

# **Convolutional Neural Networks**

**Elisa Bassignana**

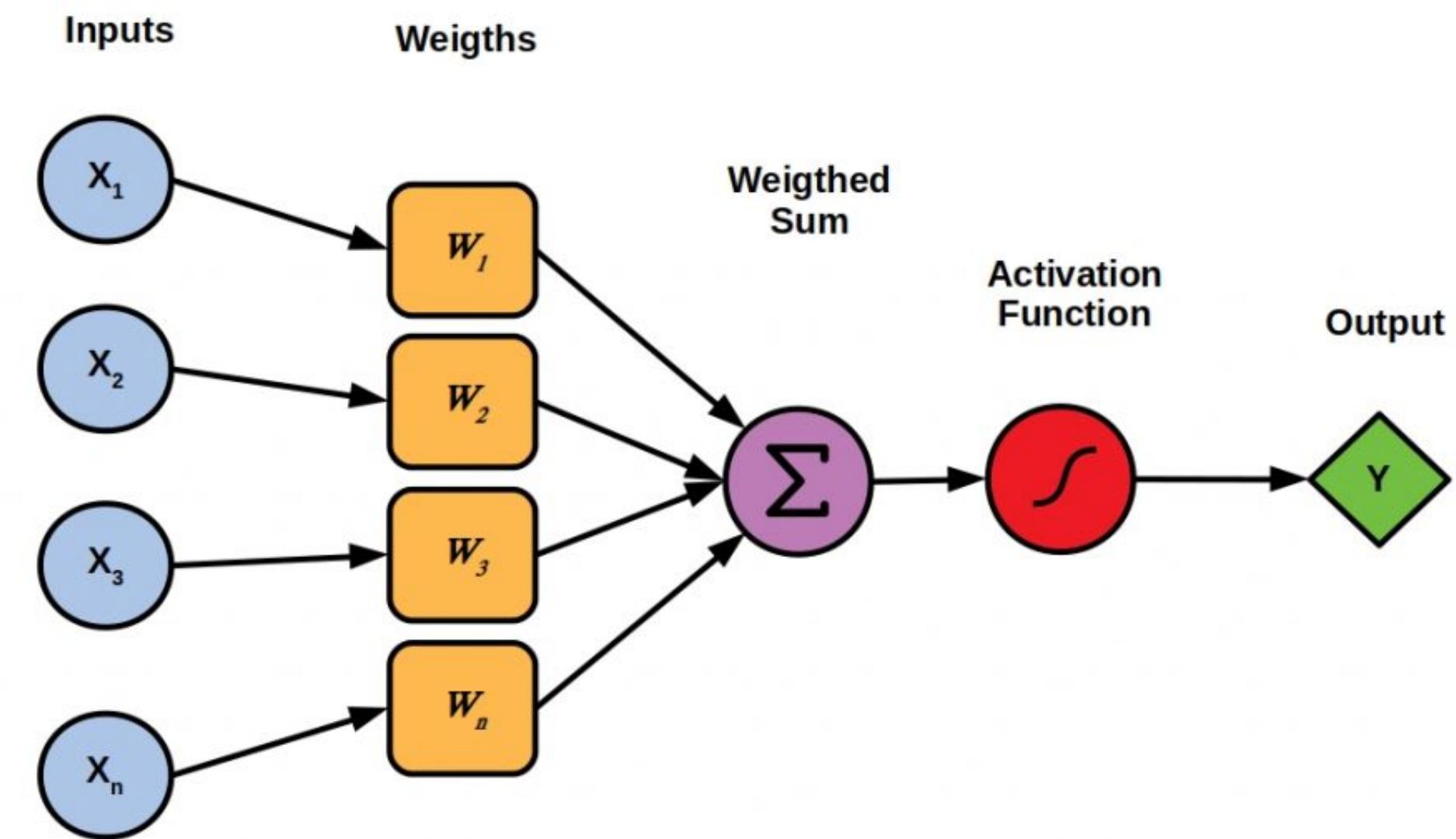
**ITU - 22/02/2023**

# **Plan of the day**

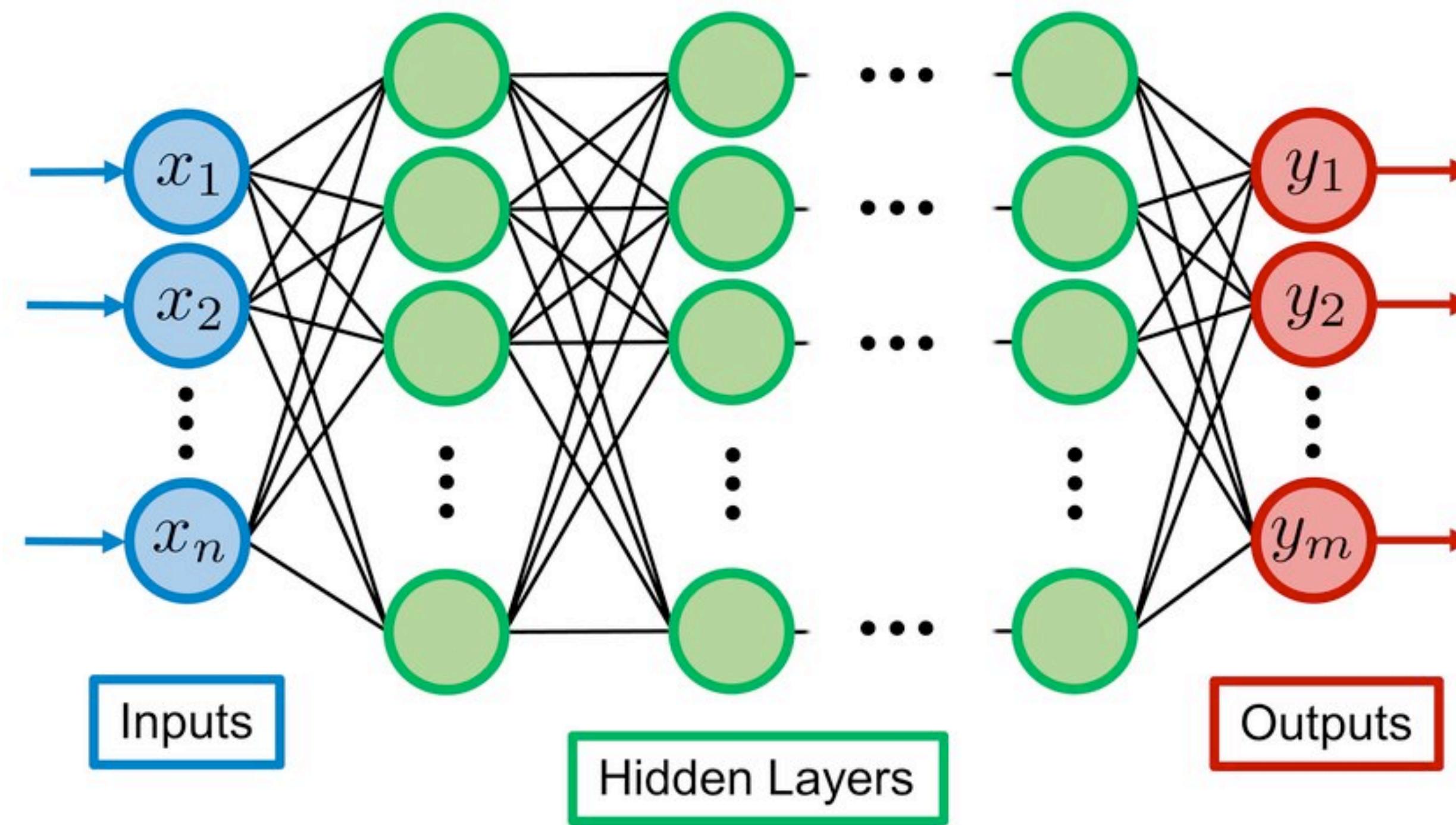
- 1. Recap of FFNN and their limitations**
- 2. Introduction to Convolutional Neural Networks (CNNs)**
- 3. CNNs for NLP**
- 4. Case study with CNNs**

# Recap: Feed-Forward Neural Networks

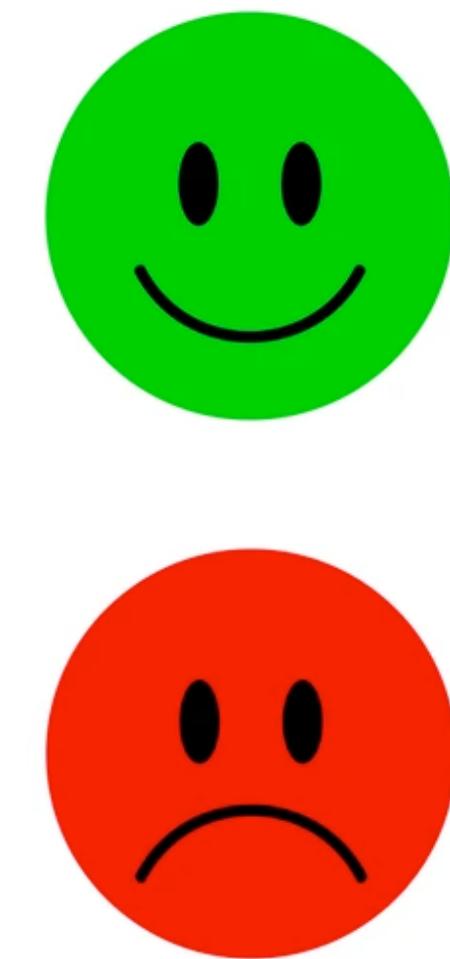
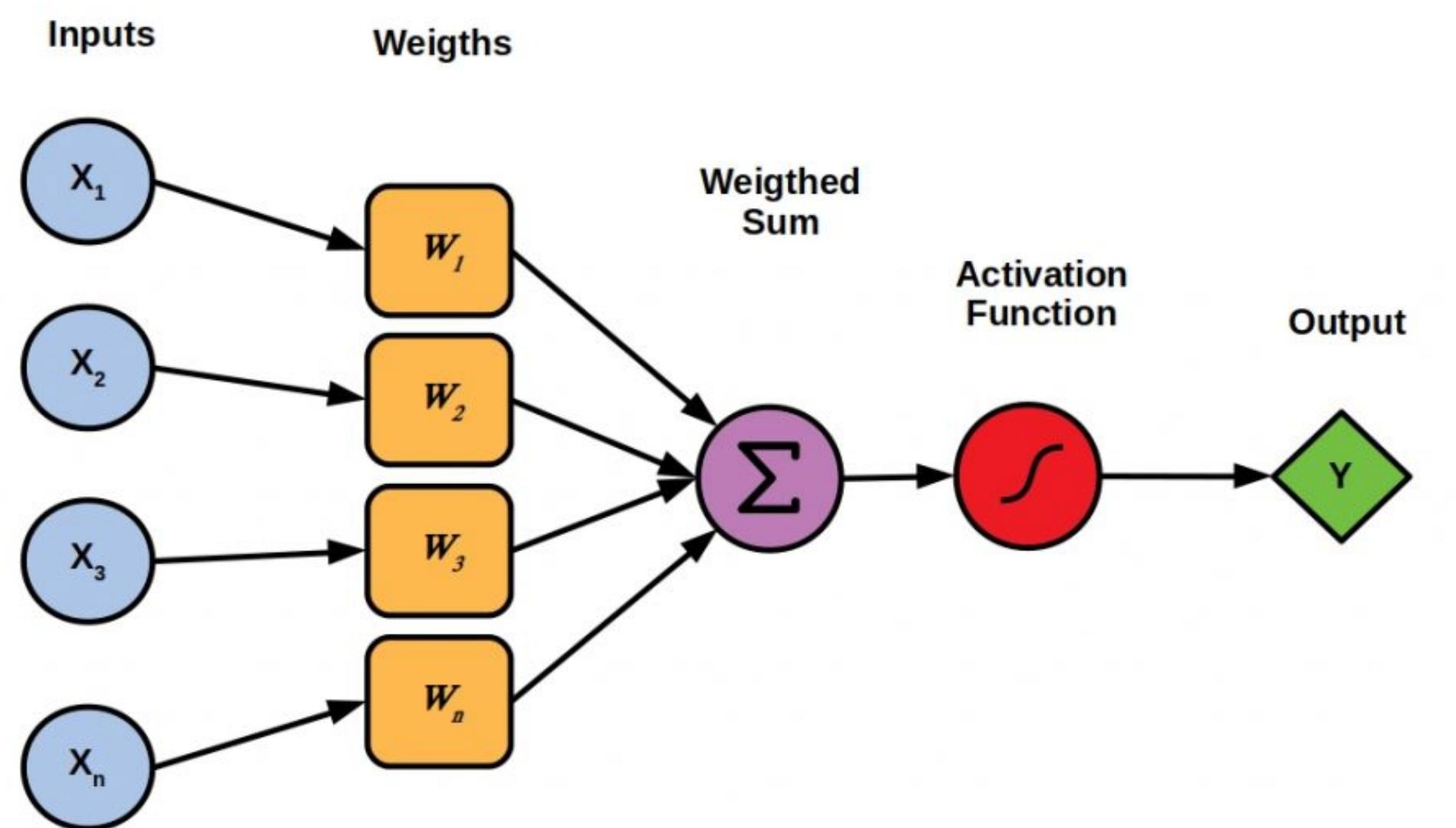
## Neural unit



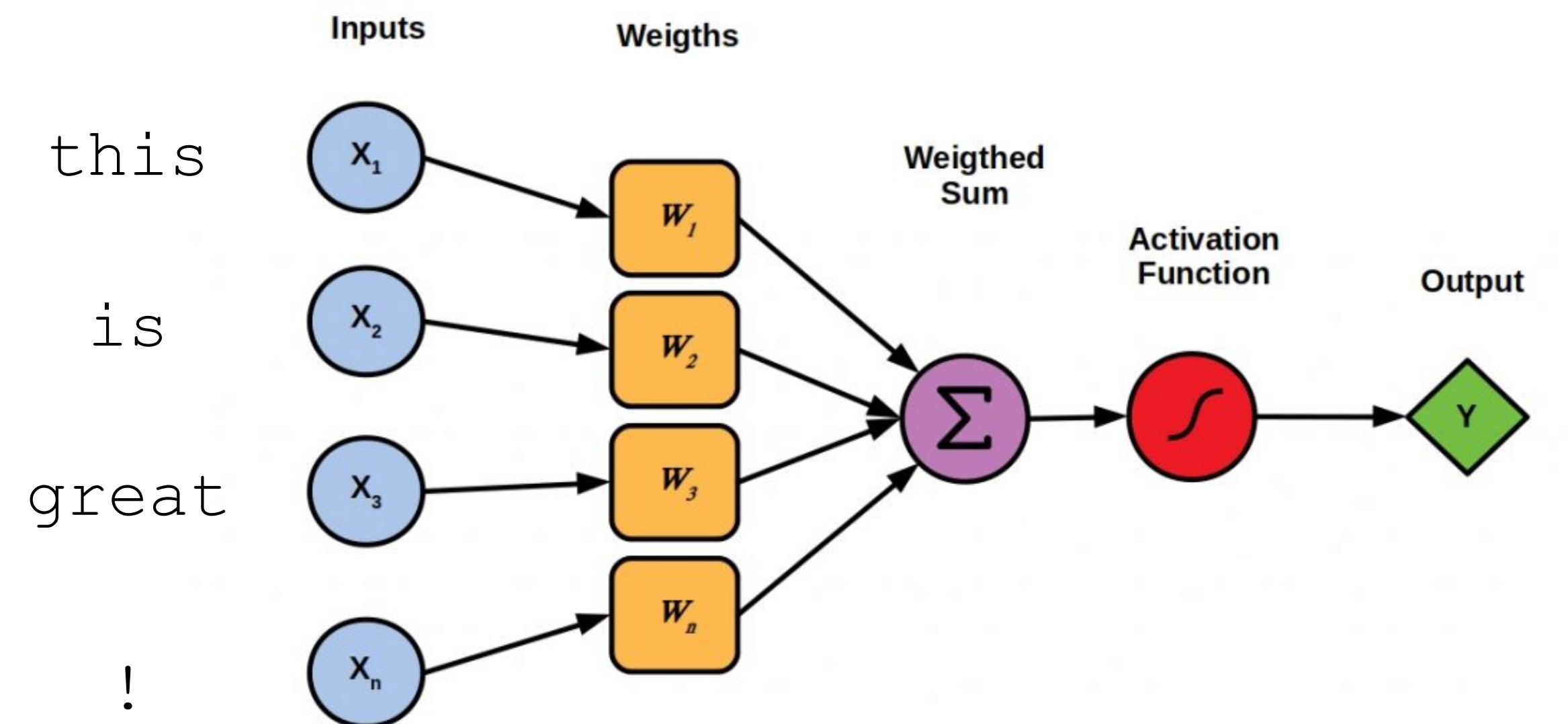
# Recap: Feed-Forward Neural Networks



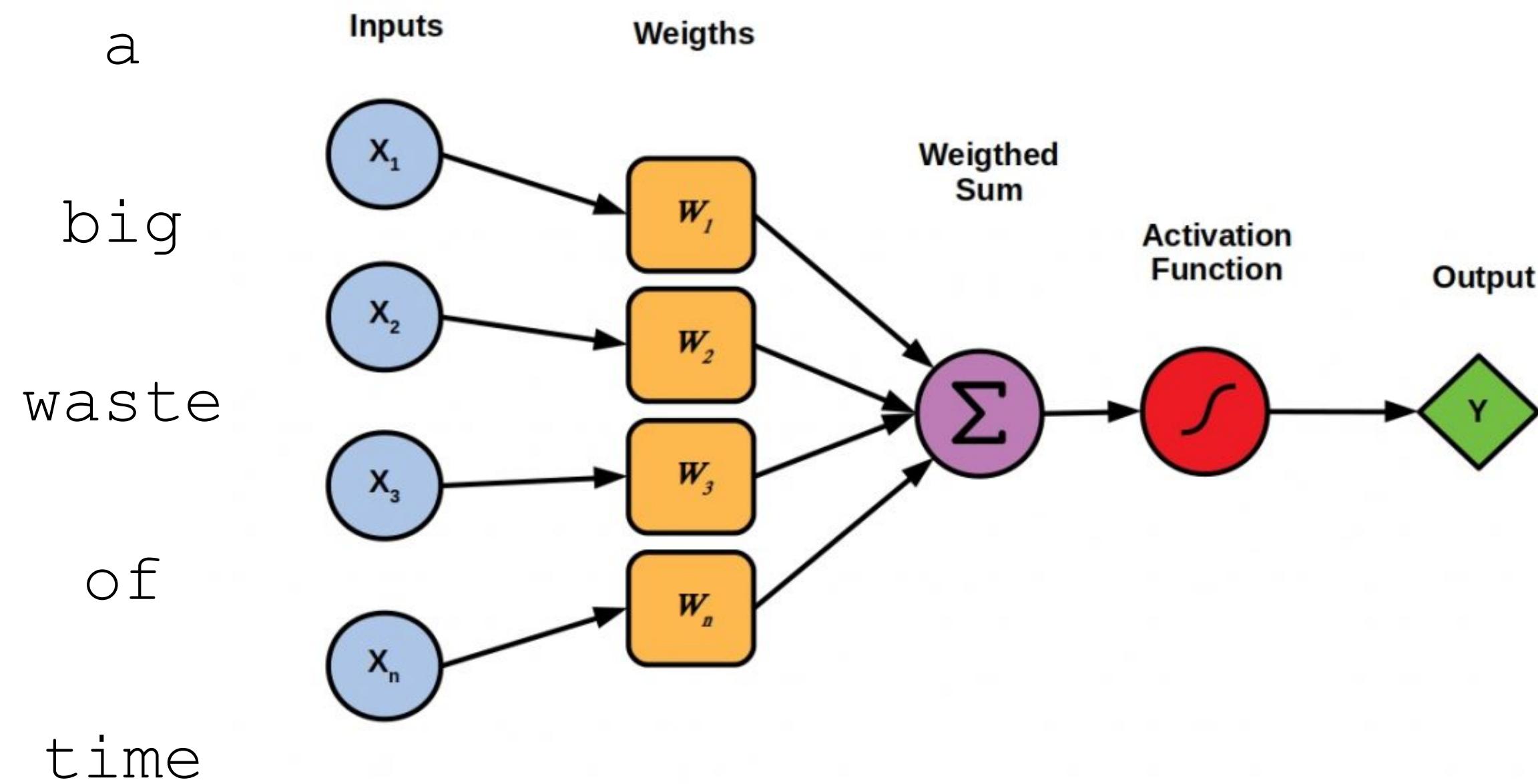
# FFNN: Any problems?



# FFNN: Any problems?

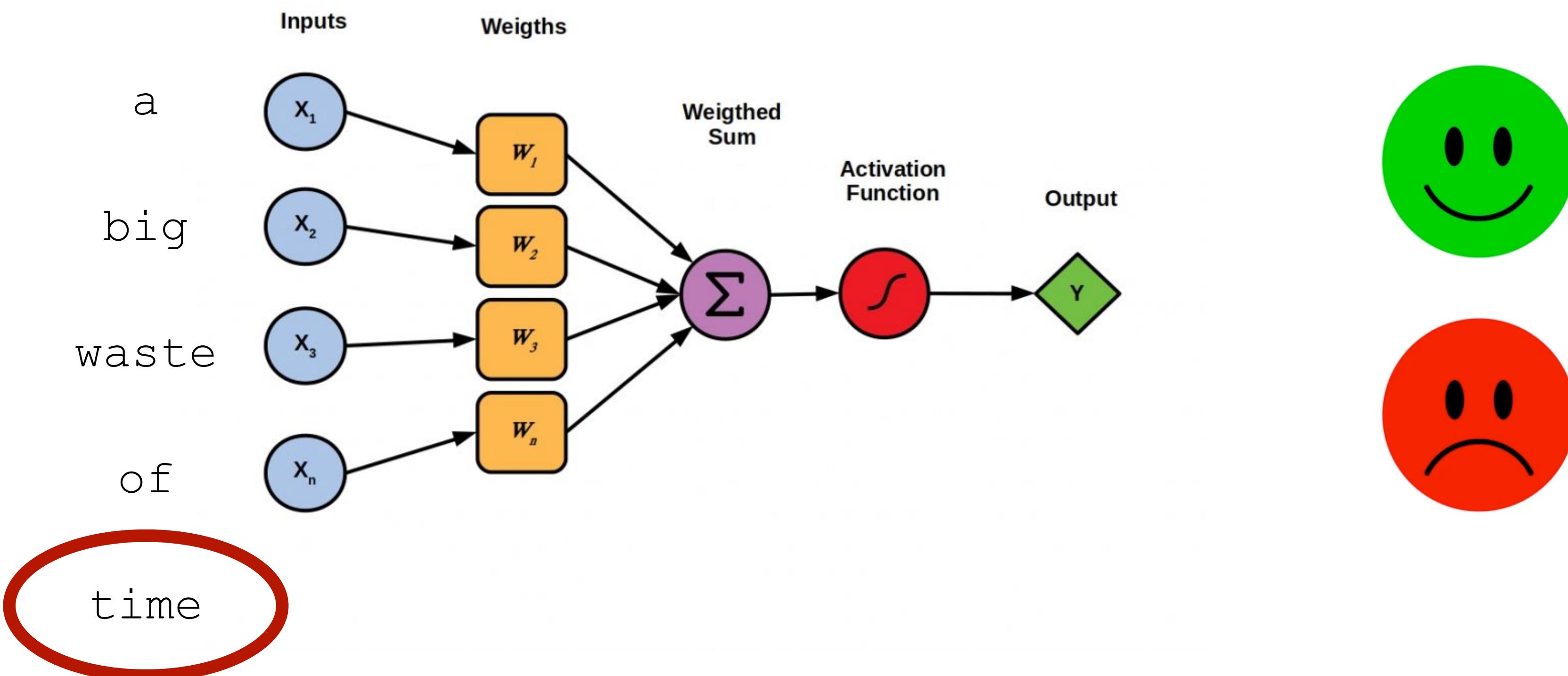


# FFNN: Any problems?



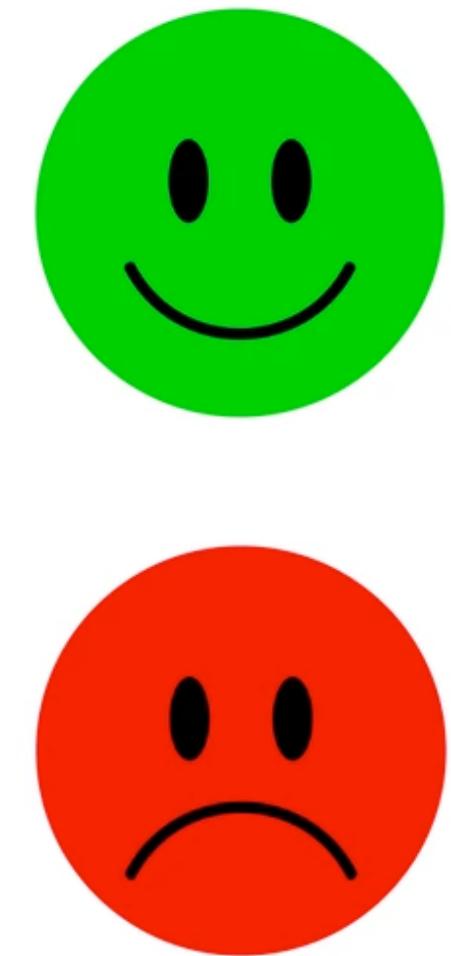
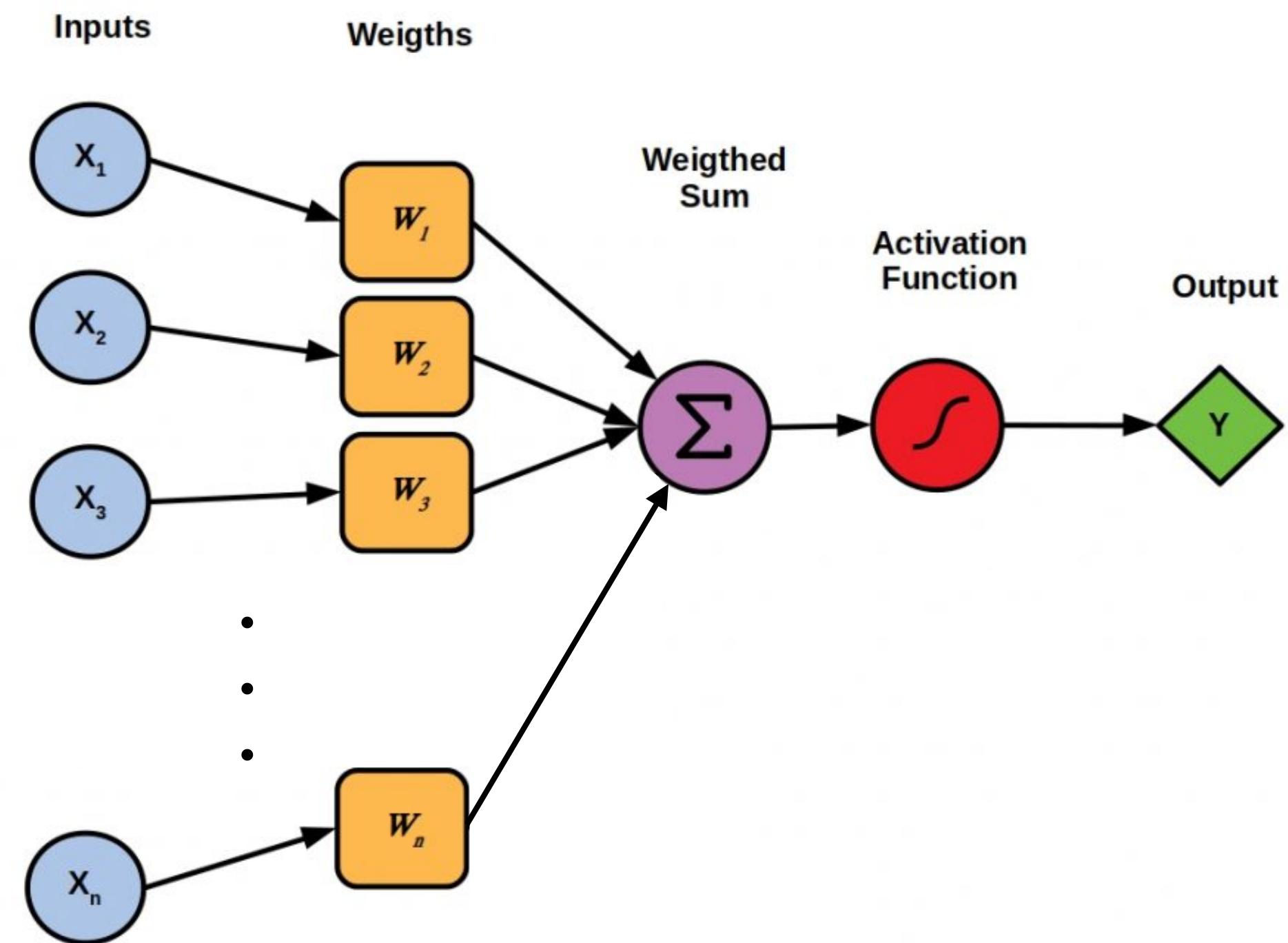
# FFNN: Any problems?

PROBLEM:  
we need fixed input size!



# FFNN: What is still missing?

it  
was  
not  
good  
, it  
was  
actually  
quite  
bad

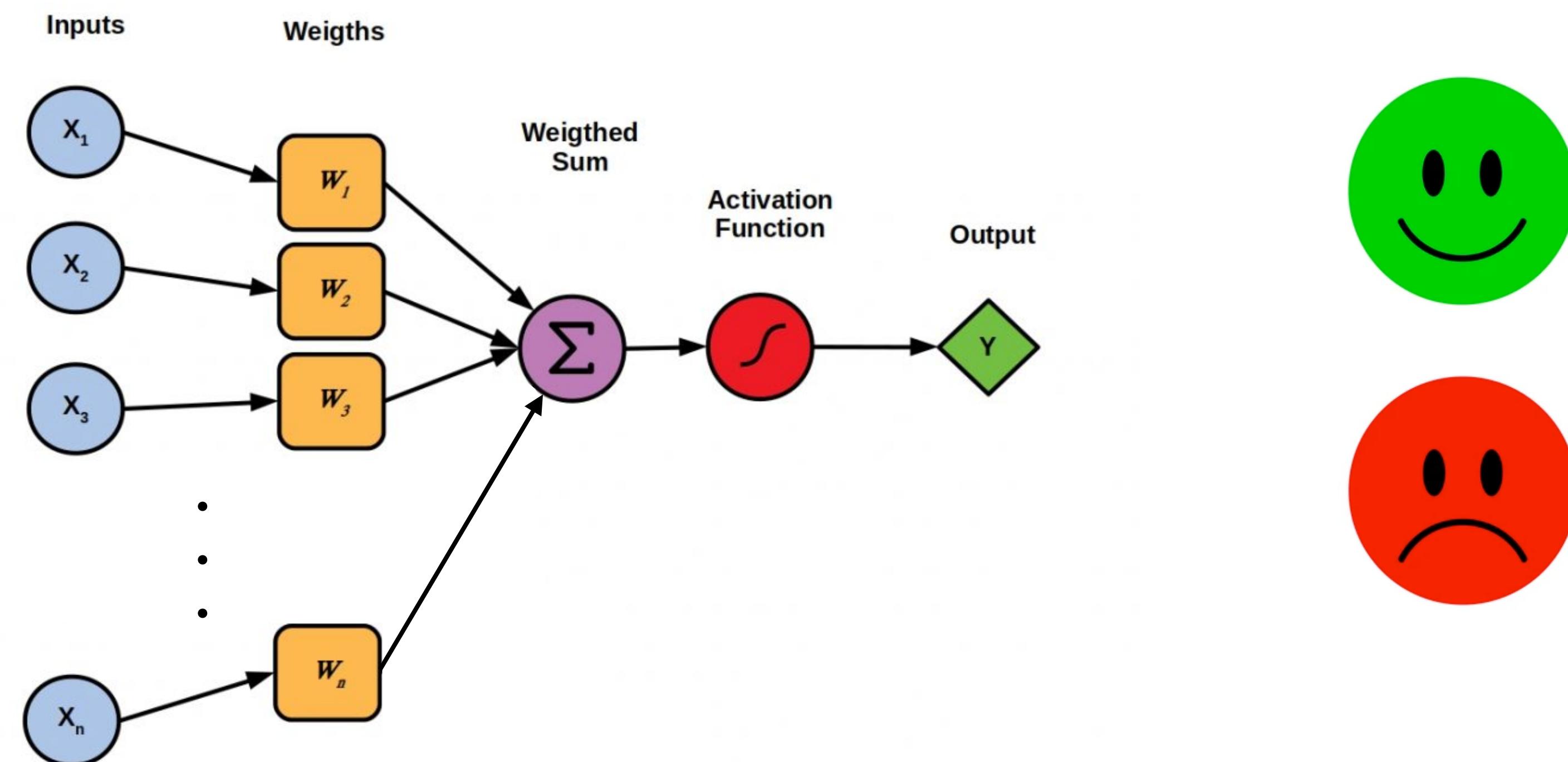


# FFNN: What is still missing?

PROBLEM:  
it does not consider the order!

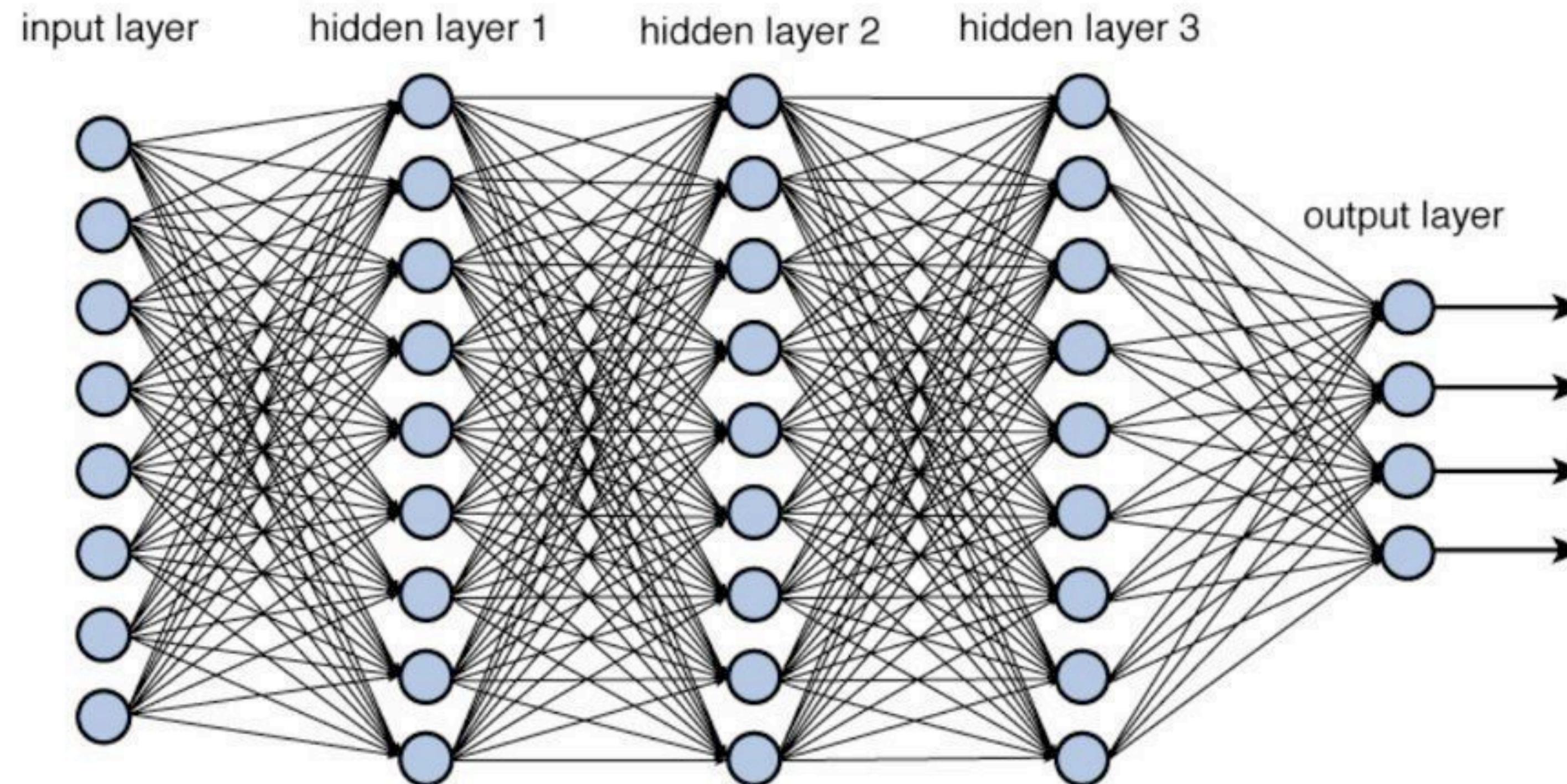
it  
was  
not  
good  
,

it  
was  
actually  
quite  
bad



# Moreover...

**PROBLEM:**  
**FFNN can became very complex:  
many parameters to tune!**



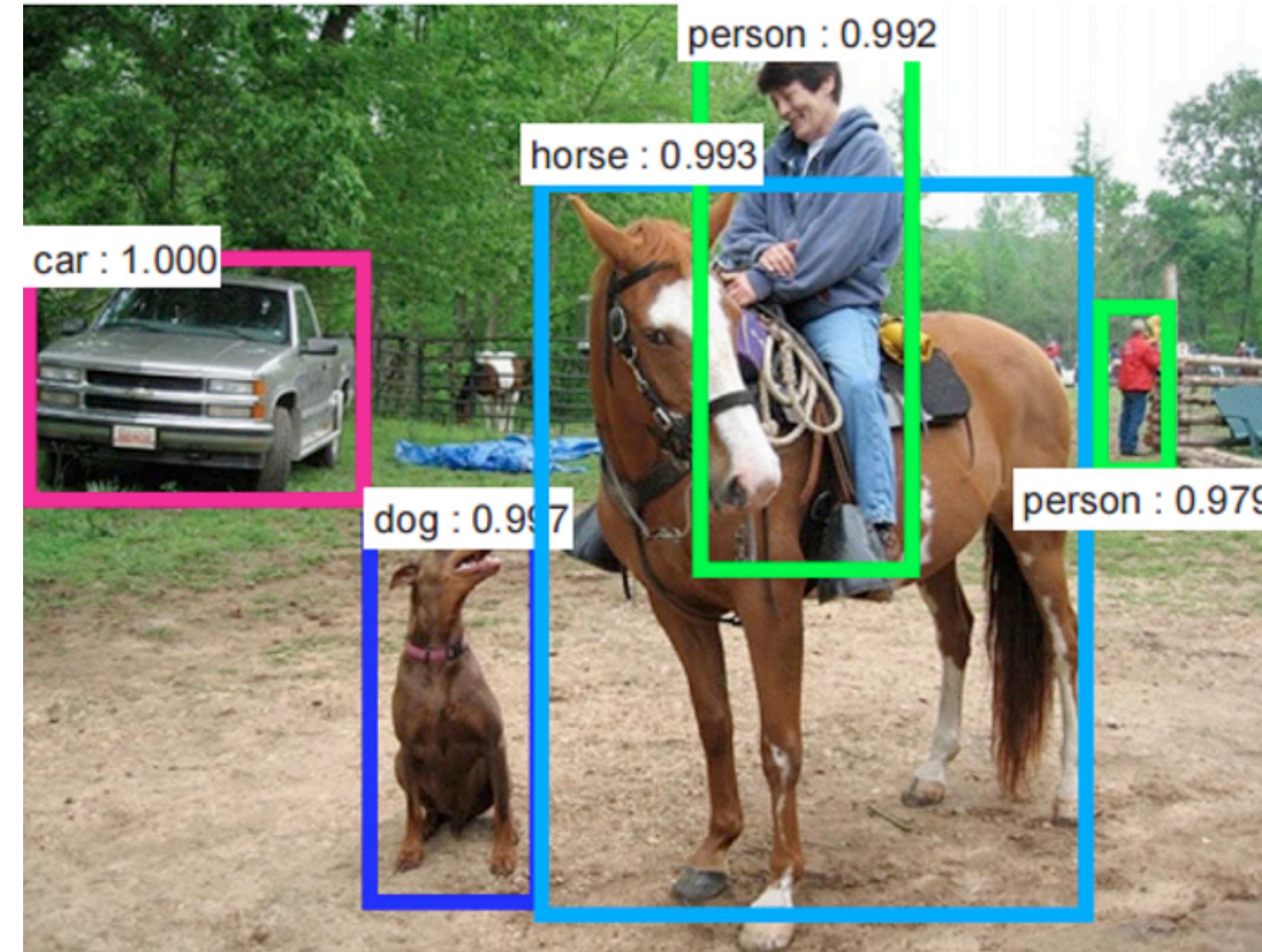
**Any questions so far?**

# **Convolutional Neural Networks (CNNs or convnets)**

- From the computer vision field
- Work with 2-dimensional data (matrix/grid): images
- Can handle arbitrary-size input and reduce them down to a fixed size vector representation
- Core ideas:
  - *Local processing*
  - *Parameter sharing*

# Intuition

- **Invariance** in the data  
(we want to find an object regardless of its position in the image)



# What are CNNs?

From a technical point of view there are 3 core concepts:

1. CNNs use **convolutions** over the input to compute the output
2. Each convolution applies different **filters**
3. Results of the convolutions are combined together via **pooling**

# Filter (or kernel)

- Sliding window of learned weights

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

input

1	0	1
0	1	0
1	0	1

filter

# Convolution

- Mathematical operation: **composition function**
- Sliding a filter over the input, generating an output (**feature map**)
- Parameters of the composition function are the same (shared)
- For each position of the filter, we multiply the overlapping values of the filter and of the image together and add up the results, to produce the output

1	0	1
0	1	0
1	0	1

filter

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

input

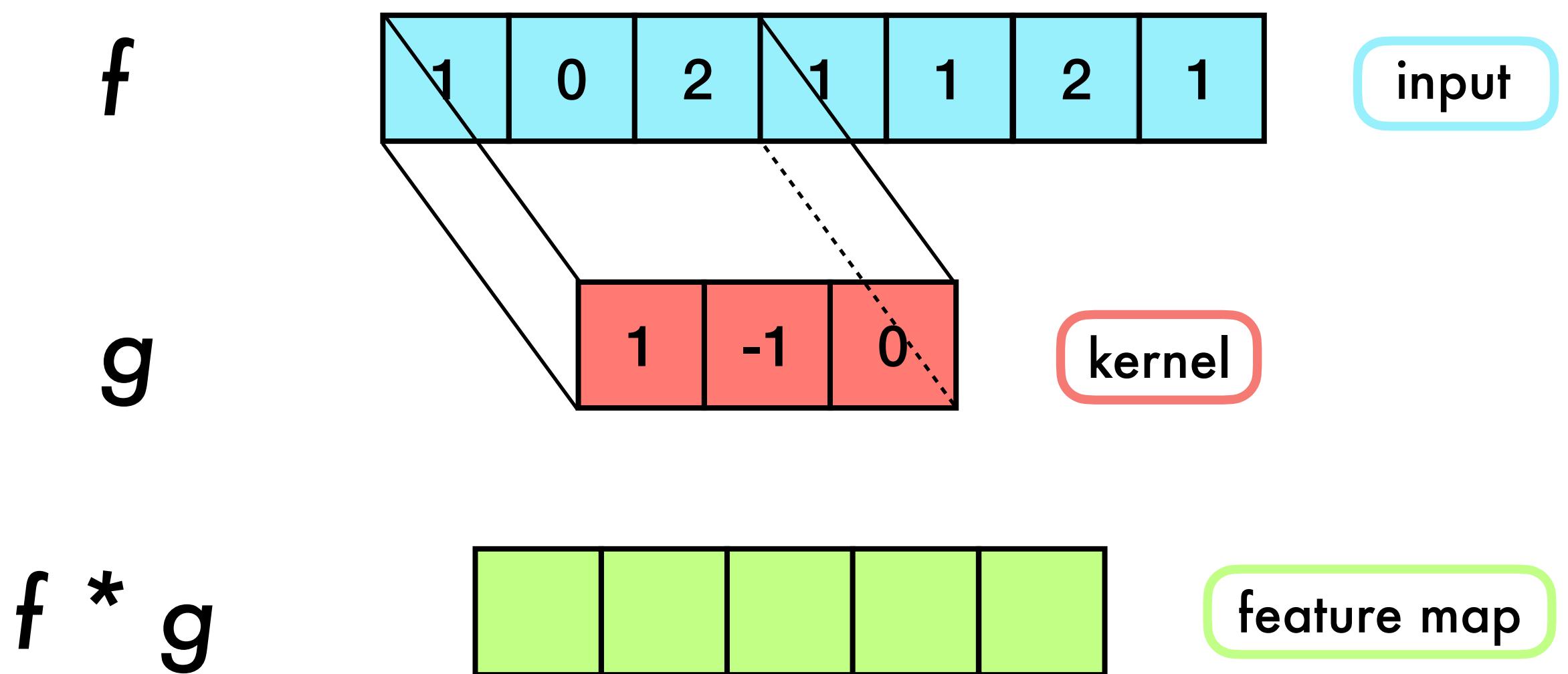
4		

feature map

# Example

$f$  is a 1-d input vector of length  $n$

$g$  is a 1-d filter (kernel) of length  $m$



# Example

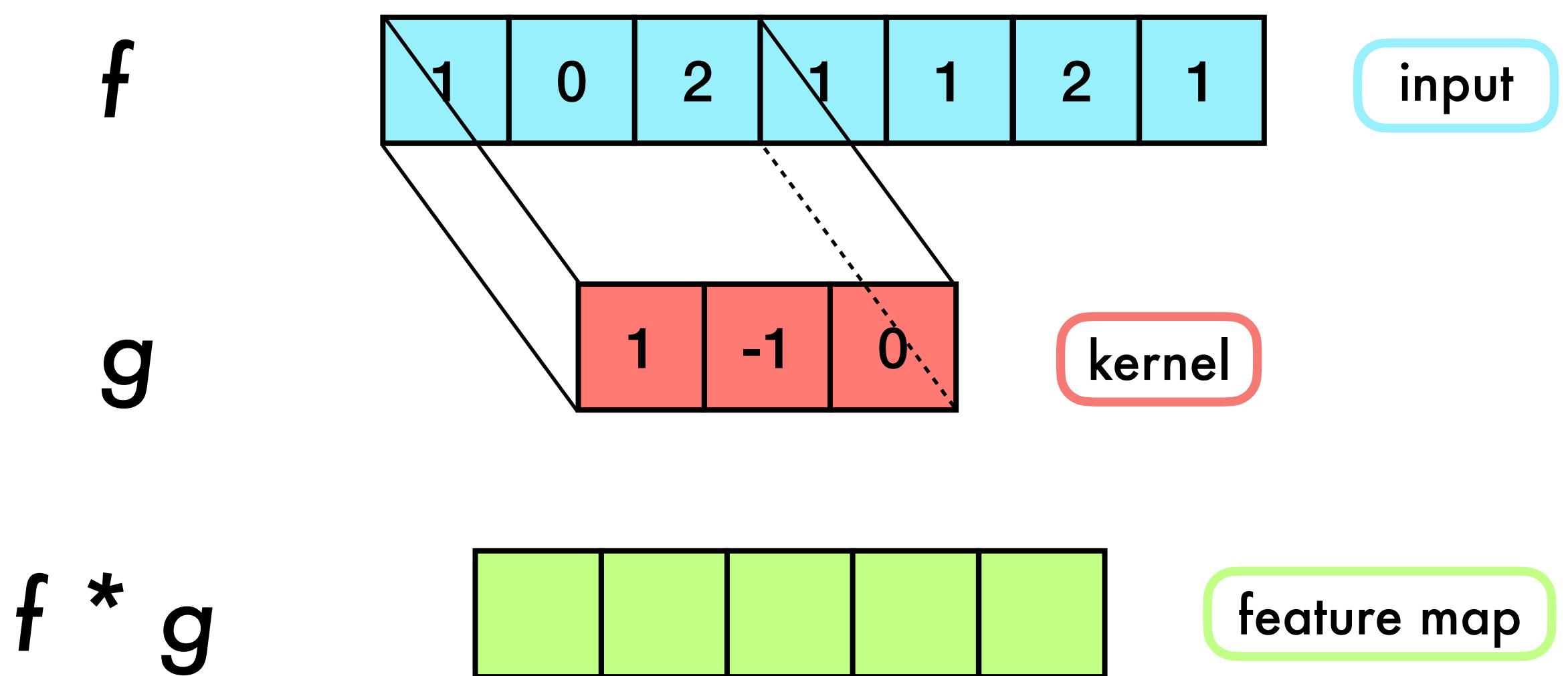
$f$  is a 1-d input vector of length  $n$

$$(f * g)(i) = \sum_{j=1}^m g(j) \cdot f(i + j - 1)$$

$g$  is a 1-d filter (kernel) of length  $m$

$$n = 7$$

$$m = 3$$



# Example

$f$  is a 1-d input vector of length  $n$

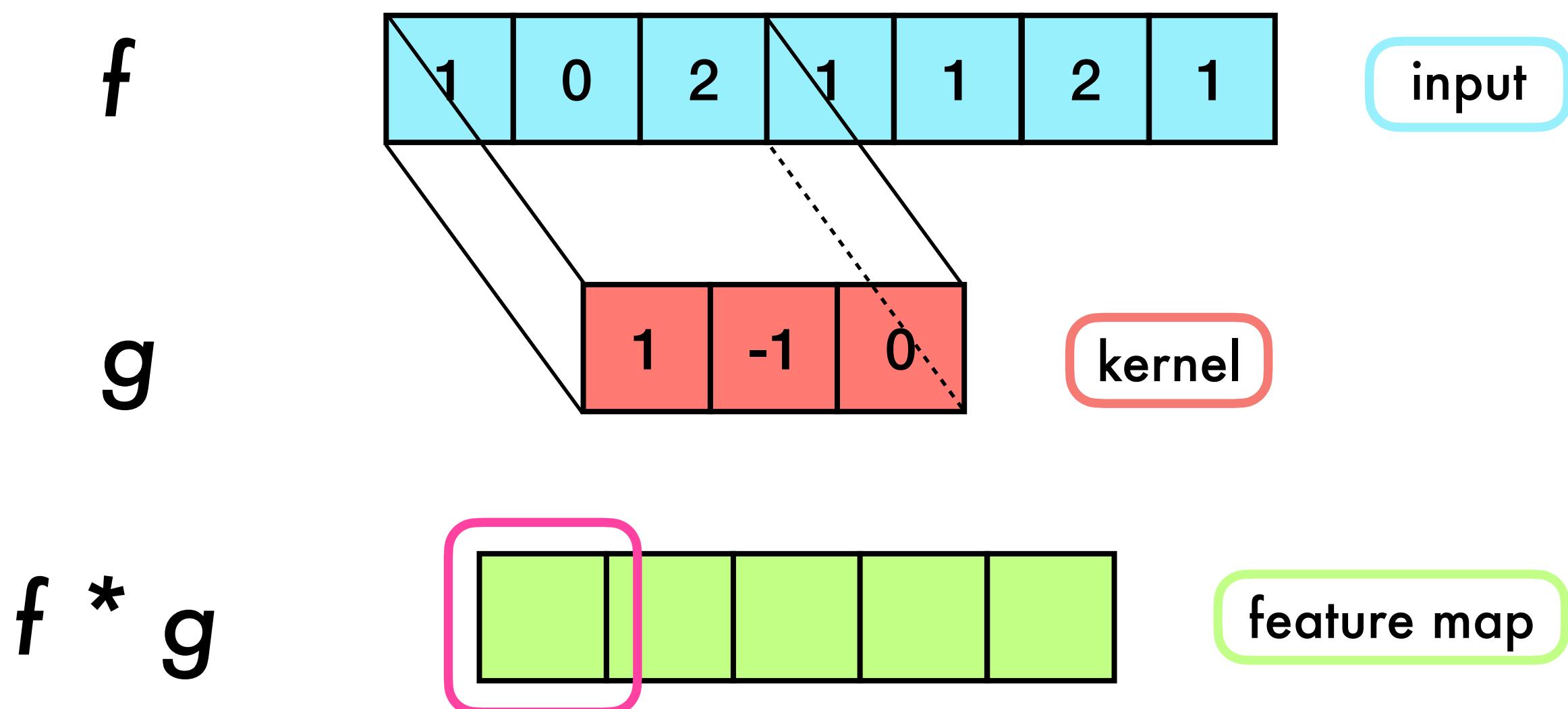
$$(f * g)(i) = \sum_{j=1}^m g(j) \cdot f(i + j - 1)$$

$g$  is a 1-d filter (kernel) of length  $m$

$$n = 7$$

$$m = 3$$

$$i = \boxed{1} \ 2 \ 3 \ 4 \ 5$$



# Example

$f$  is a 1-d input vector of length  $n$

$$(f * g)(i) = \sum_{j=1}^m g(j) \cdot f(i + j - 1)$$

$g$  is a 1-d filter (kernel) of length  $m$

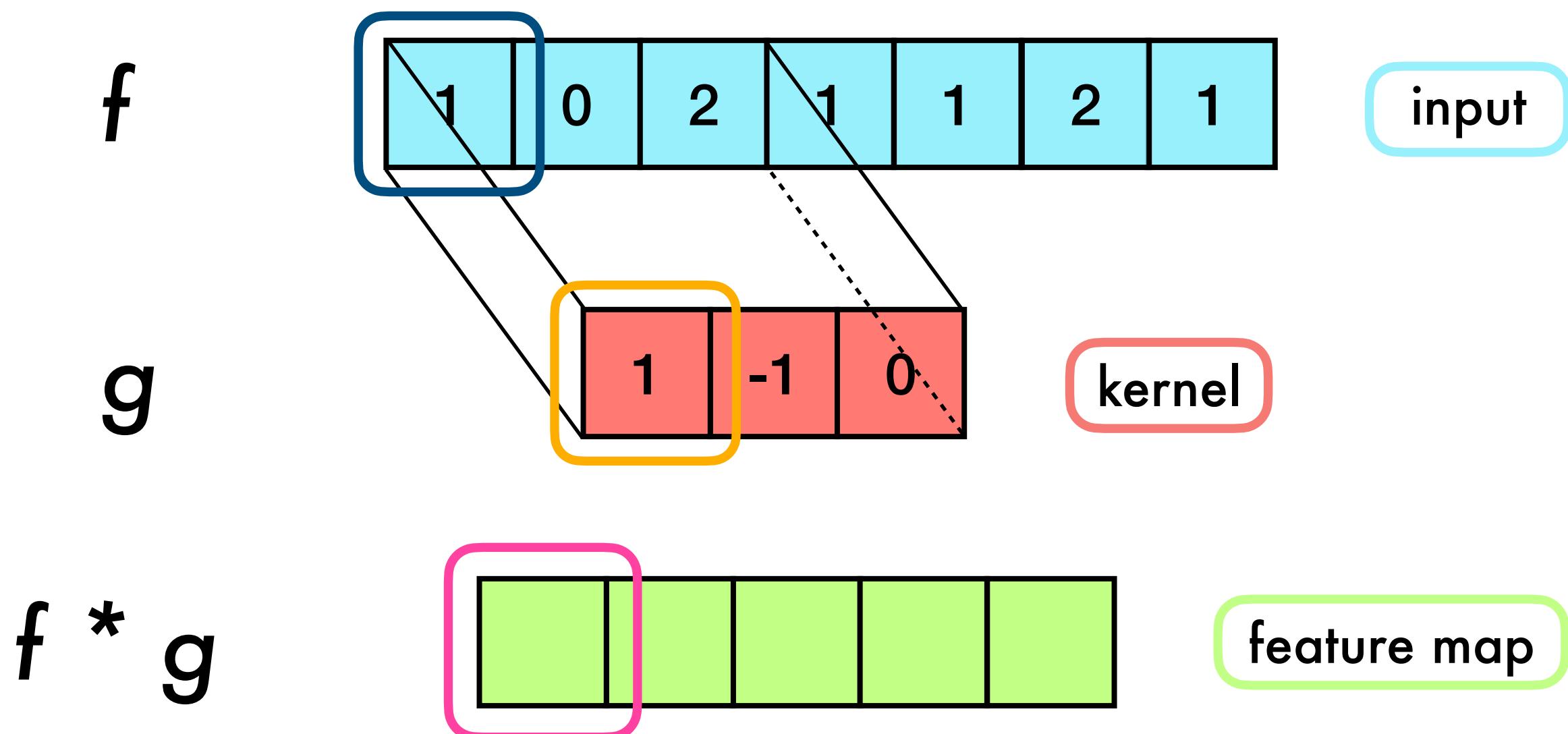
$$n = 7$$

$$m = 3$$

$$i = 1 \ 2 \ 3 \ 4 \ 5$$

$$j = 1 \ 2 \ 3$$

$$1 * 1 = 1$$

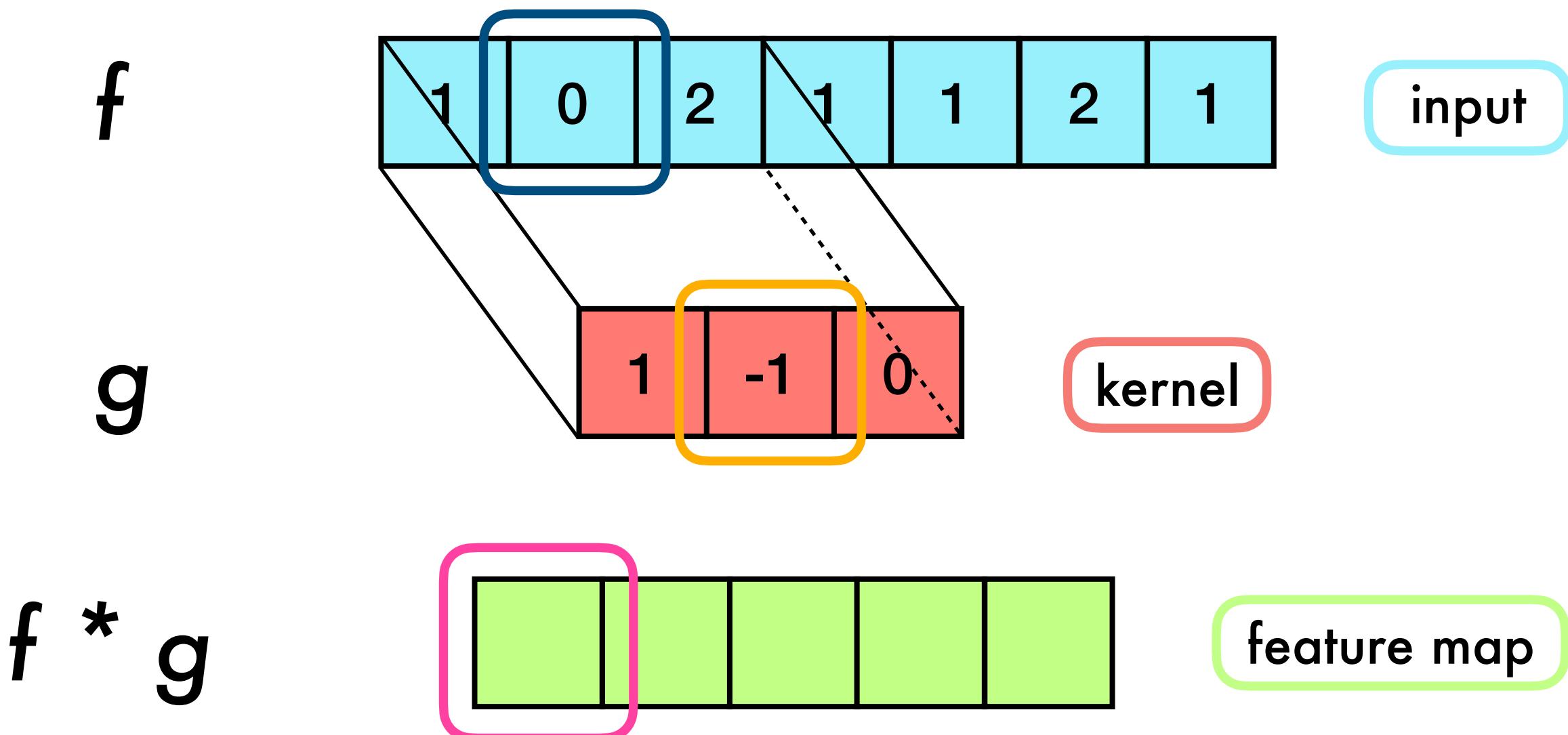


# Example

$f$  is a 1-d input vector of length  $n$

$$(f * g)(i) = \sum_{j=1}^m g(j) \cdot f(i + j - 1)$$

$g$  is a 1-d filter (kernel) of length  $m$



$$n = 7$$

$$m = 3$$

$$i = 1 \ 2 \ 3 \ 4 \ 5$$

$$i = 1 \ 2 \ 3$$

$$1 * 1 = 1$$

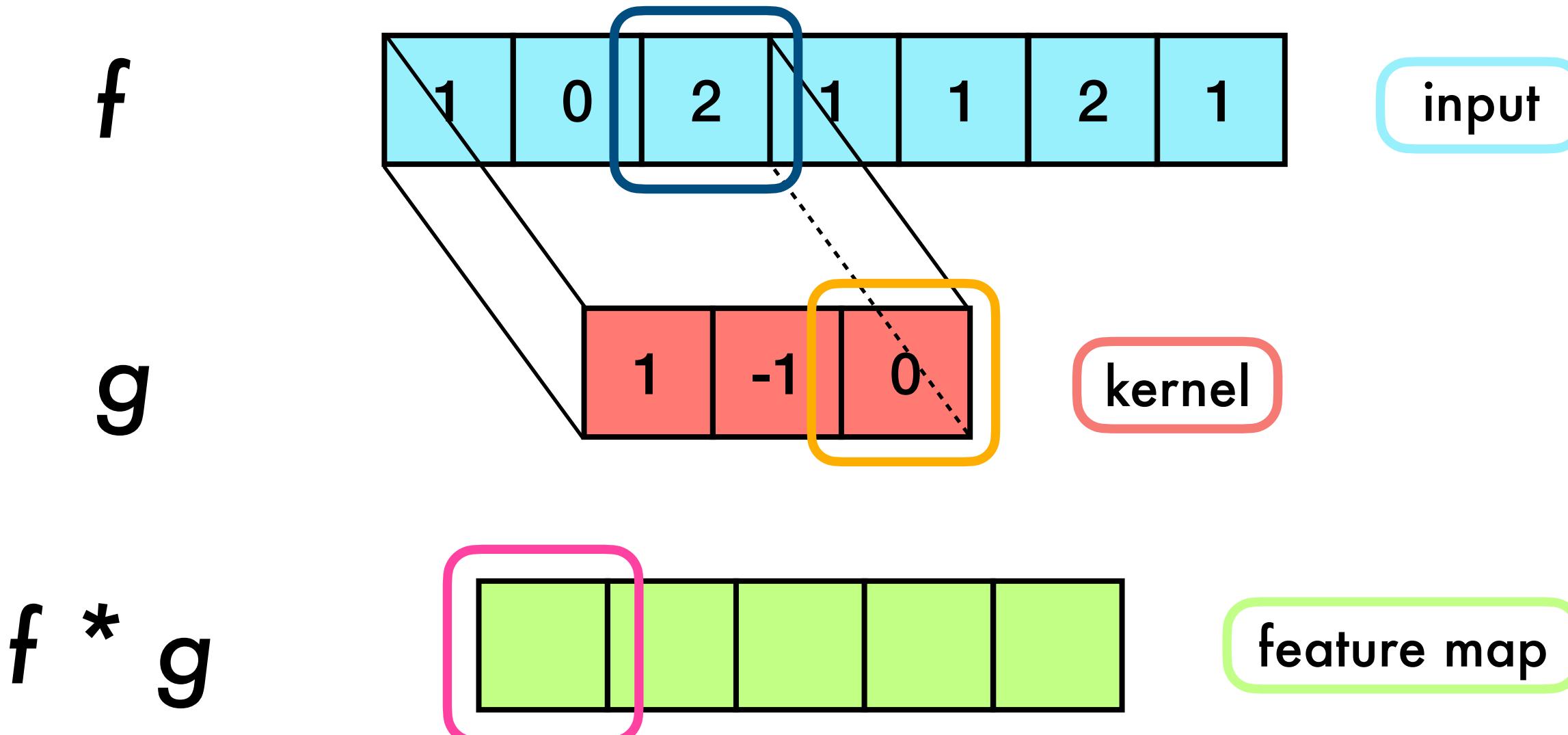
$$-1 * 0 = 0$$

# Example

$f$  is a 1-d input vector of length  $n$

$$(f * g)(i) = \sum_{j=1}^m g(j) \cdot f(i + j - 1)$$

$g$  is a 1-d filter (kernel) of length  $m$



$$n = 7$$

$$m = 3$$

$$i = 1 \ 2 \ 3 \ 4 \ 5$$

$$i = 1 \ 2 \ 3$$

$$1 * 1 = 1$$

$$-1 * 0 = 0$$

$$0 * 2 = 0$$

# Example

$f$  is a 1-d input vector of length  $n$

$$(f * g)(i) = \sum_{j=1}^m g(j) \cdot f(i + j - 1)$$

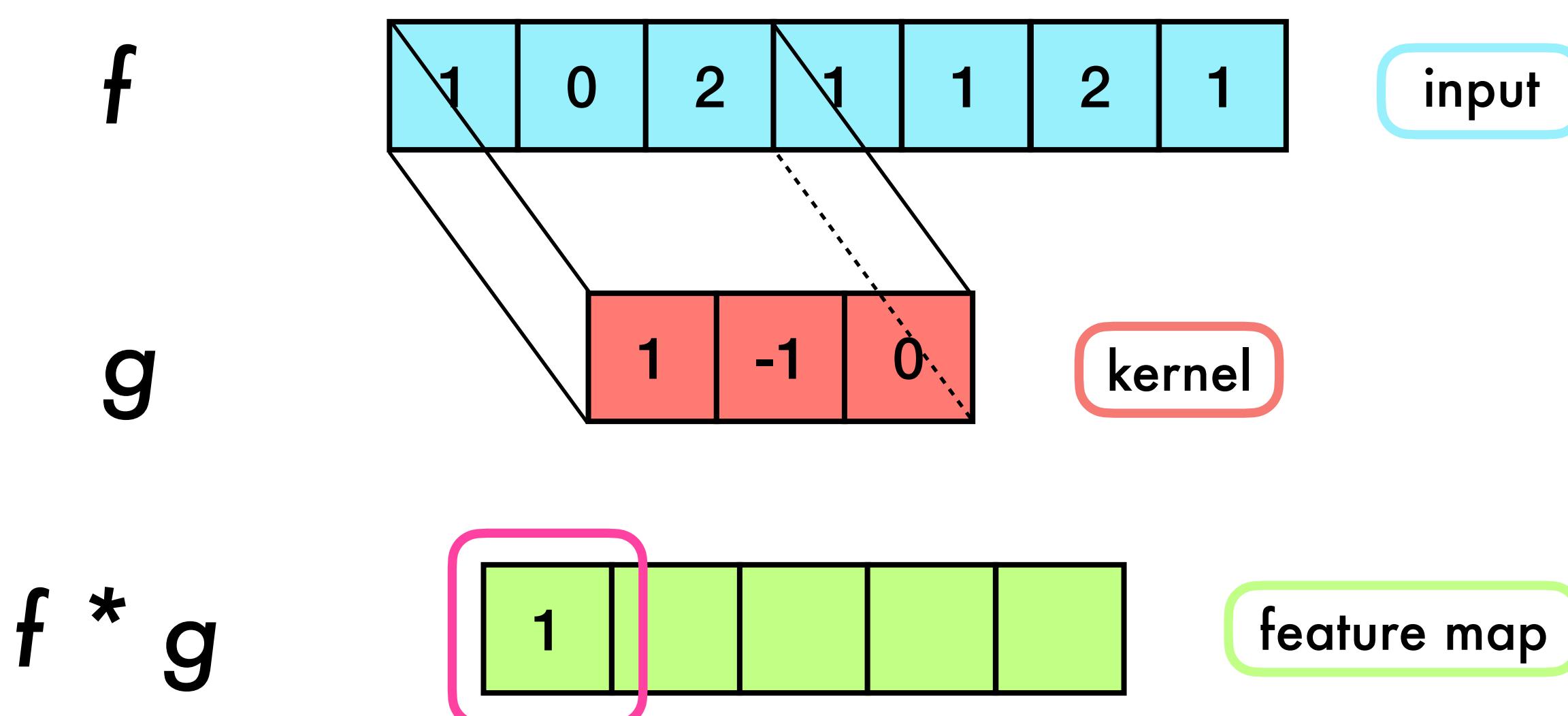
$g$  is a 1-d filter (kernel) of length  $m$

$$n = 7$$

$$m = 3$$

$$i = 1 \ 2 \ 3 \ 4 \ 5$$

$$j = 1 \ 2 \ 3$$



$$\begin{aligned} 1 * 1 &= 1 \\ -1 * 0 &= 0 \\ 0 * 2 &= 0 \end{aligned} \quad ] \quad 1 + 0 + 0$$

# Example

$f$  is a 1-d input vector of length  $n$

$$(f * g)(i) = \sum_{j=1}^m g(j) \cdot f(i + j - 1)$$

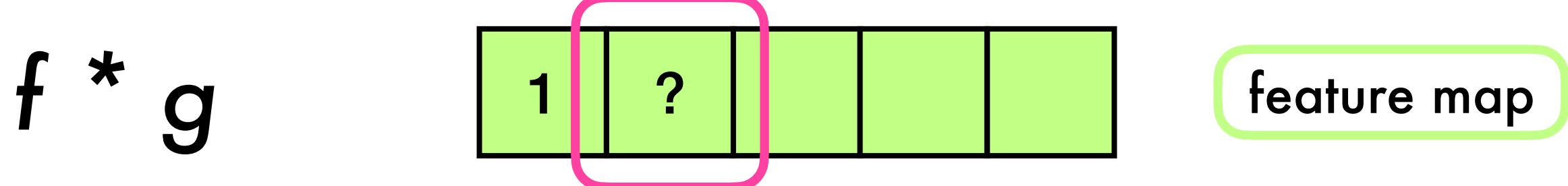
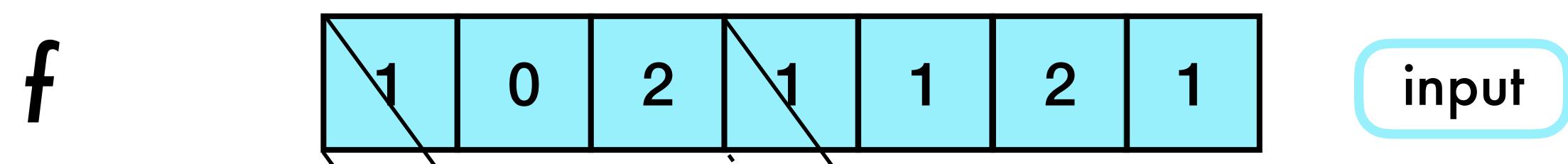
$g$  is a 1-d filter (kernel) of length  $m$

$$n = 7$$

$$m = 3$$

$$i = 1 \quad 2 \quad 3 \quad 4 \quad 5$$

$$j = 1 \quad 2 \quad 3$$



# Example

$f$  is a 1-d input vector of length  $n$

$$(f * g)(i) = \sum_{j=1}^m g(j) \cdot f(i + j - 1)$$

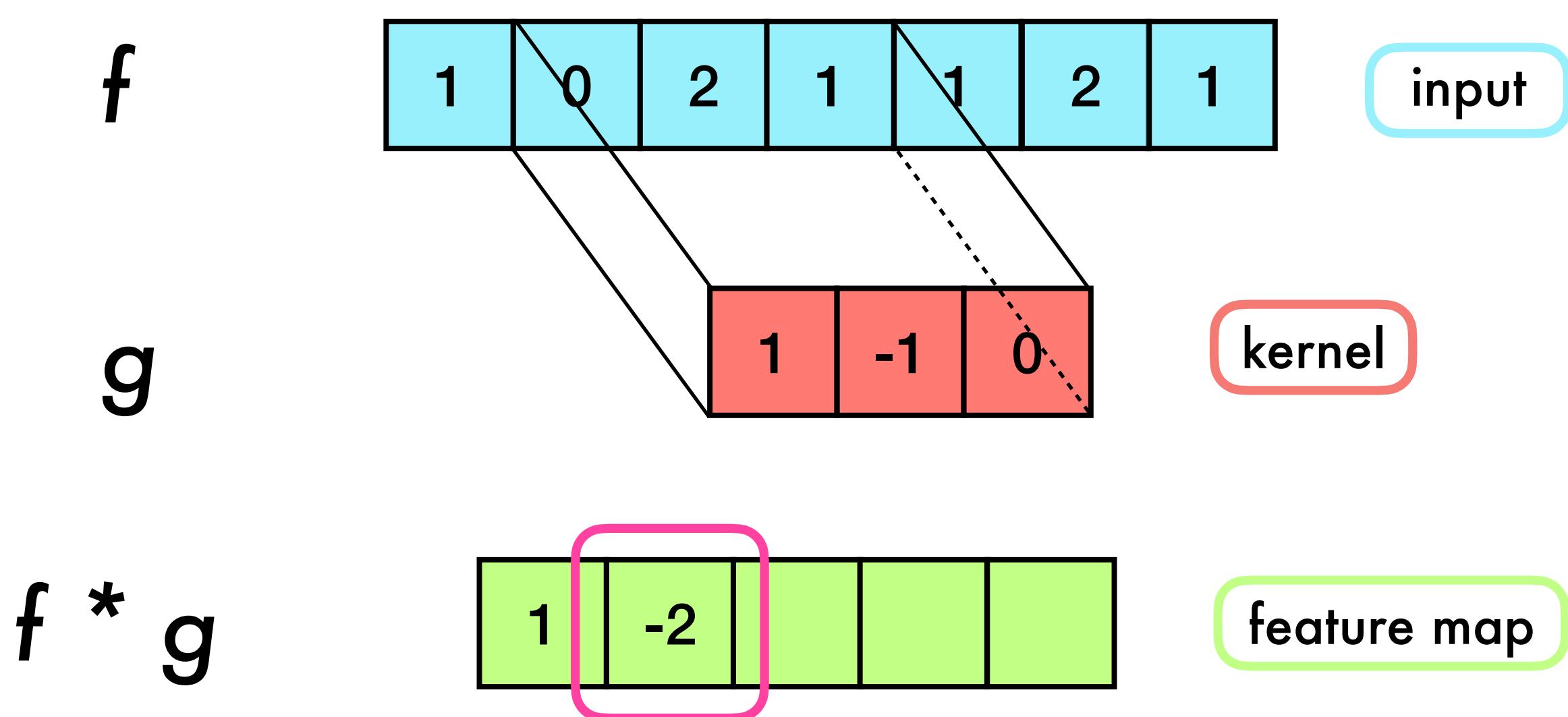
$g$  is a 1-d filter (kernel) of length  $m$

$$n = 7$$

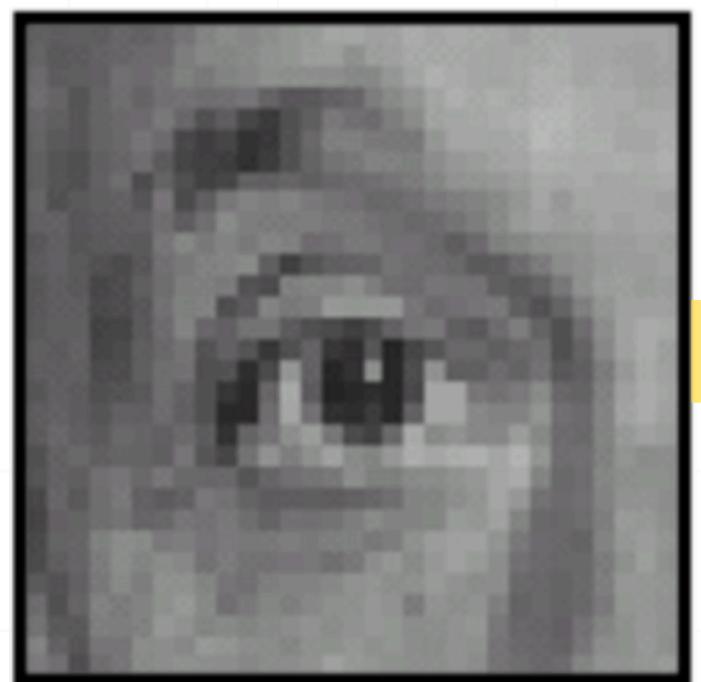
$$m = 3$$

$$i = 1 \boxed{2} 3 \ 4 \ 5$$

$$j = 1 \ 2 \ 3$$

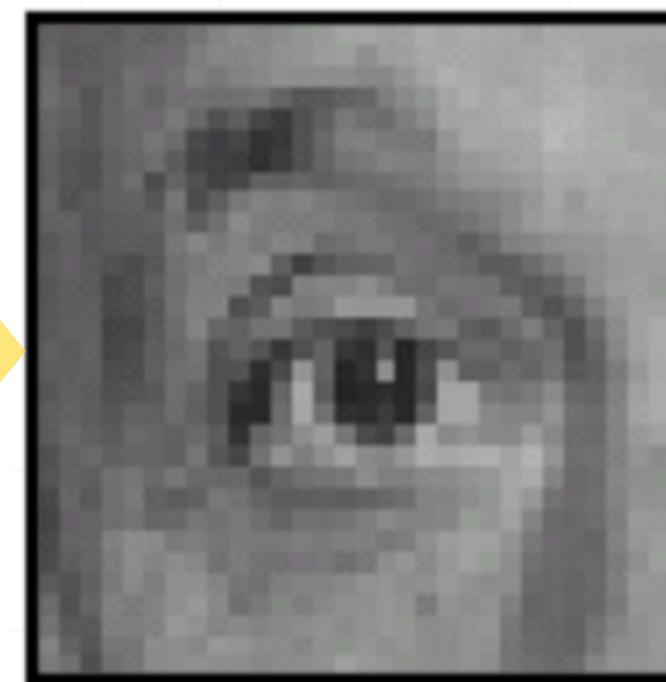


# Example of defined kernels



input

$$\begin{matrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{matrix}$$



original

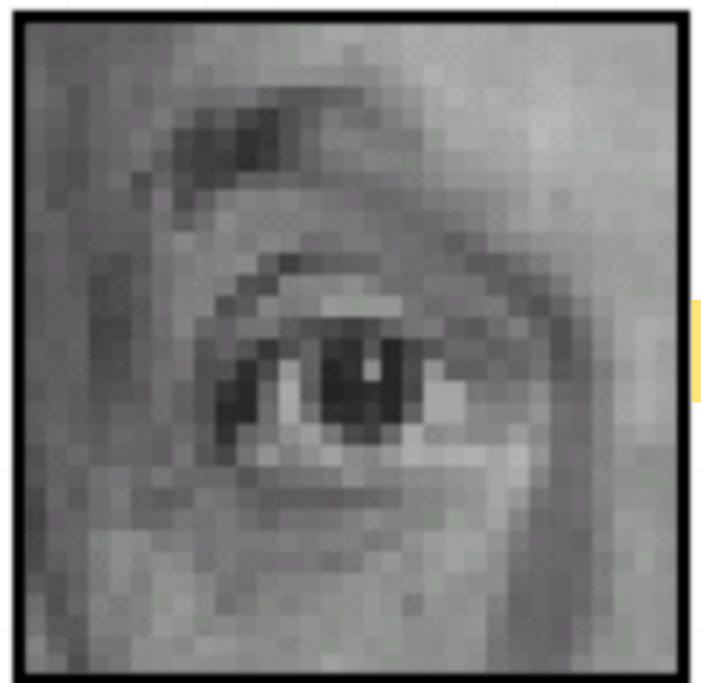


input

$$\begin{matrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{matrix}$$



negative

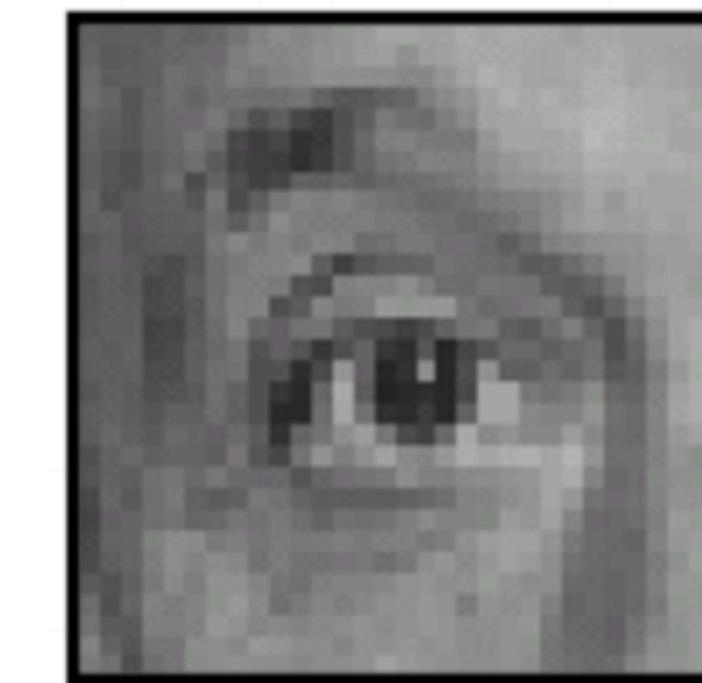


input

$$\begin{matrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$$

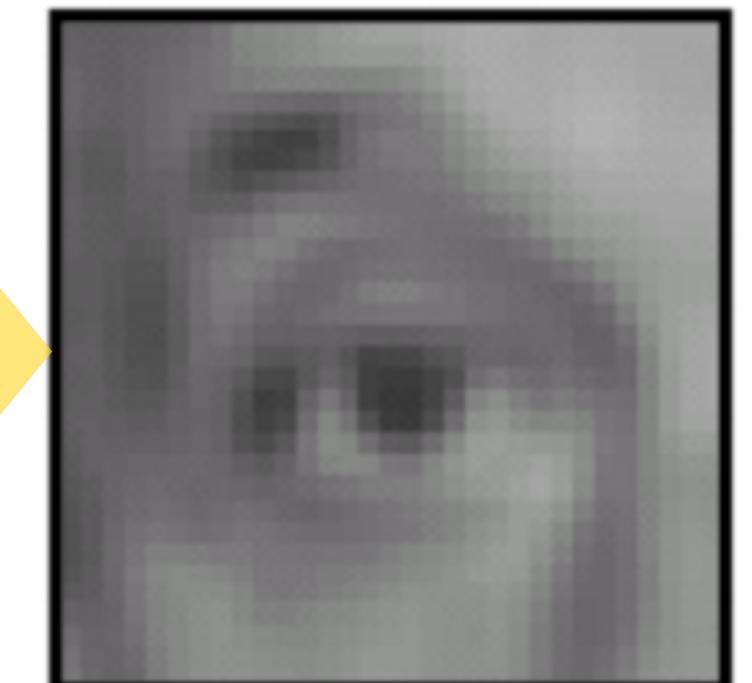


shifted left by 1 pixel



input

$$\begin{matrix} 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \end{matrix}$$



blur  
(with a mean filter)

# Two more concepts

- STRIDE: how much you shift your filter at each step

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

Step 1

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

Step 2

**stride = 1**

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

Step 1

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

Step 2

**stride = 2**

## Two more concepts

- STRIDE: how much you shift your filter at each step
- PADDING: layers of zeros added to the input image

input	filter	feature map													
<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>1</td><td>2</td></tr><tr><td>3</td><td>4</td><td>5</td></tr><tr><td>6</td><td>7</td><td>8</td></tr></table>	0	1	2	3	4	5	6	7	8	$\ast$	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>1</td></tr><tr><td>2</td><td>3</td></tr></table>	0	1	2	3
0	1	2													
3	4	5													
6	7	8													
0	1														
2	3														
	=	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>19</td><td>25</td></tr><tr><td>37</td><td>43</td></tr></table>	19	25	37	43									
19	25														
37	43														

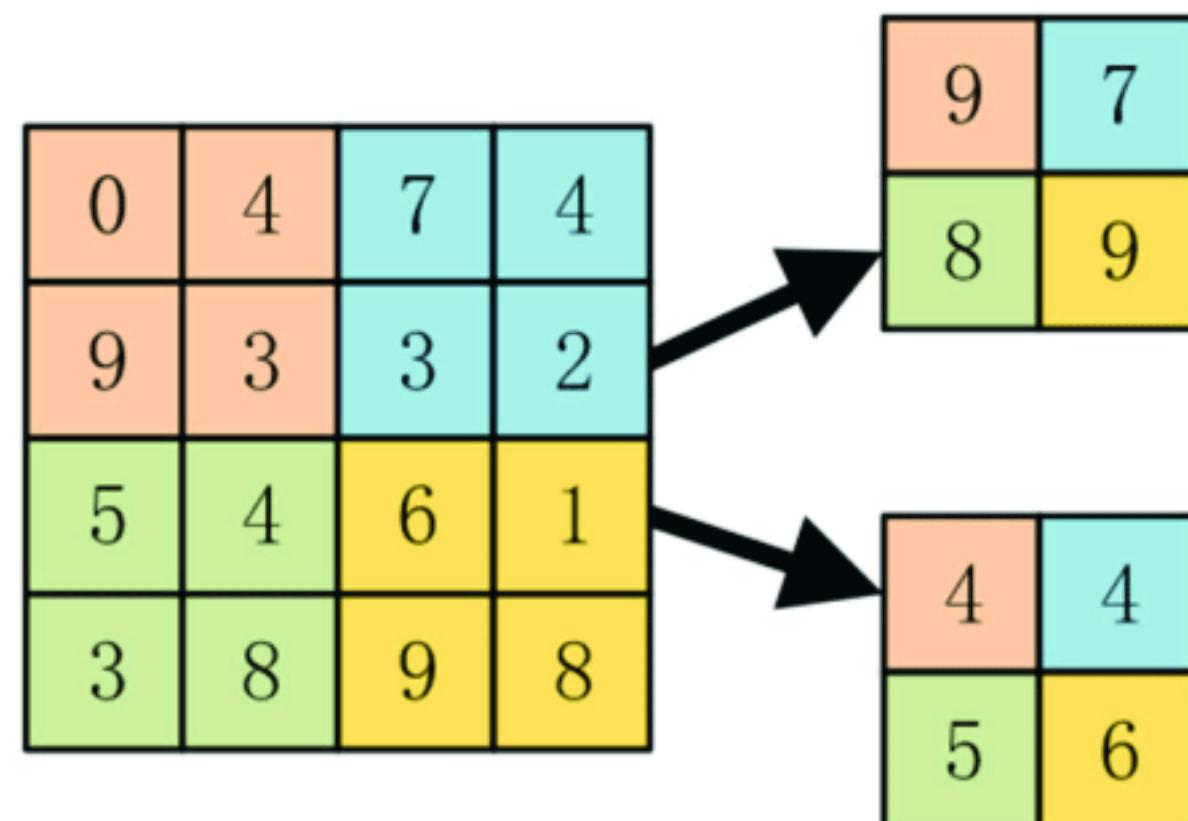
padding = 0

<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>1</td><td>2</td><td>0</td></tr><tr><td>0</td><td>3</td><td>4</td><td>5</td><td>0</td></tr><tr><td>0</td><td>6</td><td>7</td><td>8</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr></table>	0	0	0	0	0	0	0	1	2	0	0	3	4	5	0	0	6	7	8	0	0	0	0	0	0	$\ast$	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>1</td></tr><tr><td>2</td><td>3</td></tr></table>	0	1	2	3
0	0	0	0	0																											
0	0	1	2	0																											
0	3	4	5	0																											
0	6	7	8	0																											
0	0	0	0	0																											
0	1																														
2	3																														
	=	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>3</td><td>8</td><td>4</td></tr><tr><td>9</td><td>19</td><td>25</td><td>10</td></tr><tr><td>21</td><td>37</td><td>43</td><td>16</td></tr><tr><td>6</td><td>7</td><td>8</td><td>0</td></tr></table>	0	3	8	4	9	19	25	10	21	37	43	16	6	7	8	0													
0	3	8	4																												
9	19	25	10																												
21	37	43	16																												
6	7	8	0																												

padding = 1

# Pooling

- Reduces the spatial size of the feature maps
- Decrease the computational power required to process the data through dimensionality reduction by **extracting dominant features**



output of the convolution

## MAX POOLING

- Returns the **maximum** value from the portion covered by the kernel
- “Did you see this feature anywhere in the image?”

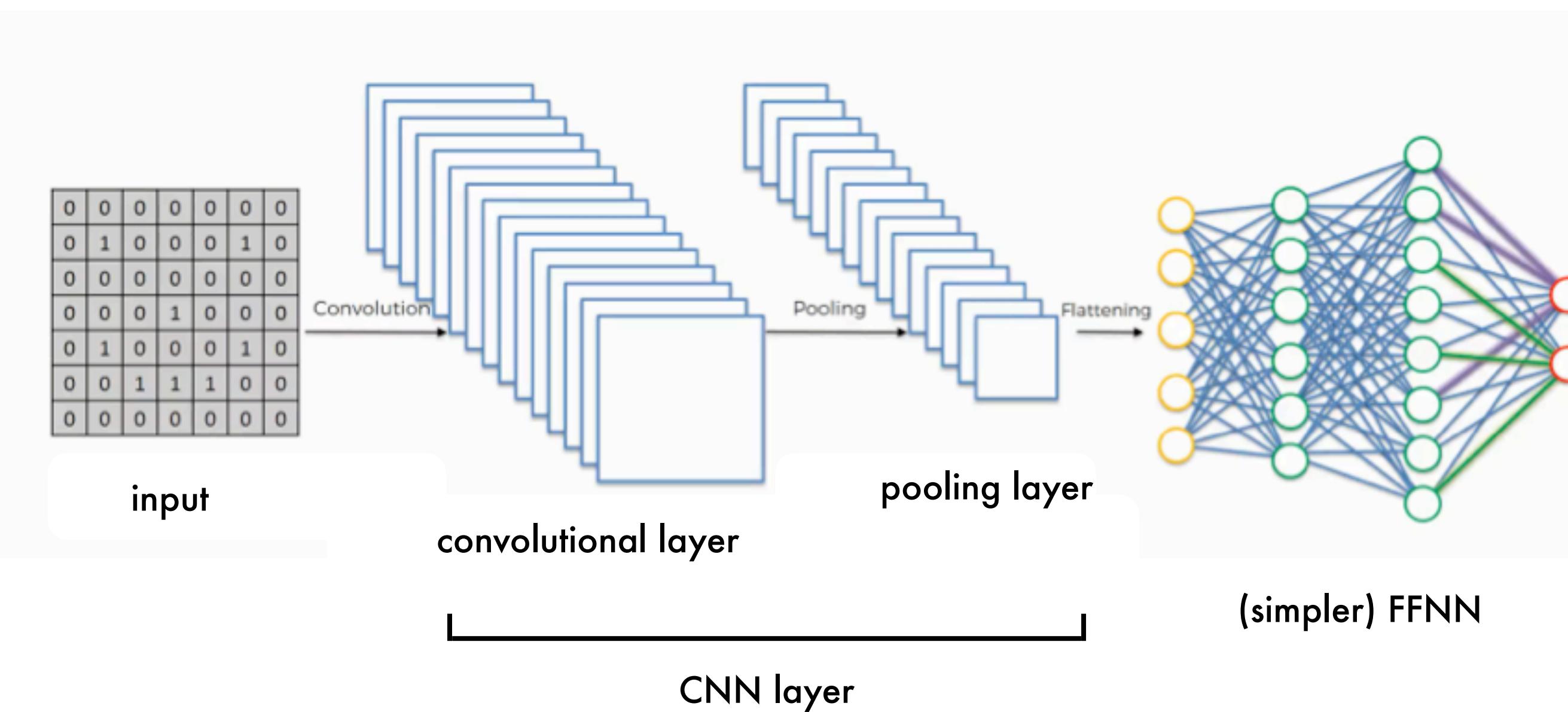
## AVERAGE POOLING

- Returns the **average** of all values from the portion covered by the kernel
- “How prevalent is this feature over the entire range?”

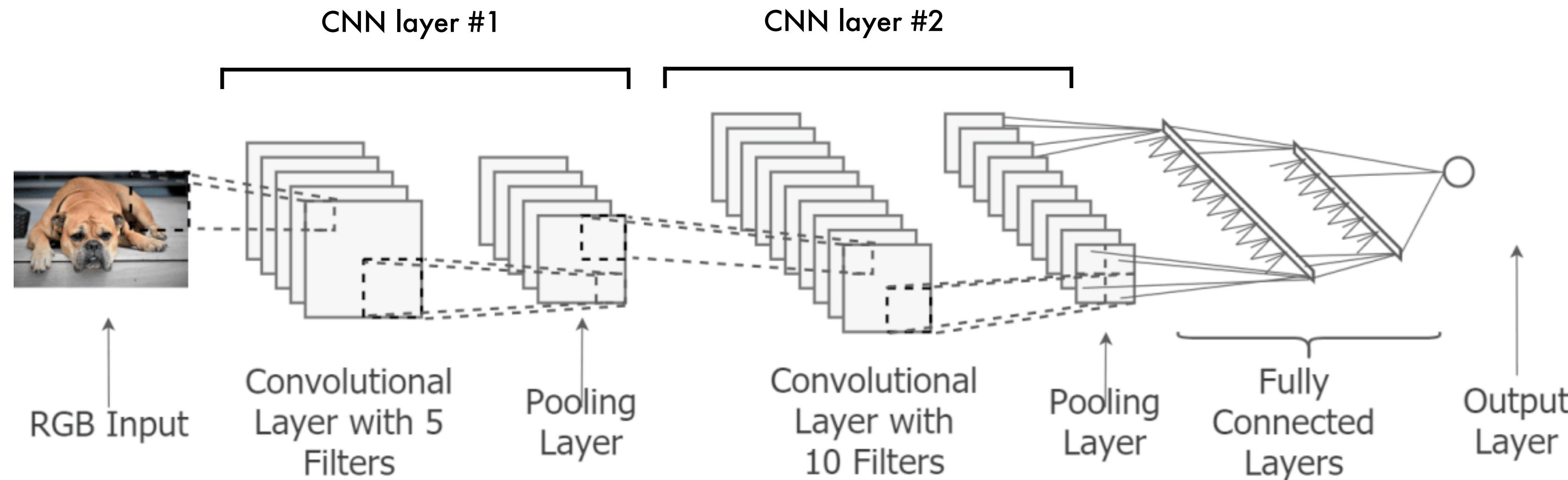
# Pooling

- Pooling layers do not add parameters to the network
- Often, max-pooling performs better than average pooling
- A **Convolutional Layer** + a **Pooling Layer** form the  $i$ -th layer of a CNN

# Overview of a CNN



# Multiple convolutional layers



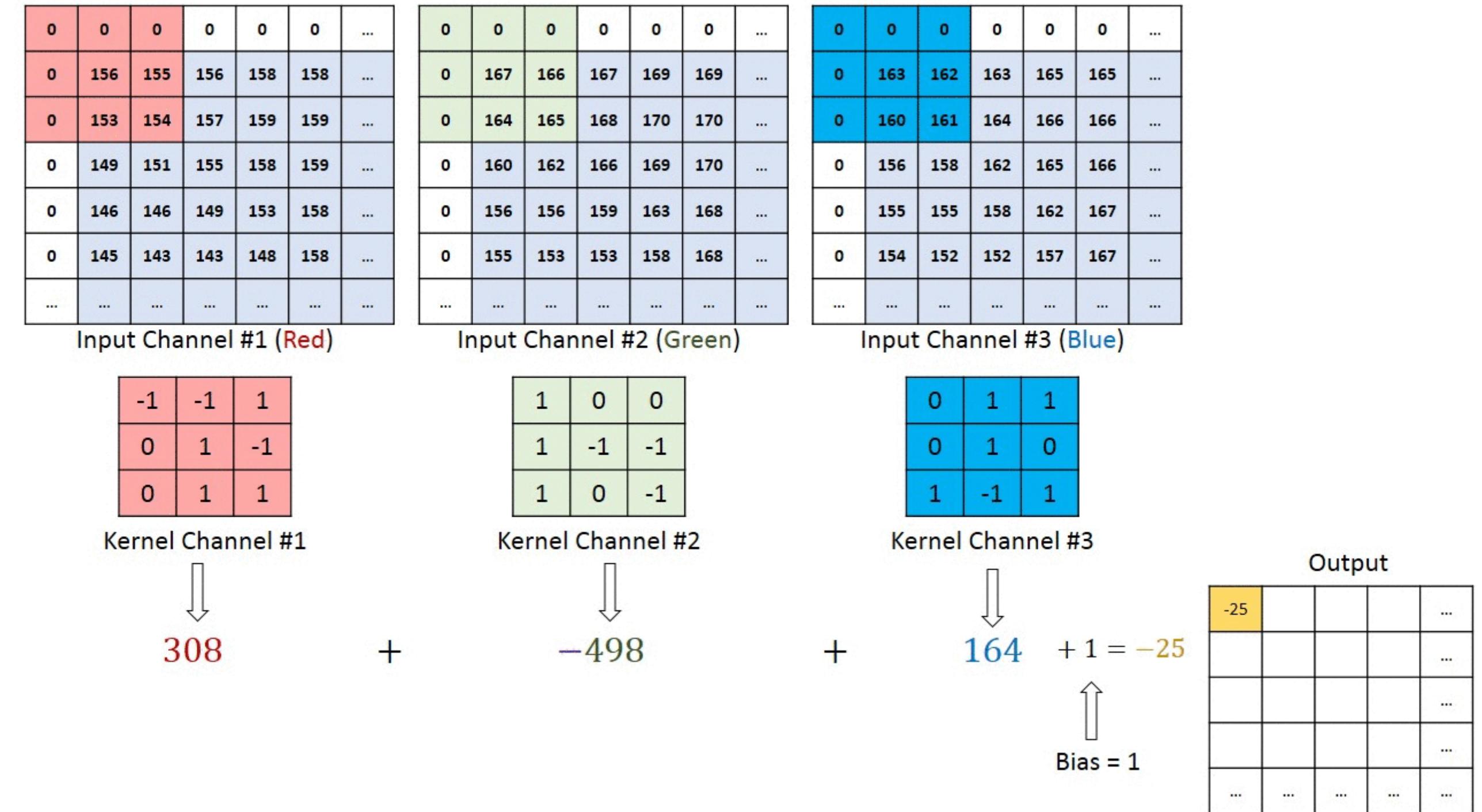
- Convolutional layers can be applied both to the input data and to the output of other layers
- The stacking of convolutional layers allows a hierarchical decomposition of the input

# Multiple channels

2-channels input      kernels

$$\begin{matrix} 0 & 1 & 2 \\ 3 & 4 & 5 \\ 6 & 7 & 8 \end{matrix} \times \begin{matrix} 1 & 2 \\ 4 & 5 \end{matrix} = \begin{matrix} 56 & 72 \\ 104 & 120 \end{matrix}$$

$$\begin{matrix} 0 & 1 & 2 \\ 3 & 4 & 5 \\ 6 & 7 & 8 \end{matrix} \times \begin{matrix} 3 & 4 \\ 7 & 8 \end{matrix} = \begin{matrix} 56 & 72 \\ 104 & 120 \end{matrix}$$



- Different “views” of the input

**Questions about CNNs components/structure?**

**Okay...but we are working with text!**

# CNNs for text

## word2vec

the →	0.2	0.4	-0.1
good →	0.7	-0.5	0.3
movie →	0.1	0.2	0.6

# CNNs for text

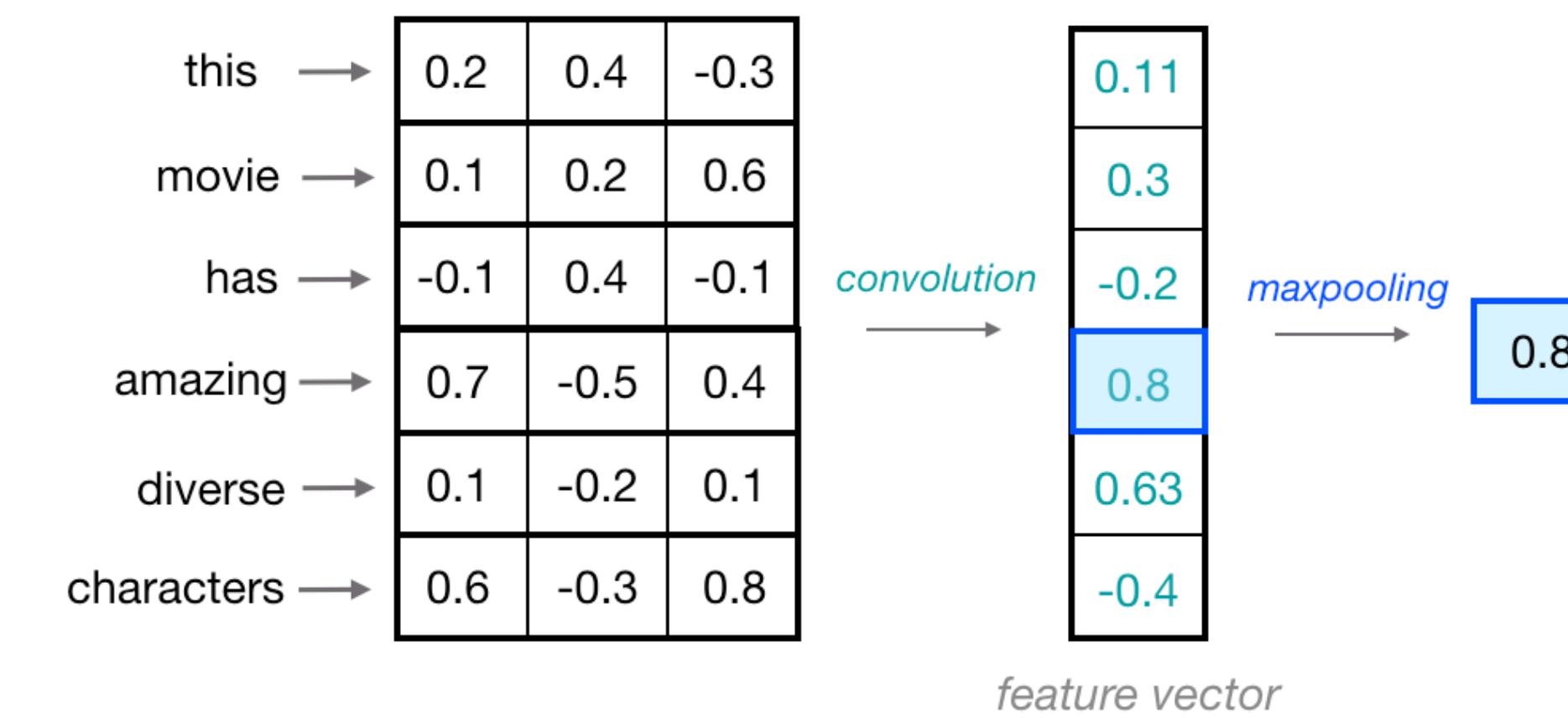
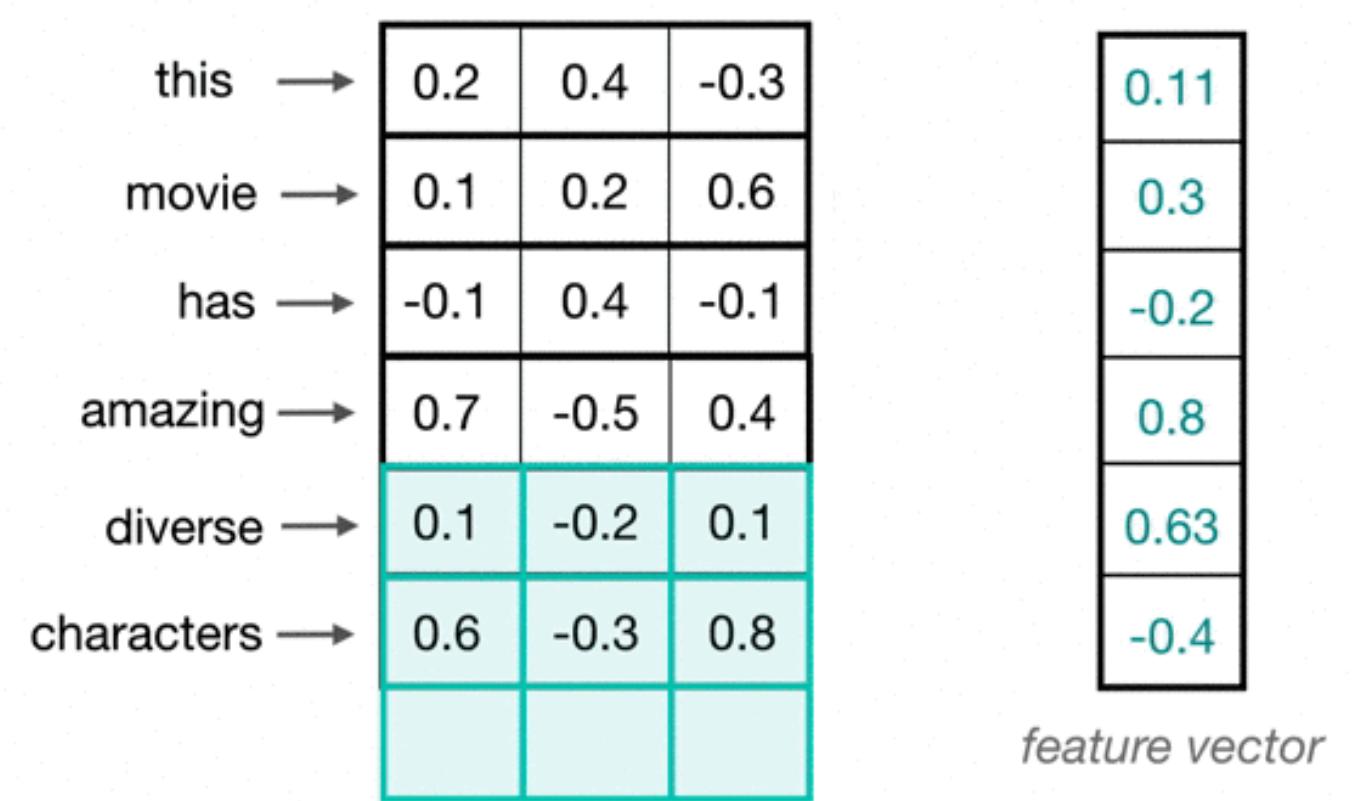
this →	0.2	0.4	-0.3
movie →	0.1	0.2	0.6
has →	-0.1	0.4	-0.1
amazing →	0.7	-0.5	0.4
diverse →	0.1	-0.2	0.1
characters →	0.6	-0.3	0.8

# CNNs for text

this →	0.2	0.4	-0.3
movie →	0.1	0.2	0.6
has →	-0.1	0.4	-0.1
amazing →	0.7	-0.5	0.4
diverse →	0.1	-0.2	0.1
characters →	0.6	-0.3	0.8



# CNNs for text



# CNNs for text

- 1-direction convolution: analysis of the text over time
- kernel width = embedding length

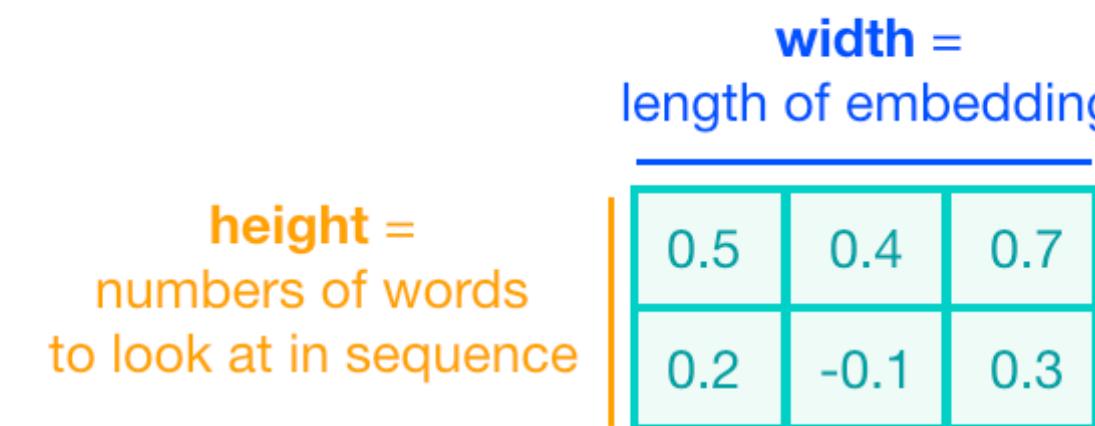
width = length of embedding		
0.5	0.4	0.7
0.2	-0.1	0.3

this →	0.2	0.4	-0.3
movie →	0.1	0.2	0.6
has →	-0.1	0.4	-0.1
amazing →	0.7	-0.5	0.4
diverse →	0.1	-0.2	0.1
characters →	0.6	-0.3	0.8

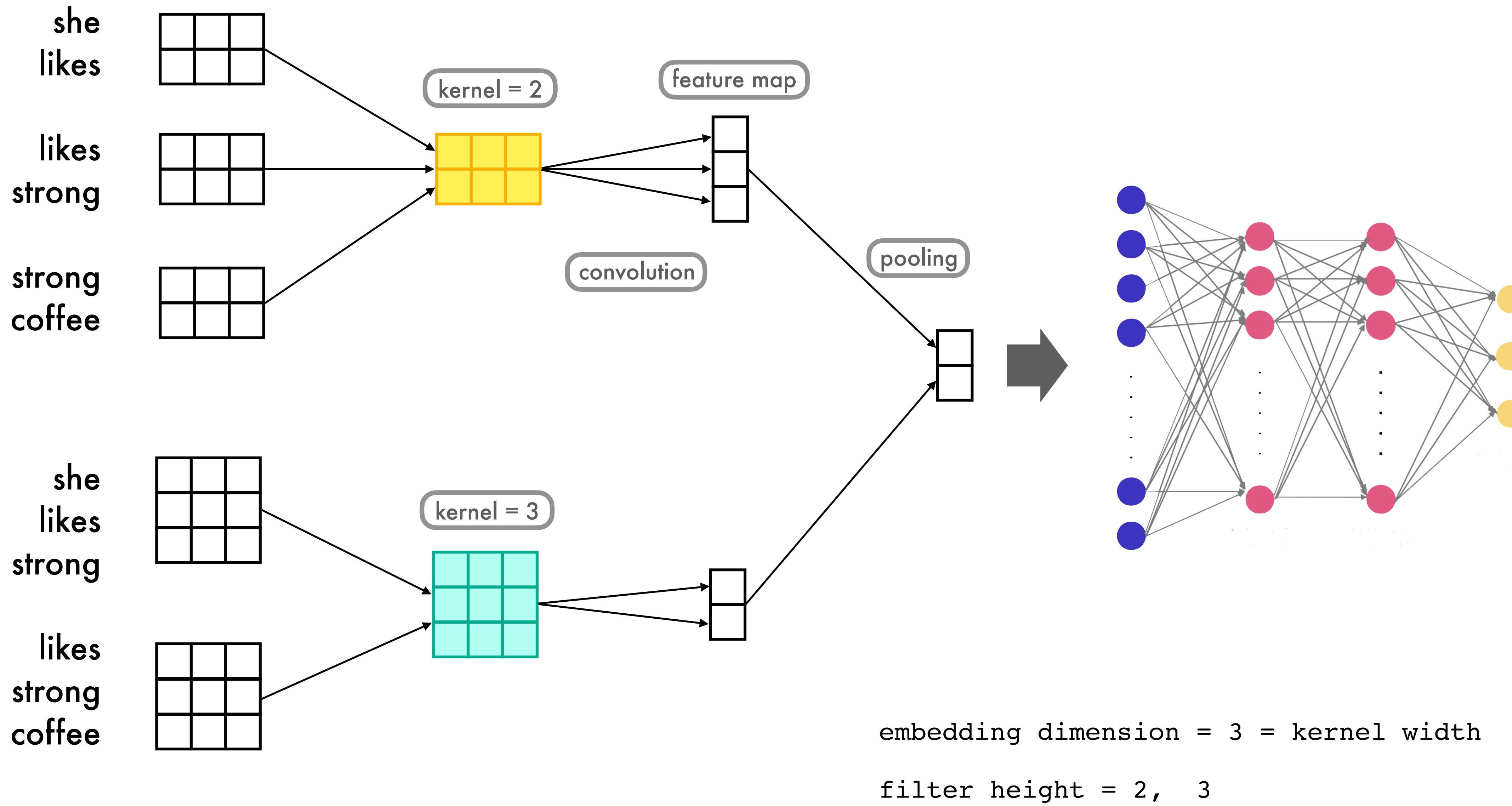


# CNNs for text

- CNNs were introduced in NLP by Collobert et al. (2011) and later by Kim (2014) and Kalchbrenner et al. (2014)
- The intention is to let the network focus on the **most important “features” in the text**, regardless of their location
- The idea is to apply a non-linear (learned) function over each instantiation of a **k-word sliding window over the text**



## Sample: She likes strong coffee



# Word representations (embeddings)

Initialisation type:

- Traditional count-based methods
- word2vec: skip-gram, cbow
- Randomly initialised

# Word representations (embeddings)

Initialisation type:

- Traditional count-based methods
- word2vec: skip-gram, cbow
- Randomly initialised

Embedding type:

- Static
- Fine-tuned (non-static)

**Questions about CNNs for NLP?**

# Relation Extraction

In 2003, the **Stade de France** was the primary site of the **2003 World Championship in Athletics**.

**LOCATION**

**PHYSICAL**

**EVENT**

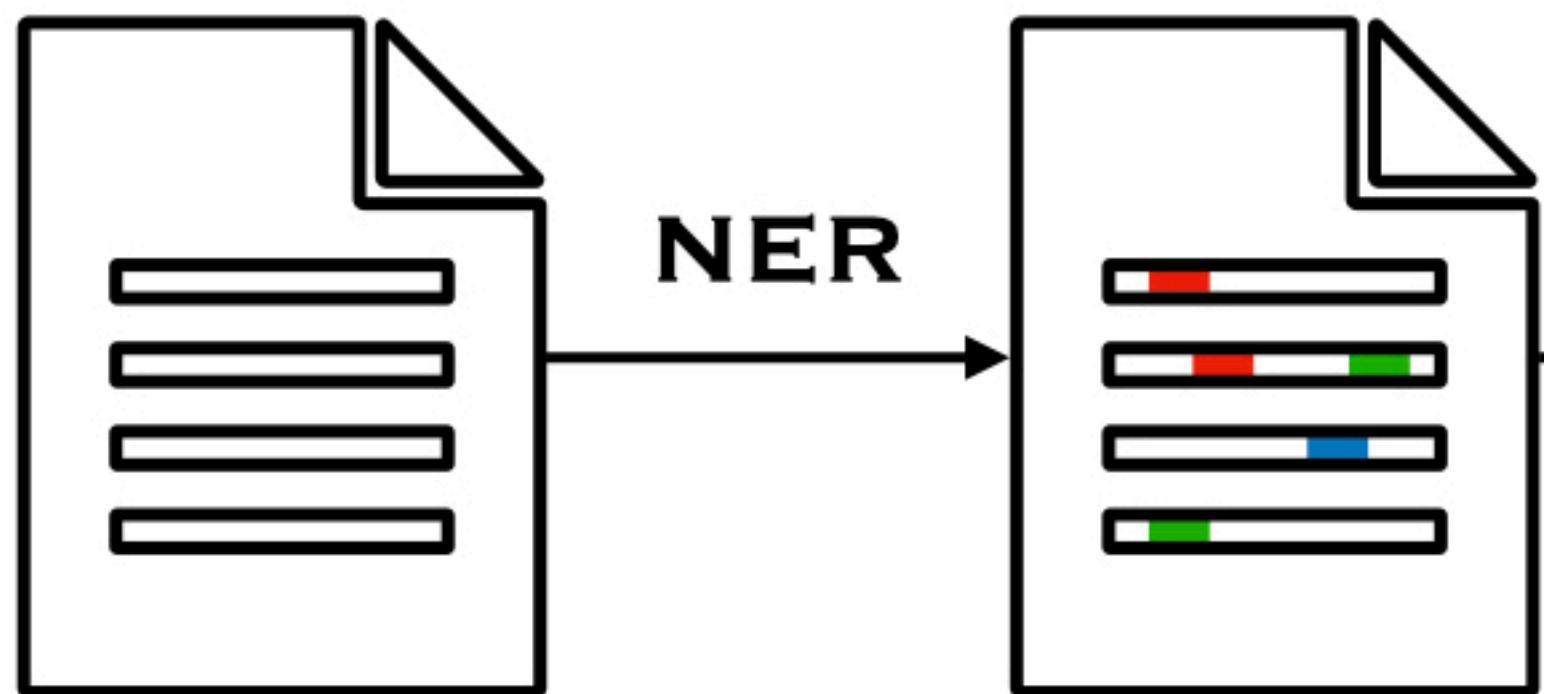
# Relation Extraction



**RAW TEXT**

*In 2003, the Stade de France was the primary site of the 2003 World Championship in Athletics.*

# Relation Extraction



**RAW TEXT**

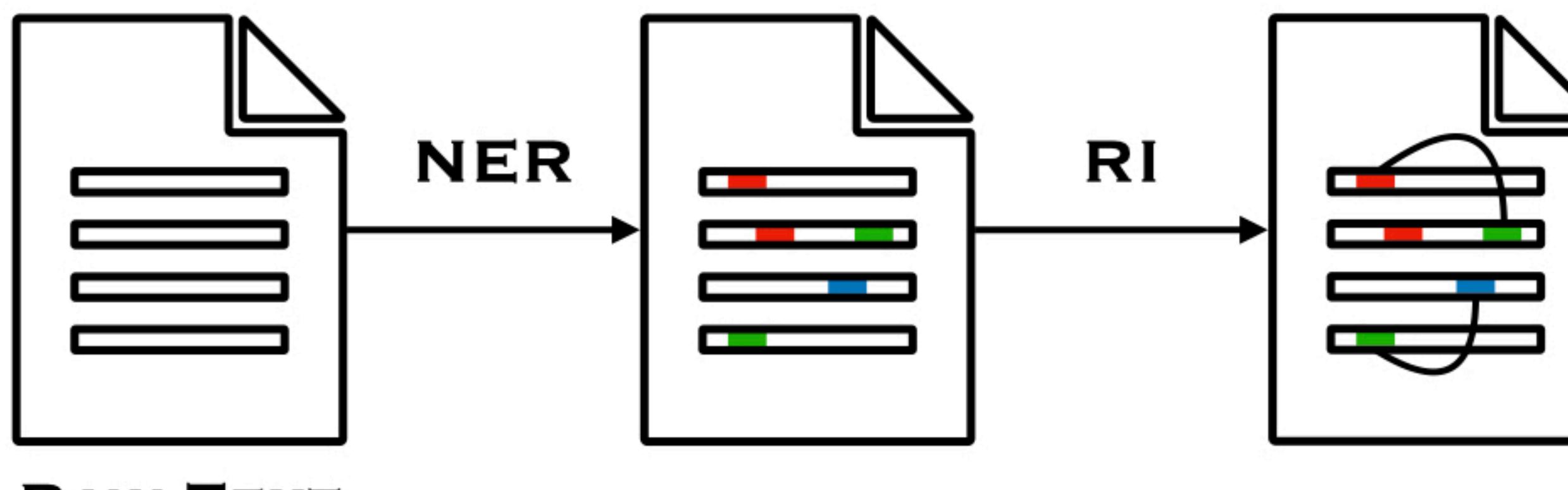
NAMED ENTITY RECOGNITION

*In 2003, the Stade de France was the primary site of the 2003 World Championship in Athletics.*

**LOCATION**

**EVENT**

# Relation Extraction



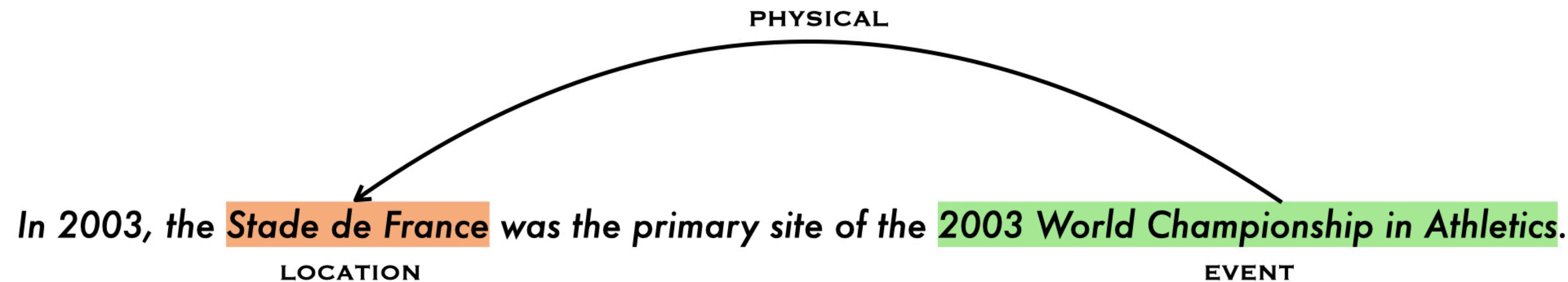
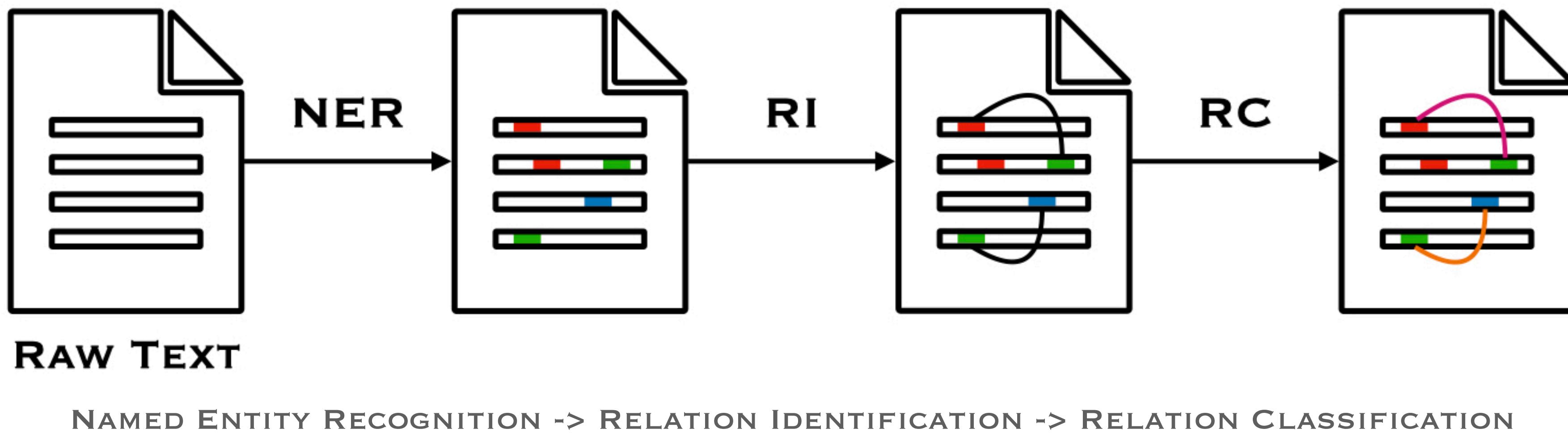
NAMED ENTITY RECOGNITION -> RELATION IDENTIFICATION

In 2003, the Stade de France was the primary site of the 2003 World Championship in Athletics.

LOCATION

EVENT

# Relation Extraction



# How to evaluate Relation Extraction?

A predicted relation is considered correct if ...

- *strict evaluation:*
- *boundaries evaluation:*
- *relaxed evaluation:*

... and the predicted relation type is correct

# How to evaluate Relation Extraction?

A predicted relation is considered correct if ...

- *strict evaluation*: the boundaries of the two entity spans, as well as the entity types are correct
  - *boundaries evaluation*:
  - *relaxed evaluation*:
- ... and the predicted relation type is correct

# How to evaluate Relation Extraction?

A predicted relation is considered correct if ...

- ***strict evaluation:*** the boundaries of the two entity spans, as well as the entity types are correct
- ***boundaries evaluation:*** the boundaries of the two spans are correct
- ***relaxed evaluation:***  
... and the predicted relation type is correct

# How to evaluate Relation Extraction?

A predicted relation is considered correct if ...

- ***strict evaluation***: the boundaries of the two entity spans, as well as the entity types are correct
- ***boundaries evaluation***: the boundaries of the two spans are correct
- ***relaxed evaluation***: at least one token per entity span is correctly typed

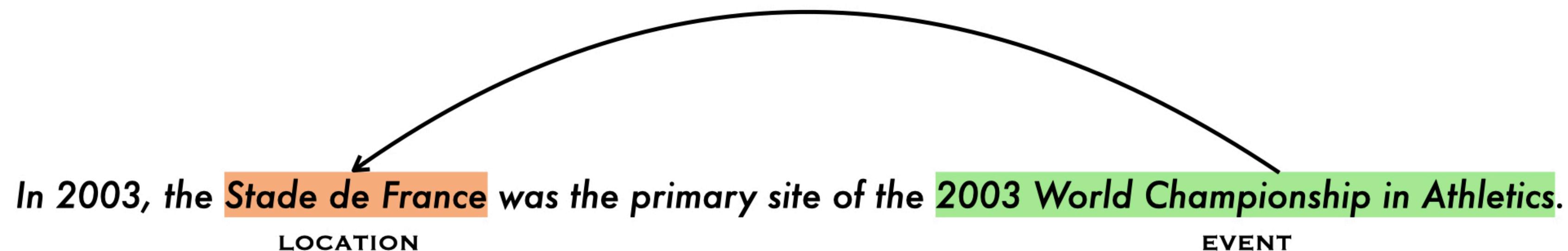
... and the predicted relation type is correct

# Metrics

- *macro-F1*: how much the model learn of each class
- *micro-F1*: how the model perform in general on the task

# Relation Classification

- Input information:
  - Pair of ordered entities
- Output:
  - Relation type from a given label set {PHYSICAL, TEMPORAL, PART-OF, TYPE-OF}



# Relation Classification

In 2003 , the Stade de France was the primary site of the 2003 World Championship in Athletics .

In [REDACTED]

2003 [REDACTED]

[REDACTED]

[REDACTED]

the [REDACTED]

Stade de France [REDACTED]

was [REDACTED]

the [REDACTED]

primary [REDACTED]

site [REDACTED]

of [REDACTED]

the [REDACTED]

2003 World Championship in Athletics [REDACTED]

.

[REDACTED]

In 2003, the Stade de France was the primary site of the 2003 World Championship in Athletics.

LOCATION

EVENT

# Relation Classification

In 2003, the Stade de France was the primary site of the 2003 World Championship in Athletics.

LOCATION

EVENT

In [REDACTED]  
2003  
,  
the  
**Stade de France**  
was [REDACTED]  
the [REDACTED]  
primary [REDACTED]  
site [REDACTED]  
of [REDACTED]  
the [REDACTED]  
**2003 World Championship in Athletics**  
.

# Relation Classification

In [REDACTED]  
2003  
,  
the  
**Stade de France**  
was  
the  
primary  
site  
of  
the  
**2003 World Championship in Athletics**  
.

In 2003, the **Stade de France** was the primary site of the **2003 World Championship in Athletics**.

LOCATION

EVENT

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

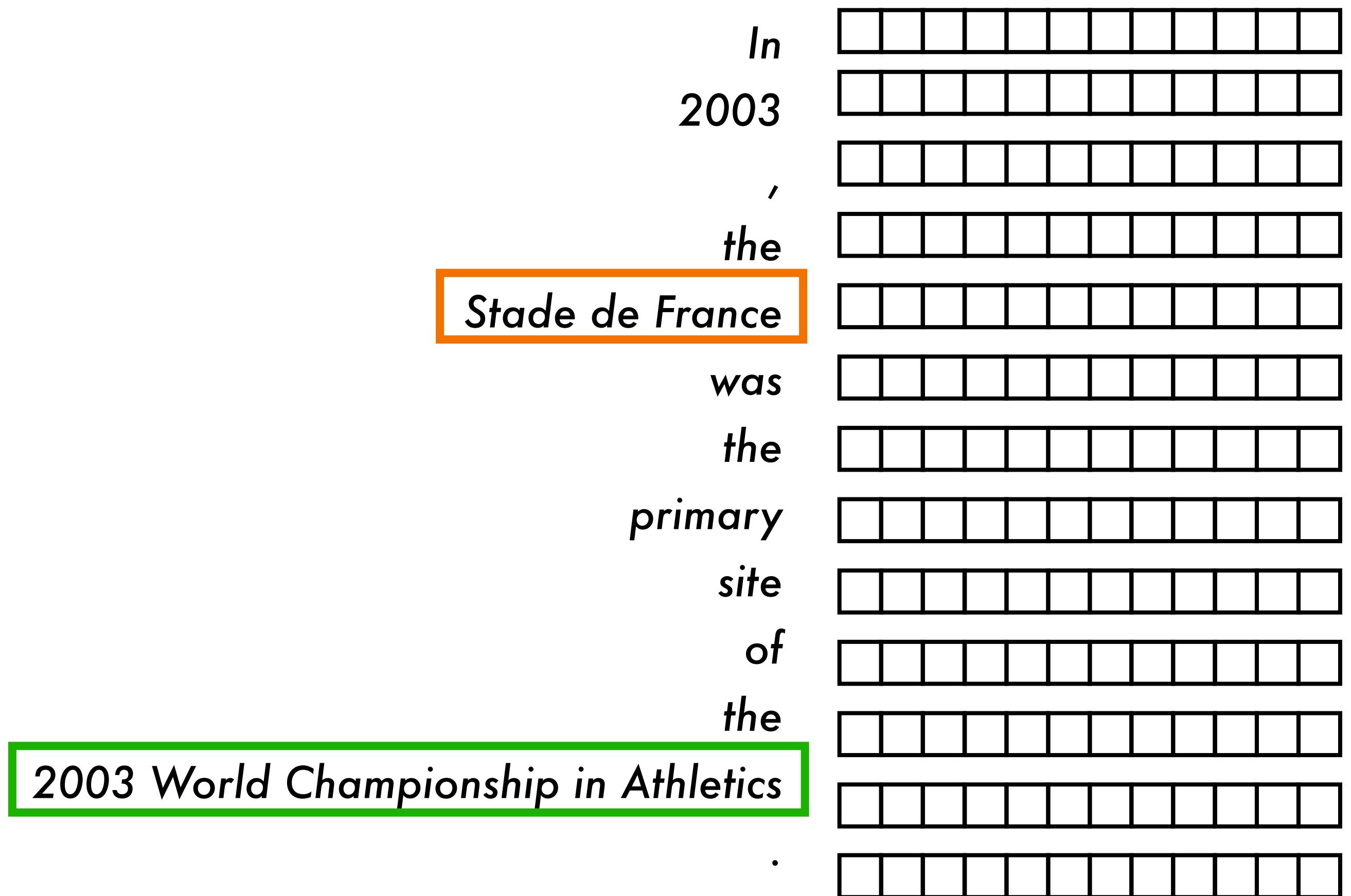
Information about which  
are the entities

# Relation Classification

In 2003, the Stade de France was the primary site of the 2003 World Championship in Athletics.

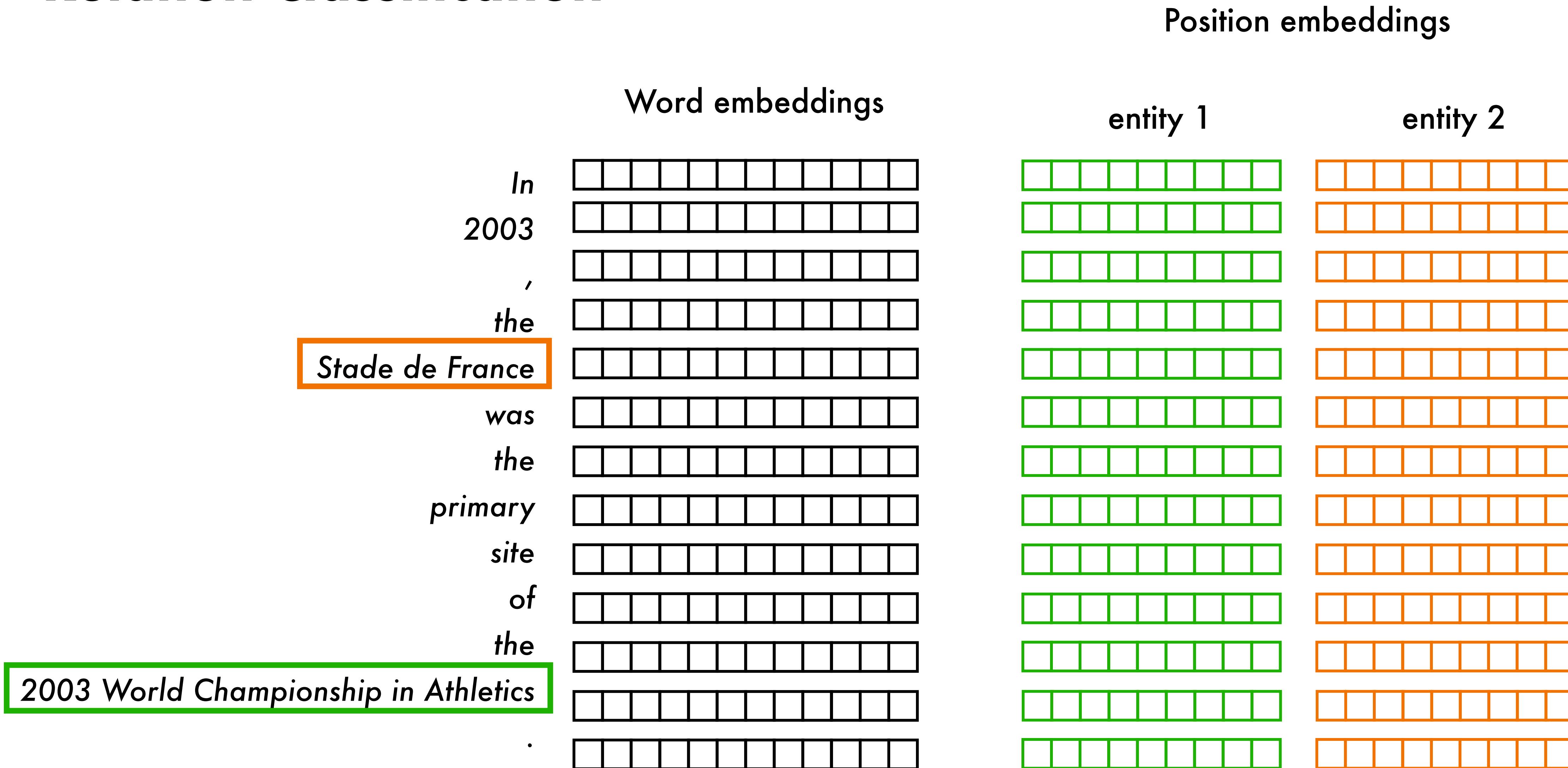
LOCATION

EVENT



More fine-grained:  
Relative distance of each  
token from the 2 entities!

# Relation Classification

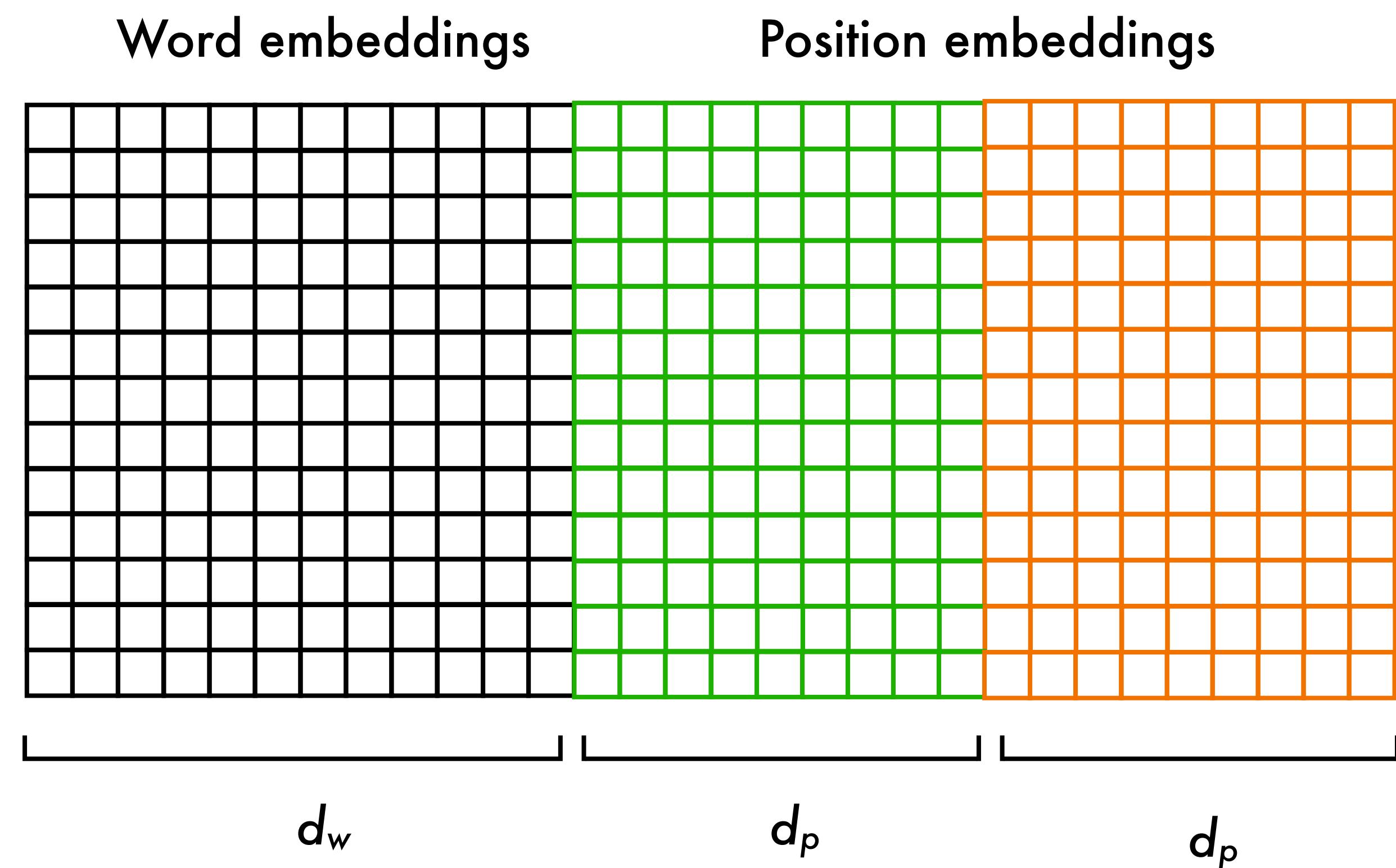


# Relation Classification

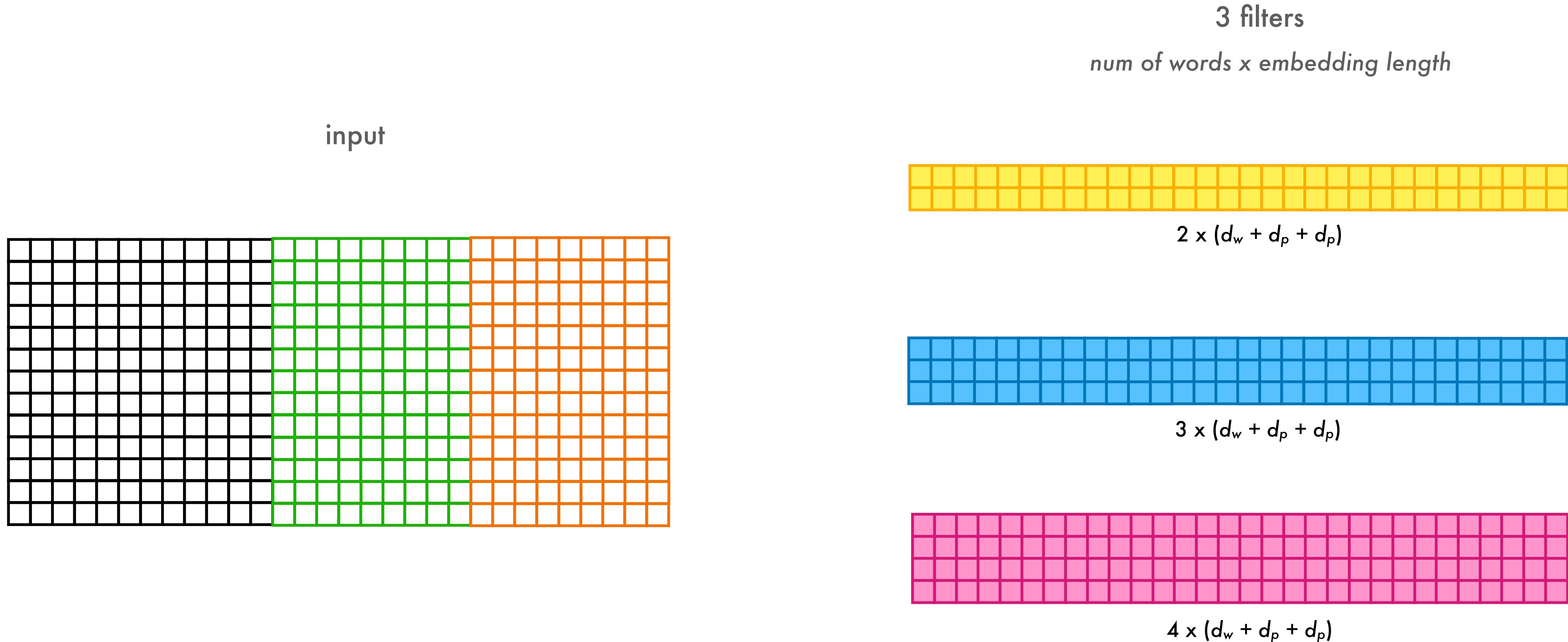
In 2003, the Stade de France was the primary site of the 2003 World Championship in Athletics.

LOCATION EVENT

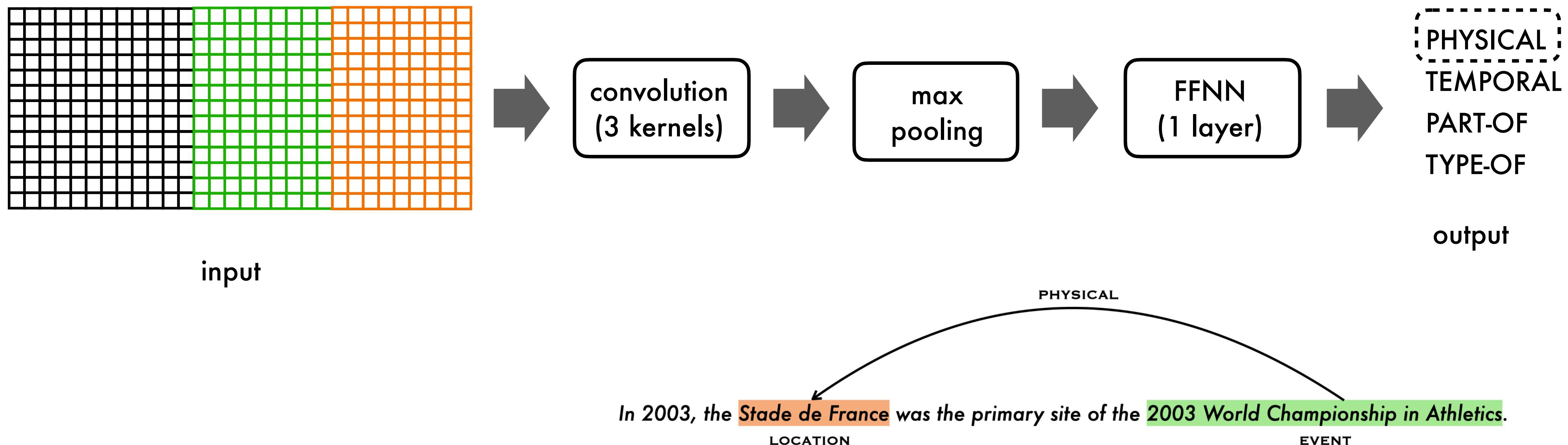
In  
2003  
'  
the  
**Stade de France**  
was  
the  
primary  
site  
of  
the  
**2003 World Championship in Athletics**  
.



# Relation Classification



# Relation Classification model example



**Questions?**

# **Sum up of the lecture**

1. Recap of FFNN and their limitations
2. Introduction to Convolutional Neural Networks (CNNs)
3. CNNs for NLP
4. Case study with CNNs

# References

- Yoav Goldberg (Primer), 2015. A Primer on Neural Networks for Natural Language Processing. <http://u.cs.biu.ac.il/~yogo/nlp.pdf>
- Kim 2014. Convolutional Neural Networks for Sentence Classification. <https://www.aclweb.org/anthology/D14-1181/>
- Cezanne Camacho blogpost. [https://cezannec.github.io/CNN\\_Text\\_Classification/](https://cezannec.github.io/CNN_Text_Classification/)