

10-403 Recitation (1/17/20): Introduction to Deep Learning, CNN & RNN

Checklist for logistics

- Can you log into Piazza and Gradescope?
 - If you are not added, email to amritsin@andrew.cmu.edu with **your name** and **Andrew ID**
- How do you know when the first assignment will be released?

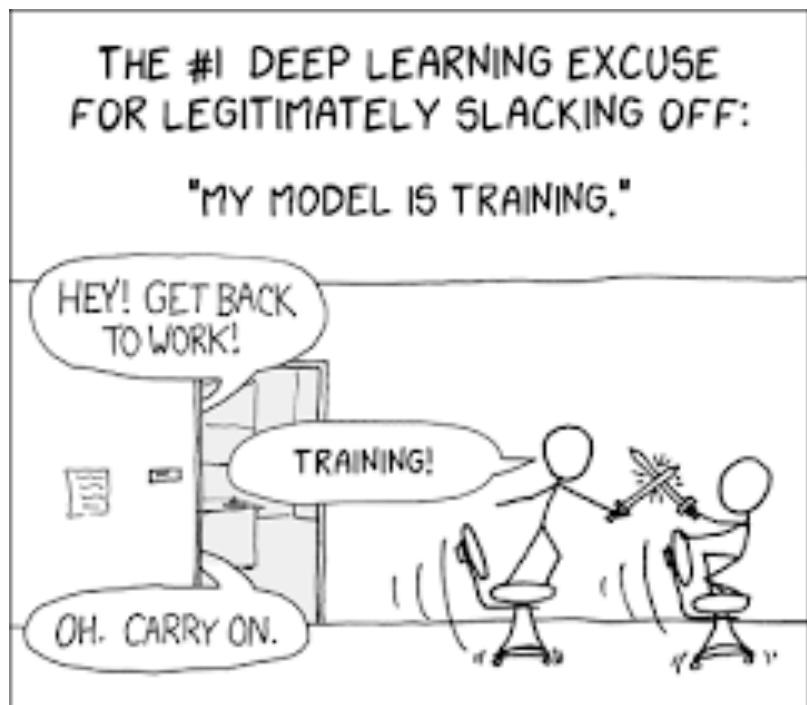
Questions to be answered in this recitation

- Why do we need deep learning?
- How do we tailor neural network to each data domain?
- Why are CNN and RNN important in real-world application?
- What is intuition behind CNN and RNN?

Disclaimer: Some of the material were borrowed from
10-707 (Deep Learning)

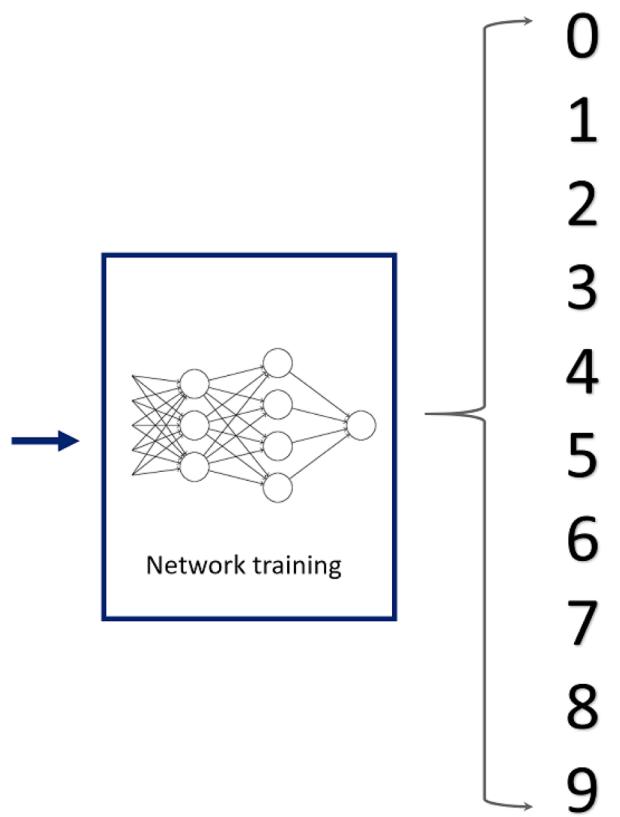
Deep Learning

- What is it?
- Why the heck is it so popular?



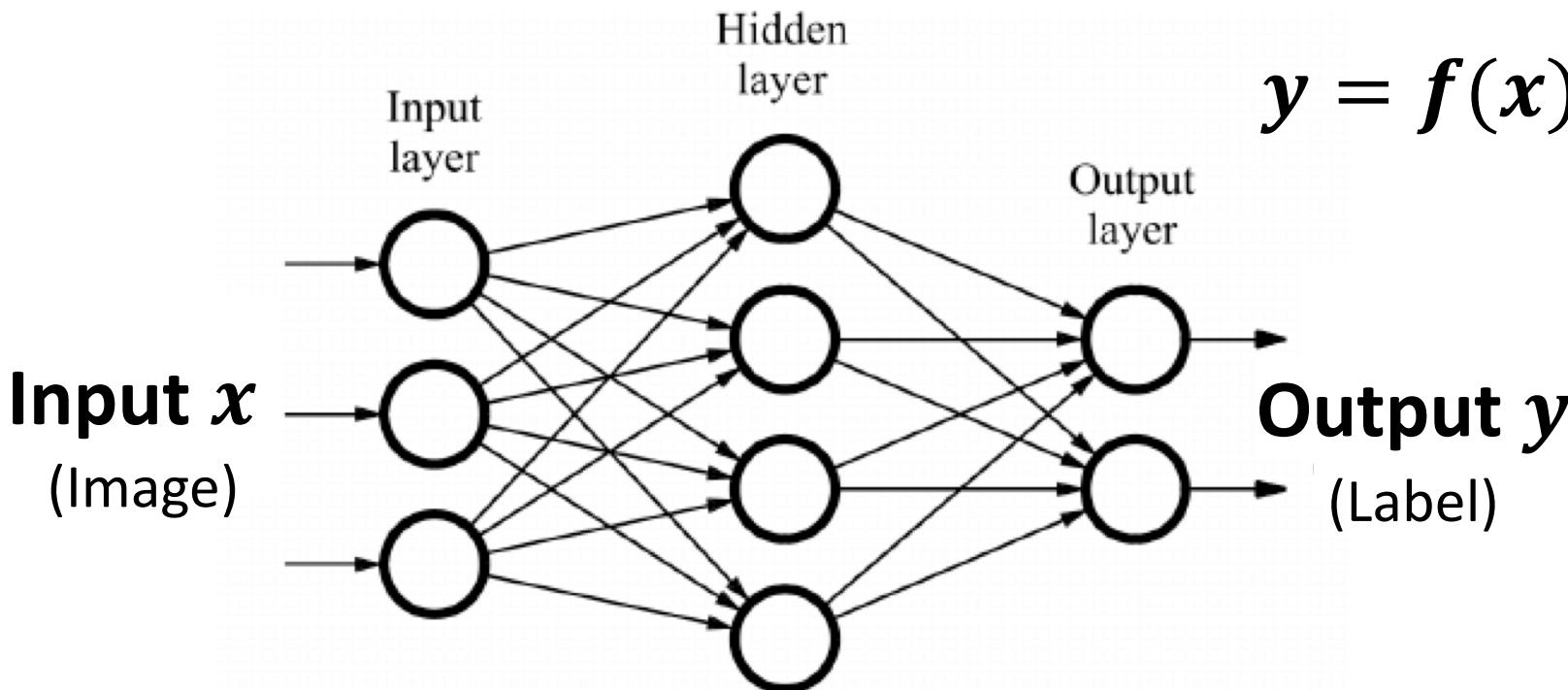
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9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

Data & Labels



Feed Forward Neural Network

- Feed-forward Neural Network = Multi-level perceptron = Fully-connected layer
- Hidden layer = Hidden units = Hidden nodes



$$\mathbf{y} = f(\mathbf{x})$$

$$a = b_h + W_{xh}x$$

$$h = \tanh(a)$$

$$o = b_o + W_{ho}h$$

$$y = \text{softmax}(o)$$

Deep Learning (= Deep Neural Network)

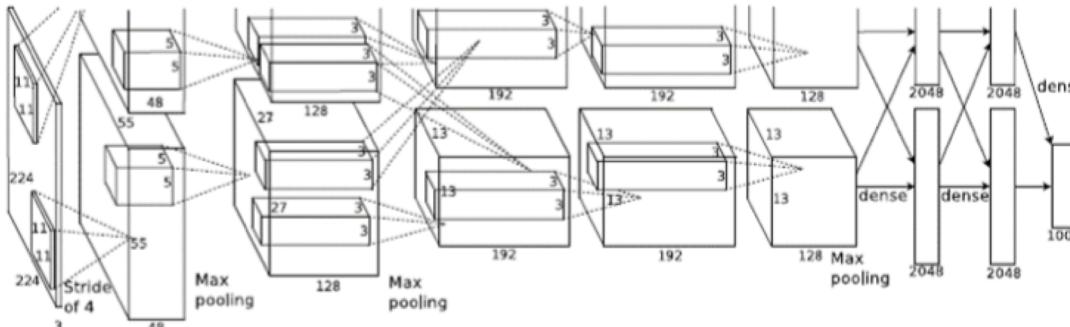
- Does feed forward neural network works for all data domain?
- What is key question to design neural network for each domain?

AI APPLICATIONS



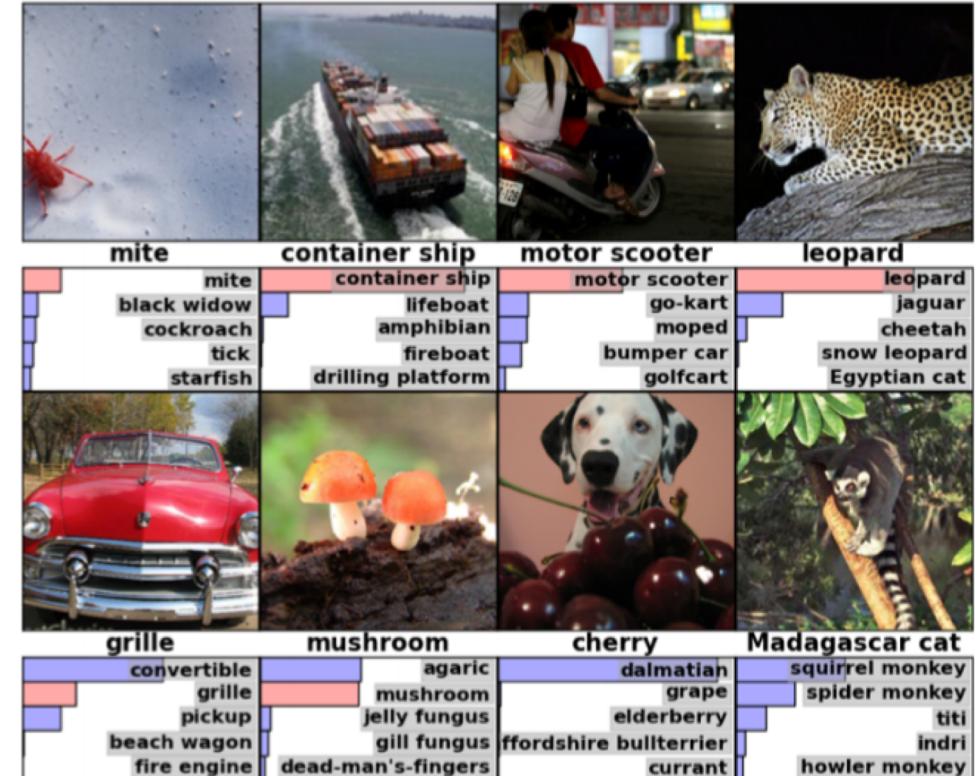
Convolutional Neural Network

Application: Computer Vision

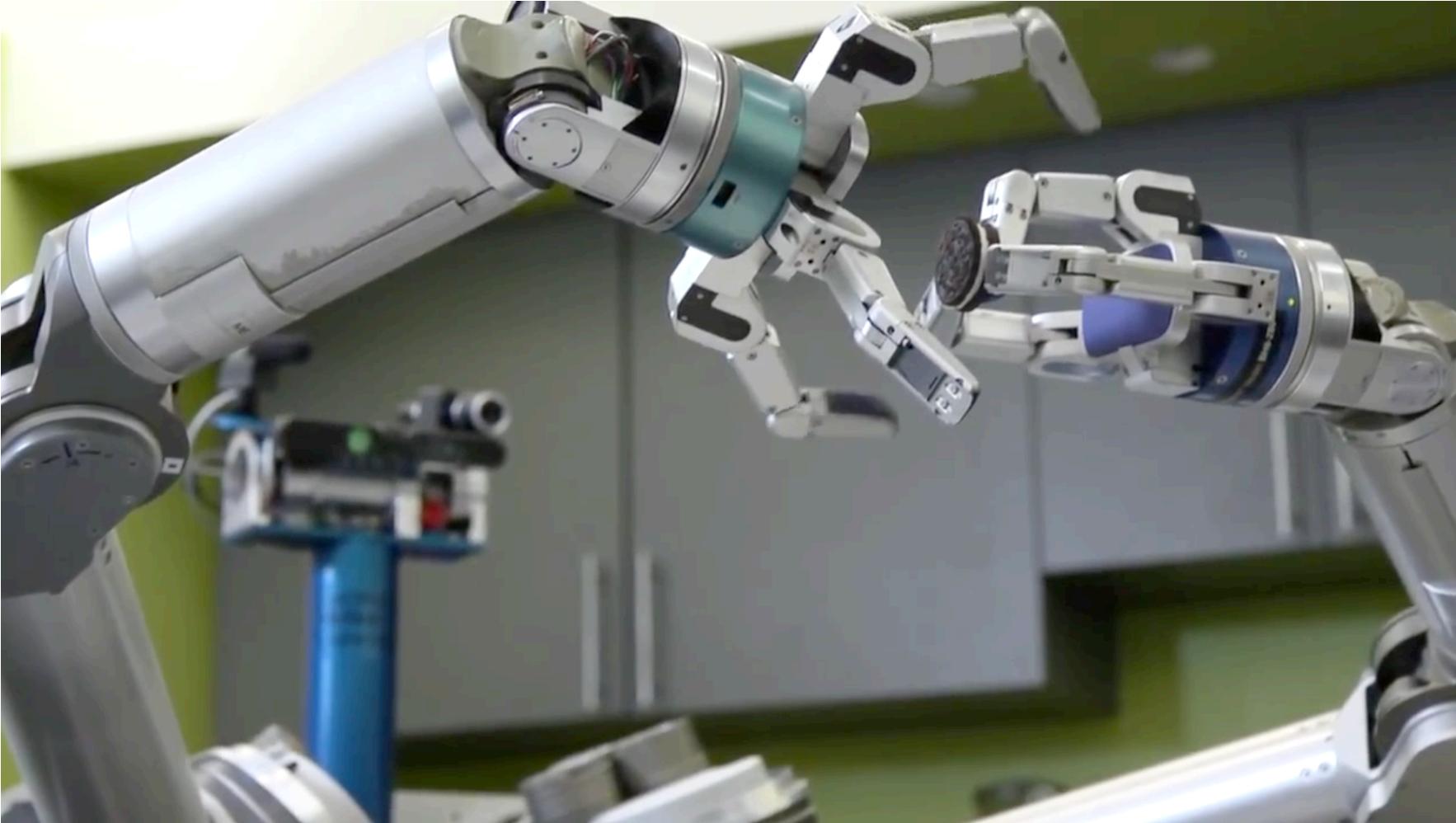


IMAGENET

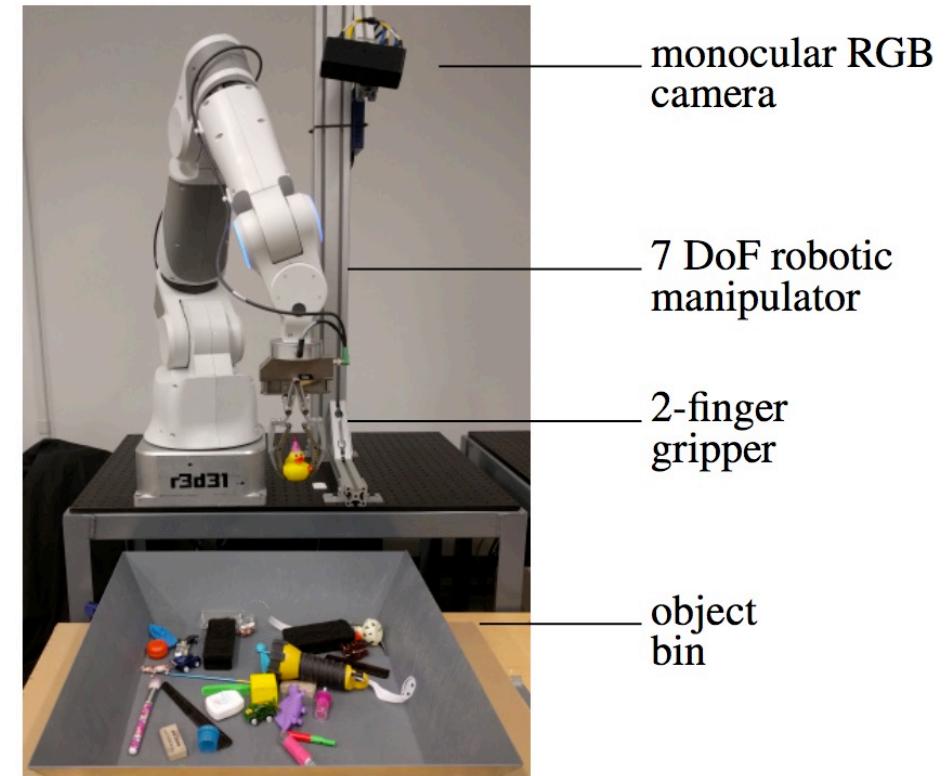
1.2 million training images
1000 classes



Application: Computer Vision



Application: Computer Vision



Shared Objective

- Design algorithms that can process **visual** data to accomplish a given task (e.g. object recognition)



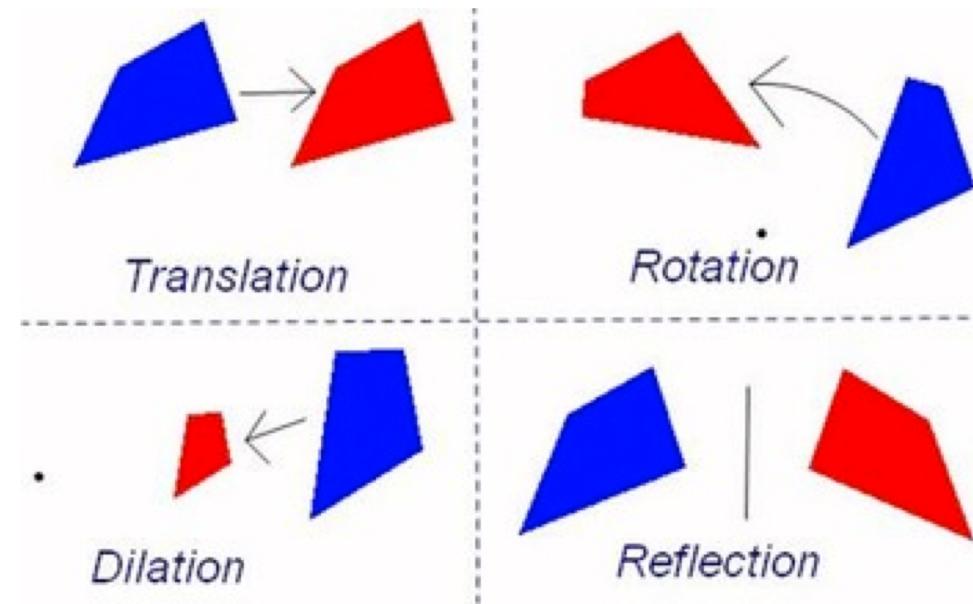
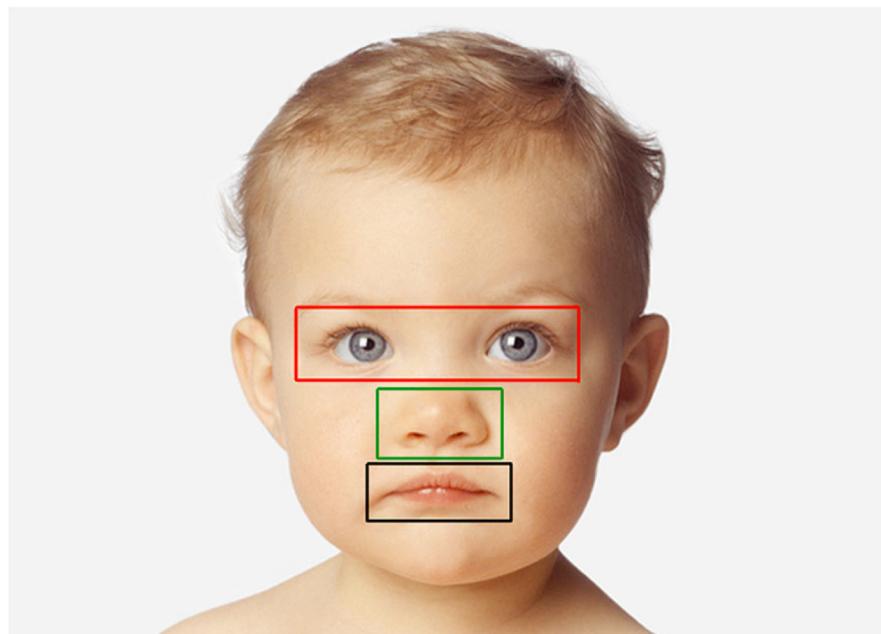
Shared Objective

- Design a **neural network** that can process **visual** data to accomplish a given task (e.g. object recognition)



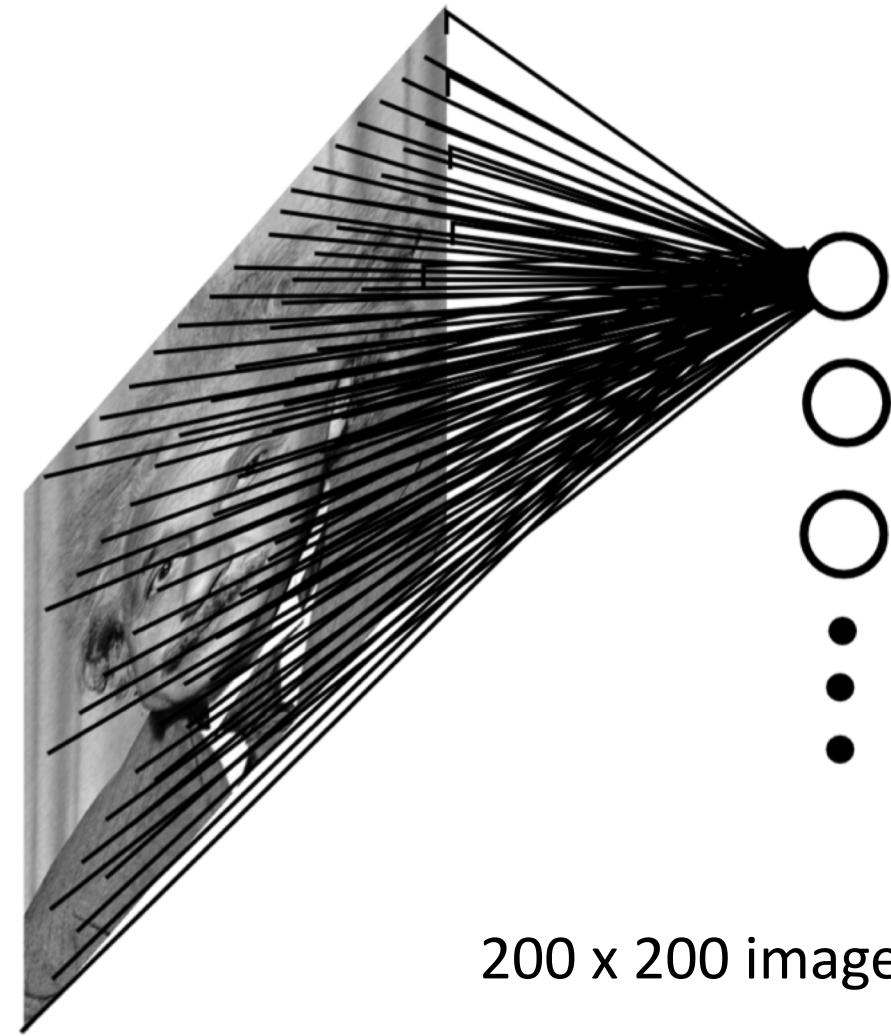
Characteristics of Visual Data

- **High-dimensional** inputs: $150 \times 150 \text{ pixels} = 22500 \text{ inputs}$ ($\times 3$ if RGB)
- 2D or 3D **Topology** of pixels (Spatial correlation)
- **Invariant properties** to certain variations: translation, rotation, etc.



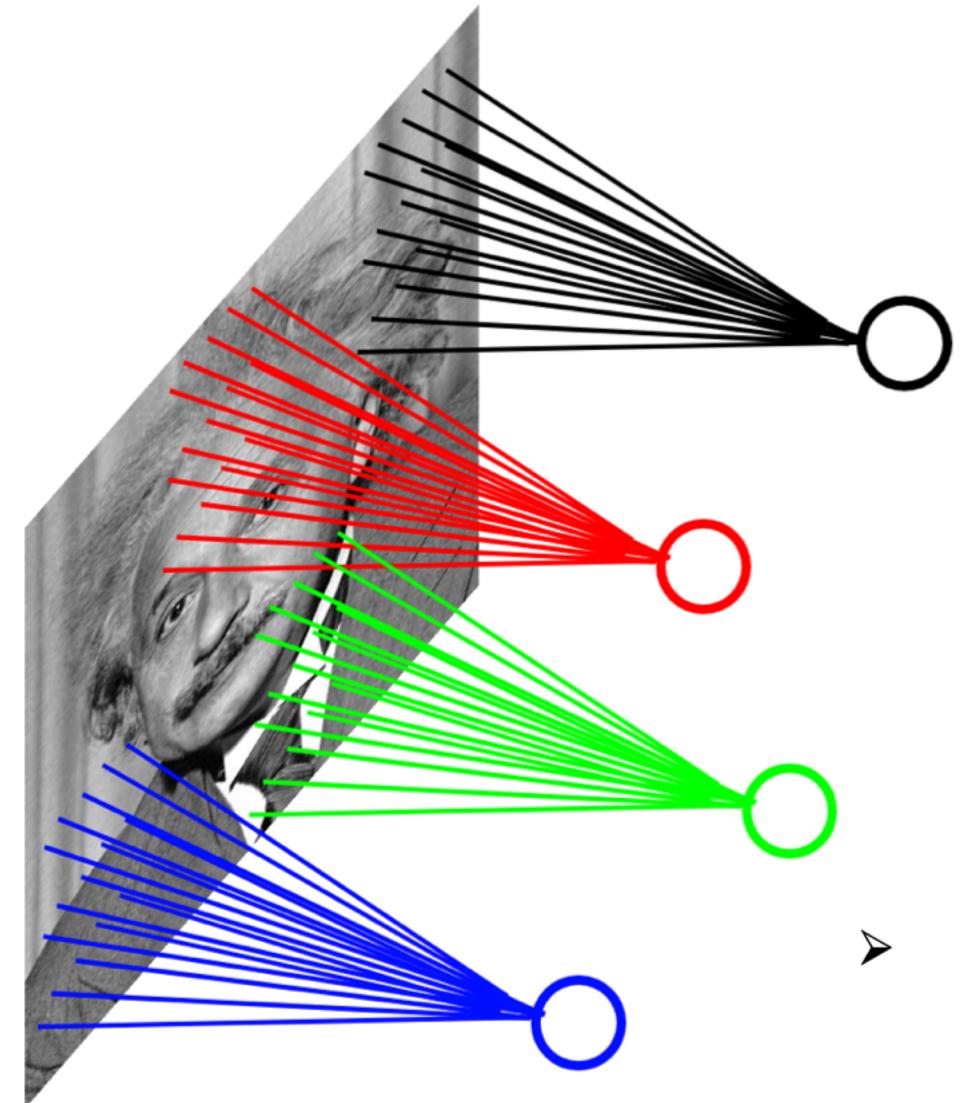
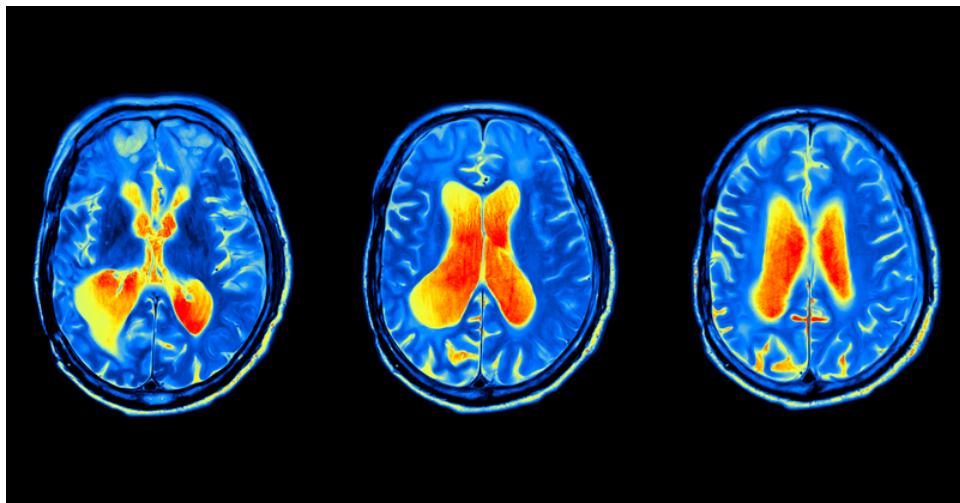
Naïve approach

- Assume grayscale image and 40K hidden units
- How many parameters?
- What's the problem?



Local Connectivity

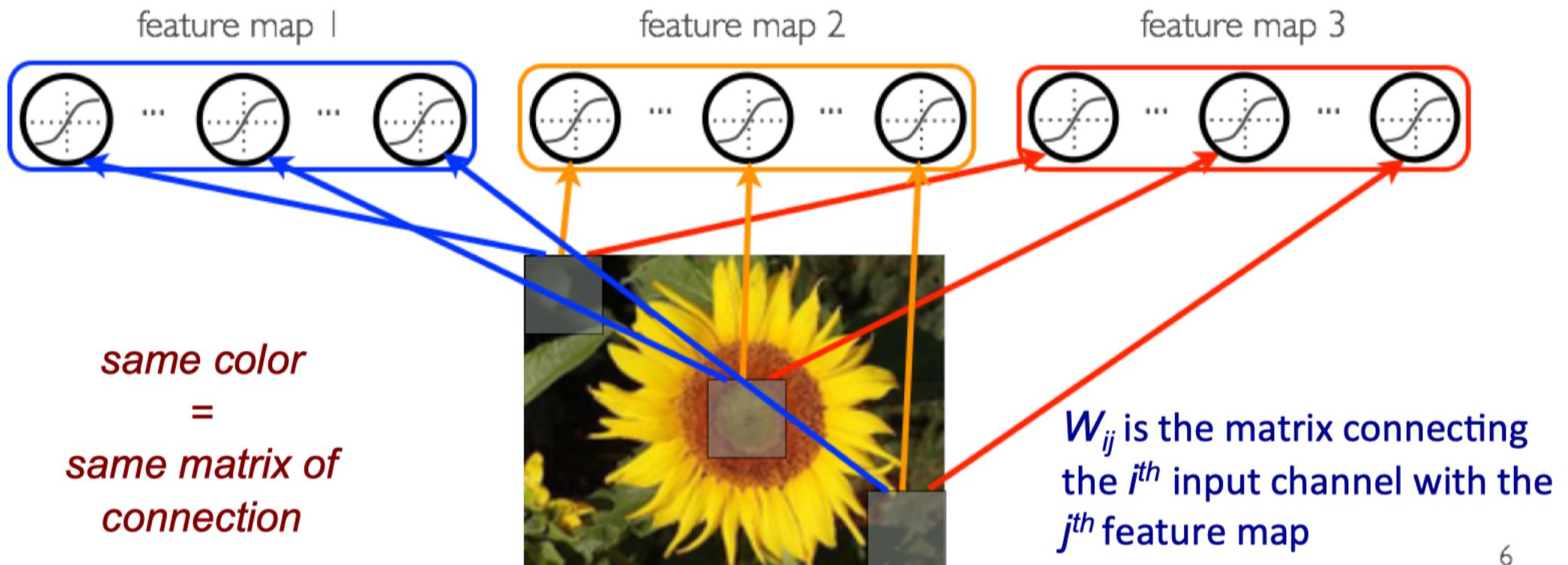
- Assume cell (= filter) size of 10x10 in a grid (20 x 20) and 40K hidden units
- How many parameters?
- What's the problem?



200 x 200 image

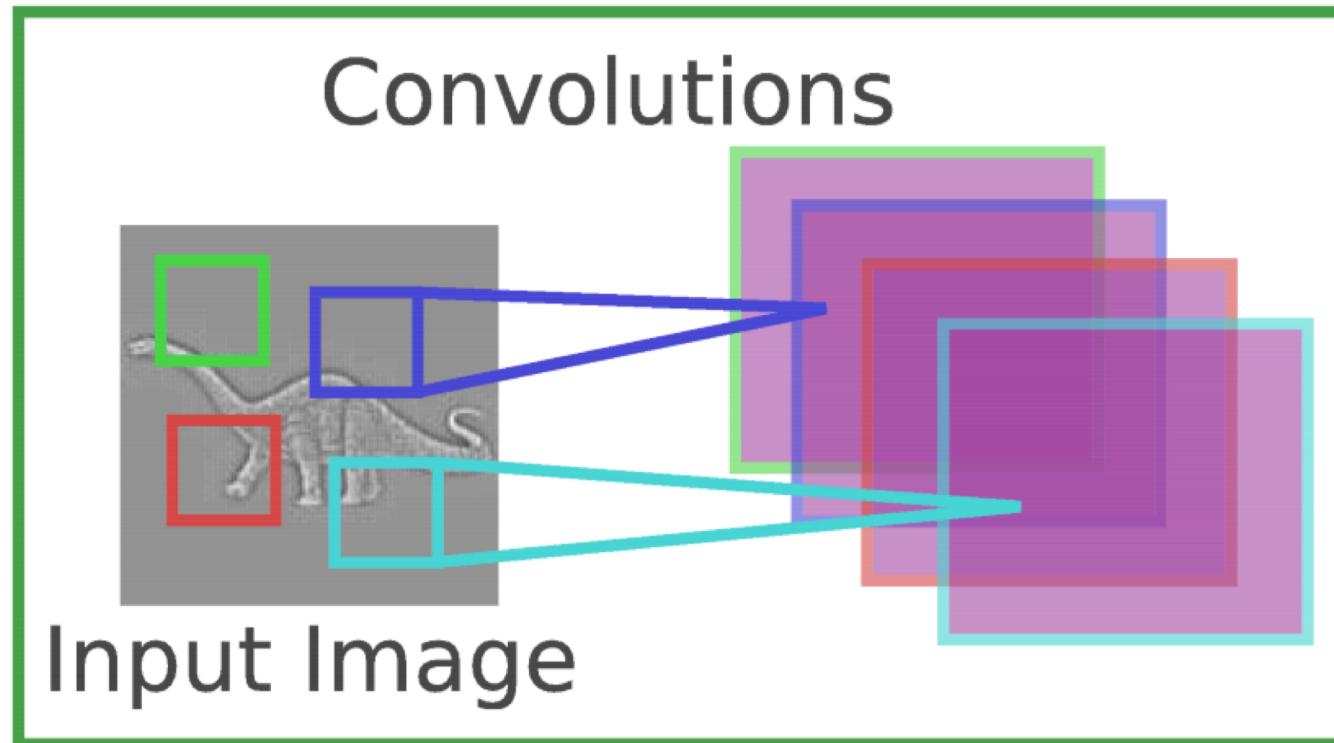
Parameter sharing

- Feedforward NN where hidden units are organized into multiple groups called **feature map** and hidden units in each feature map **share parameters**



Convolution

- Each feature map forms a 2D grid of features computed through **convolution** with a **sliding** filter over an image



Discrete Convolution

- The convolution of an image x with a kernel k is computed as follows:

$$(x * k)_{ij} = \sum_{pq} x_{i+p, j+q} k_{r-p, r-q}$$

- Example:

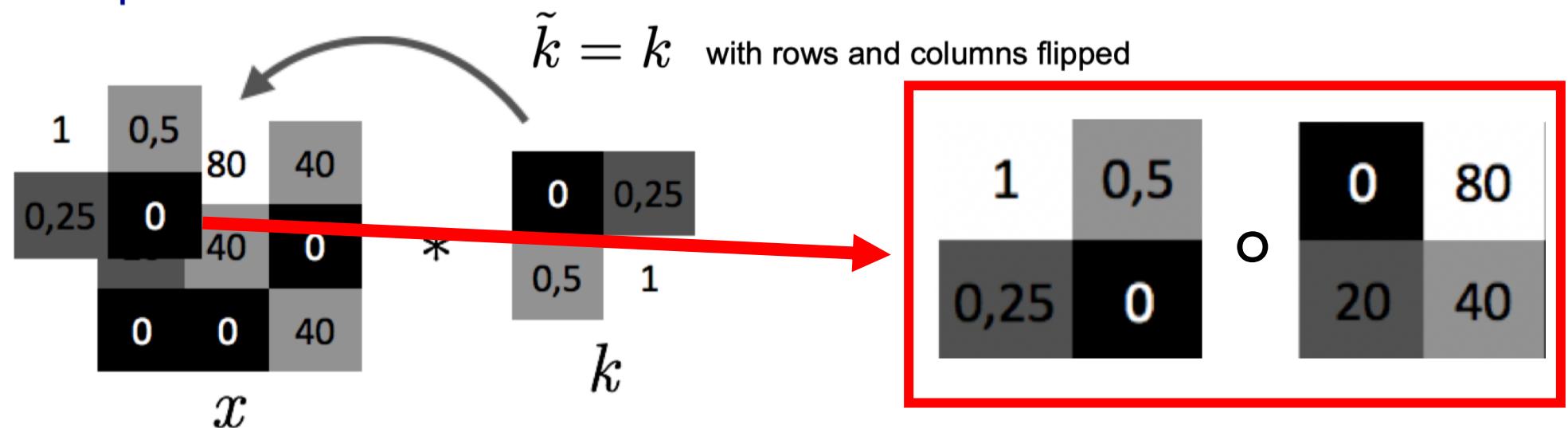
$$\begin{array}{ccc|c} 0 & 80 & 40 & \\ \hline 20 & 40 & 0 & \\ 0 & 0 & 40 & \\ \hline x & & & \end{array} * \begin{array}{cc} 0 & 0,25 \\ \hline 0,5 & 1 \\ k & \end{array} =$$

Discrete Convolution

- The convolution of an image x with a kernel k is computed as follows:

$$(x * k)_{ij} = \sum_{pq} x_{i+p, j+q} k_{r-p, r-q}$$

- Example:

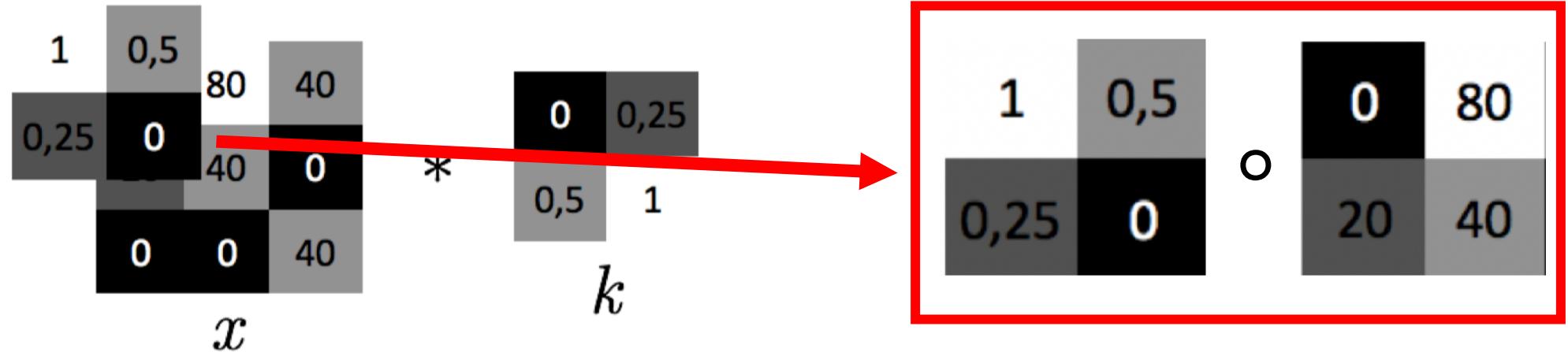


Discrete Convolution

- The convolution of an image x with a kernel k is computed as follows:

$$(x * k)_{ij} = \sum_{pq} x_{i+p,j+q} k_{r-p,r-q}$$

- Example: $1 \times 0 + 0.5 \times 80 + 0.25 \times 20 + 0 \times 40 =$



Discrete Convolution

- The convolution of an image x with a kernel k is computed as follows:

$$(x * k)_{ij} = \sum_{pq} x_{i+p,j+q} k_{r-p,r-q}$$

- Example:** $1 \times 0 + 0.5 \times 80 + 0.25 \times 20 + 0 \times 40 = 45$

The diagram shows the convolution of a 3x3 input image x with a 2x2 kernel k . The input x has values [1, 0.5, 80; 0.25, 0, 40; 0, 0, 40]. The kernel k has values [0, 0.25; 0.5, 1]. The result of the convolution is 45.

$$\begin{matrix} 1 & 0,5 & 80 \\ 0,25 & 0 & 40 \\ 0 & 0 & 40 \end{matrix} \quad * \quad \begin{matrix} 0 & 0,25 \\ 0,5 & 1 \end{matrix} = 45$$

x k

Discrete Convolution

- The convolution of an image x with a kernel k is computed as follows:

$$(x * k)_{ij} = \sum_{pq} x_{i+p, j+q} k_{r-p, r-q}$$

- Example:** $1 \times 80 + 0.5 \times 40 + 0.25 \times 40 + 0 \times 0 = 110$

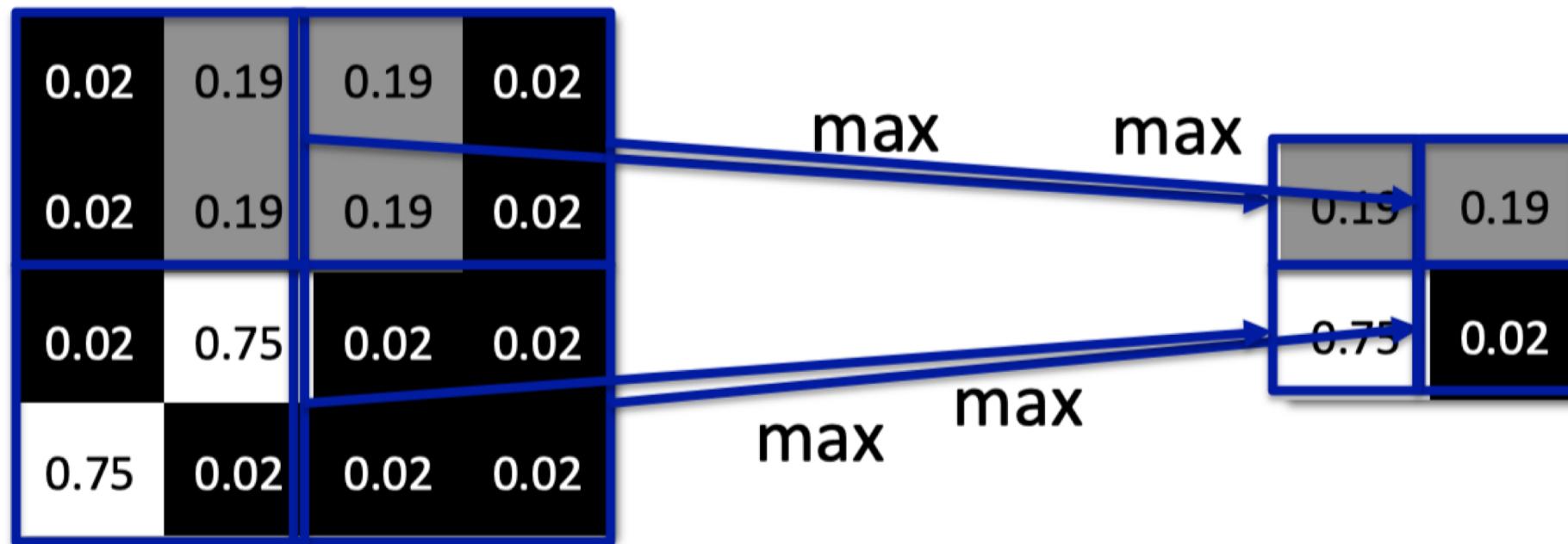
The diagram shows the convolution process between two 3x3 matrices. The image x (left) has values: top row [1, 0.5, 40], second row [0.25, 0, 0], bottom row [0, 0, 40]. The kernel k (right) has values: top row [0, 0.25], second row [0.5, 1]. The result of the convolution is shown as a 2x2 matrix: [45, 110]. The asterisk (*) indicates the convolution operation, and the equals sign (=) indicates the result.

$$\begin{matrix} 1 & 0,5 & 40 \\ 0,25 & 0 & 0 \\ 0 & 0 & 40 \end{matrix} \quad * \quad \begin{matrix} 0 & 0,25 \\ 0,5 & 1 \end{matrix} \quad = \quad \begin{matrix} 45 & 110 \end{matrix}$$

x k

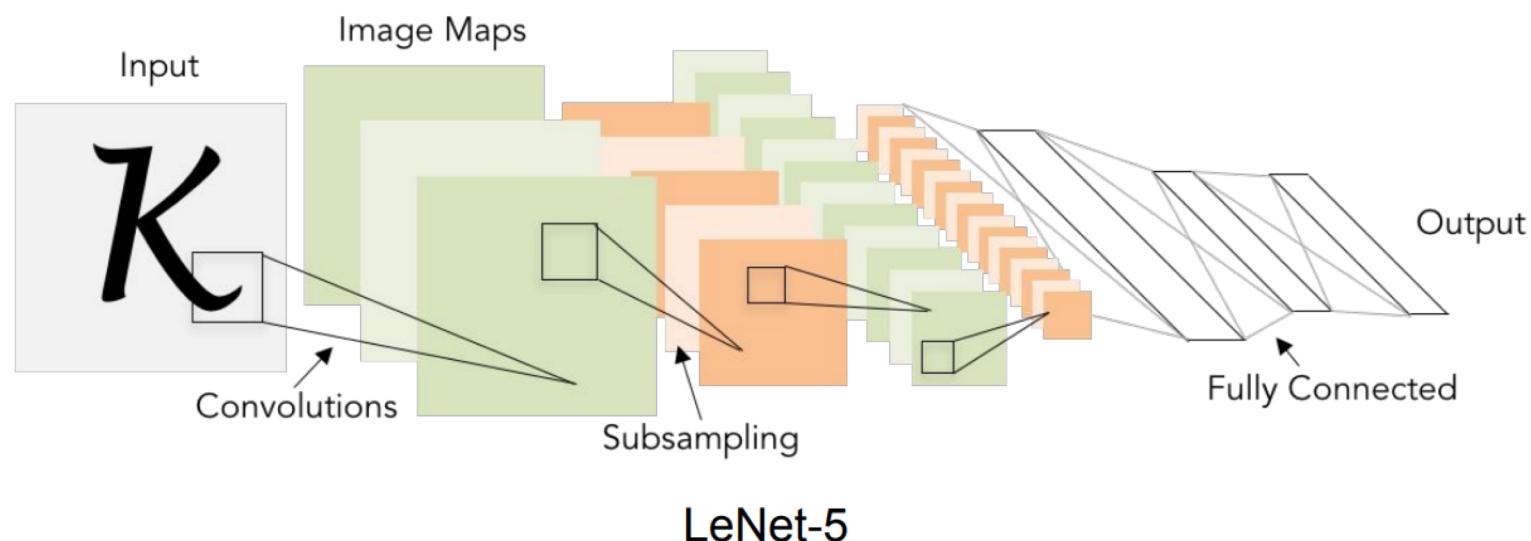
Pooling/Subsampling

- Pooling is performed in non-overlapping neighborhoods
- What are the advantages?



Convolutional Neural Network in a Nutshell

- Local connectivity
- Parameter sharing
- Convolution
- Pooling / subsampling hidden units



How does it work in mathematical detail?

- Deep learning textbook, Chapter 9
- http://cs231n.stanford.edu/slides/2019/cs231n_2019_lecture05.pdf

Recurrent Neural Network

Application: Sequential Data (Stock Market)

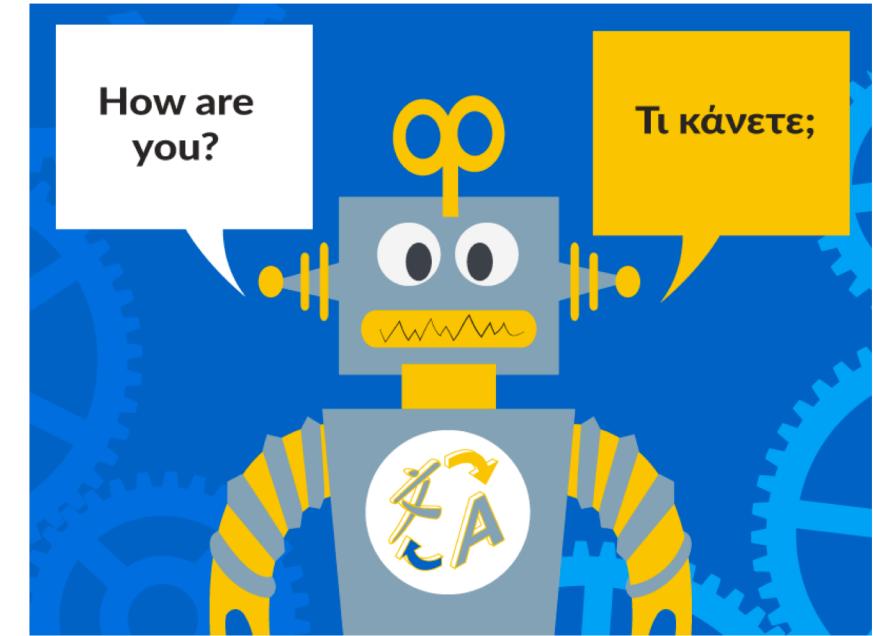


Application: Sequential Data (Language)

"I love this movie.
I've seen it many times
and it's still awesome."



"This movie is bad.
I don't like it at all.
It's terrible."



Shared Objective

- Design algorithms that can process **sequential** data to accomplish a given task (e.g. sentiment classification)

"I love this movie.
I've seen it many times
and it's still awesome."



"This movie is bad.
I don't like it at all.
It's terrible."



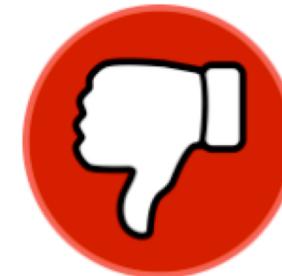
Shared Objective

- Design a **neural network** that can process **sequential** data to accomplish a given task (e.g. sentiment classification)

"I love this movie.
I've seen it many times
and it's still awesome."

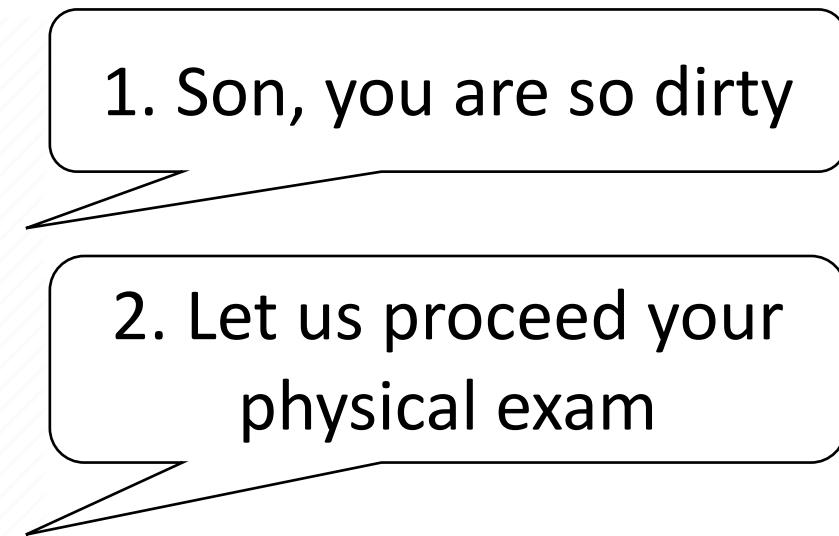


"This movie is bad.
I don't like it at all.
It's terrible."



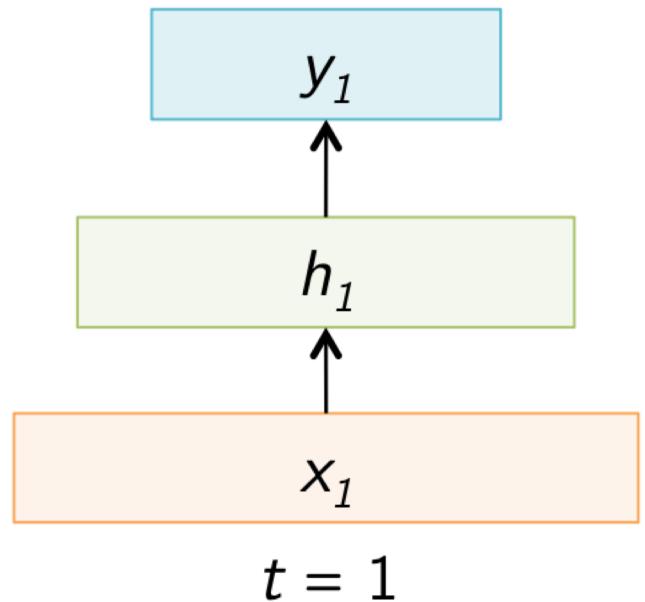
Characteristics of Sequential Data

- **Variable-length** input: “I like you” vs. “I think I like you”
- **Context** information: “Take your clothes off”
- **Variable-length** context window

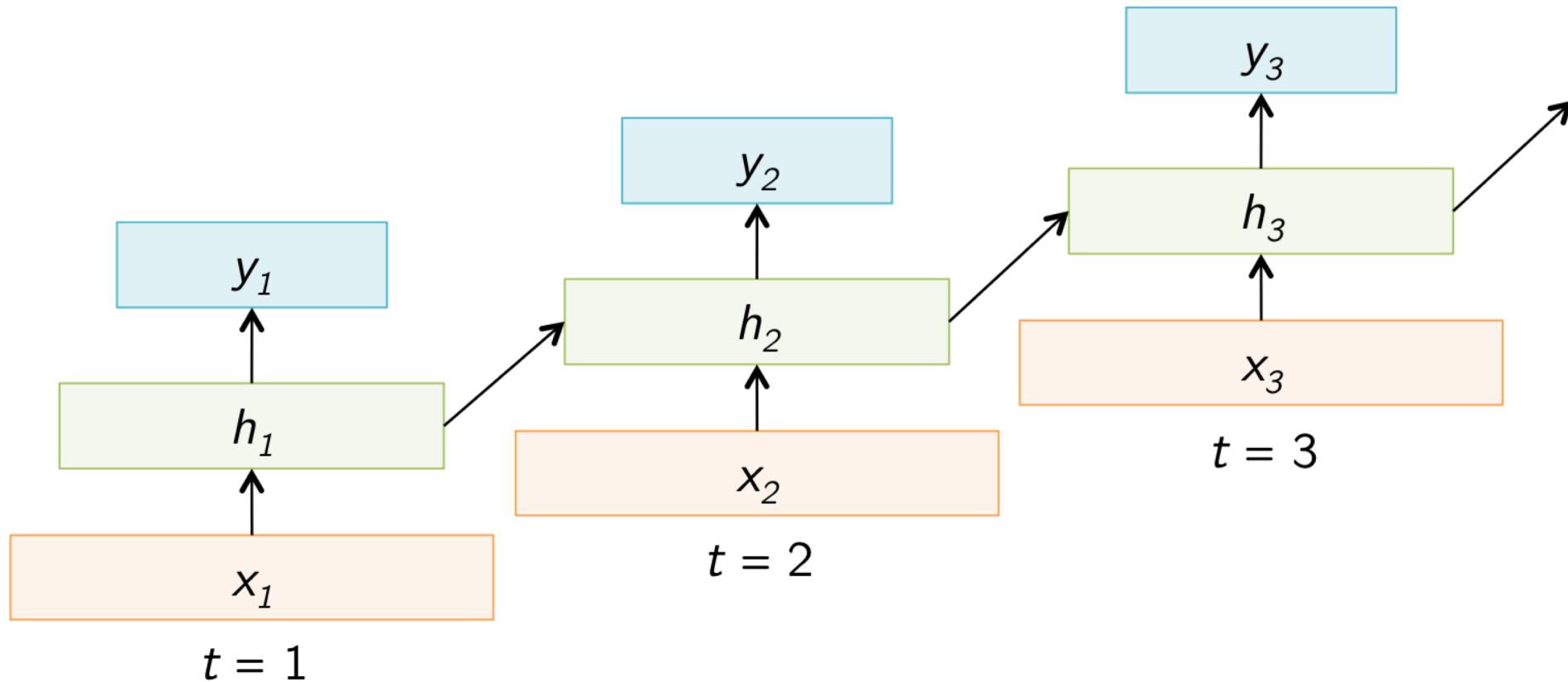


Naïve approach

- What's the problem?



Recurrent Neural Network (RNN)



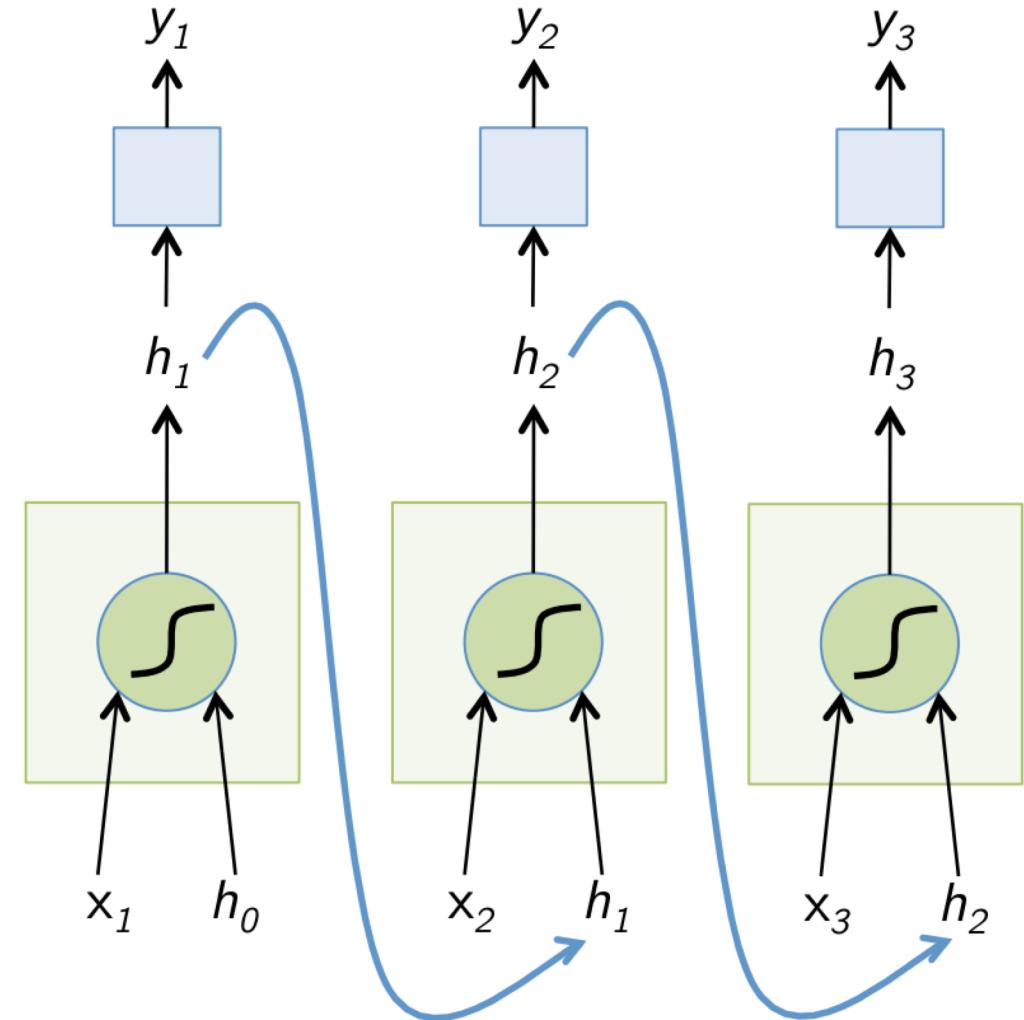
Recurrent Neural Network (RNN)

$$a_t = b_h + W_{hh} h_{t-1} + W_{xh} x_t$$

$$h_t = \tanh(a_t)$$

$$o_t = b_o + W_{hy} h_t$$

$$y_t = \text{softmax}(o_t)$$



Example: Sentiment Classification

- Input?
- Output?

"I love this movie.
I've seen it many times
and it's still awesome."



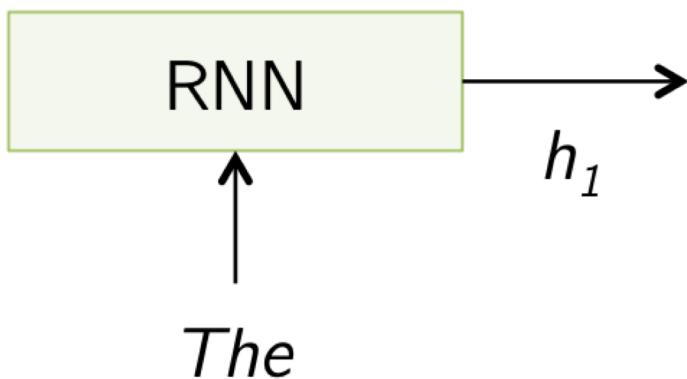
"This movie is bad.
I don't like it at all.
It's terrible."



Example: Sentiment Classification

Input : “The food is really good”

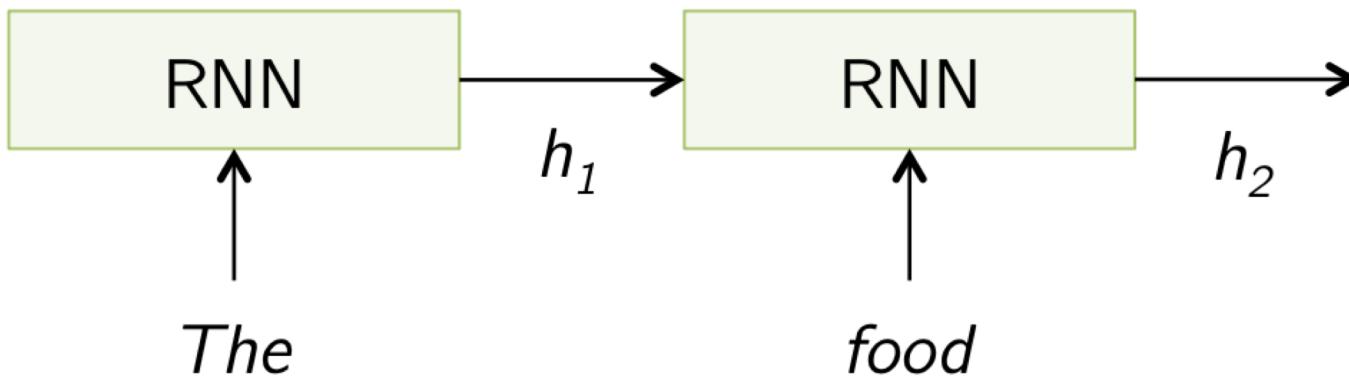
Output: 0 or 1 (Negative / Positive)



Example: Sentiment Classification

Input : “The food is really good”

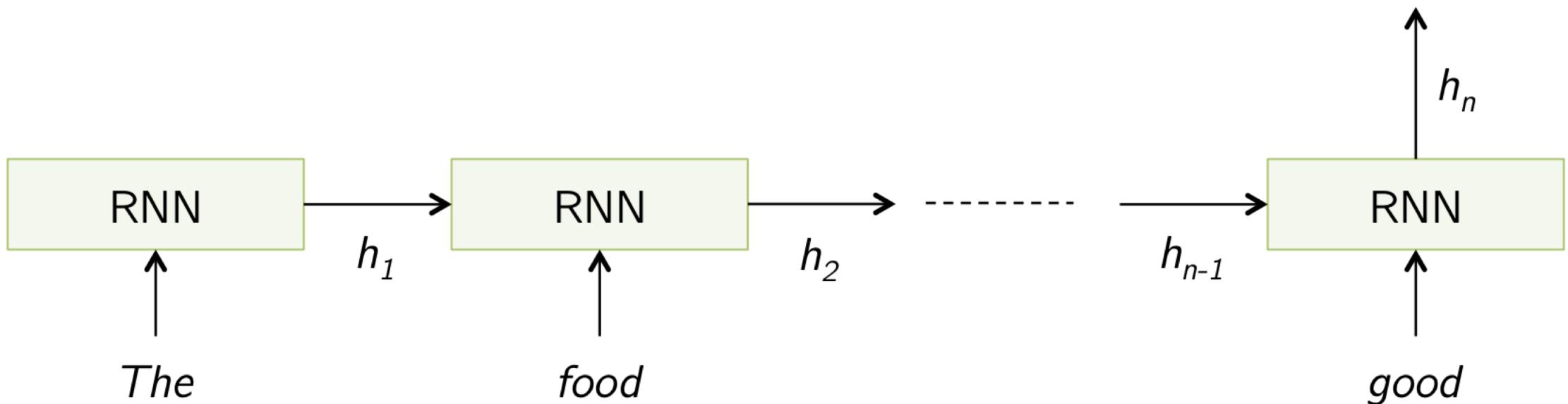
Output: 0 or 1 (Negative / Positive)



Example: Sentiment Classification

Input : “The food is really good”

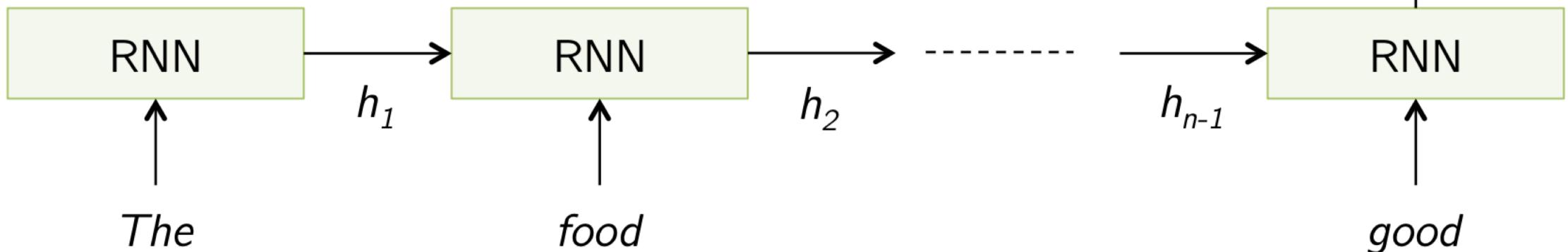
Output: 0 or 1 (Negative / Positive)



Example: Sentiment Classification

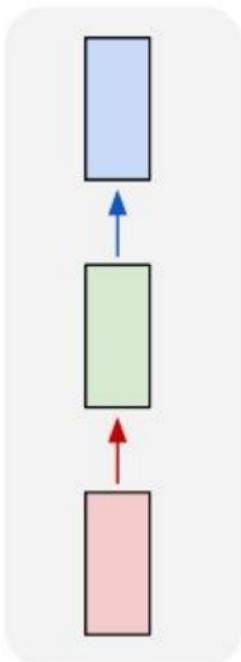
Input : “The food is really good”

Output: 0 or 1 (Negative / Positive)

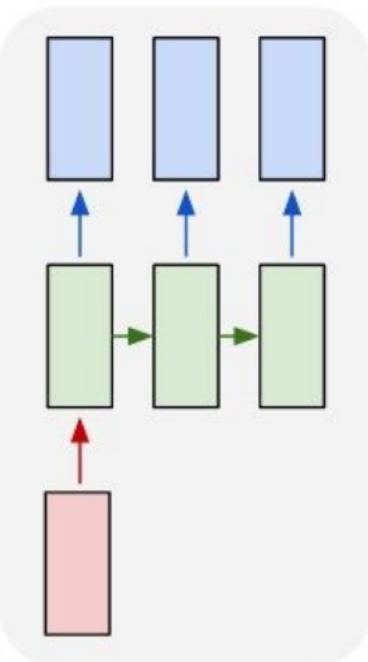


Variants of RNN

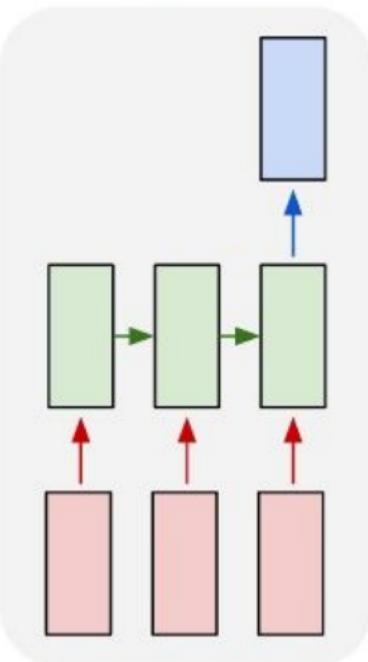
one to one



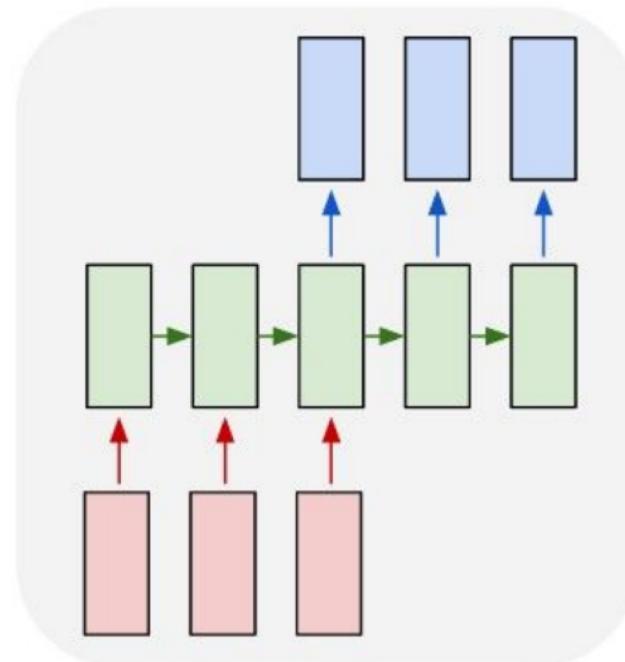
one to many



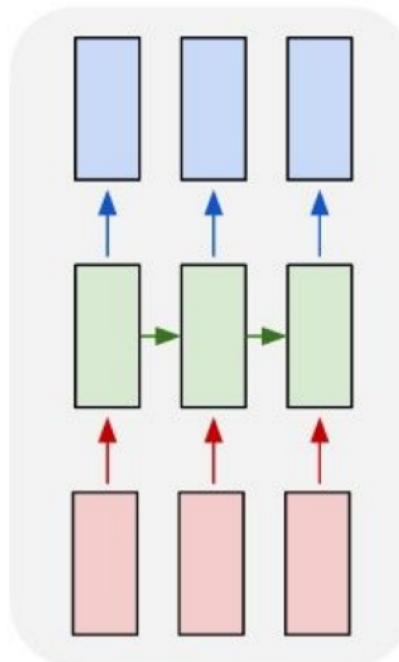
many to one



many to many



many to many



How does it work in mathematical detail?

- Deep learning textbook, Chapter 10
- [http://cs231n.stanford.edu/slides/2017/cs231n 2017 lecture10.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture10.pdf)
- <https://www.cs.ubc.ca/labs/lci/mlrg/slides/rnn.pdf>

Take-home Messages

- Key question to design neural network solution for domain: How do we parametrize neural networks based on domain characteristic?
 - In other words, how do we share parameters effectively?
- CNN & RNN are amazing representation learning tool changing the world
 - Visual data → CNN, Sequential data → RNN
- CNN processes grid-like topology by employing specialized operations such as convolution and pooling.
- RNN processes variable-length inputs and outputs by maintaining state information across time steps