

Natural Language Processing

Lecture 21

Dirk Hovy

dirk.hovy@unibocconi.it

 @dirk_hovy

Goals for Today

- Learn about **recurrent neural network architectures**
- Understand the difference to convolutional networks
- Learn about **different architecture**
- Understand the **attention mechanism**

Recurring Matters

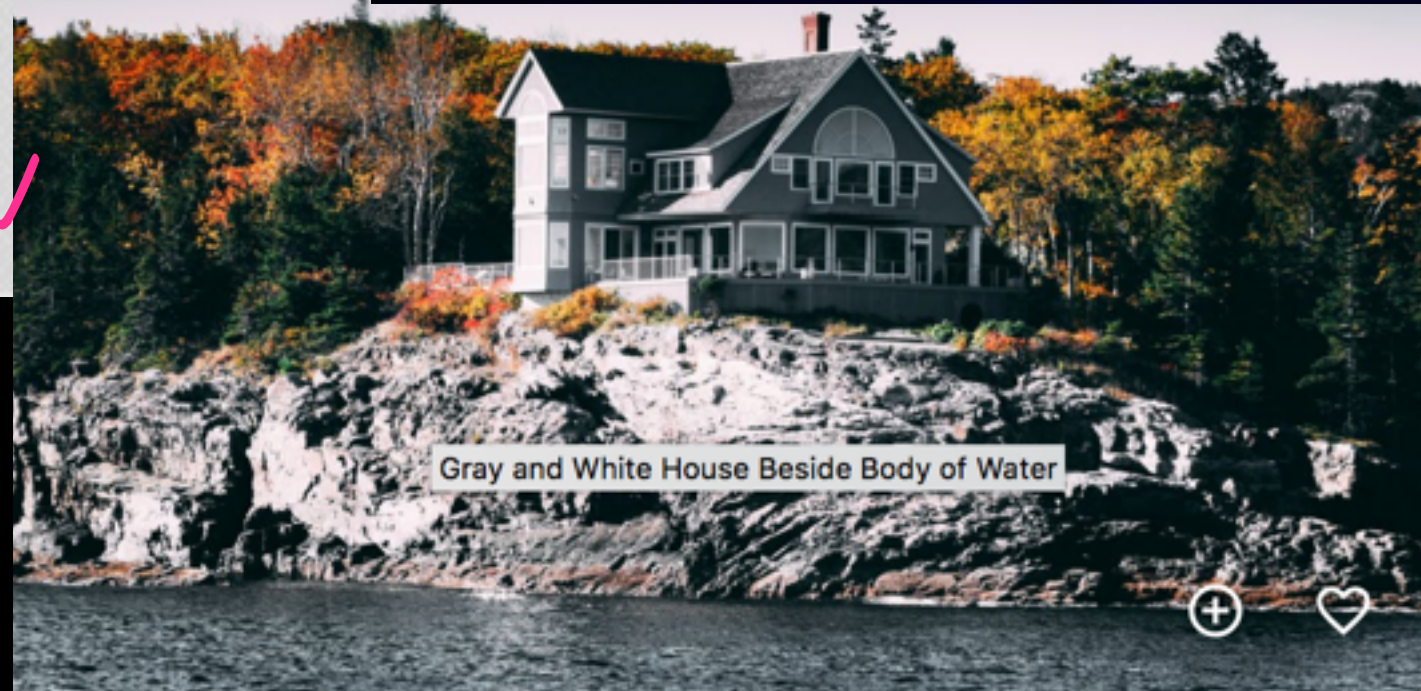


IMAGE CAPTIONS

Long-Term Trouble

SUBJECT

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

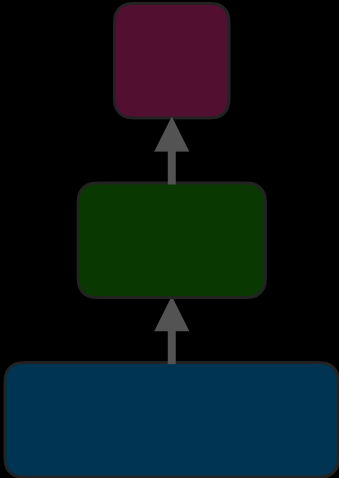
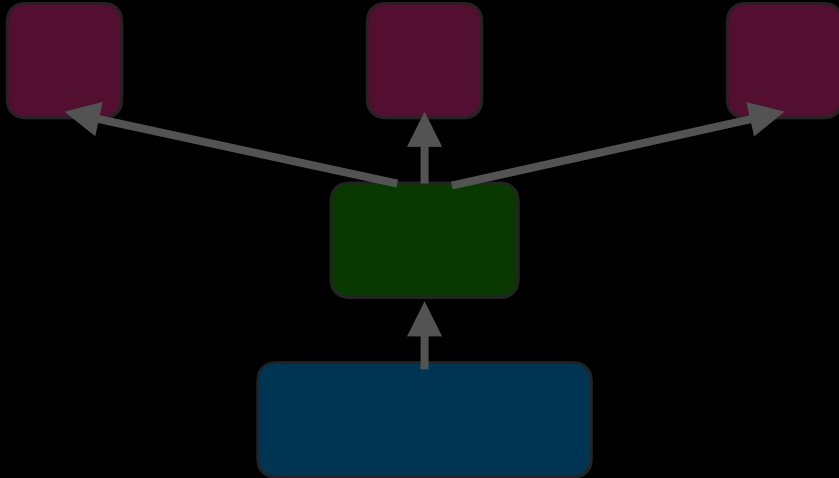
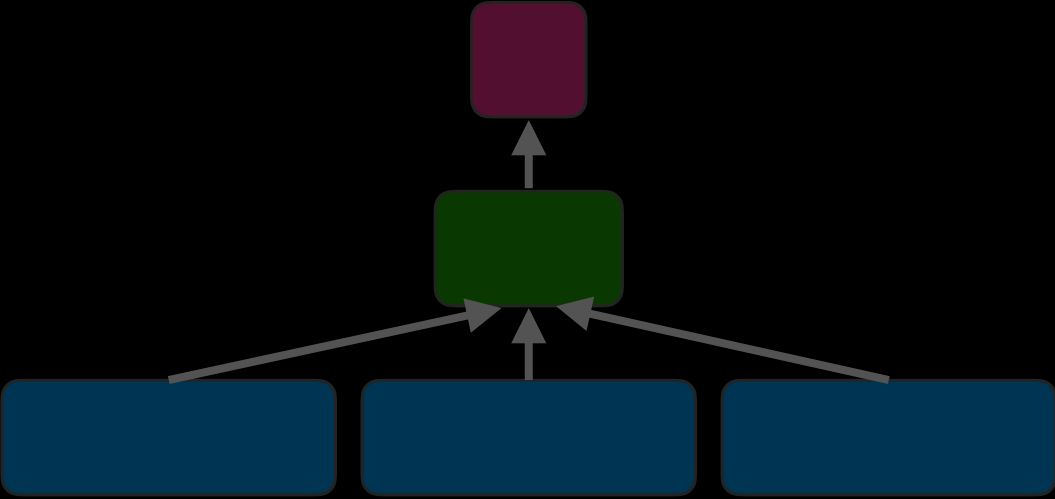
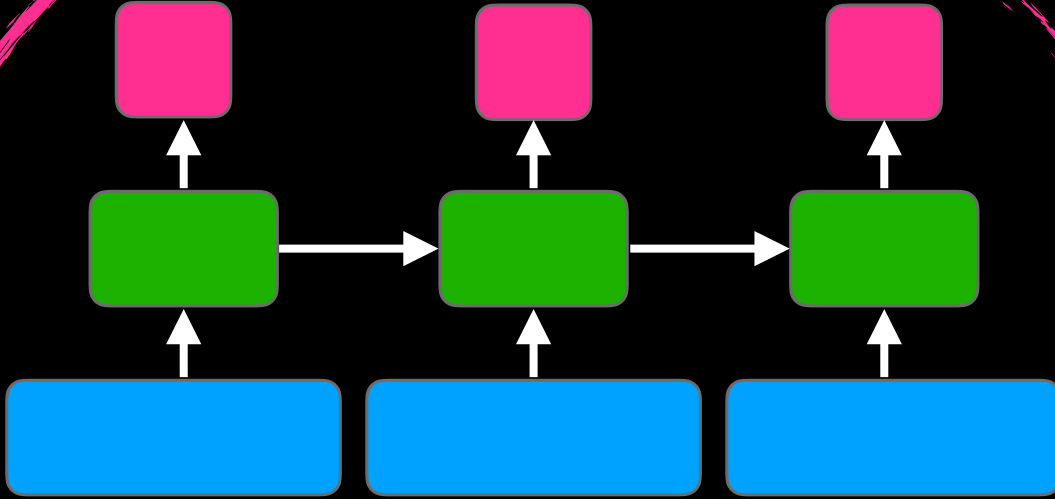
show {**VERB**, **NOUN**}

PART **show**
show
PRON **show**
show

DET **show**
show
show
ADJ **show**

Structured prediction: depends on the POS of a previous word

Types of Text Classification

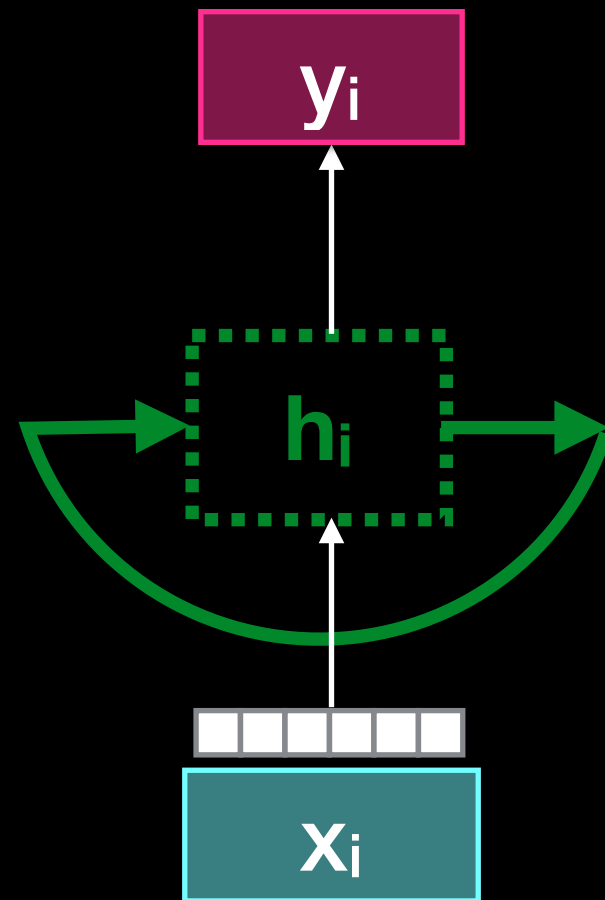
	Fixed length	Variable length
Fixed length	 <p>Logistic Regression, Perceptron, Feed-Forward Network, Random Forest, Naive Bayes, SVM, ...</p>	 <p>Multitask Learning, Decoder</p>
Variable length	 <p>Convolutional Neural Networks (CNN)</p>	 <p>Recurrent Neural Networks (RNN), Hidden Markov Models (HMM), Conditional Random Fields</p>

Recurrent Networks

Recurrence

$$y_i = f(h_i)$$

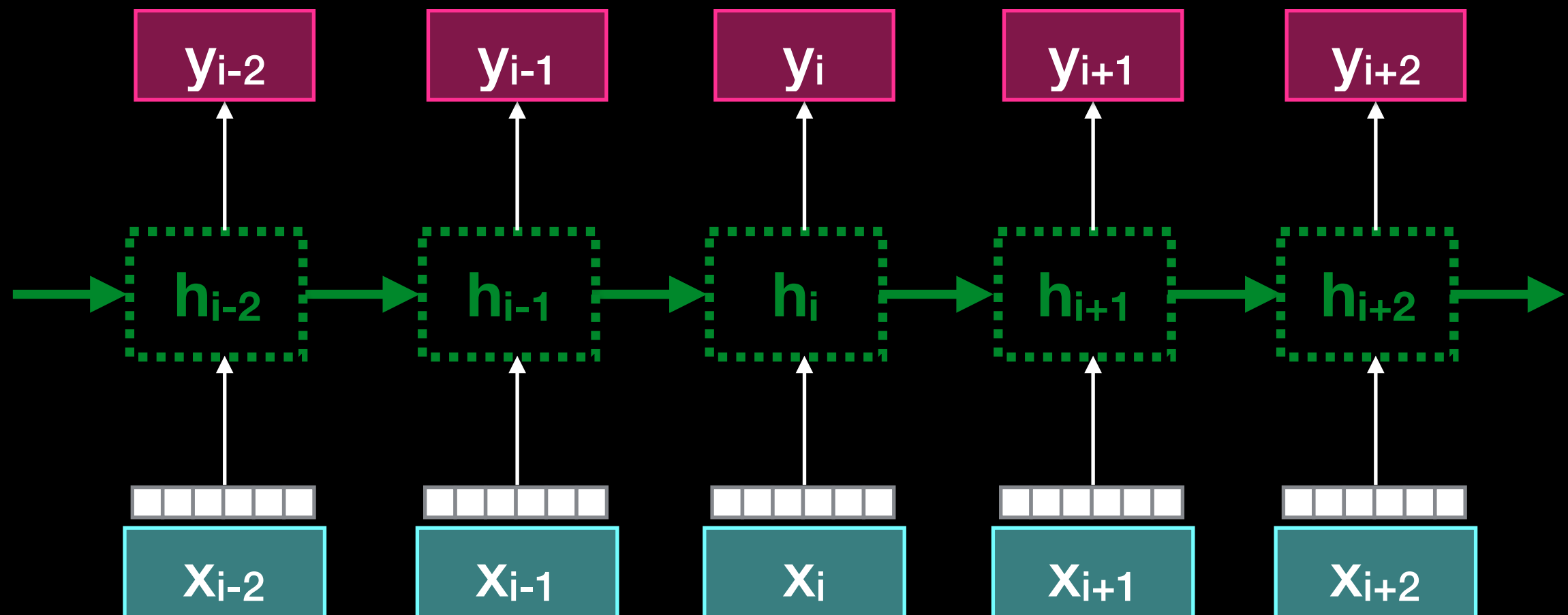
$$h_i = s(h_{i-1}, x_i)$$



...Unrolled

$$y_i = f(h_i)$$

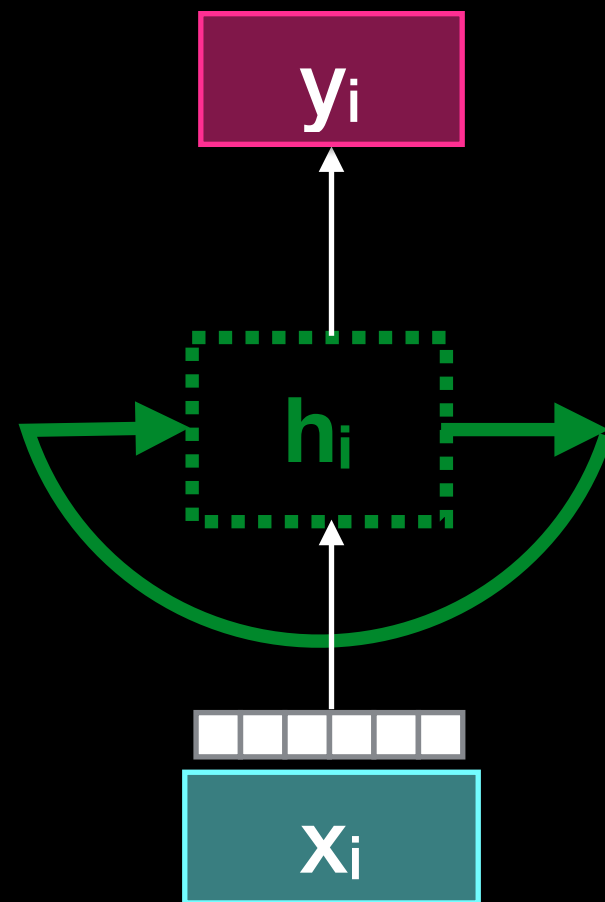
$$h_i = s(h_{i-1}, x_i)$$



Concretely

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

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



SIMPLE RNN

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**

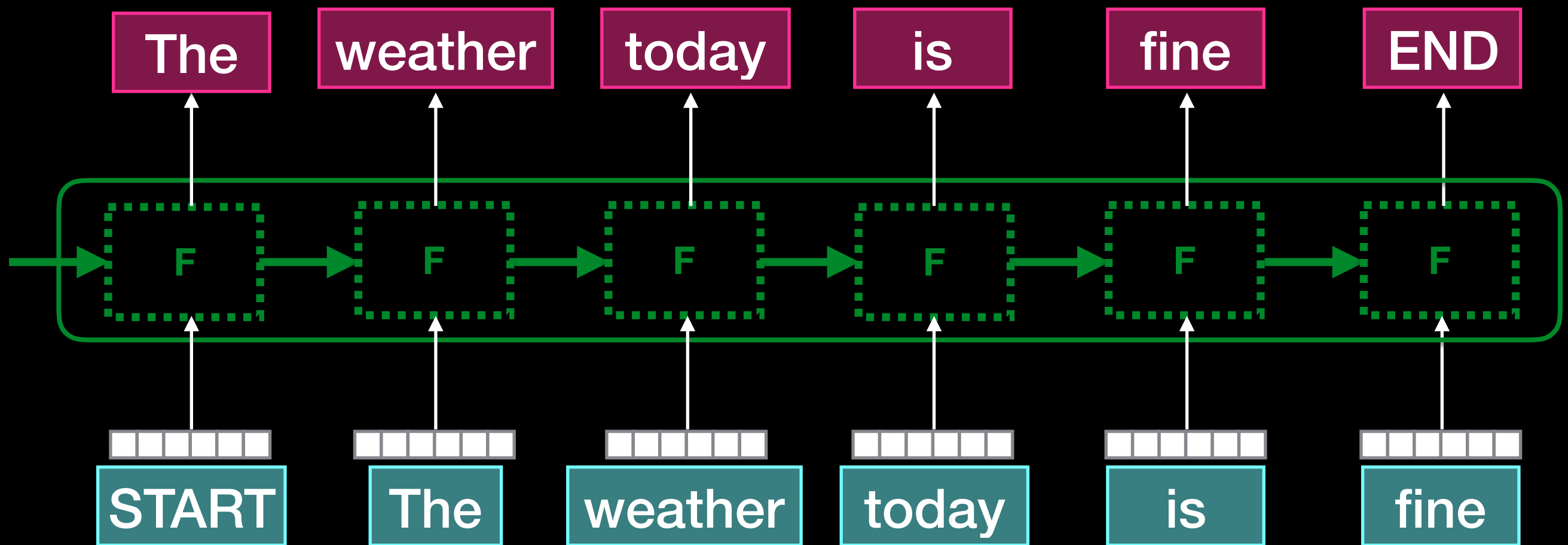
$$\begin{aligned} P(S) = P(w_1, \dots, w_n) = & P(\text{The} | * *) \\ & \times P(\text{weather} | * \text{The}) \\ & \times P(\text{today} | \text{The weather}) \\ & \times P(\text{is} | \text{weather today}) \\ & \times P(\text{fine} | \text{today is}) \\ & \times P(\text{STOP} | \text{is fine}) \end{aligned}$$

CHAIN RULE

Neural LMs

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

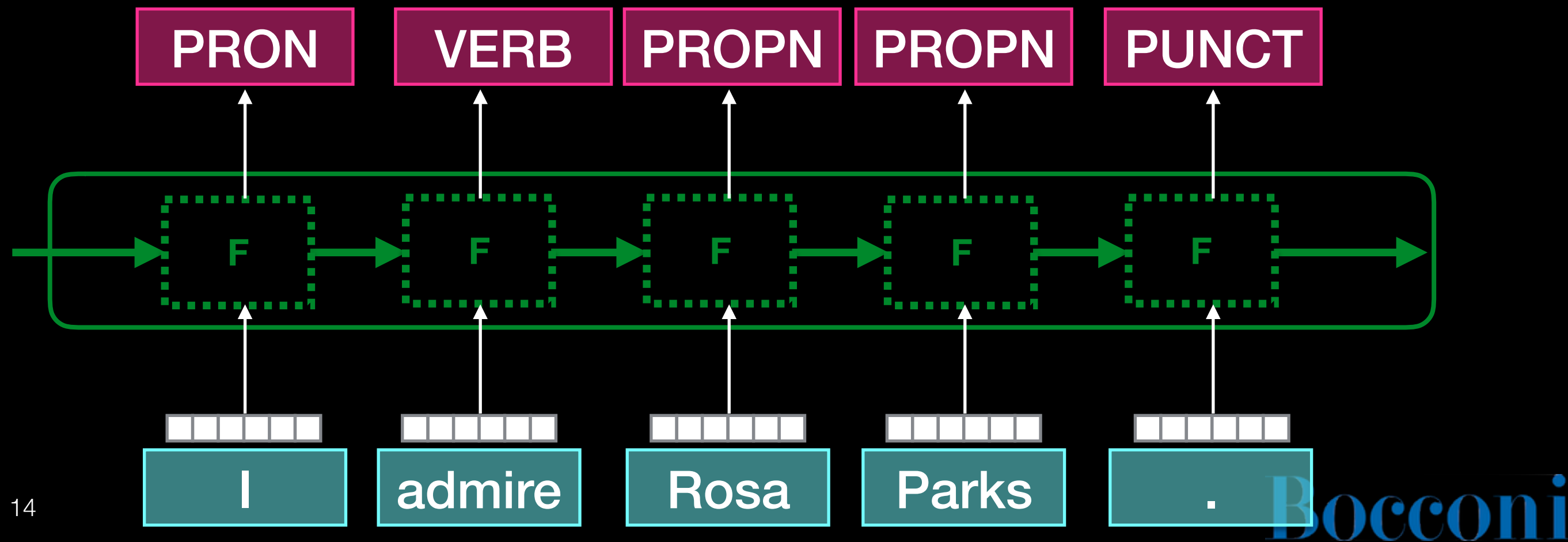
INFINITE MODEL



PREDICT NEXT WORD GIVEN HISTORY

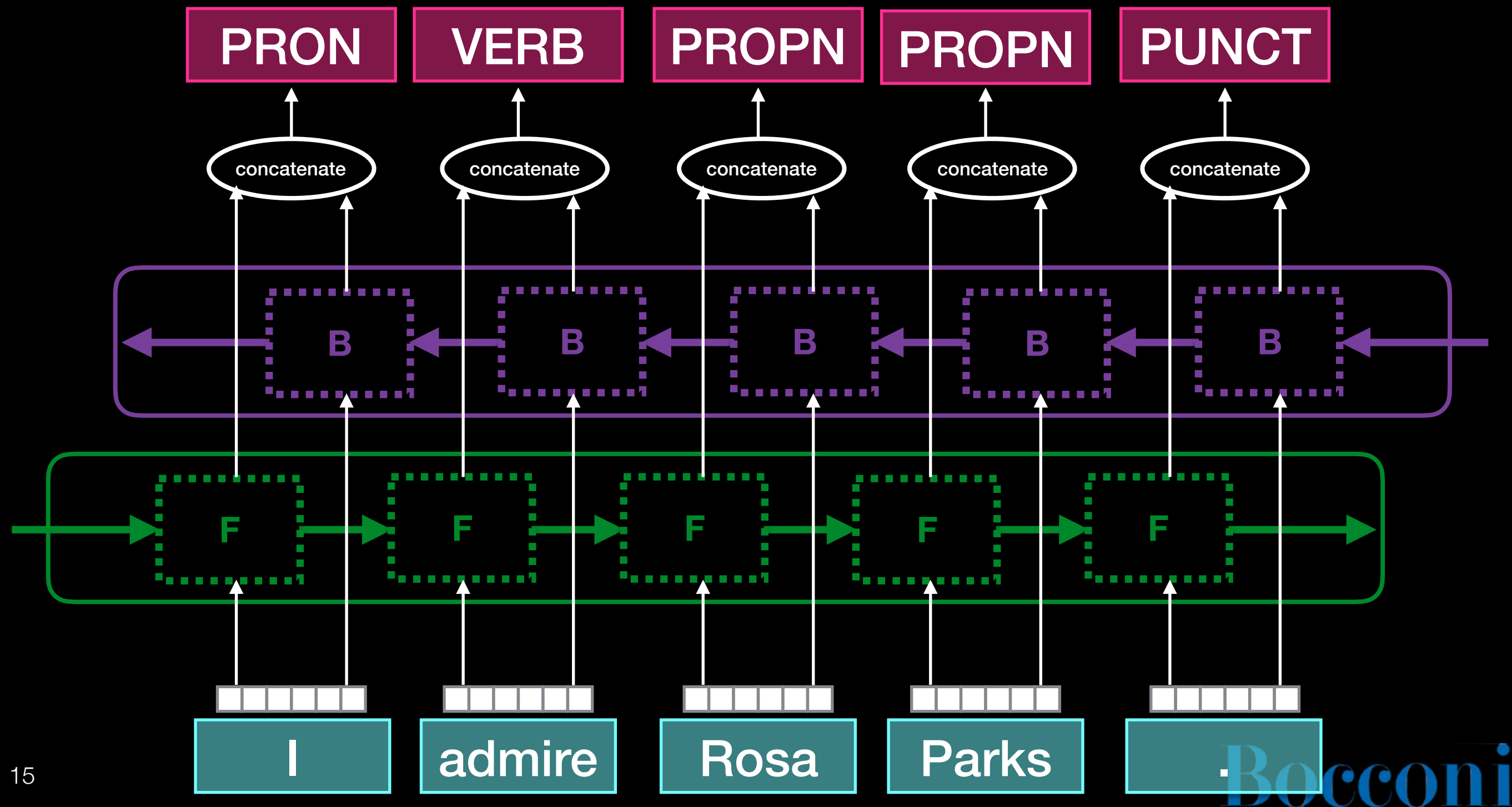
RNN Tagging

STRUCTURED PREDICTION



Bidirectional-RNN

STRUCTURED PREDICTION



Special Recurrent Networks

Vanishing Memory

WHERE WERE YOU MARCH 3, 2016?

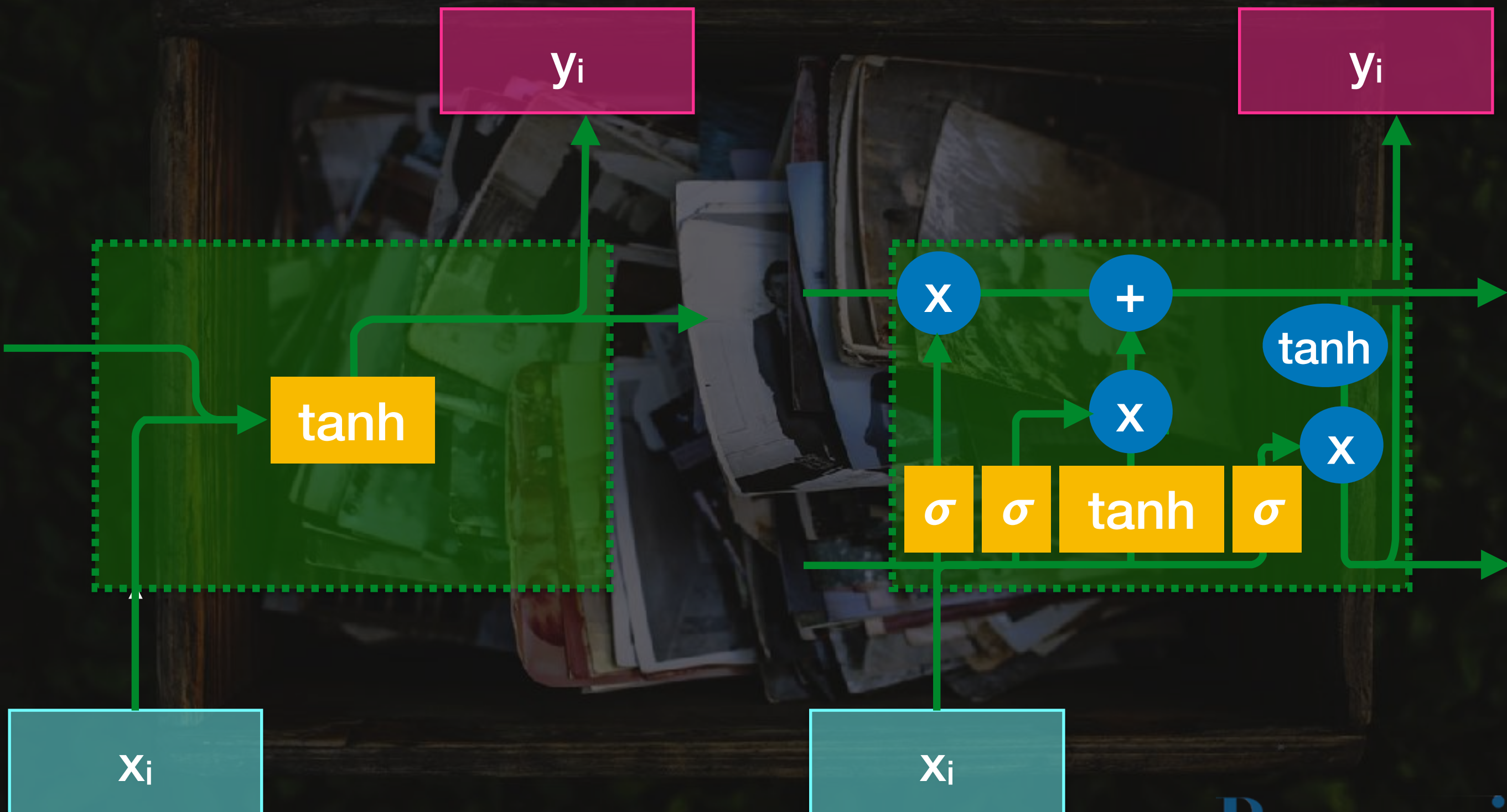


PROBLEM WITH LONG SEQUENCES

Selective Forgetting

SIMPLE RNN

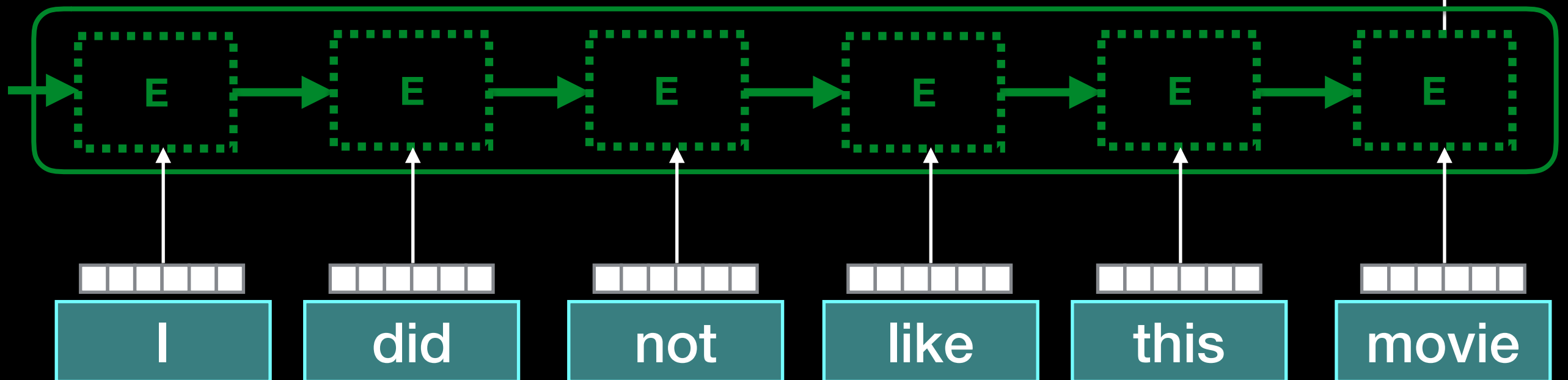
LSTM



Acceptor

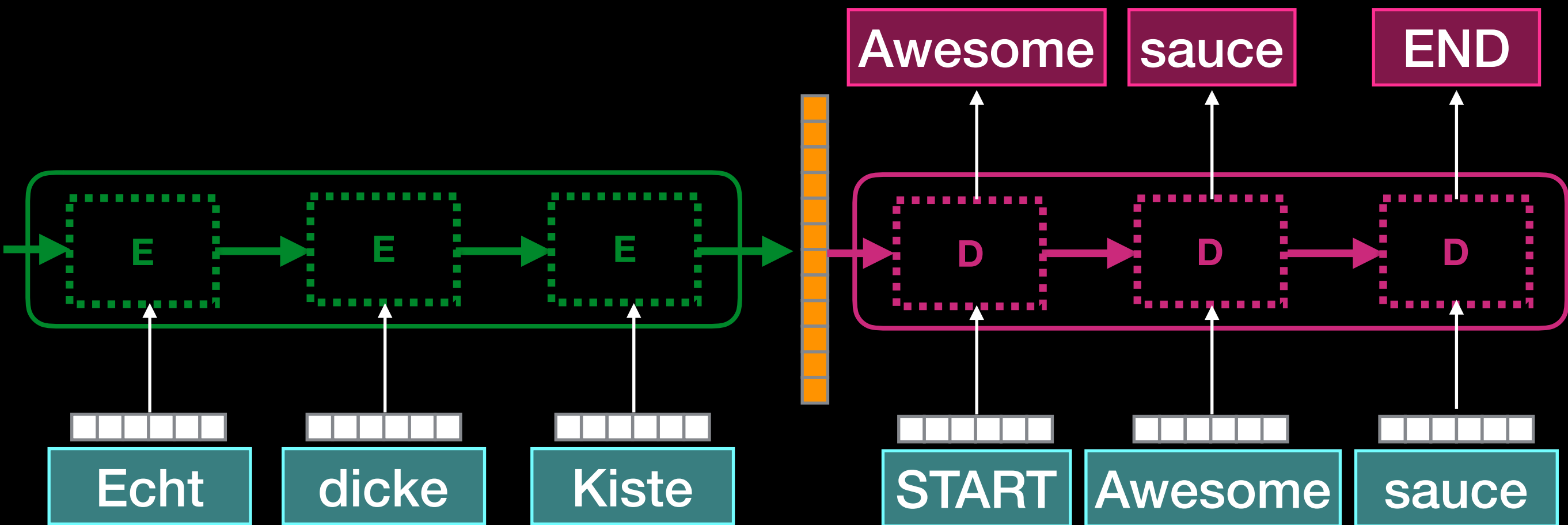
PREDICT OFF OF FINAL STATE

NEGATIVE



Encoder-Decoder

*...AND GENERATE
OUTPUT FROM IT*



*GOBBLE UP SEQUENCE
INTO A VECTOR...*

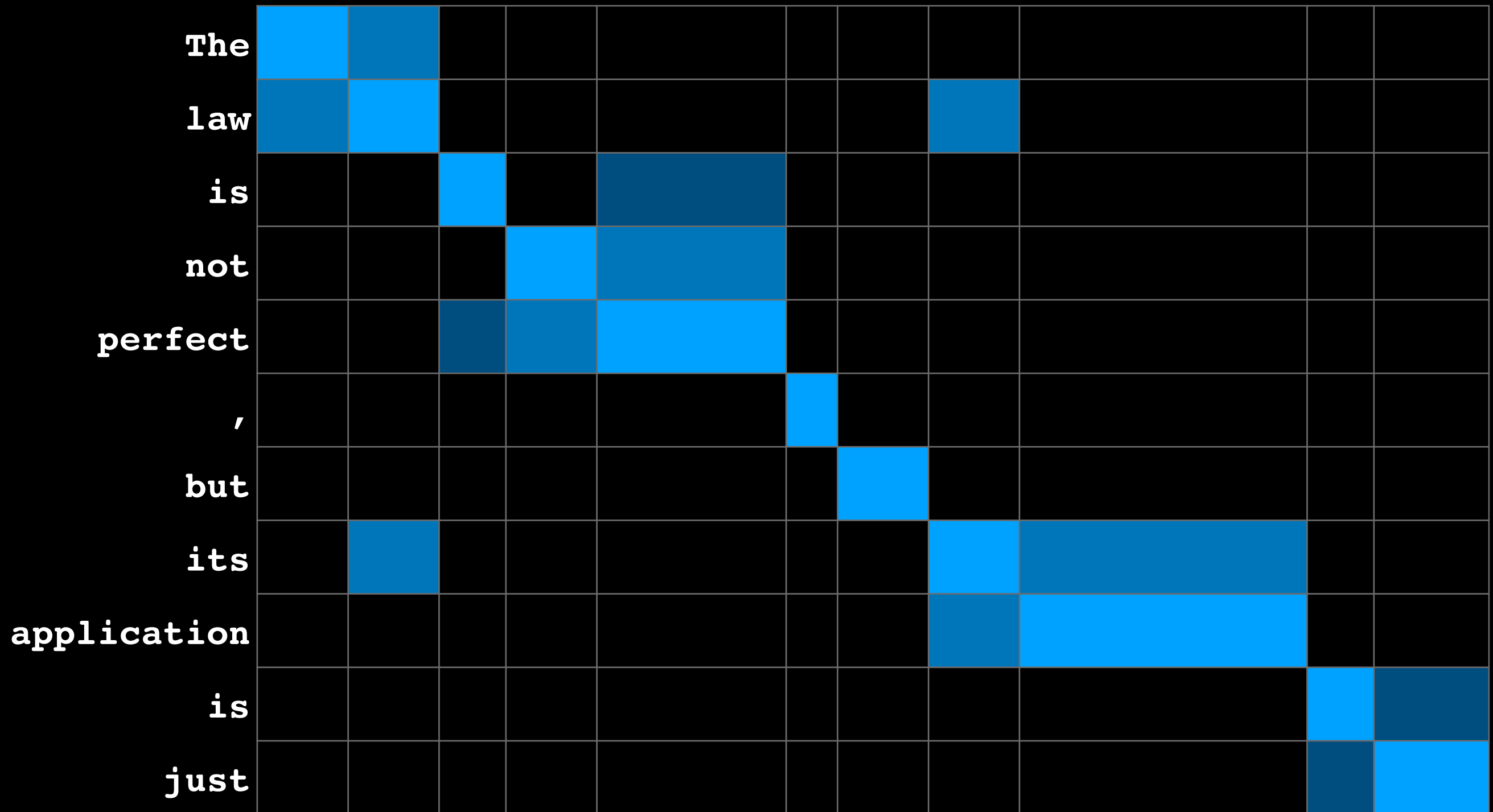
The Attention Mechanism

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)

CNN with Attention

The law is not perfect , but its application is just



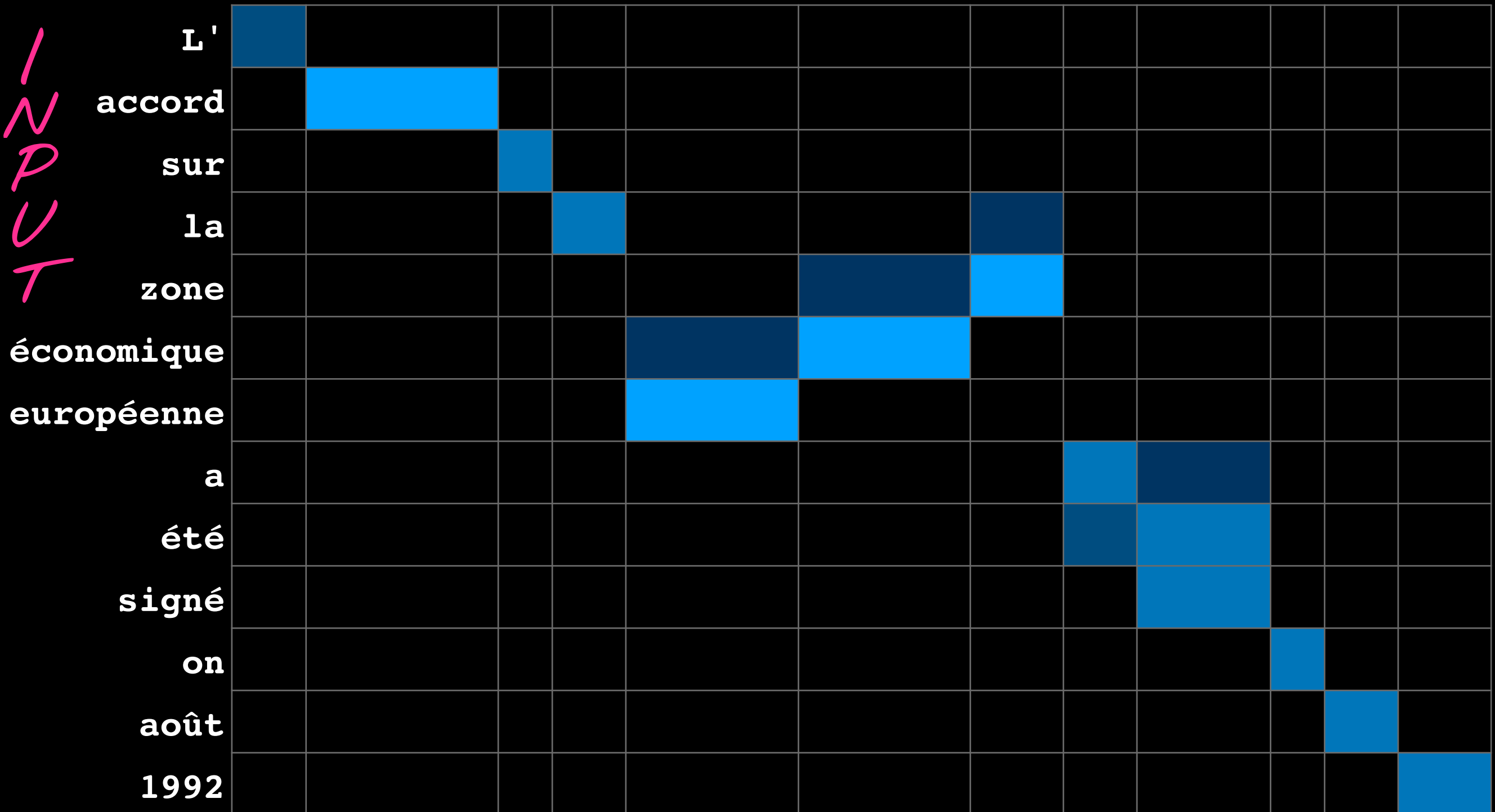
FIND LONG-RANGE DEPENDENCIES

RNN with Attention

OUTPUT

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

INPUT



LEARN REORDERING

Bocconi

RNN with Attention

OUTPUT

Economic growth has slowed down in recent years

INPUT

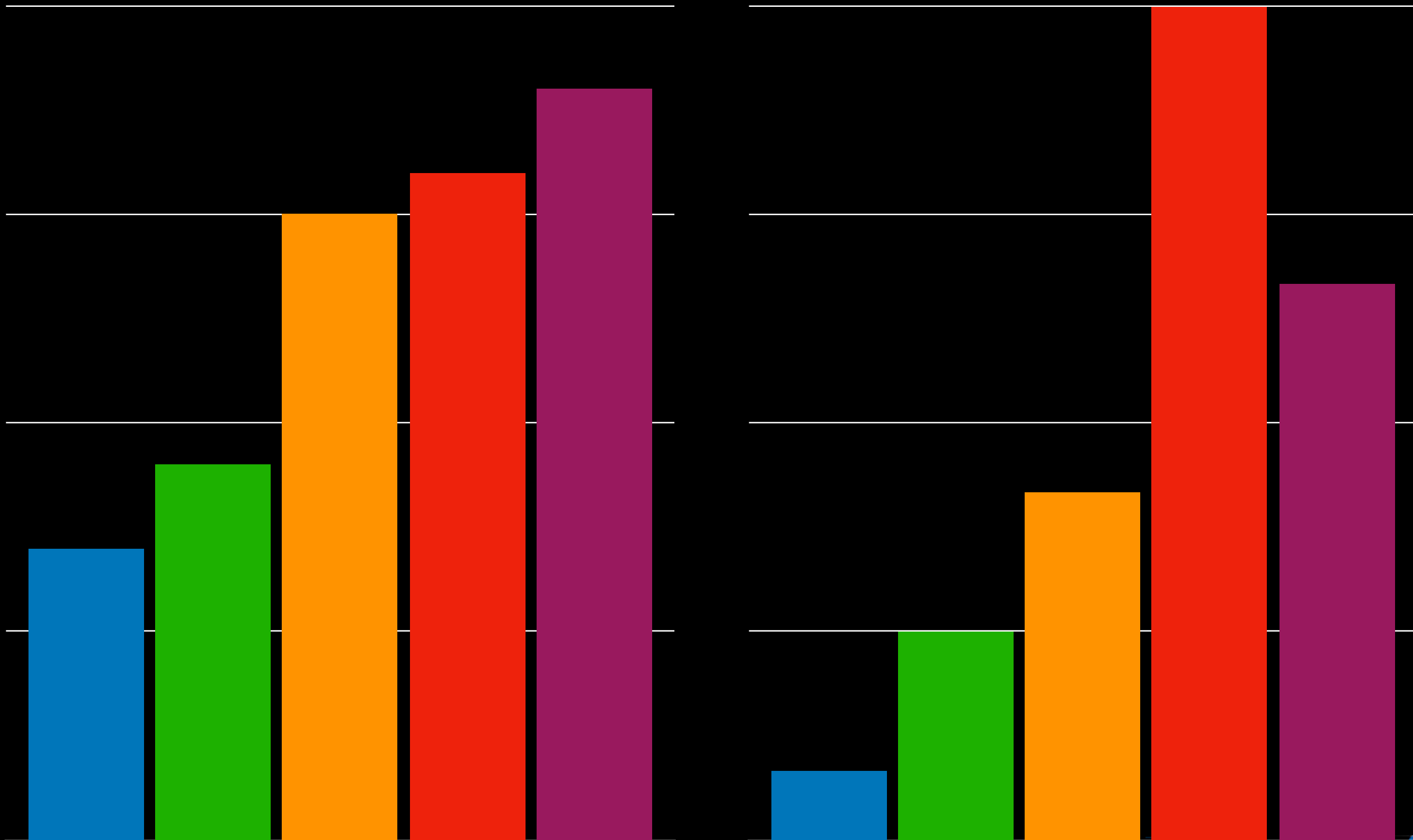
Das
Wirtschafts-
wachstum
hat
sich
in
den
letzten
Jahren
verlangsamt

LEARN REORDERING

Bocconi

Performance

Logistic Regression Feed-Forward CNN RNN CNN+Attention



performance

parameters

Wrapping up

Take Home Points

- **Recurrent Neural Nets** address long-range dependencies
- Condition each word on all previous ones:
 - better **LMs**
 - better **sequence labels**
- **Bidirectional RNNs** condition on following words
- **LSTMs** learn to forget useless input
- **Attention** improves coherence