

# Natural Language Processing

Lecture 19

Dirk Hovy

[dirk.hovy@unibocconi.it](mailto:dirk.hovy@unibocconi.it)

 @dirk\_hovy

Bocconi

# Goals for Today

- Learn about **recurrent neural network architectures**
- Learn about **convolutional neural network architectures**
- Understand the concept of **convolution** and **pooling**
- Understand the **difference between recurrent and convolutional networks**
- Understand the **attention mechanism**

# Recurring Matters



IMAGE CAPTIONS

Bocconi



# Long-Term Trouble

SUBJECT

"Wenn er aber auf der Strasse der in Sammt und Seide gehüllten jetzt sehr ungenirt nach der neusten Mode gekleideten Regierungsräthen begegnet."

VERB

Mark Twain, *The Awful German Language*

# Long-Term Trouble

## *NEGATION*

This is **not** in any sense of the word a **funny** movie.

# Sequence Tagging

PRON	VERB	ADP	DET	???	PUNCT
I	went	to	the	show	.

show {VERB, NOUN}

PART show

show

PRON show

show

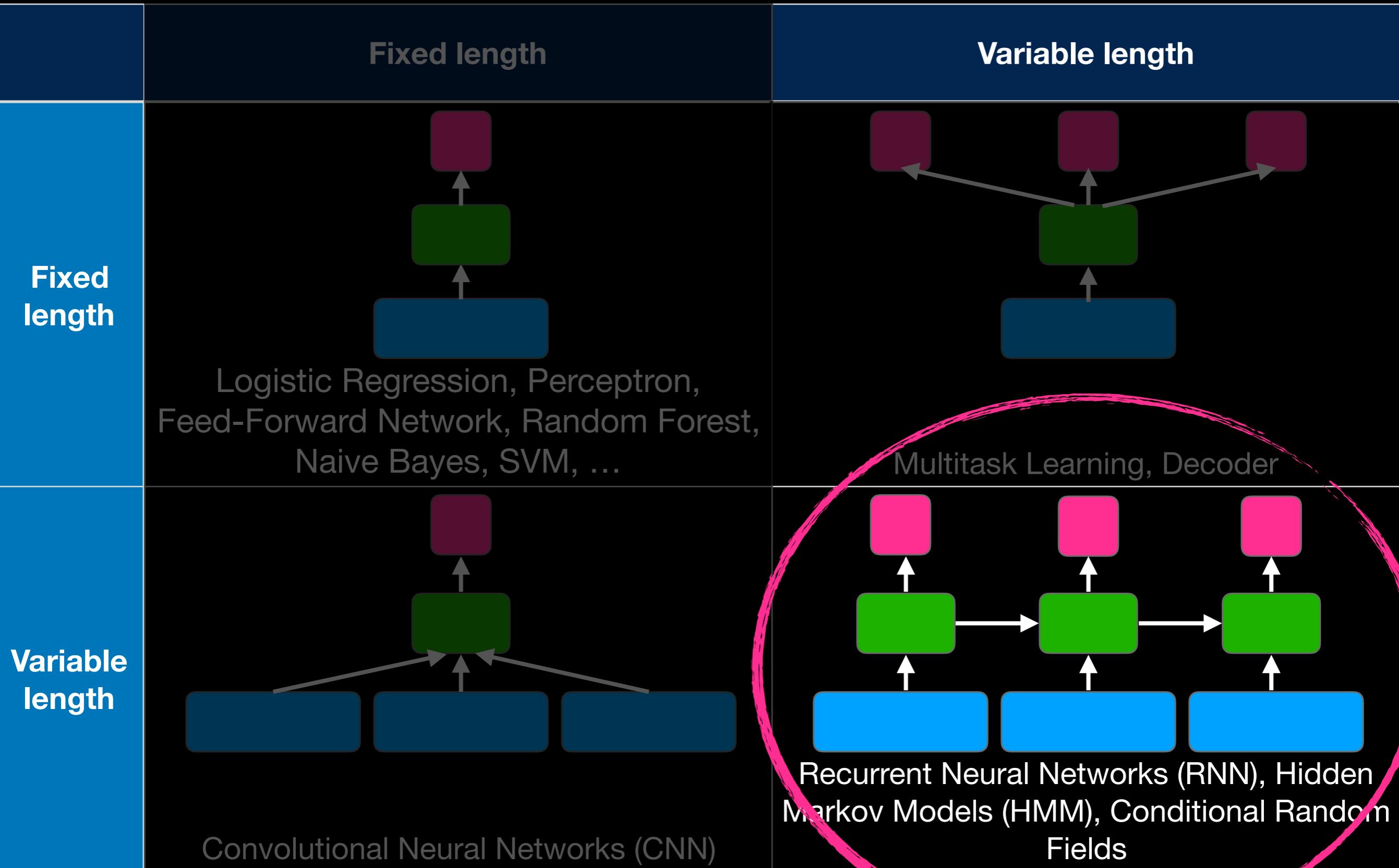
show  
show  
show  
show

DET

ADJ

Structured prediction: depends on the POS of a previous word

# Types of Text Classification

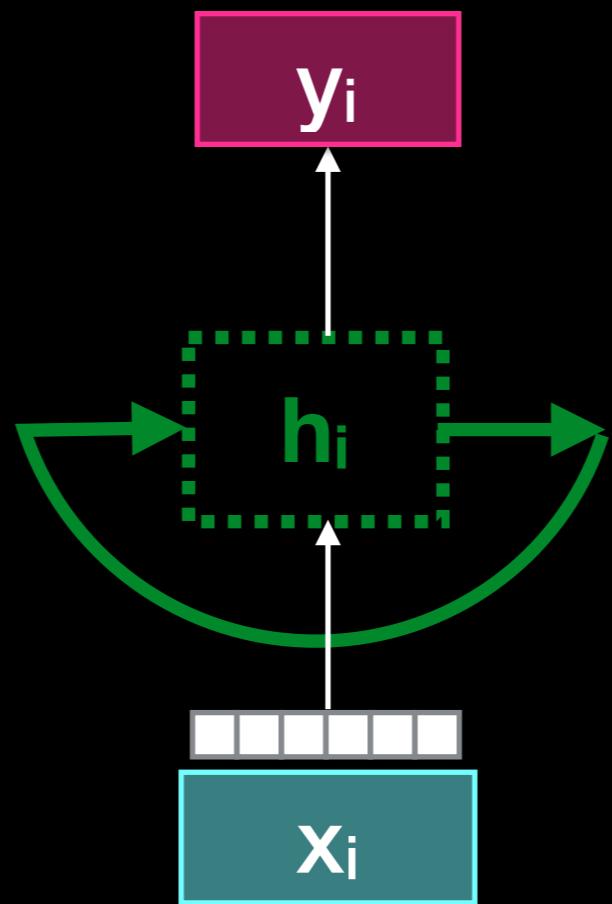


# Recurrent Networks

# Recurrence

$$y_i = f(h_i)$$

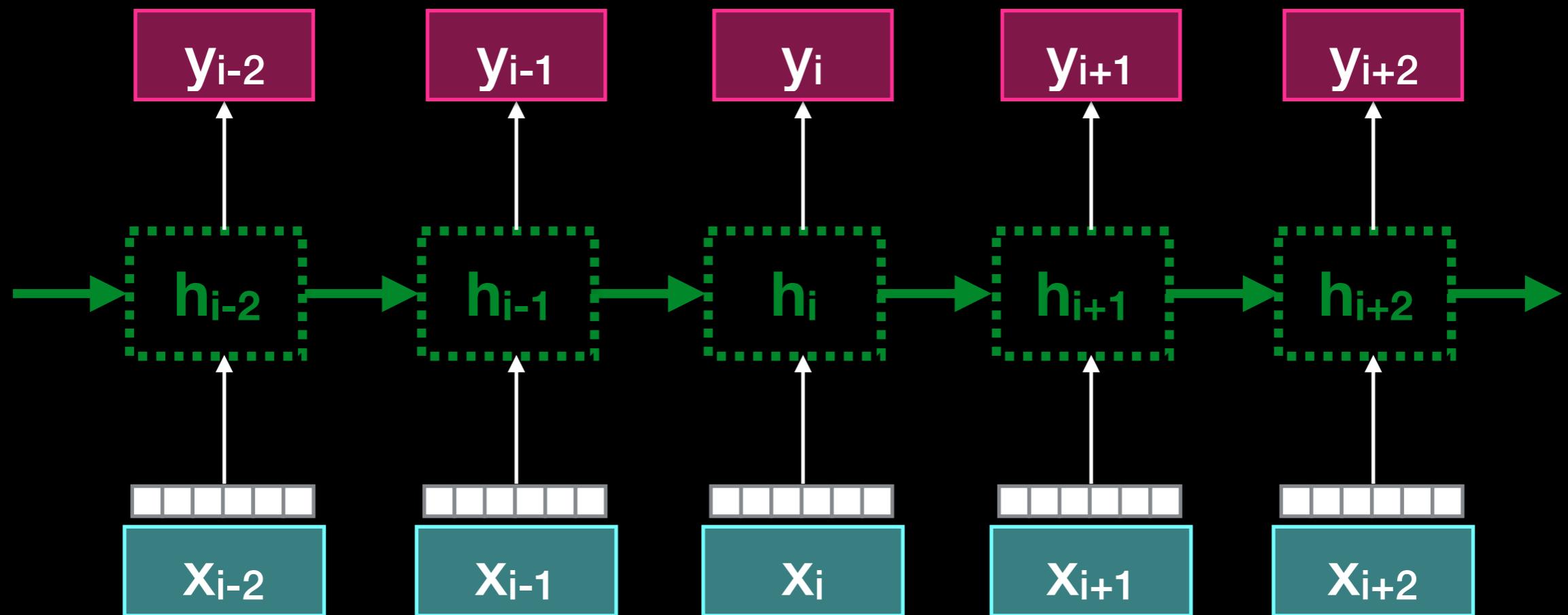
$$h_i = s(h_{i-1}, x_i)$$



# ...Unrolled

$$y_i = f(h_i)$$

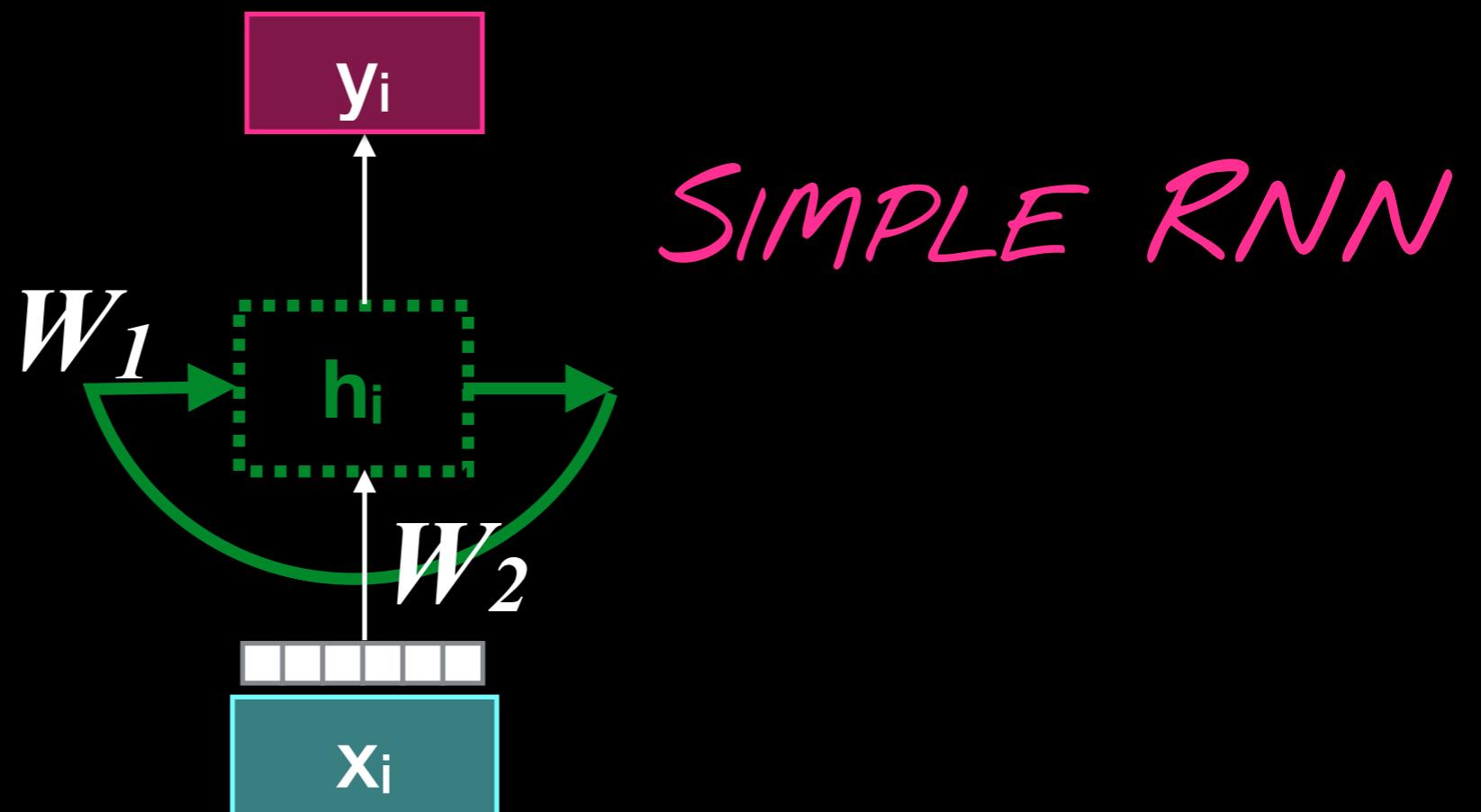
$$h_i = s(h_{i-1}, x_i)$$



# Concretely

$$y_i = f(h_i) = h_i$$

$$h_i = s(h_{i-1}, x_i) = \tanh(W_1 h_{i-1} + W_2 x_i + b)$$



# Recap: LMs

$$P(w_1, w_2, \dots, w_n) \approx \prod_{i=1}^N P(w_i | w_{i-2}, w_{i-1})$$

*TRIGRAM MODEL*

\* \* The weather today is fine STOP

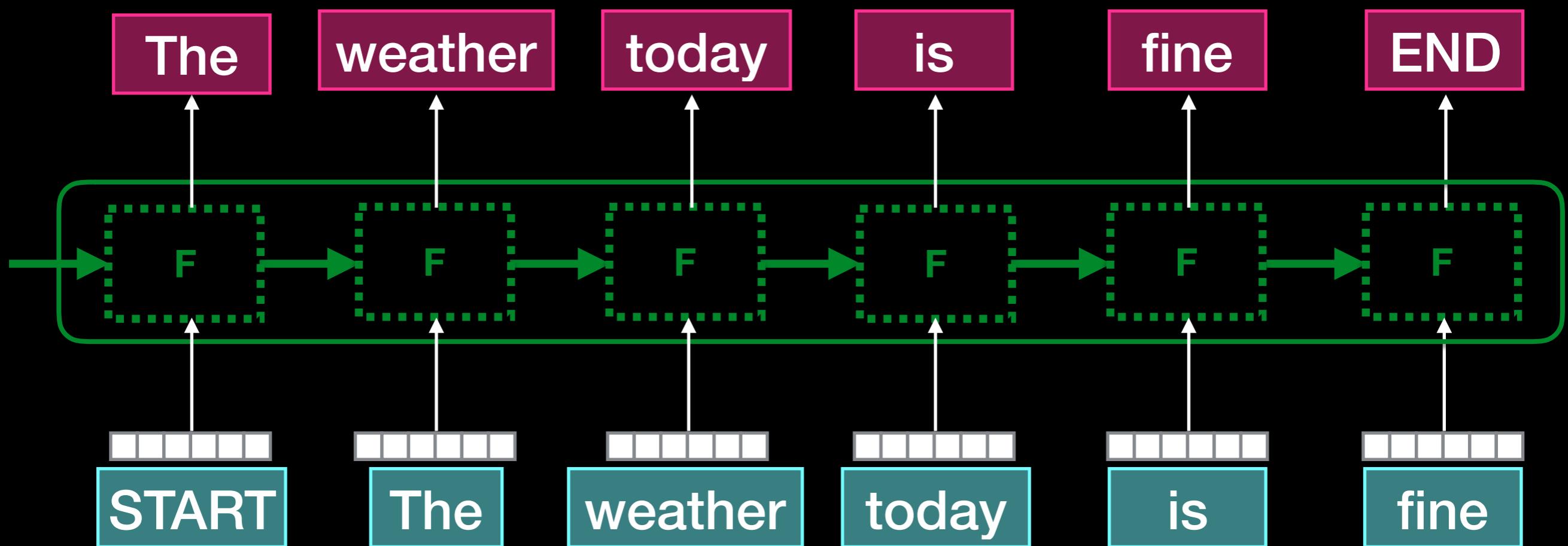
$$\begin{aligned} P(S) = P(w_1, \dots, w_n) &= P(\text{The} | * *) \\ &\quad \times P(\text{weather} | * \text{ The}) \\ &\quad \times P(\text{today} | \text{The weather}) \\ &\quad \times P(\text{is} | \text{weather today}) \\ &\quad \times P(\text{fine} | \text{today is}) \\ &\quad \times P(\text{STOP} | \text{is fine}) \end{aligned}$$

*CHAIN RULE*

# Neural LMs

$$P(w_1, w_2, \dots, w_n) \approx \prod_{i=1}^N P(w_i | w_1 \dots w_{i-1})$$

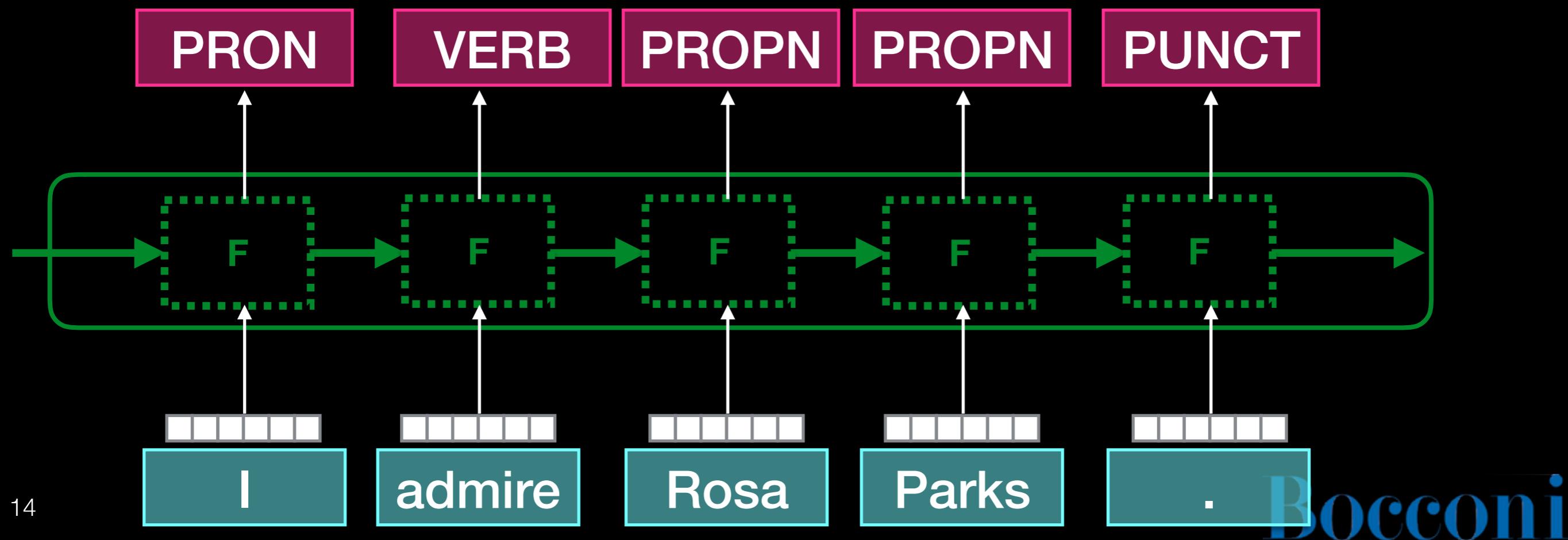
INFINITE MODEL



PREDICT NEXT WORD GIVEN HISTORY

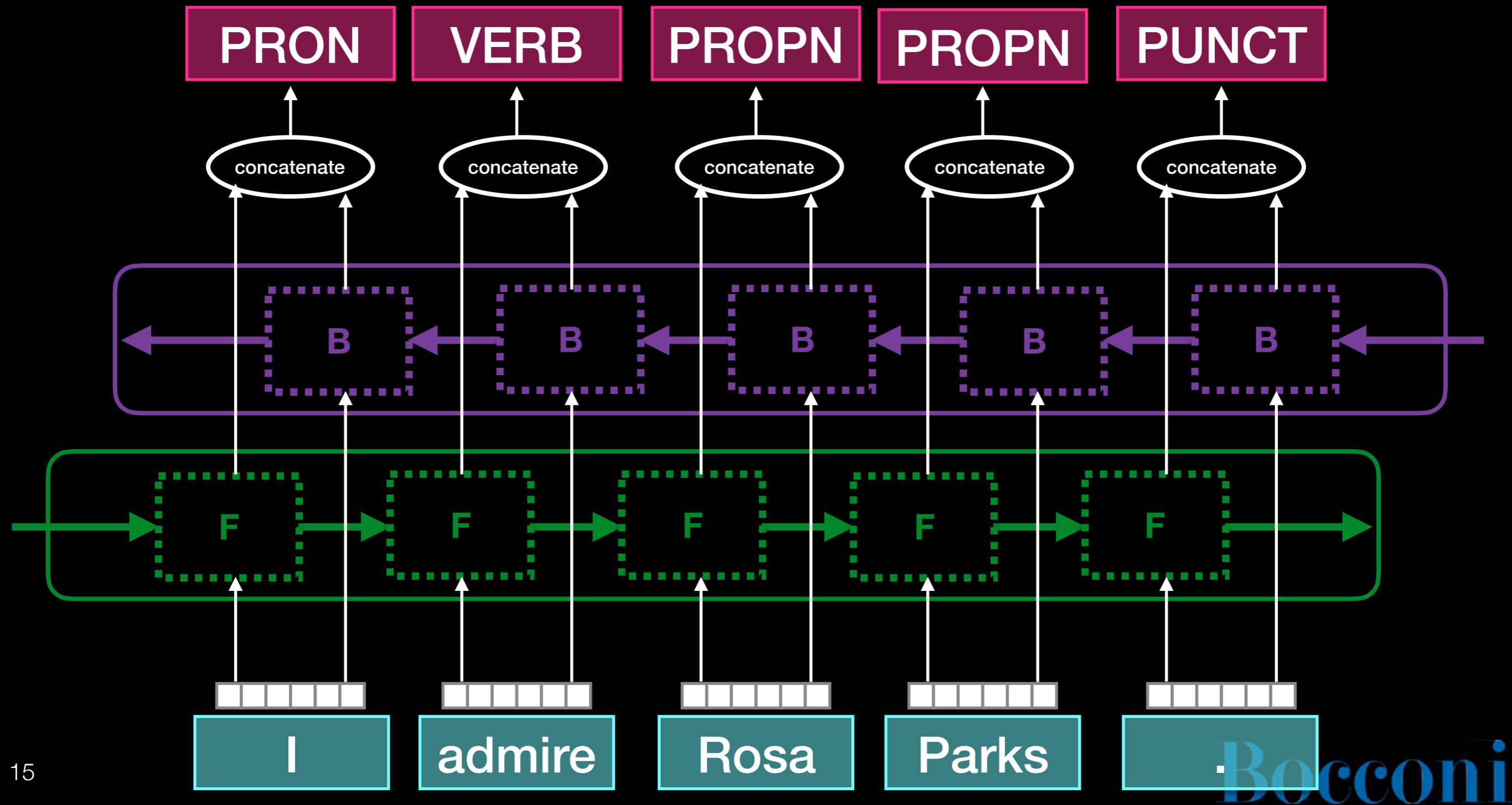
# RNN Tagging

*STRUCTURED PREDICTION*



# Bidirectional-RNN

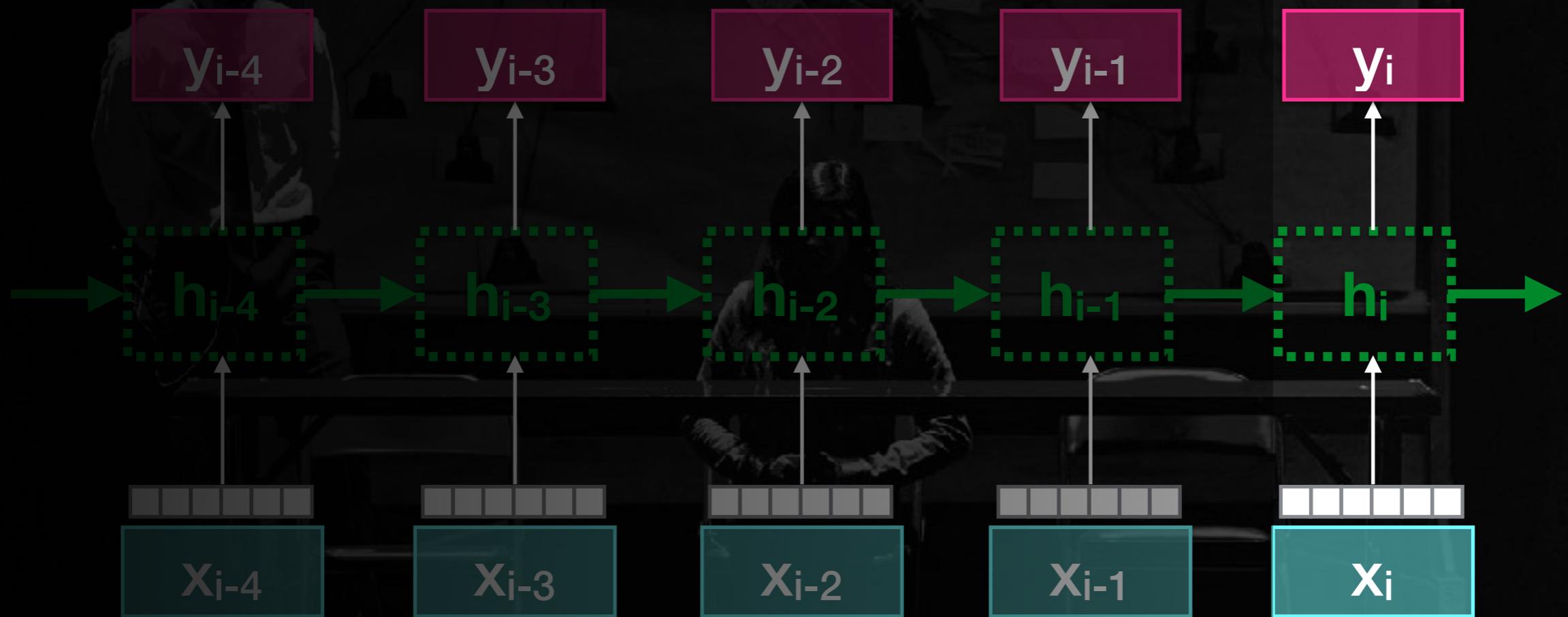
*STRUCTURED PREDICTION*



# Special Recurrent Networks

# Vanishing Memory

WHERE WERE YOU MARCH 3, 2016?



PROBLEM WITH LONG SEQUENCES

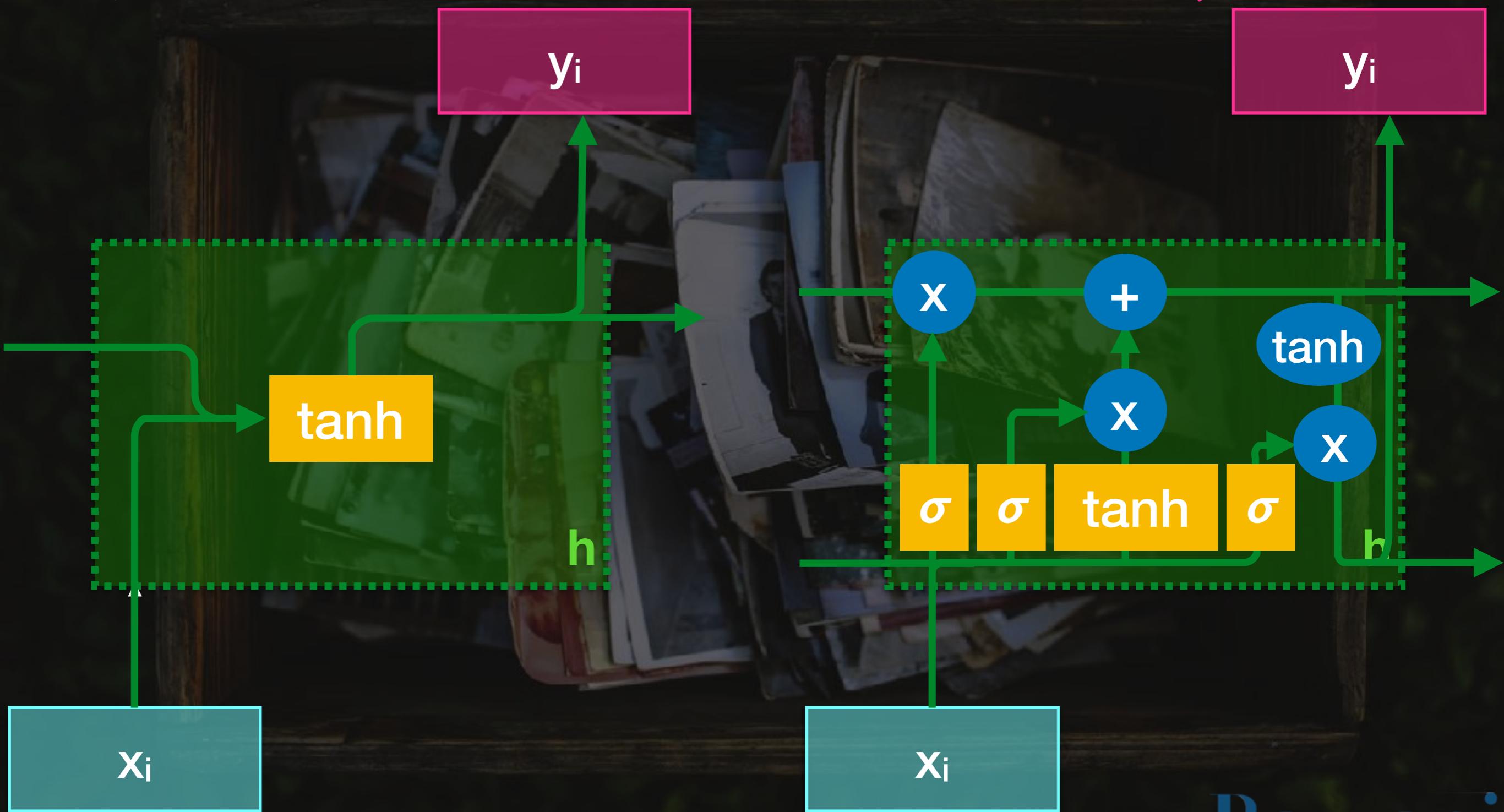
Bocconi

# Selective Forgetting

$$h_i = \tanh(W_1 h_{i-1} + W_2 x_i + b)$$

SIMPLE RNN

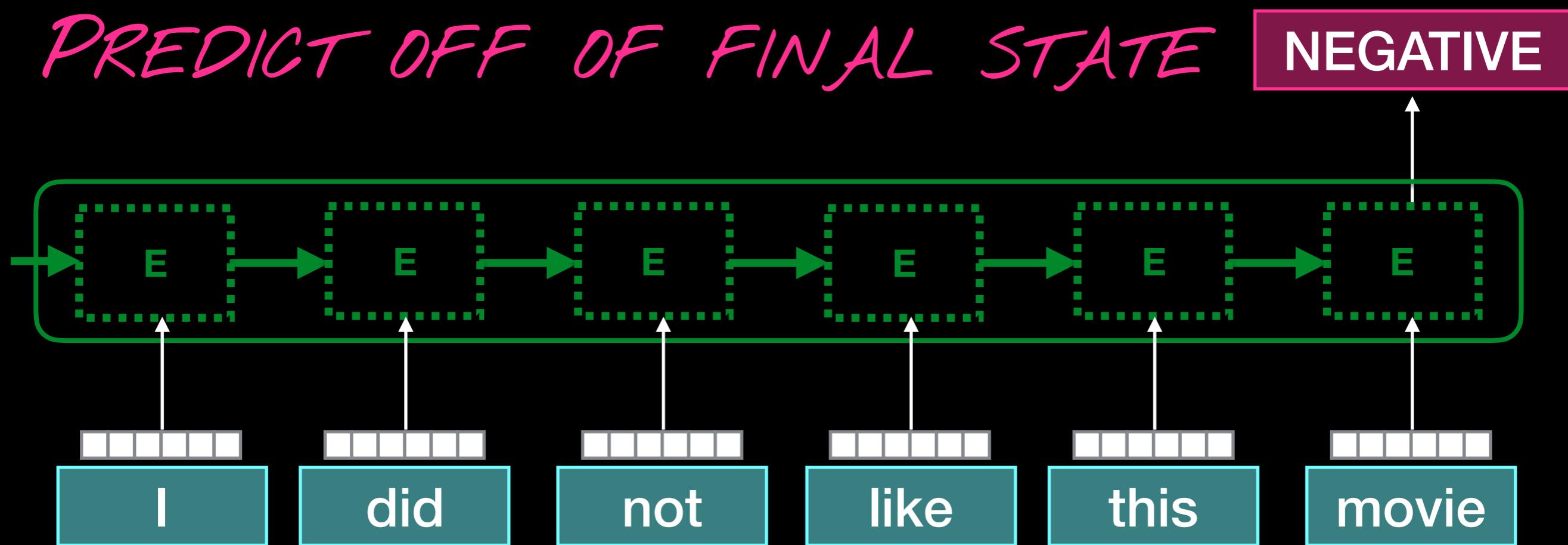
LSTM



# Acceptor

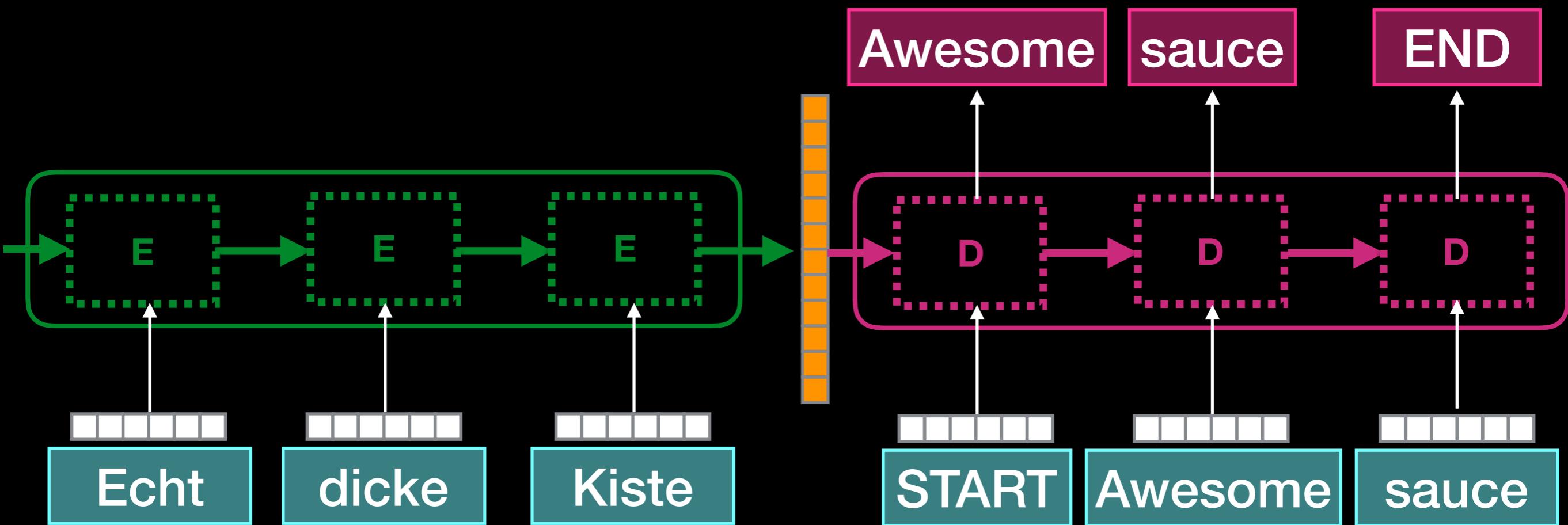
*PREDICT OFF OF FINAL STATE*

NEGATIVE



# Encoder-Decoder

...AND GENERATE  
OUTPUT FROM IT

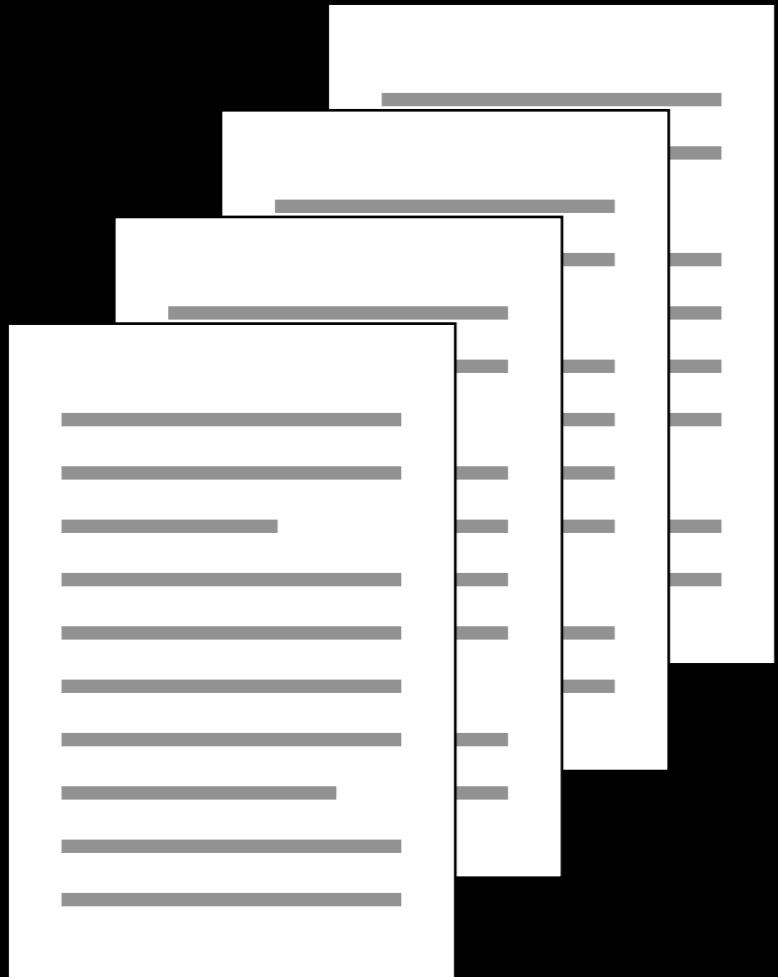


GOBBLE UP SEQUENCE

INTO A VECTOR...

# Convolution

# Convolved Matters



*TEXT SORT CLASSIFICATION*

Retail

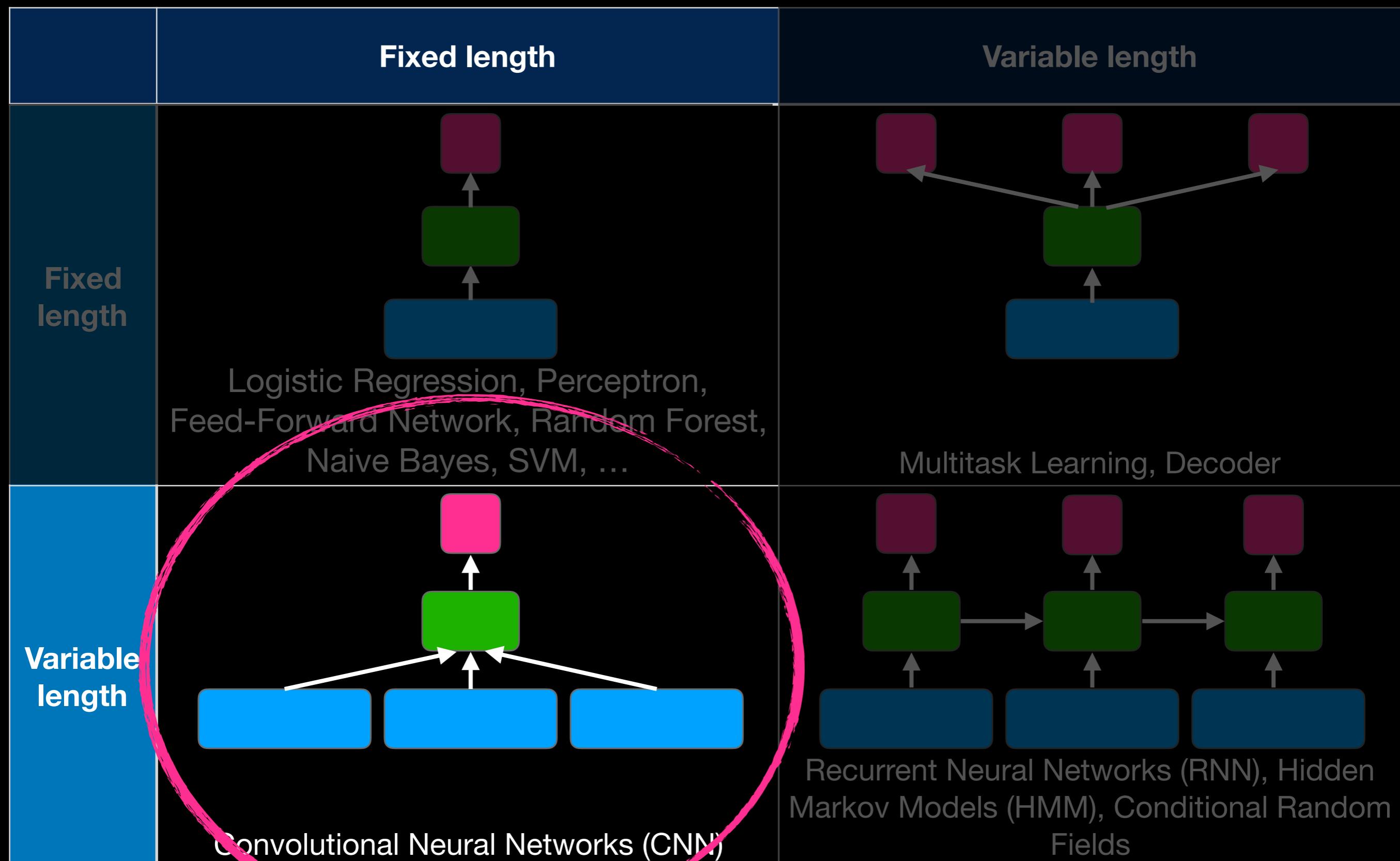
*SENTIMENT ANALYSIS*

positive

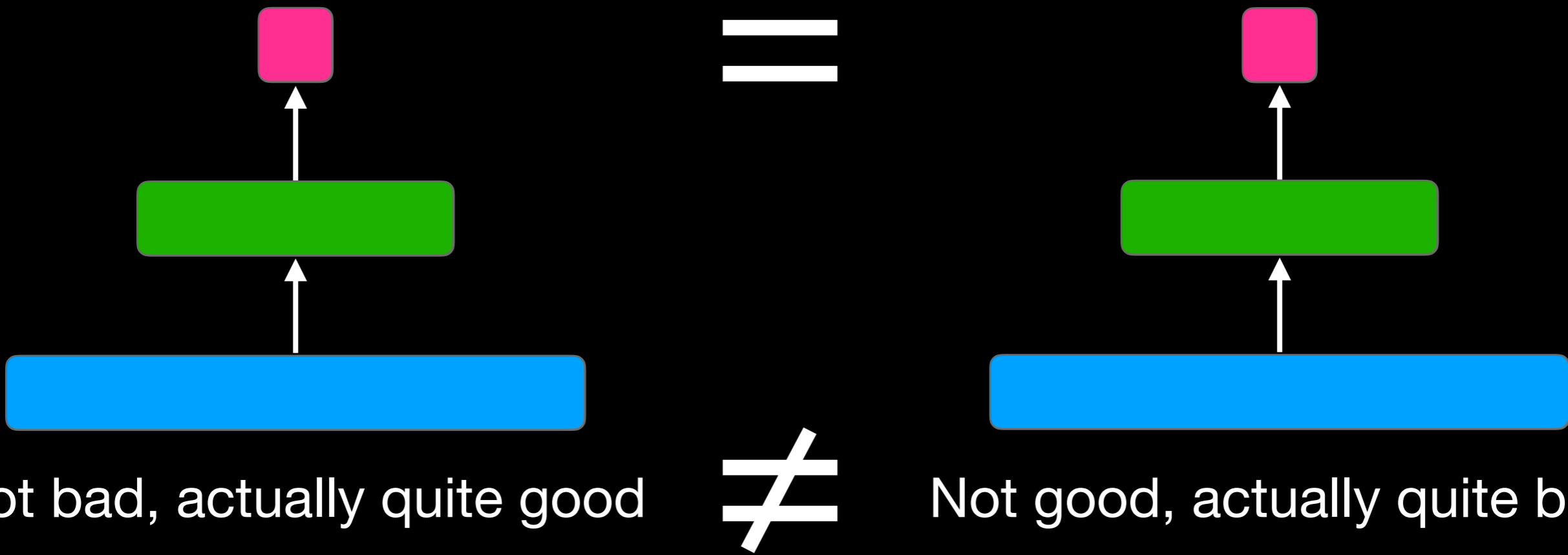
*RELATION EXTRACTION*

founded\_by(Amazon, Jeff Bezos)

# Types of Text Classification

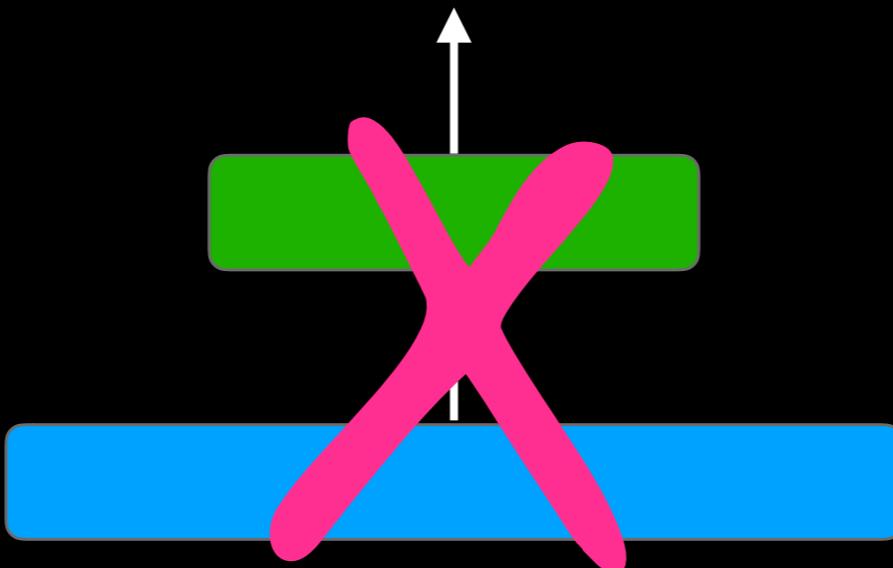


# Problems with MLPs



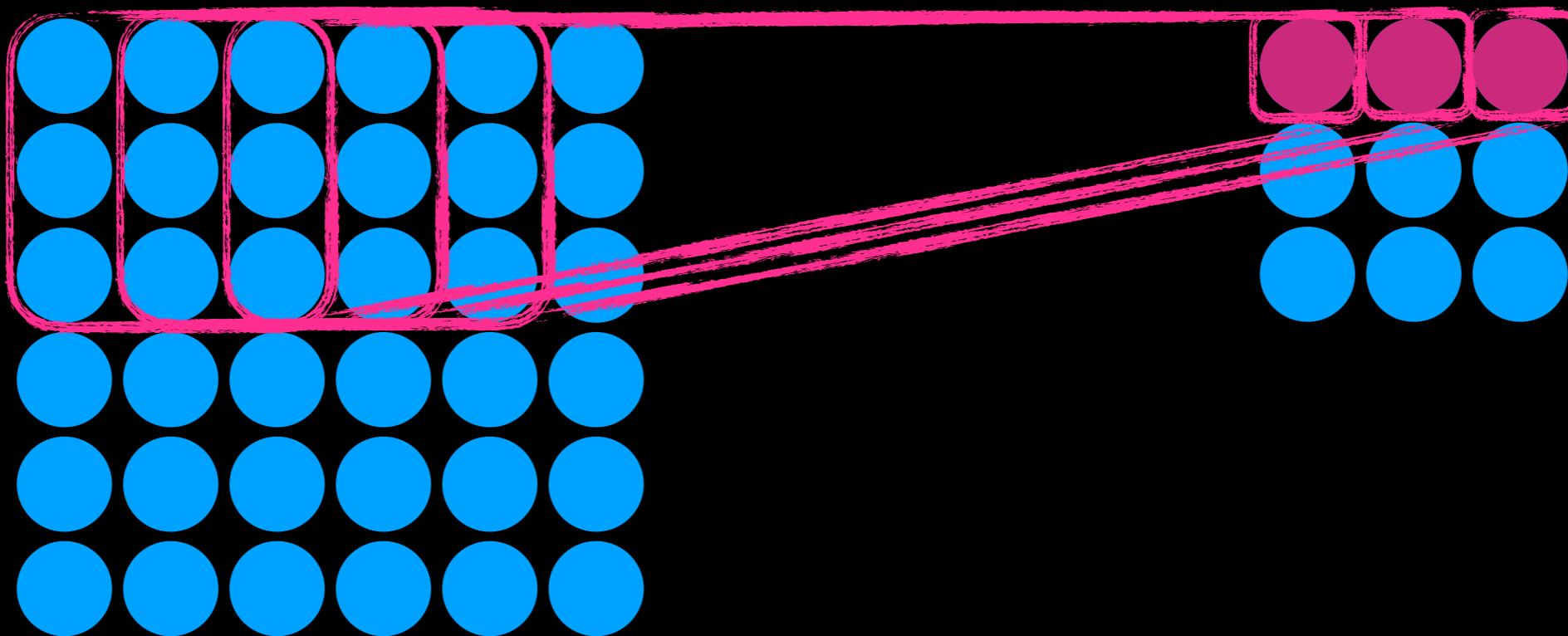
# Problems with MLPs

founded\_by(Amazon, Jeff Bezos)



Jeff Bezos, or what Dr. Evil would look like on steroids, went from book seller to billionaire when he founded Amazon in 1994.

# Convolution



# Convolution Extraction

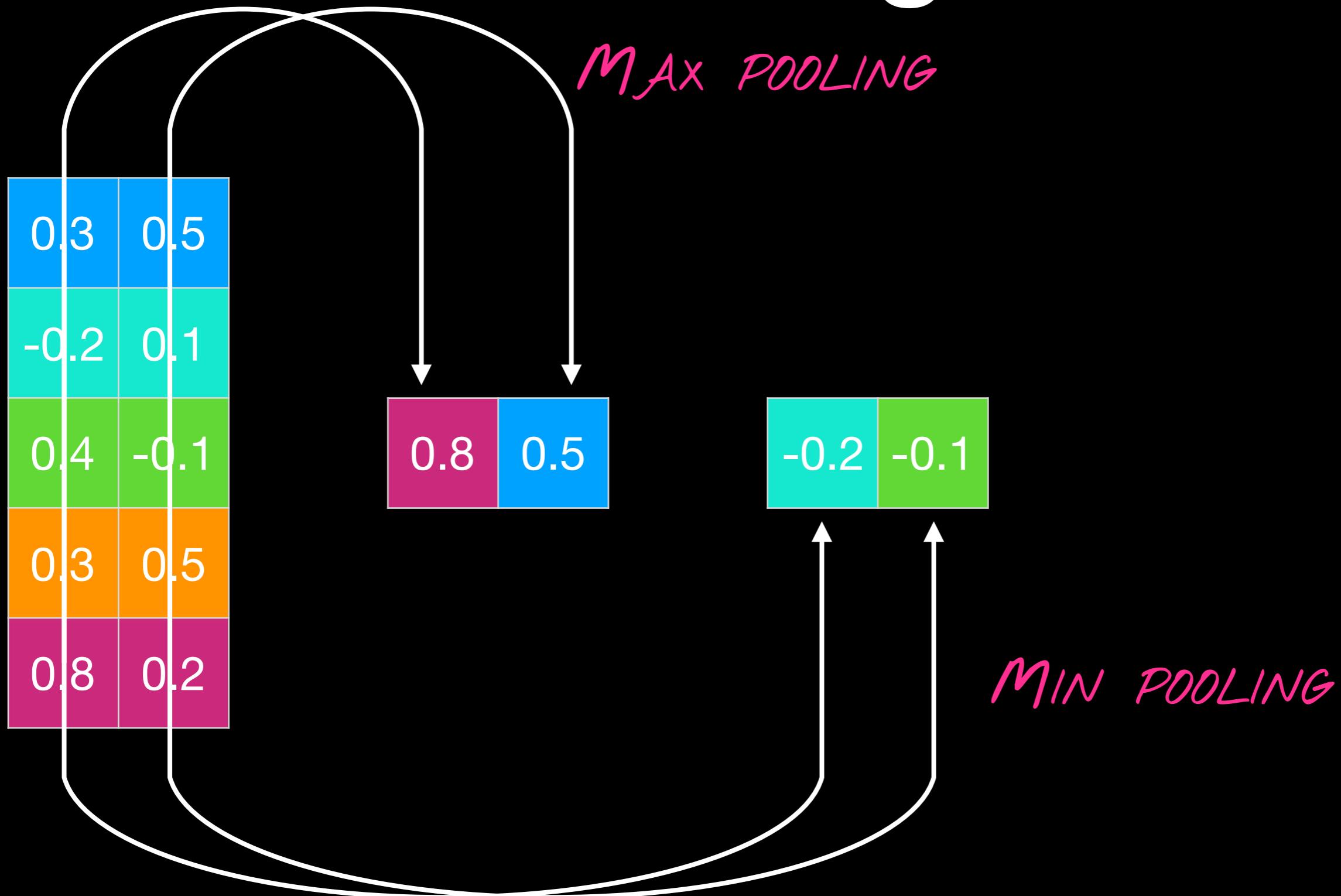
1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

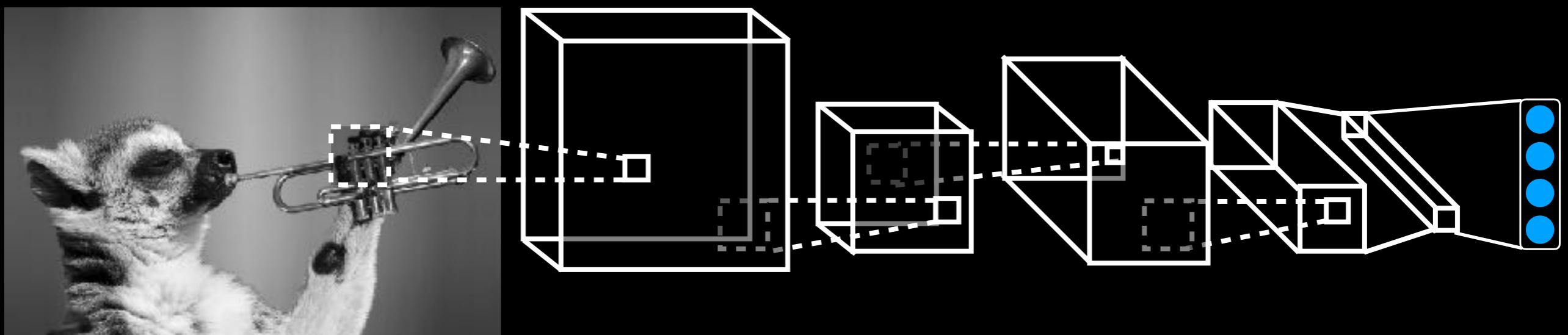
4		

Convolved  
Feature

# Pooling



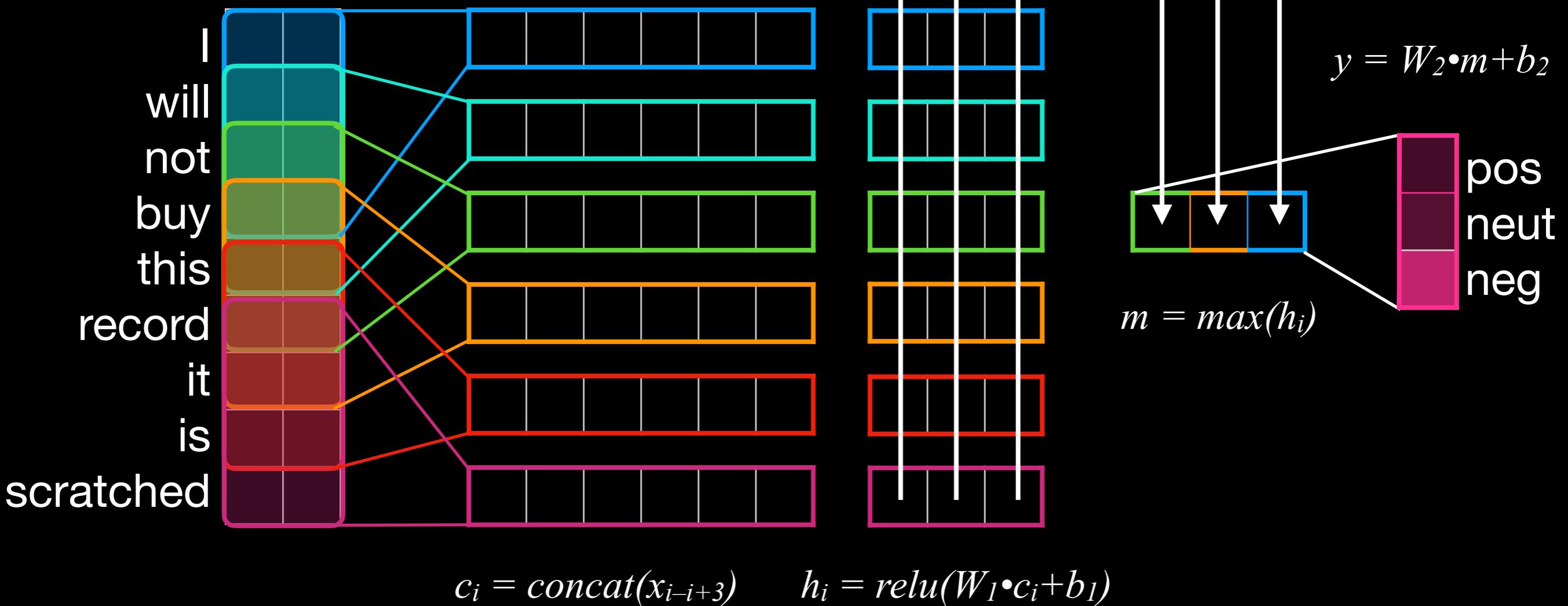
# CNNs in Images



# CNNs in Text

*CONVOLUTION WINDOW = 4*

*STEP SIZE = 1*



# The Attention Mechanism

# Attention!

- Learn syntactic and semantic relations between words in
  - the input and output (RNNs)
  - only the input (CNNs)
- Good for machine translation (word alignment) and classification (complex expressions)
- Basis of the Transformer model

# Visual Attention



Oh look, a furry ear!

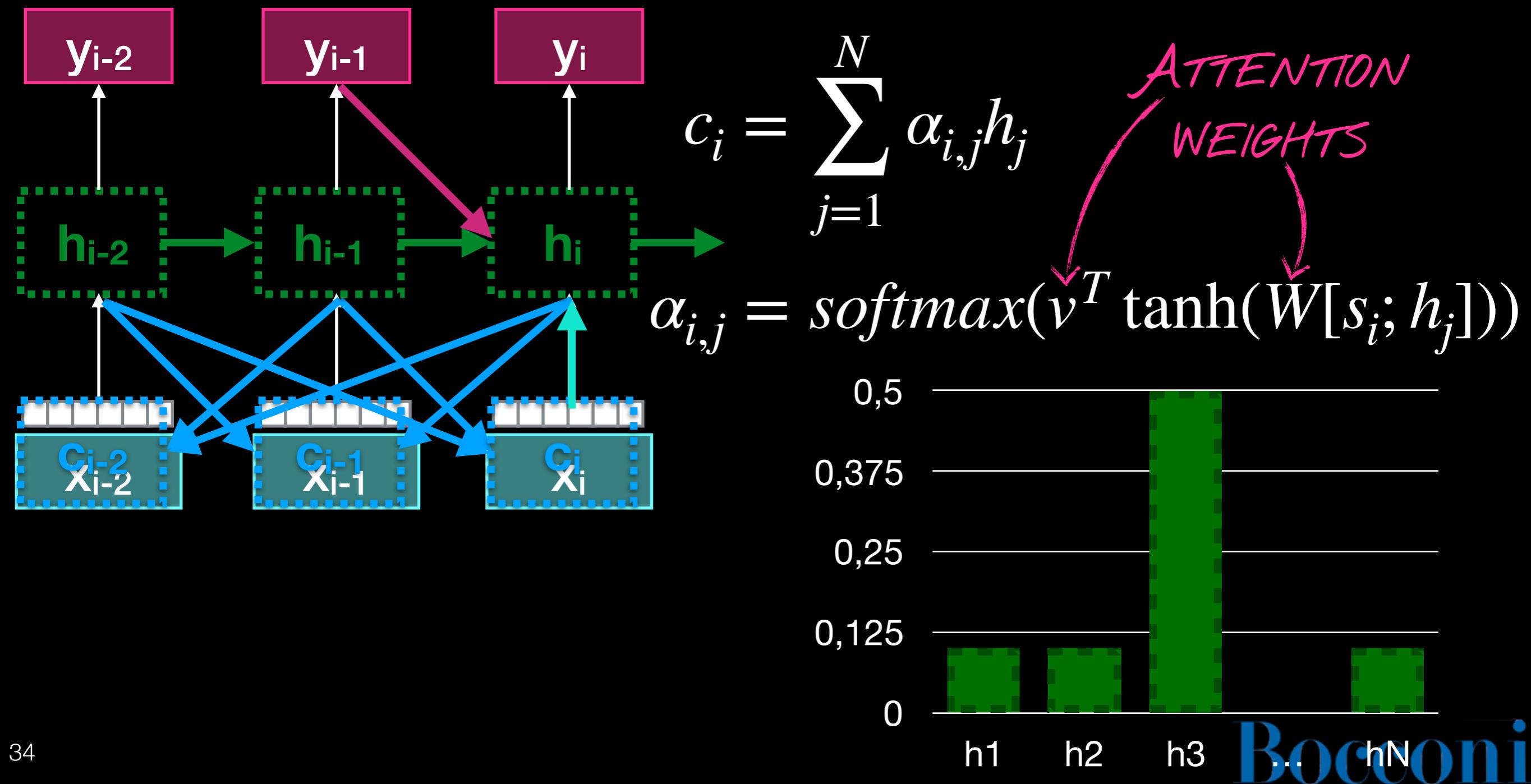
Of course, there's a doge in the picture:  
look at all the other doge parts

Ignore those other parts, they don't tell you anything about doges

see <https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html>

# Ok, but how?

$$h_i = f(h_{i-1}, y_{i-1}, c_i)$$



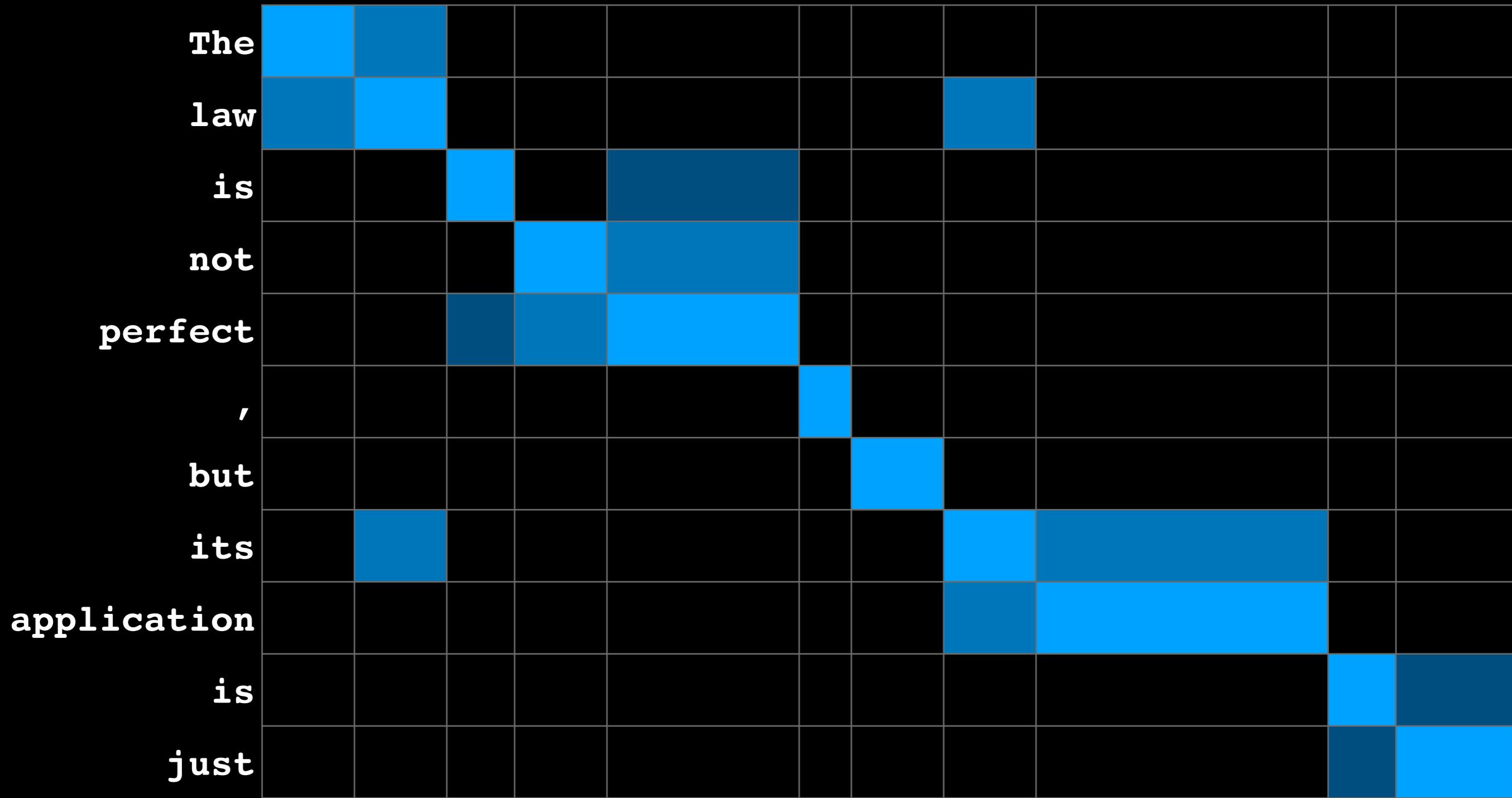
# But Why?



- RNNs can't store long sequences
- CNNs don't have a sense of location in the final vector
- Attention tells us how each part relates to each other part

# $\alpha$ CNN with Attention

The law is not perfect , but its application is just



FIND LONG-RANGE DEPENDENCIES

# RNN with Attention

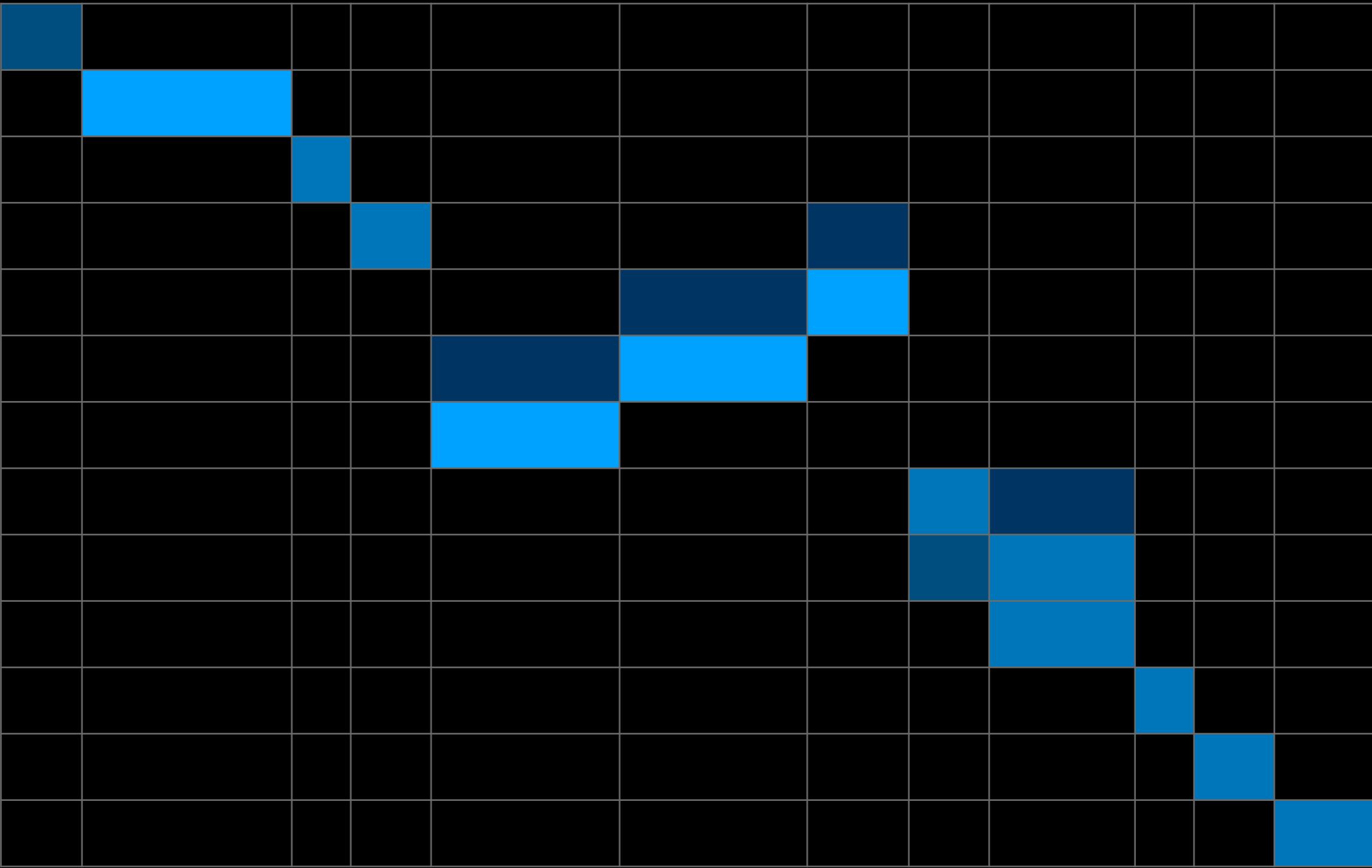
OUTPUT

$\alpha$

The agreement on the European Economic Area was signed in Aug 1992

INPUT

L'  
accord  
sur  
la  
zone  
économique  
européenne  
a  
été  
signé  
on  
août  
1992



LEARN REORDERING

Bocconi

# RNN with Attention

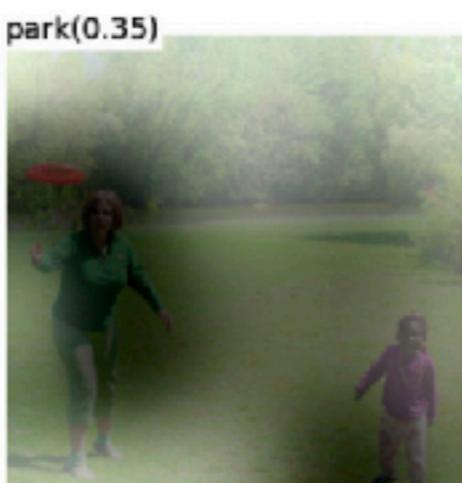


LEARN REORDERING

Bocconi

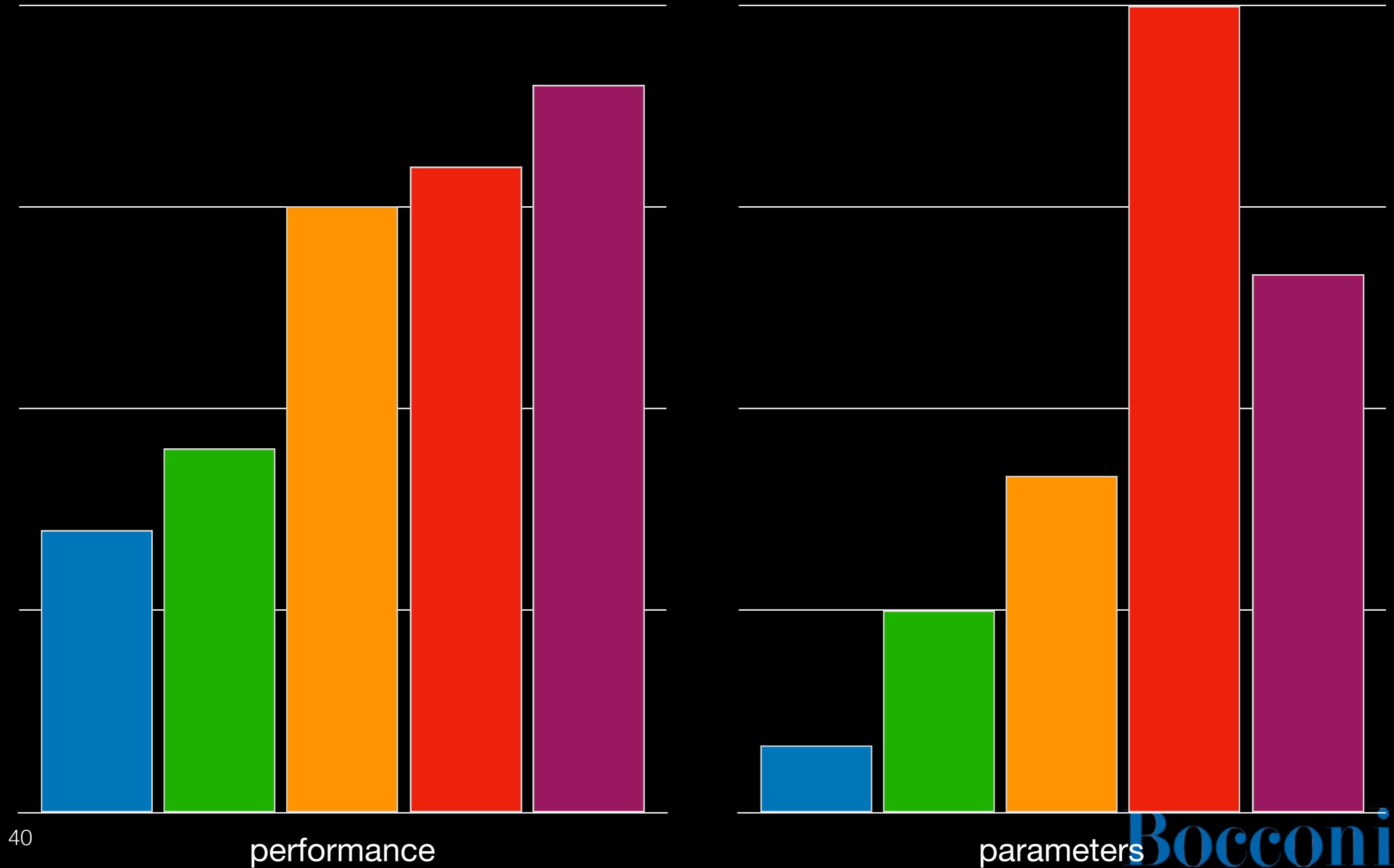
# Multimodal Attention

A woman is throwing a frisbee in the park .



# Performance vs Parameters

■ Logistic Regression ■ Feed-Forward ■ CNN ■ RNN ■ CNN+Attention



performance

parameters **Bocconi**

# Wrapping up

# Take Home Points

- **Recurrent Neural Nets** address long-range dependencies
- Condition each word on all previous ones (better for **LMs** and **sequence labels**)
- **Bidirectional RNNs** condition on following words
- **LSTMs** learn to forget useless input
- **Convolution windows** captures **different views** of input
- **Pooling** reduces dimensionality
- CNNs are often **better for text classification** than feedforward NNs
- **Attention** improves coherence and performance