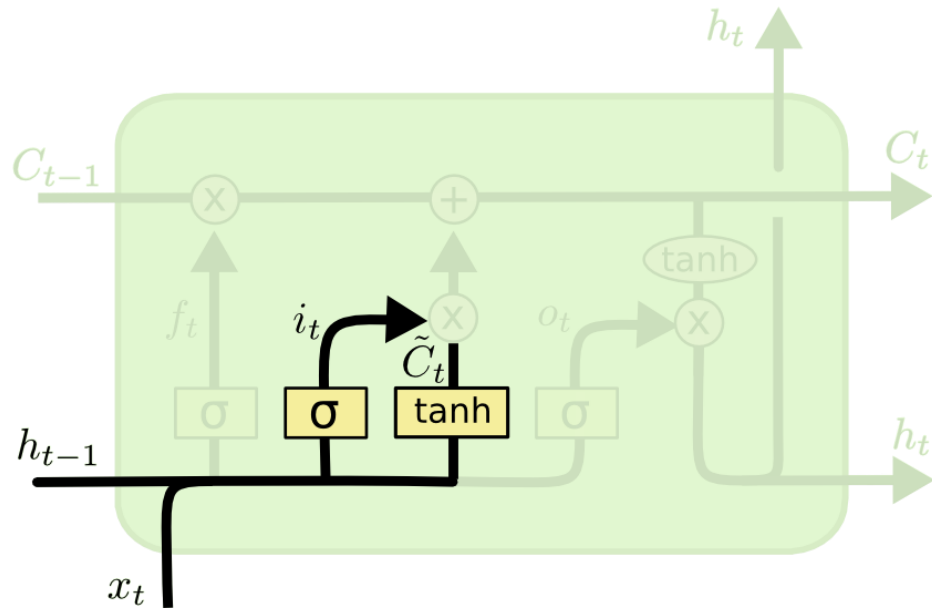
The word "Mike" is highlighted in yellow, "boat" in red, and "Maria" in pink. An arrow labeled x_t points to the word "Maria".

New subject,
Forget Mike?

An arrow points from this text box up to the word "Maria" in the sentence above.

- f_t outputs 0...1 for each element in C_{t-1} , based on h_{t-1} and x_t
- $f_t = 0 \rightarrow$ completely forget
- $f_t = 1 \rightarrow$ completely remember

2. What to update



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

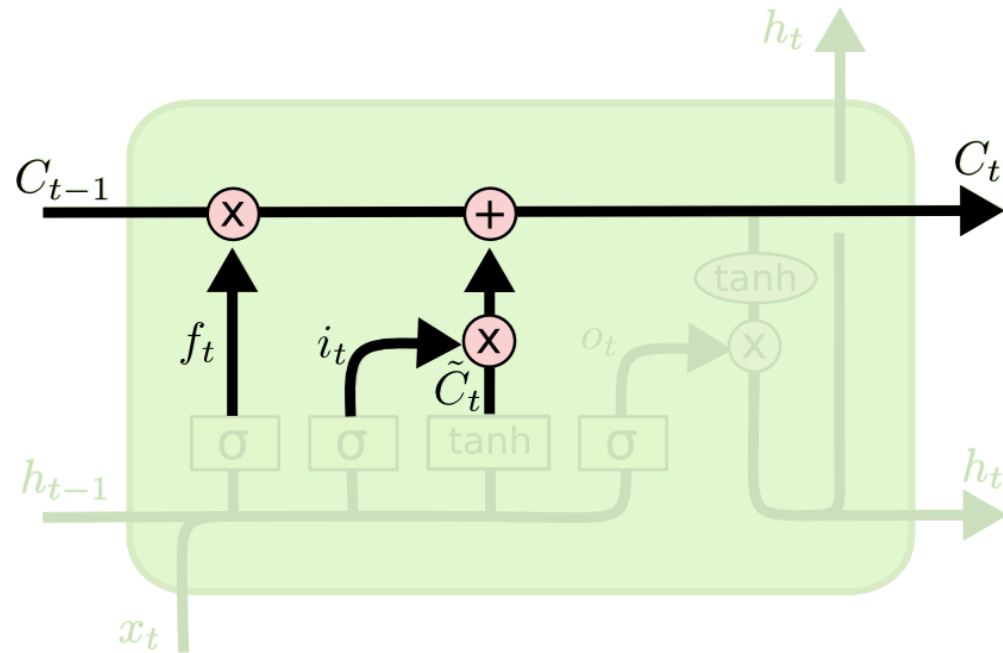
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

"Mike likes his boat, but Maria..."

Candidate for
new Subject

- tanh-layer creates new candidates for cell state C (like in regular RNN)
- σ -layer decides which of these to update (factor 0..1)

Updating cell state



- Cell state is updated elementwise! (computationally cheap)

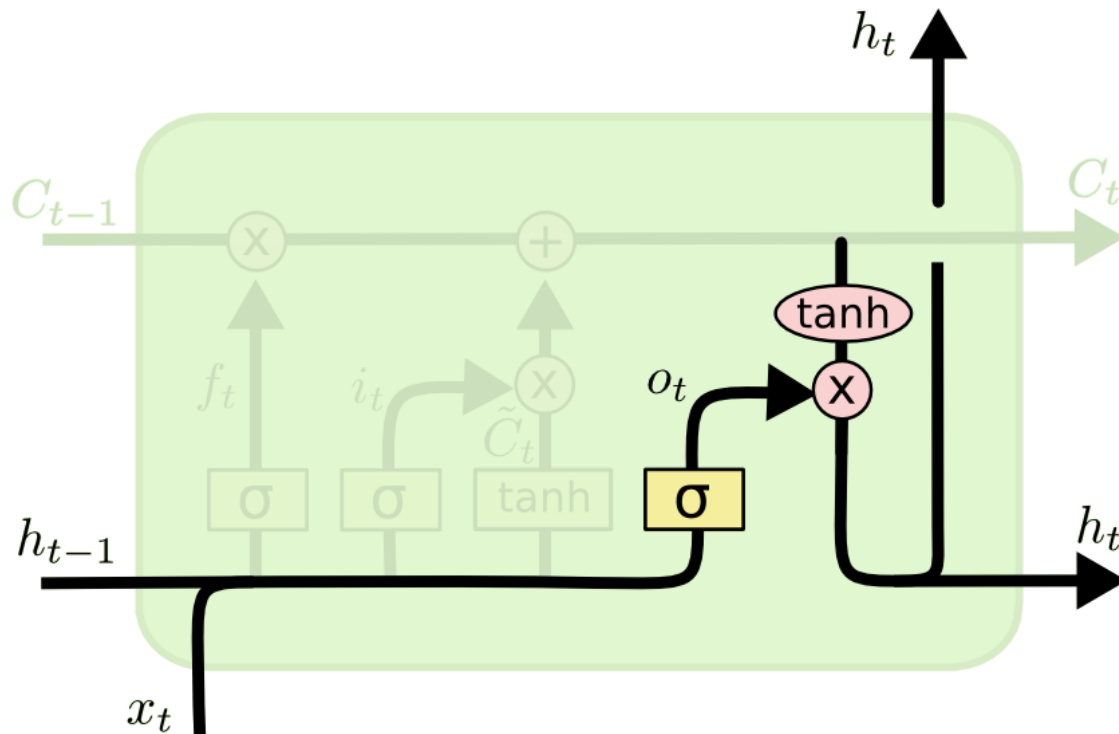
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

"Mike likes his boat, but Maria..."

Update C by
forgetting Mike

Update C with
new Subject

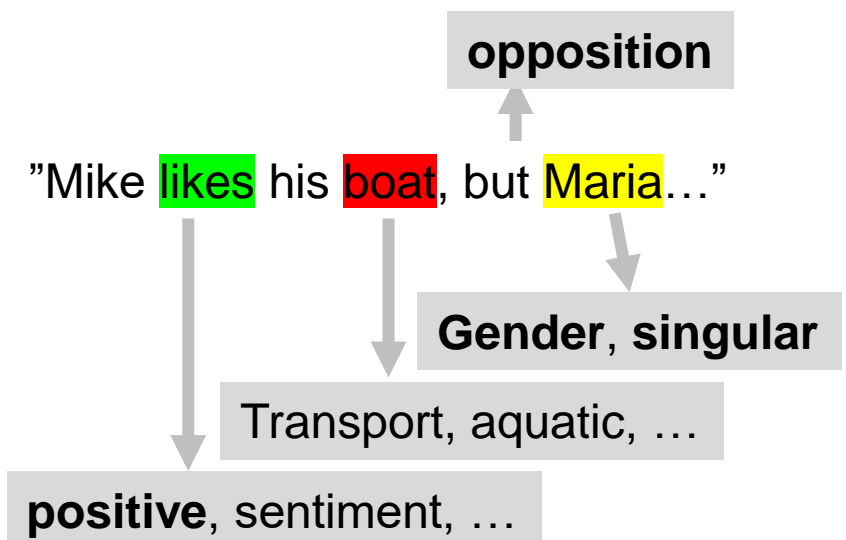
Outputting to next timestep



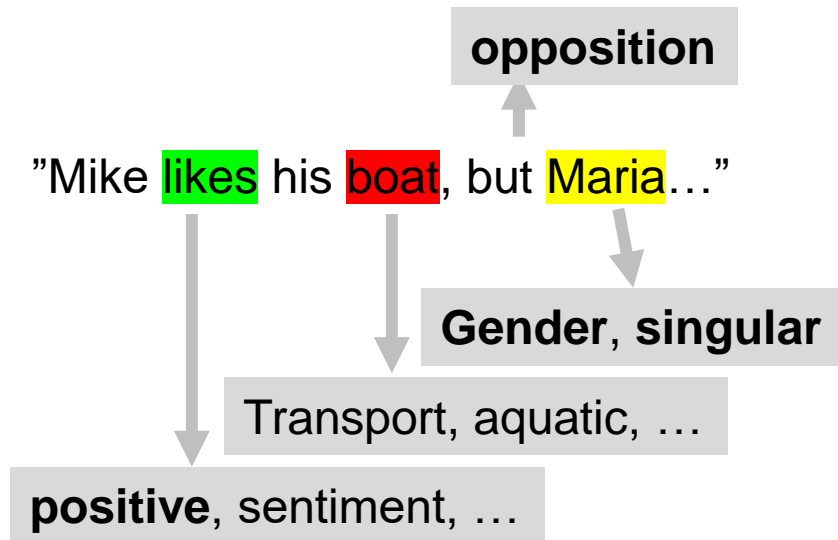
- \tanh -function bounds output between -1 and 1
- σ -layer decides what to output (factor 0..1)
 - Relevant to what should come next (verb etc)

$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$



2 min Exercise – What could come next?



Example – Music generation

The diagram illustrates a process for music generation. On the left, a table lists characters and their corresponding indices and hex values. The top section of the table is labeled "Source".

Dec	Hx	Oct	Html	Chr
88	58	130	X	X
89	59	131	Y	Y
90	5A	132	Z	Z
91	5B	133	[[
92	5C	134	\	\
93	5D	135]]
94	5E	136	^	^
95	5F	137	_	_

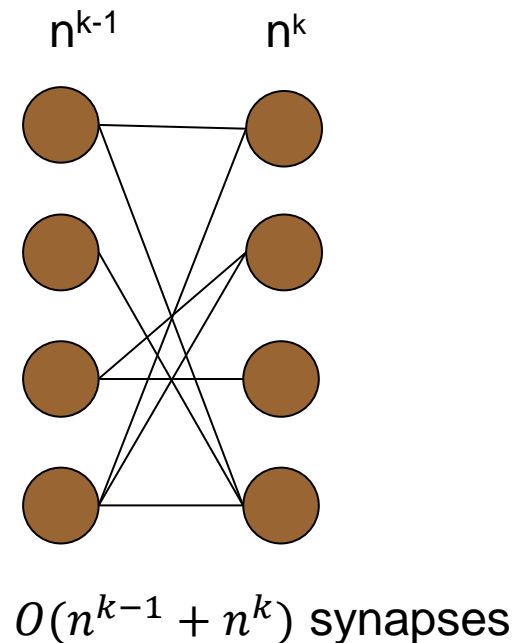
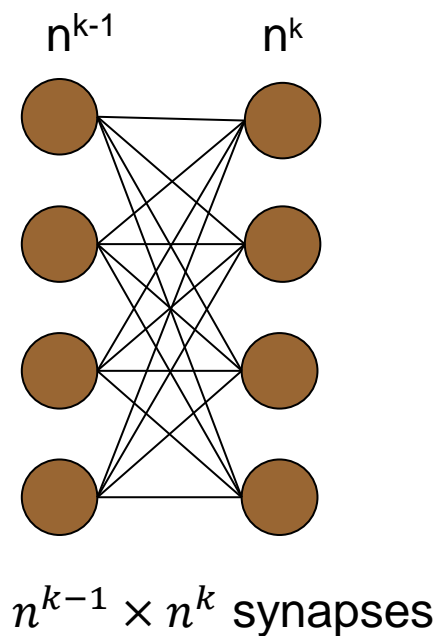
Below this, another section of the table is shown, with the first row labeled "Dec Hx Oct Html Chr".

Dec	Hx	Oct	Html	Chr
96	60	140	`	`
97	61	141	a	a
98	62	142	b	b
99	63	143	c	c
100	64	144	d	d
101	65	145	e	e
102	66	146	f	f
103	67	147	g	g
104	68	150	h	h

On the right, a word cloud features the word "Yaeh" in large black letters. Below it, the numbers 89, 97, 101, and 104 are written in blue. Blue arrows point from these numbers to the corresponding rows in the character set table. A musical staff with a treble clef is shown, with a red circle around the first four notes (A, C, E, G) and the text "MAJOR CHORD IN FIRST INVERSION!" written in red next to it.

Sparse and Locally connected networks

- Typical parameter size in Neural Network: 1-10 million
 - Work-intensive
 - Risk of overfitting
 - Scalability issues!



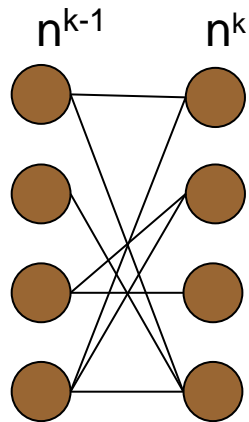
- Can layers with $O(n^{k-1} + n^k)$ connections perform as good?
- What type of connectivity would then be ideal?

Sparse Network topologies

Erdos-Renyi network (random graph)

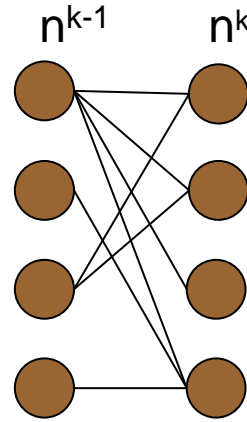
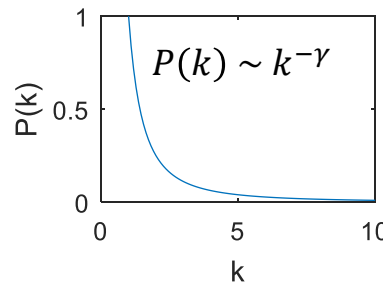
$$P(W_{ij}^k) = \frac{\varepsilon(n^k + n^{k-1})}{n^k n^{k-1}}, \varepsilon \in \mathbb{R}^+$$

$$\rightarrow |W^k| = \varepsilon(n^k + n^{k-1}) \text{ if } \varepsilon \ll n^k, n^{k-1}$$

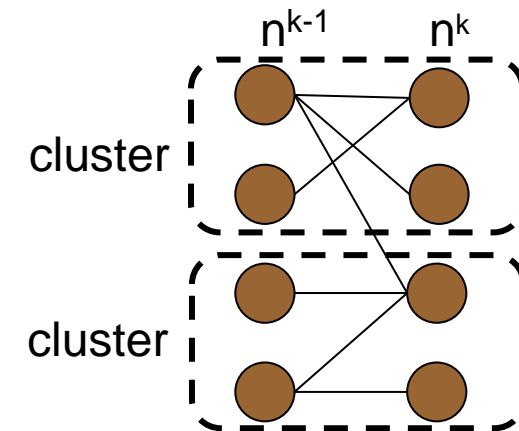
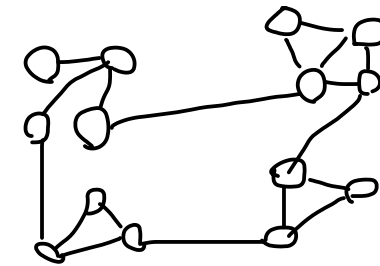


$O(n^{k-1} + n^k)$ synapses

Scale-free network (few nodes with many connections)



Small world network (most nodes are not neighbors, but connected in a few steps)



Example – Adaptive sparse connectivity

- Sparse-Evolutionary-Training (SET) algorithm

Algorithm 1: SET pseudocode

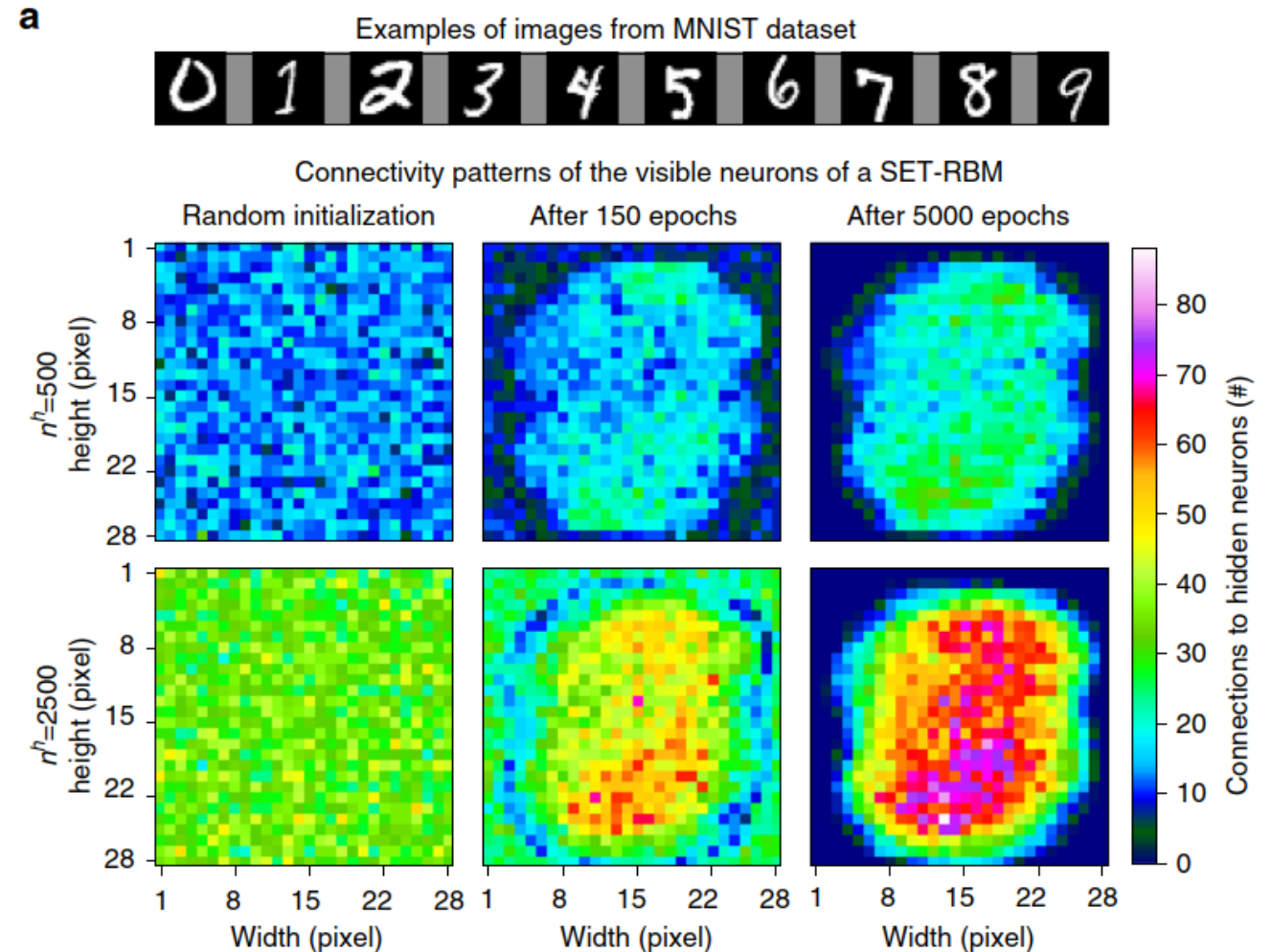
```
1  %Initialization;  
2  initialize ANN model;  
3  set  $\epsilon$  and  $\zeta$ ;  
4  for each bipartite fully-connected (FC) layer of the ANN do  
5  |   replace FC with a Sparse Connected (SC) layer having a Erdős-Rényi topology given by  $\epsilon$  and Eq.1;  
6  end  
7  initialize training algorithm parameters;  
8  %Training;  
9  for each training epoch  $e$  do  
10 |   perform standard training procedure;  
11 |   perform weights update;  
12 |   for each bipartite SC layer of the ANN do  
13 |   |   remove a fraction  $\zeta$  of the smallest positive weights;  
14 |   |   remove a fraction  $\zeta$  of the largest negative weights;  
15 |   |   if  $e$  is not the last training epoch then  
16 |   |   |   add randomly new weights (connections) in the same amount as the ones removed previously;  
17 |   |   end  
18 |   end  
19 end
```

Handwritten annotations:

- 1. points to line 5: replace FC with a Sparse Connected (SC) layer having a Erdős-Rényi topology given by ϵ and Eq.1;
- 2. points to line 10: perform standard training procedure;
- 3. points to line 13: remove a fraction ζ of the smallest positive weights;
- 4. points to line 16: add randomly new weights (connections) in the same amount as the ones removed previously;

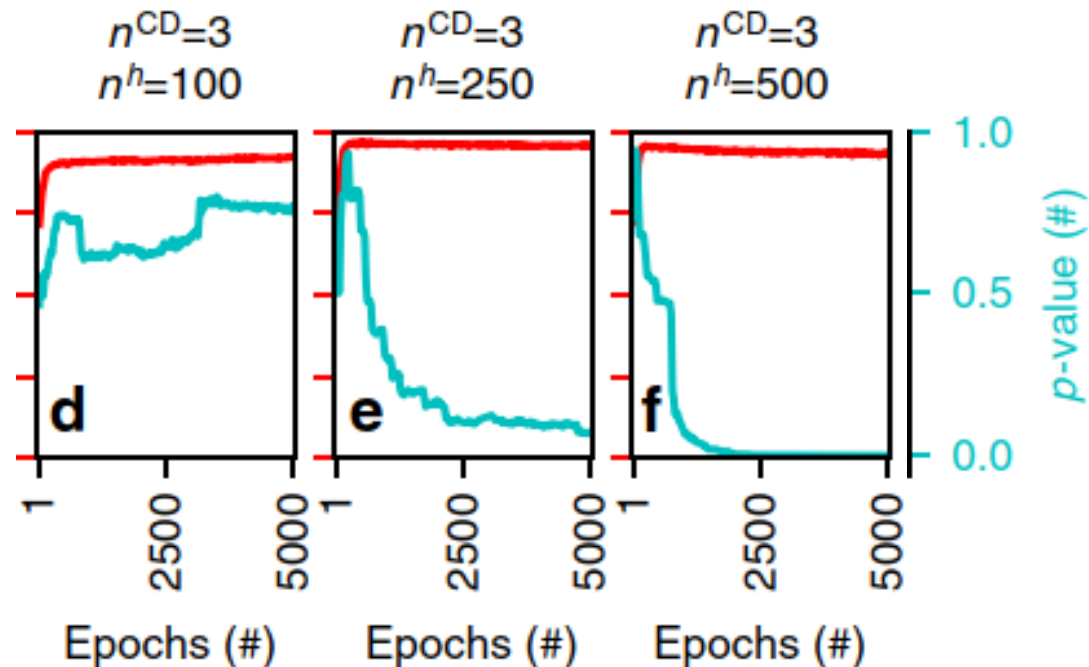
Clustering of connections

- Towards end of training:
 - Pixels in the center of the images (MNIST) have many connections
 - Pixels on the edges have few or zero connections
- Natural data-dependent clustering (scale-free network?)



Evolution towards scale-free network

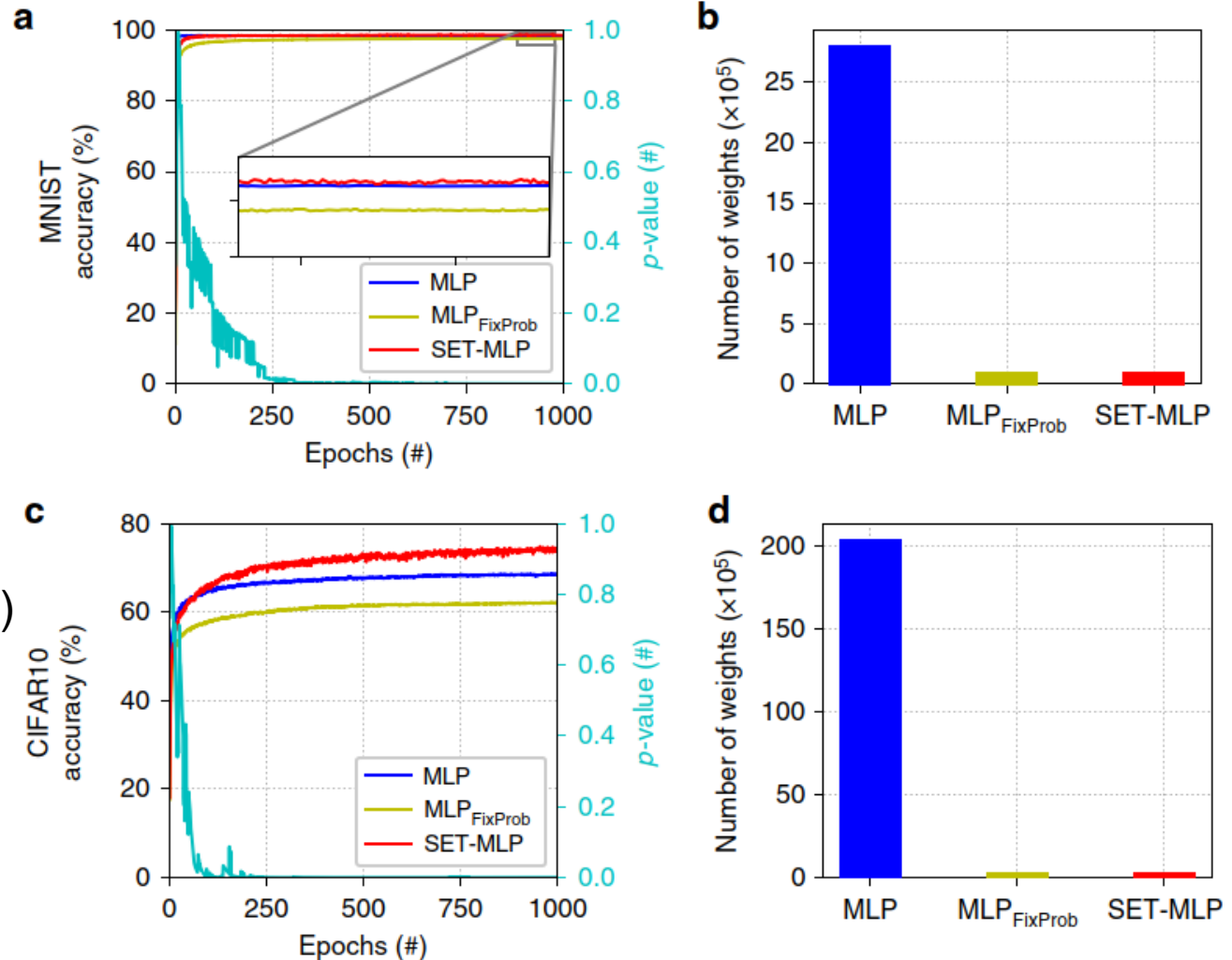
- Indeed: Adaptive sparsity evolves towards a scale-free network
- Starts from a binomial probability distribution (Erdos-Renyi)
- With enough neurons in a layer \rightarrow scale-free ($p < 0.05$)



Results compared to fully connected networks

- $\text{MLP}_{\text{FixProb}}$: Fixed sparse network (E-R topology)
- MLP: Fully connected network
- SET-MLP: Adaptively sparse network
- Fixed sparsity performs worst
 - 62% on CIFAR10
 - But with only 1.4% as many weights
- Adaptive sparsity outperforms all!
 - 75% on CIFAR10 (278K params)
 - 2nd best in literature has 74% (31M params)

	MLP	$\text{MLP}_{\text{FixProb}}$	SET-MLP
MNIST	98.55%	97.68%	98.74%
CIFAR10	68.70%	62.19%	74.84%



Possibilities with LCNs

- The sparsity enables:
 - MUCH larger models (1 billion neurons...) on supercomputers
 - Powerful tiny models trained directly on wearables
- Problem: GPU optimized for dense matrices, not for sparse..
 - More on GPU next lecture...

Summary

