

# NLP, Deep Learning & Some Examples

Reyyan Yeniterzi

# NLP, Deep Learning & Some Examples

History

Reyyan Yeniterzi

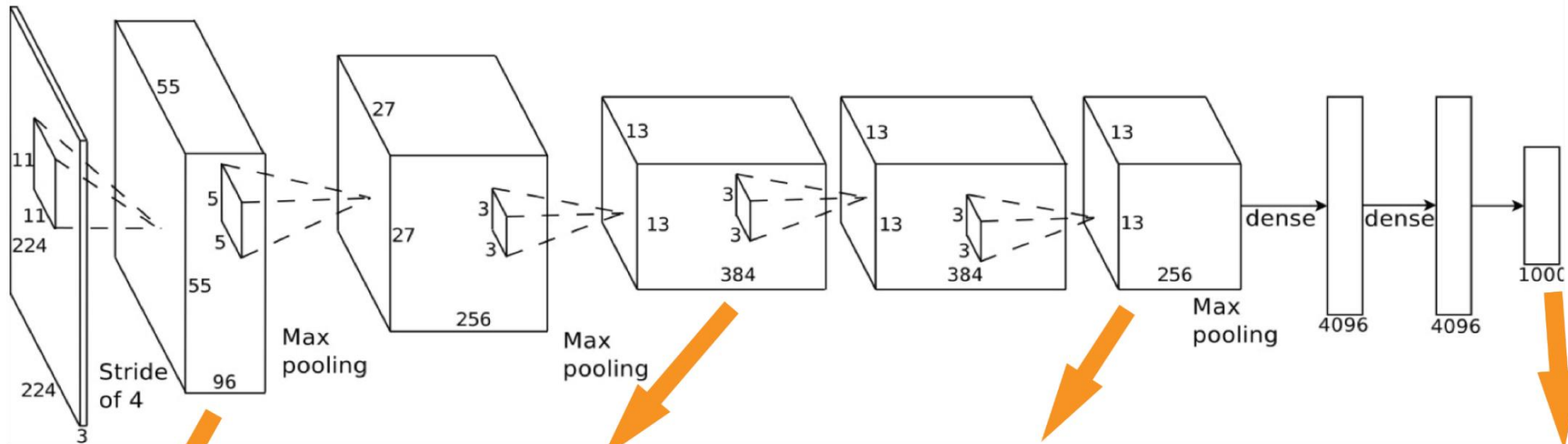
2018

the year of NLP

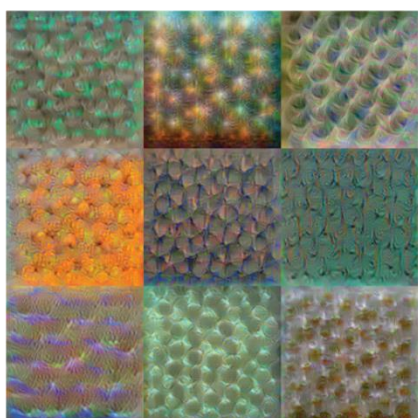
2018

the year of NLP

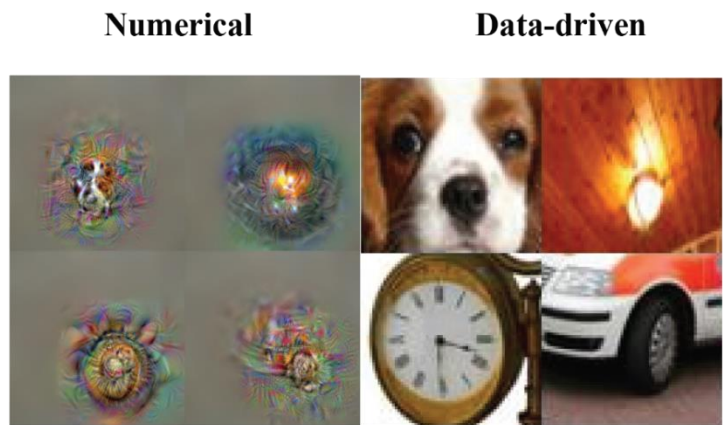
The ImageNet moment of NLP



**Conv 1: Edge+Blob**



**Conv 3: Texture**



**Conv 5: Object Parts**



**Fc8: Object Classes**



“Transfer learning  
will be the next  
driver of ML  
success.”

Andrew Ng,  
NIPS 2016 tutorial

Transfer Learning in NLP?

“You shall know a word by  
the company it keeps”

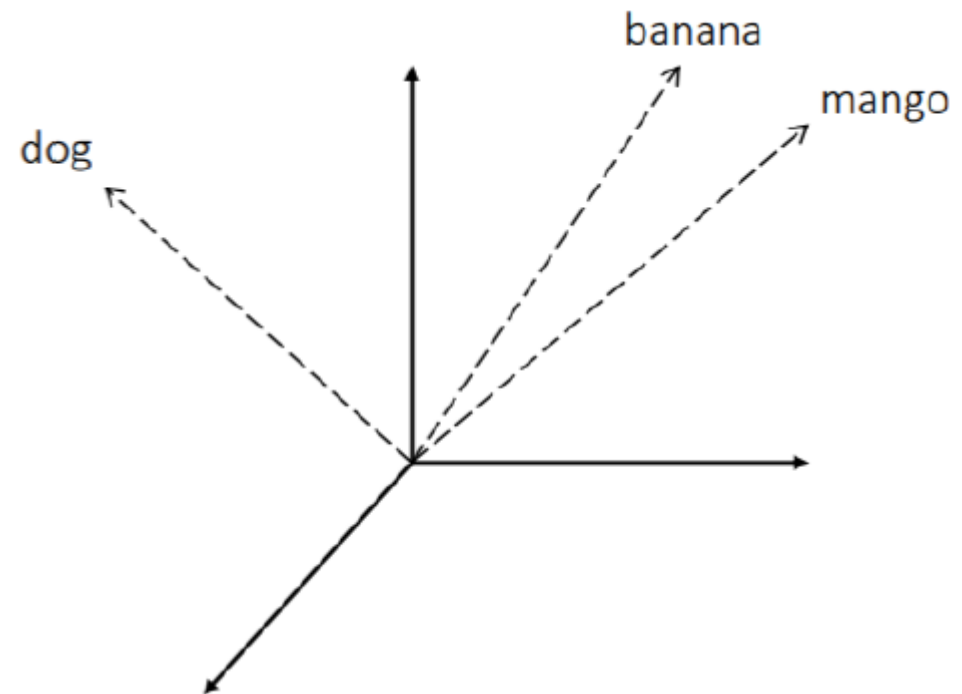
J. R. Firth, 1957

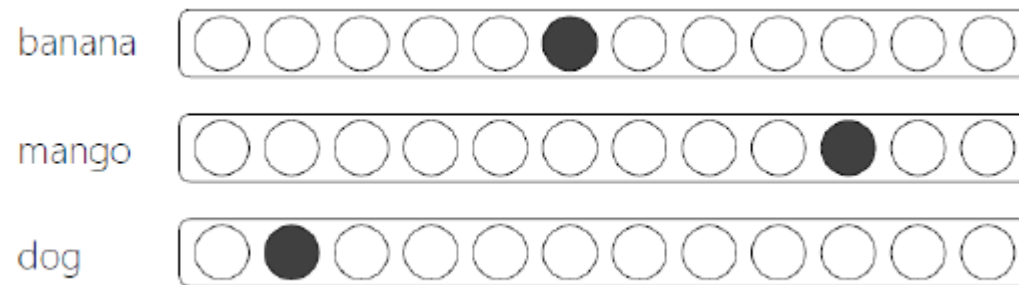
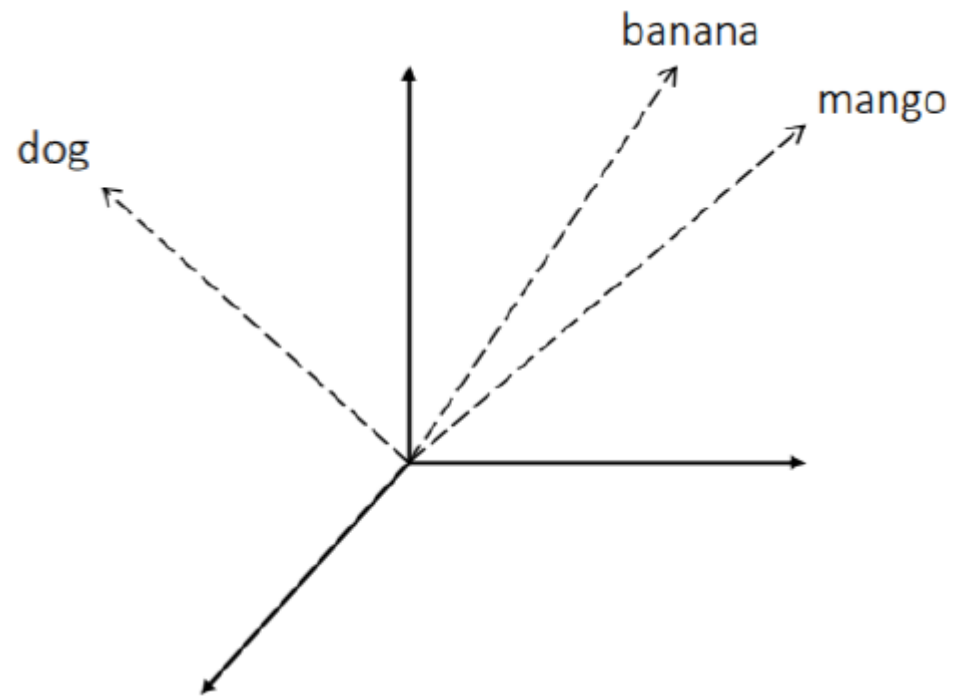


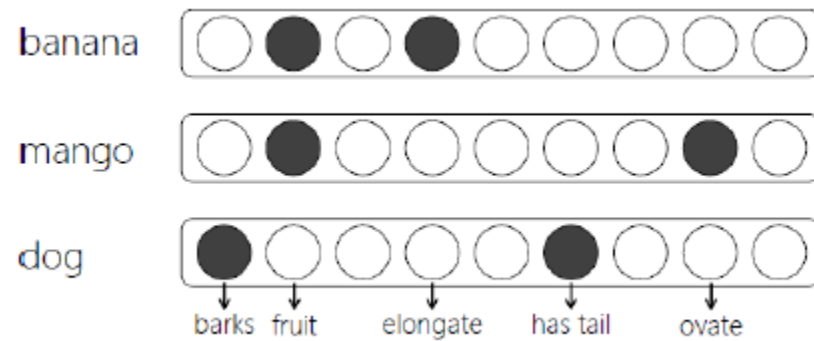
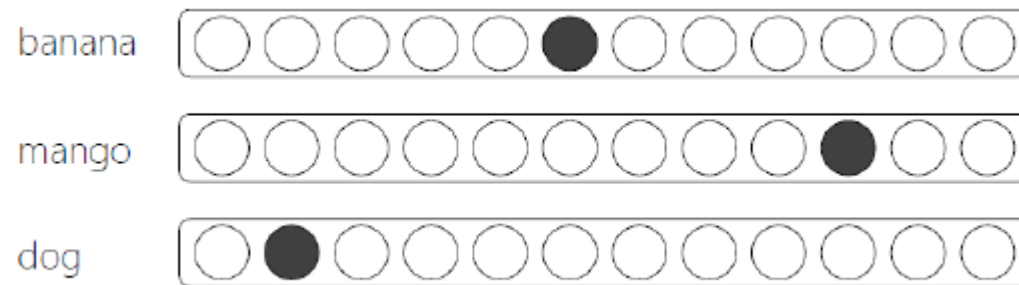
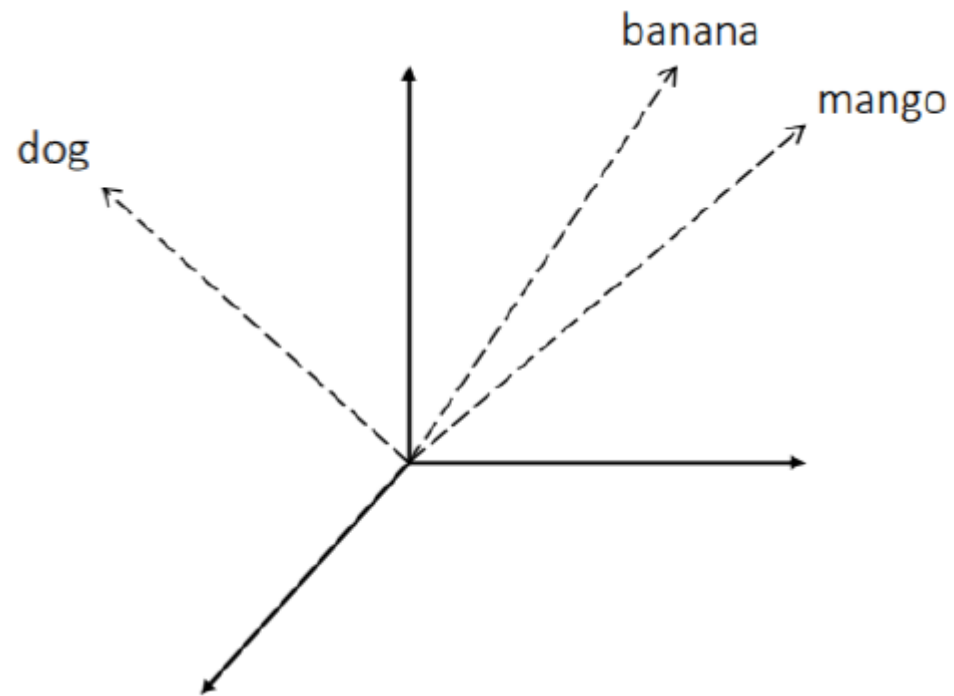
What does *wampimuk* mean?

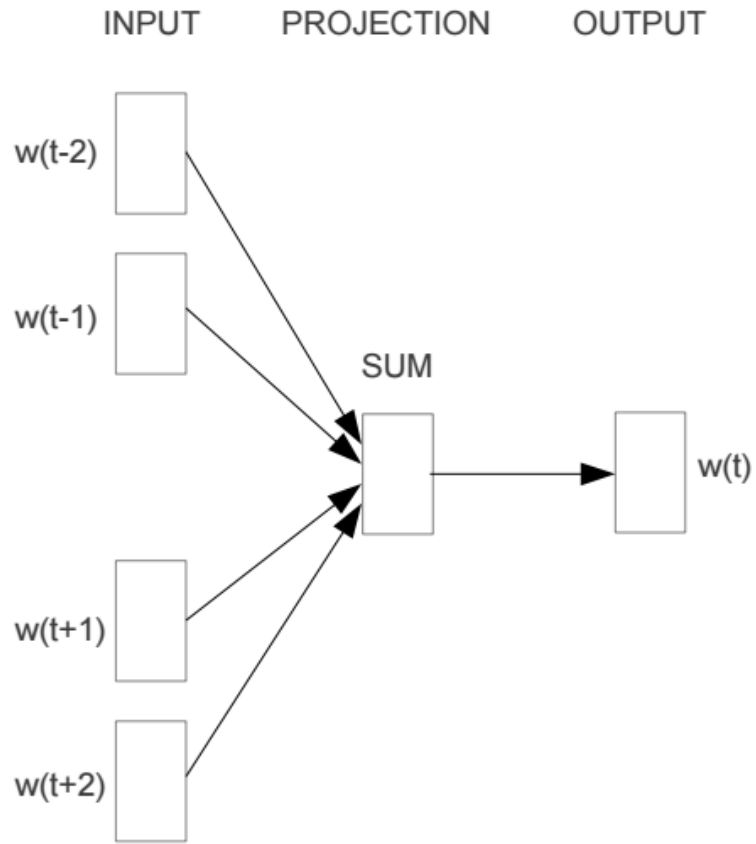
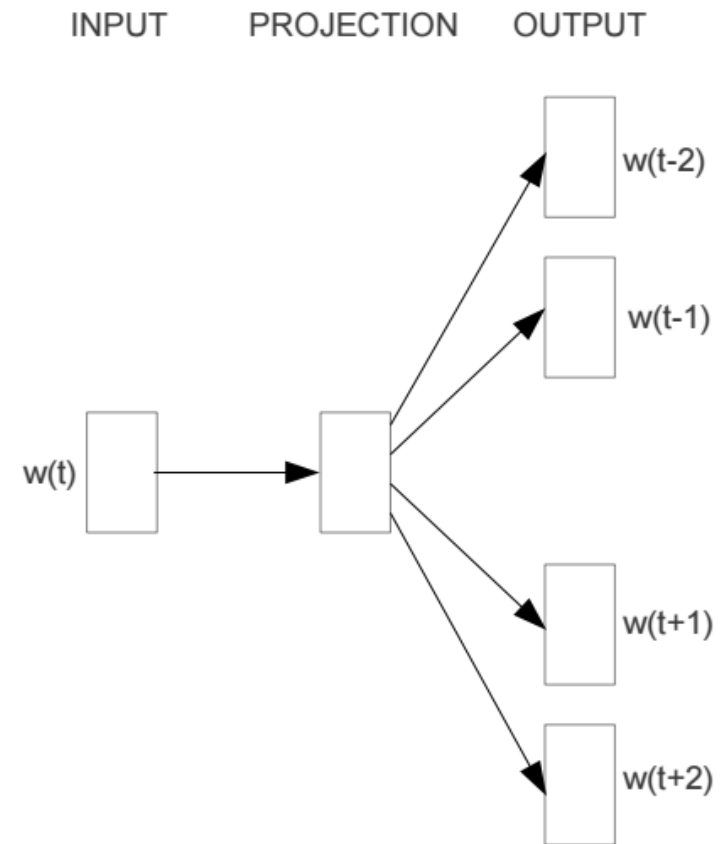
What does *wampimuk* mean?

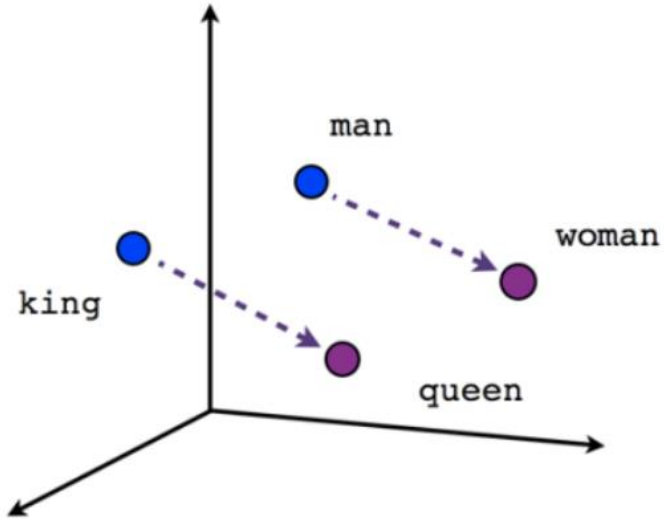
*Marco saw a hairy little wampimuk crouching behind a tree.*



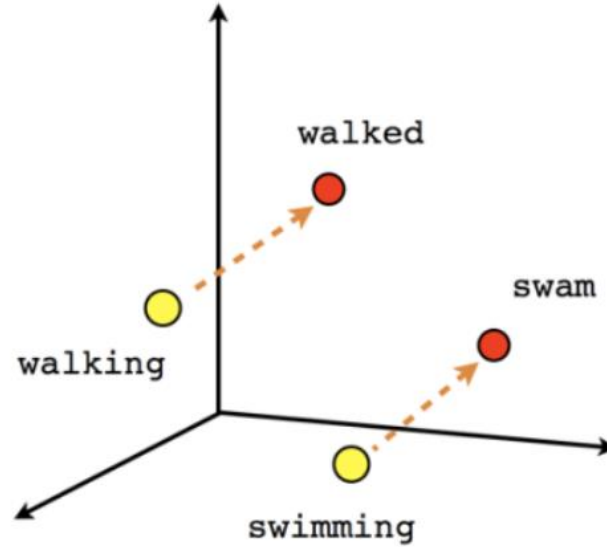




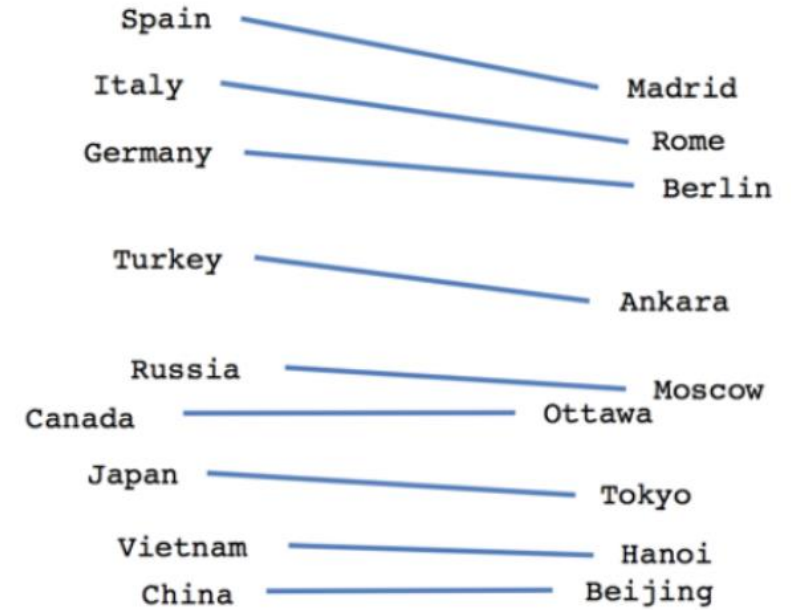
**CBOW****Skip-gram**



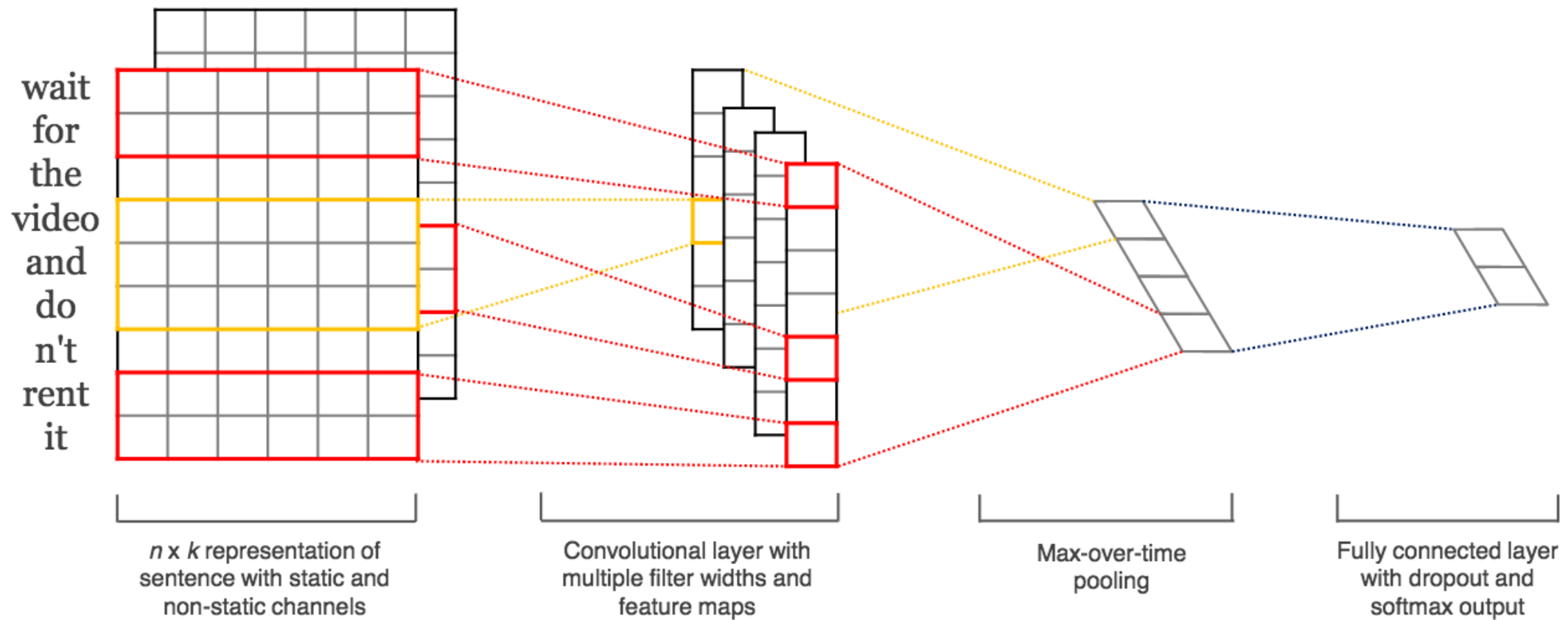
Male-Female



Verb tense

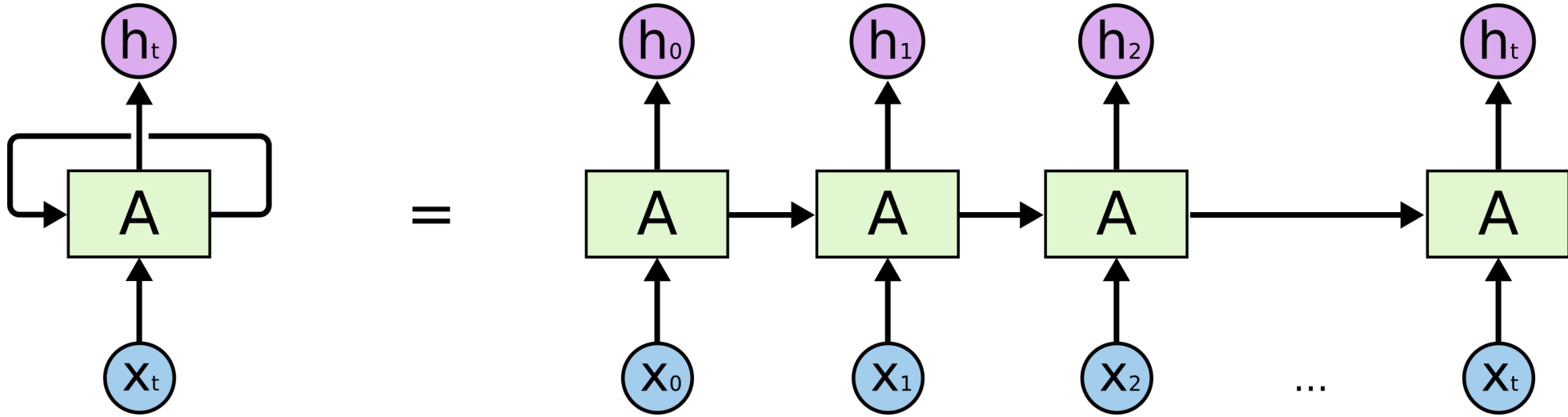


Country-Capital

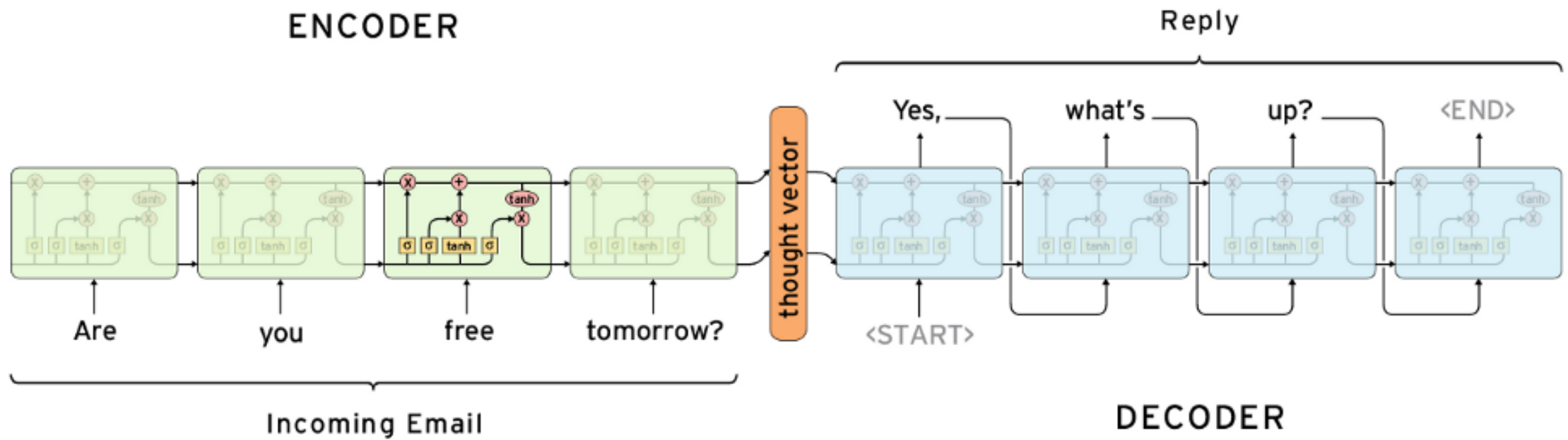




# RNN

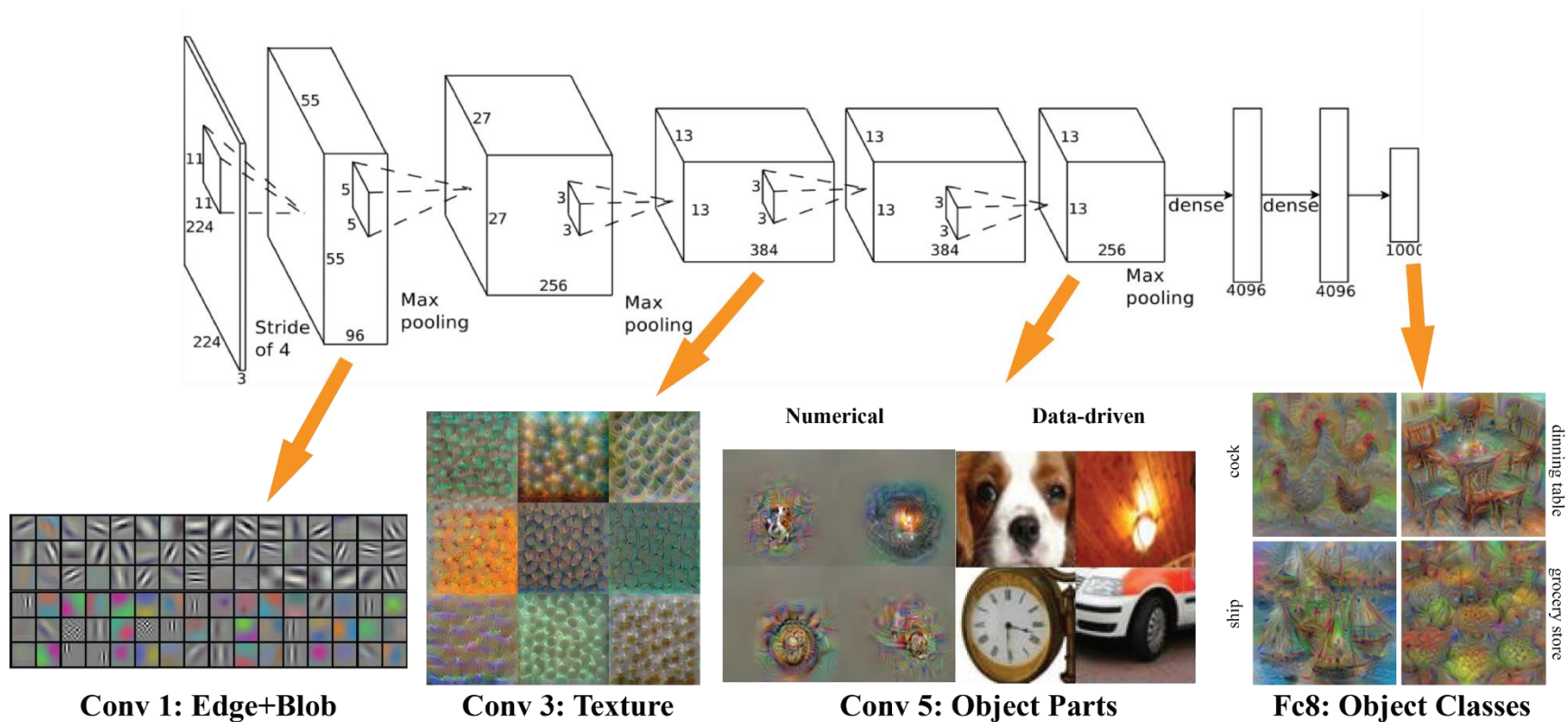


# Encoder&Decoder



# Transfer learning in NLP through the embedding layer

Shallow representation is not enough,  
still require large training sets

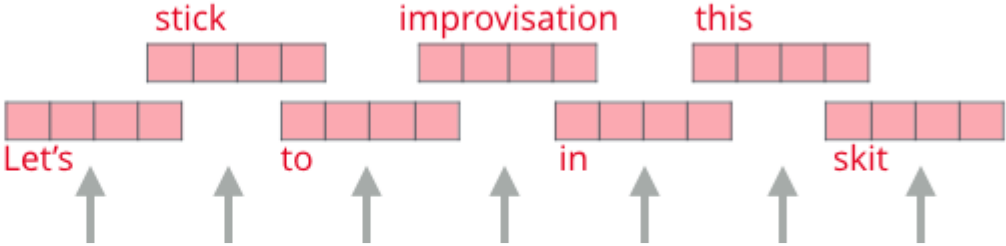


Word embeddings are cool,  
but **context matters!**

Word embeddings are cool,  
but **context matters!**

It is time for **contextual word embeddings**

# ELMo Embeddings



# Embeddings from Language Models



Words to embed

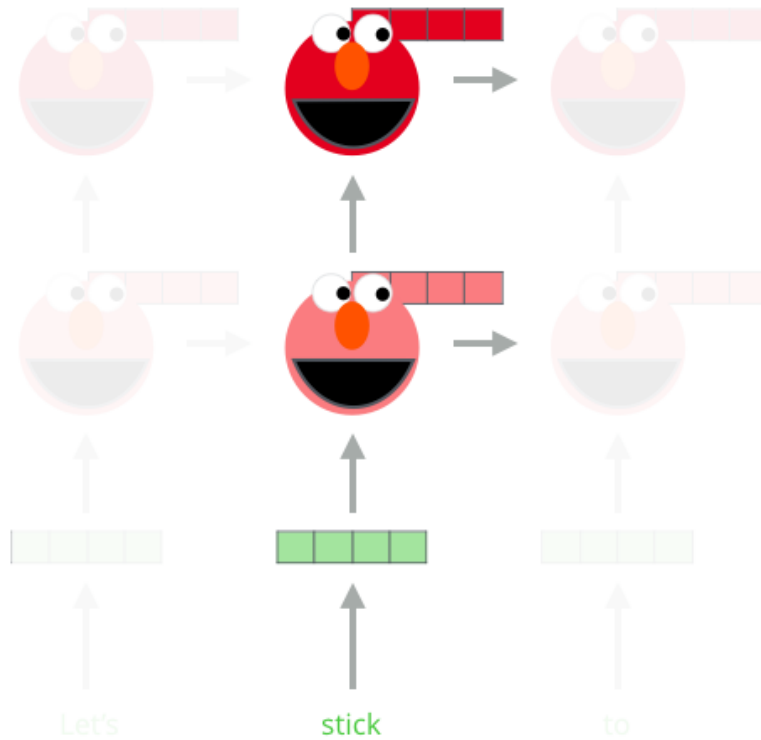


# Embedding of "stick" in "Let's stick to" - Step #2

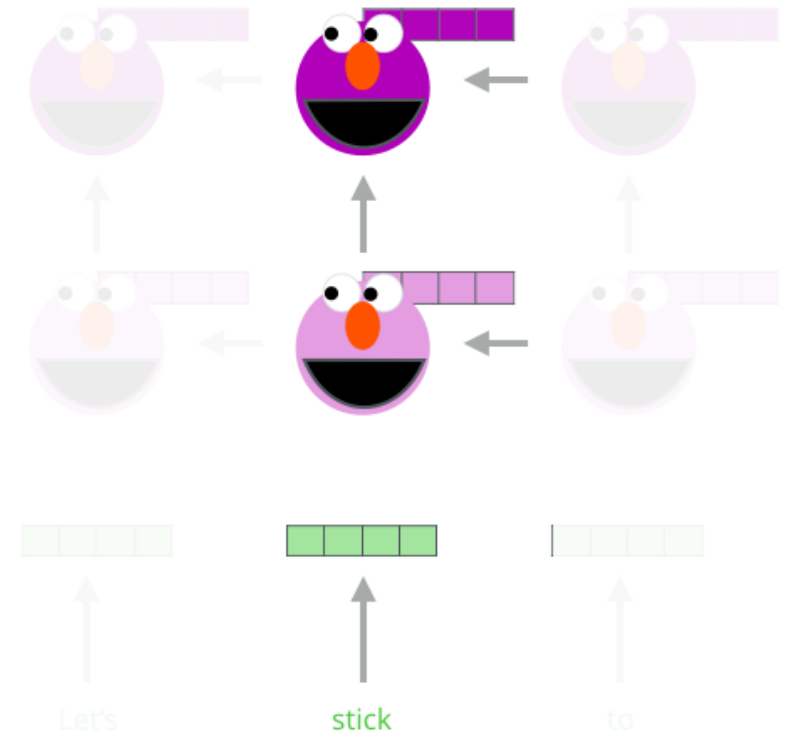
1- Concatenate hidden layers



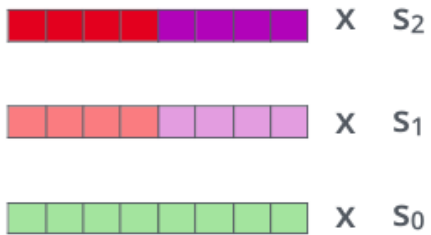
Forward Language Model



Backward Language Model



2- Multiply each vector by a weight based on the task



3- Sum the (now weighted) vectors



ELMo embedding of "stick" for this task in this context

Peters, Matthew E., et al. "Deep contextualized word representations." 2018.

<https://jalammar.github.io/illustrated-transformer/>

2018

ELMo

Contextualized embeddings

---

2013

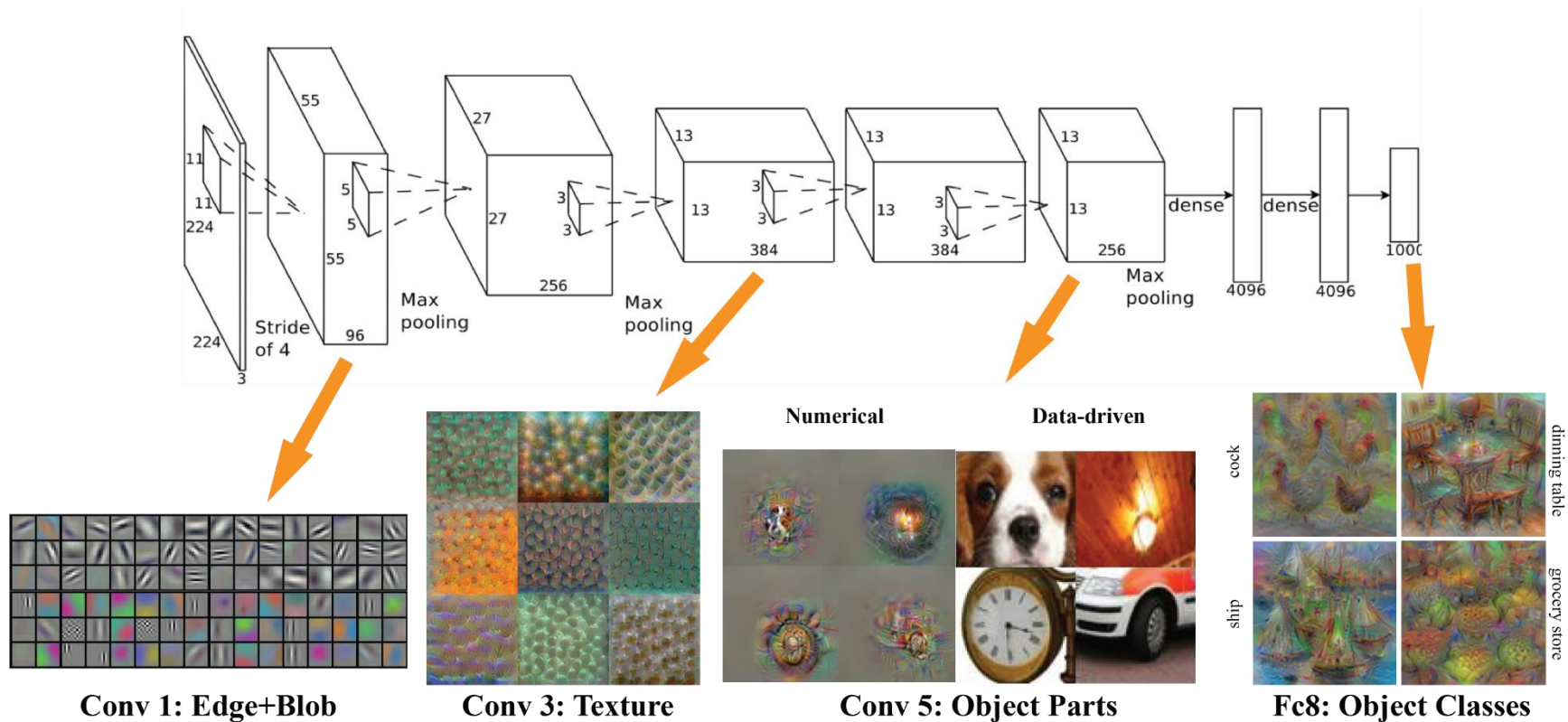
Word2Vec

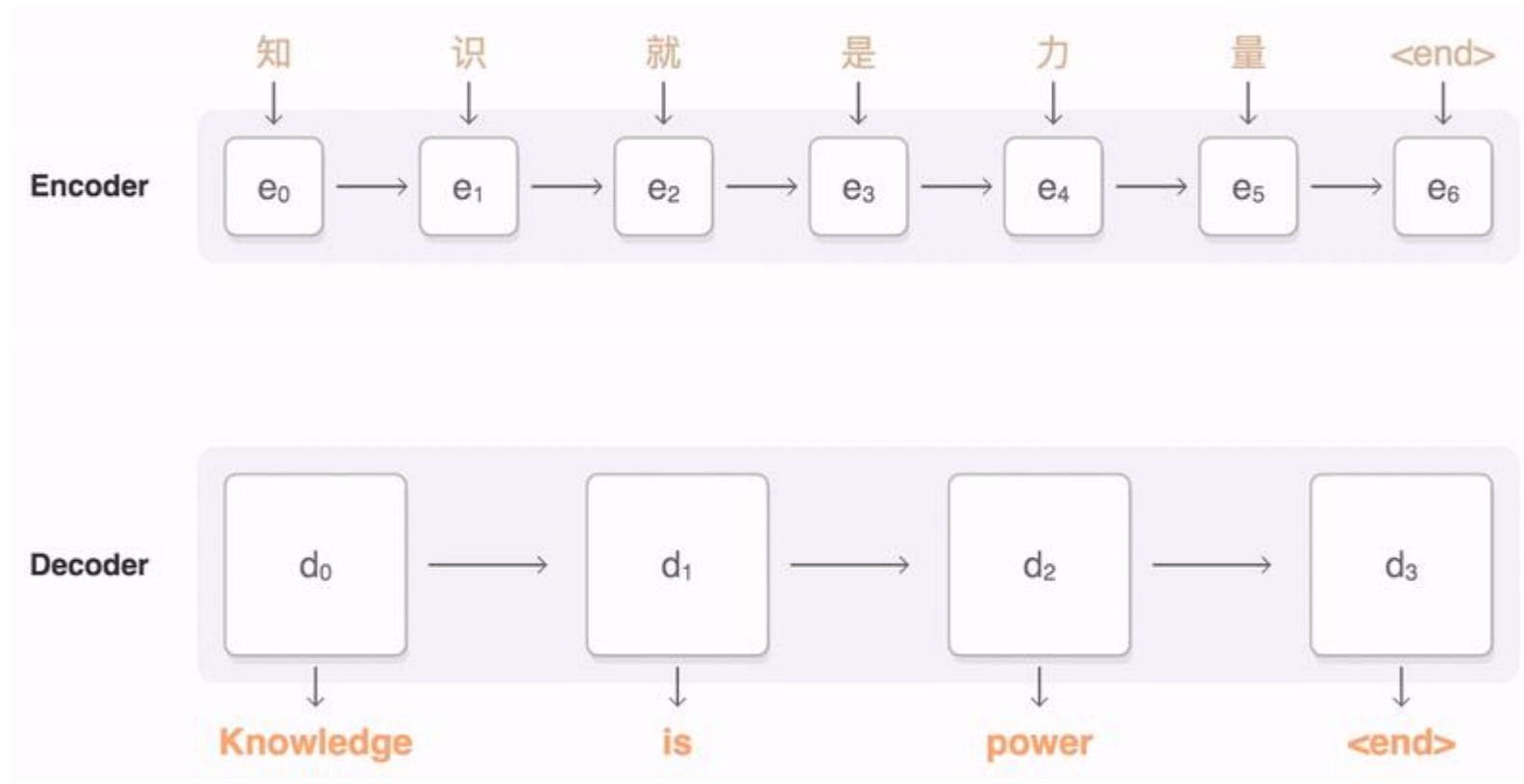
Context free embeddings



# Transfer learning in NLP through the embedding layer

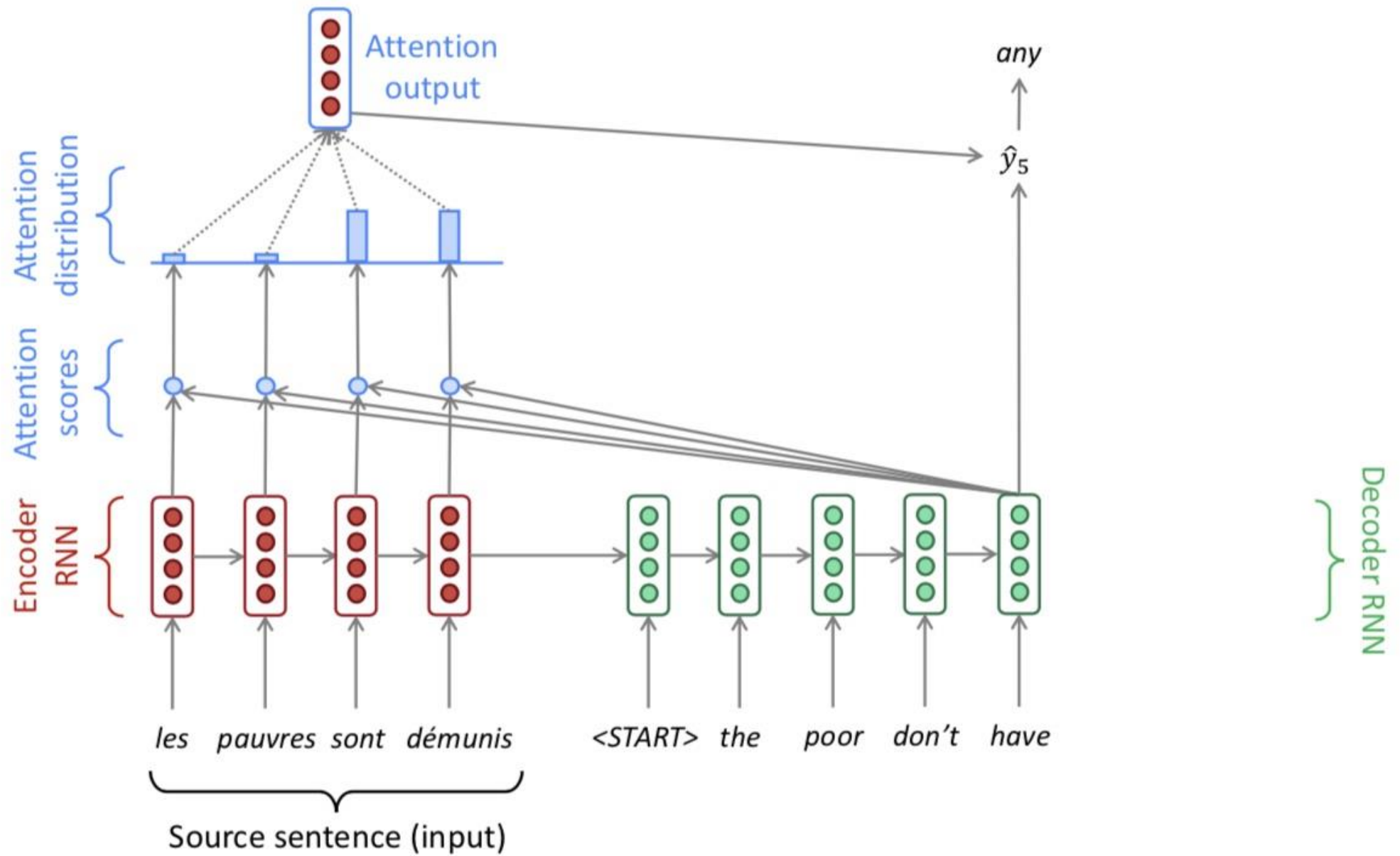
Shallow representation is not enough,  
still require large training sets

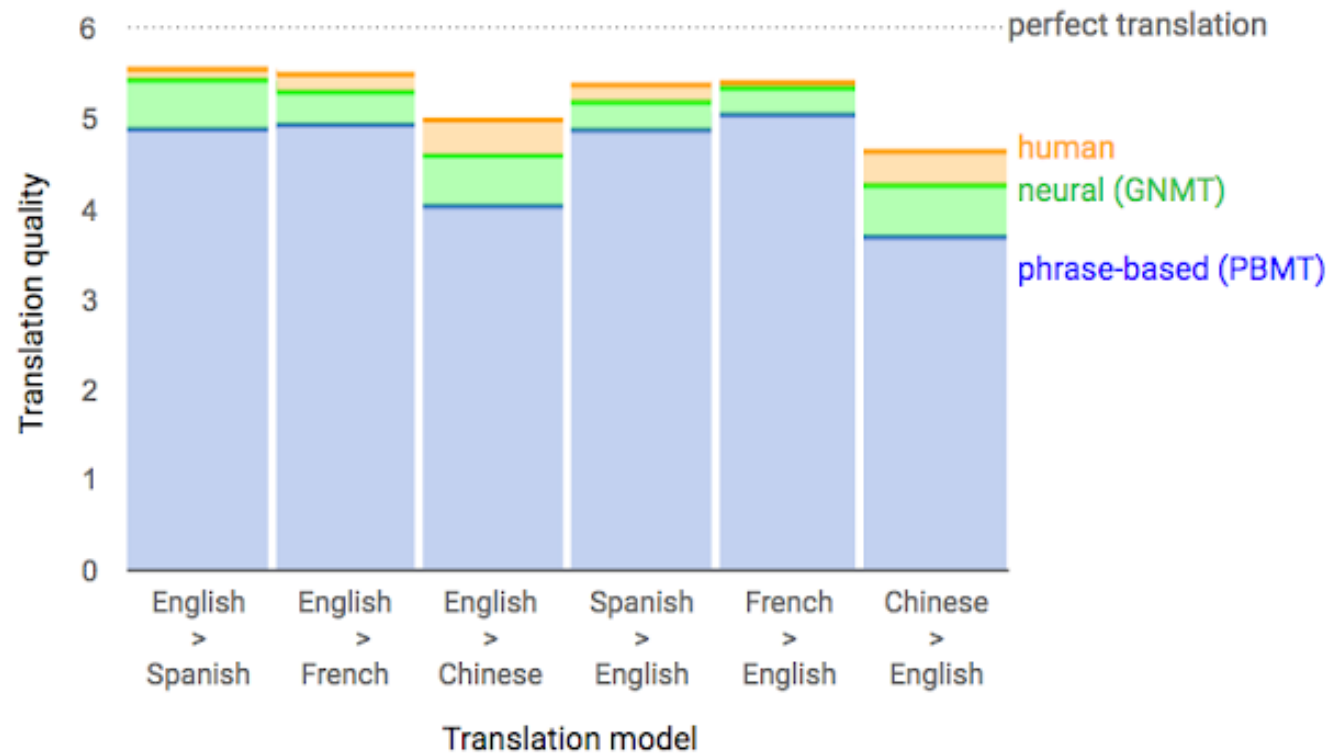


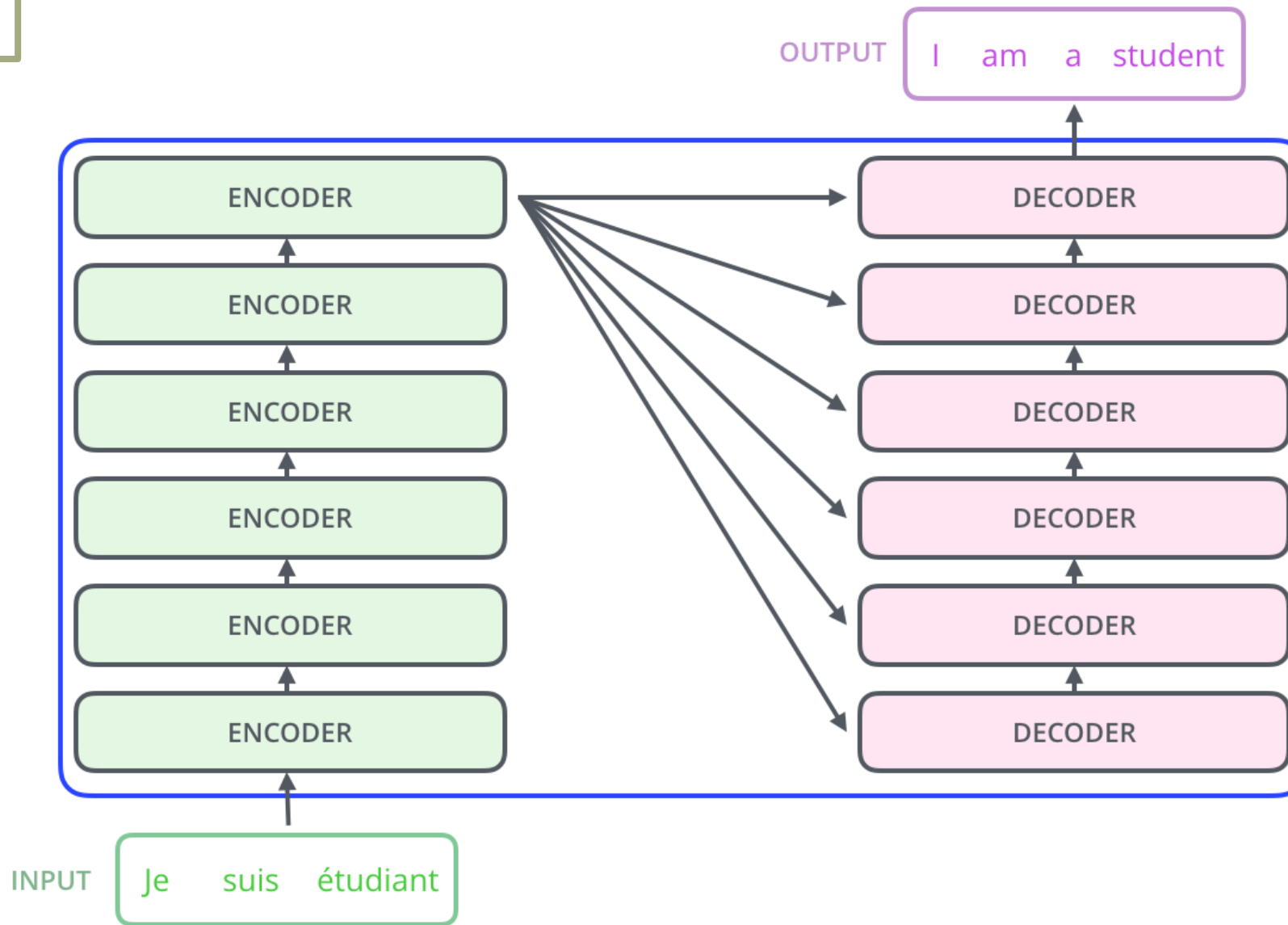


Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." 2014.

<https://ai.googleblog.com/2016/09/a-neural-network-for-machine.html>

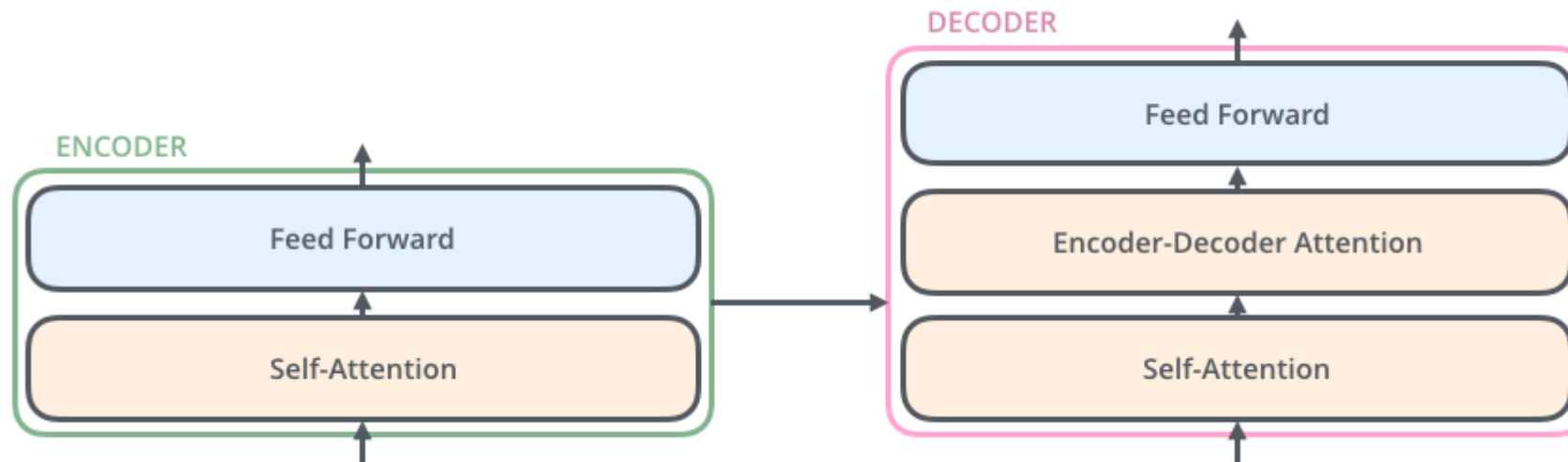






Vaswani, Ashish, et al. "Attention is all you need." *NIPS*. 2017.

<https://jalamar.github.io/illustrated-transformer/>



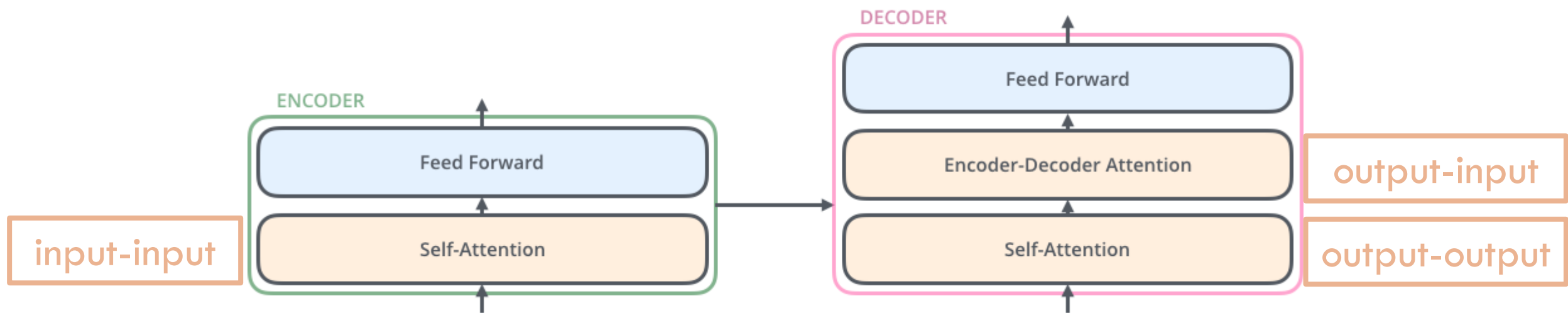
Vaswani, Ashish, et al. "Attention is all you need." *NIPS*. 2017.

<https://jalammar.github.io/illustrated-transformer/>

Transformer

Self-Attention

2017

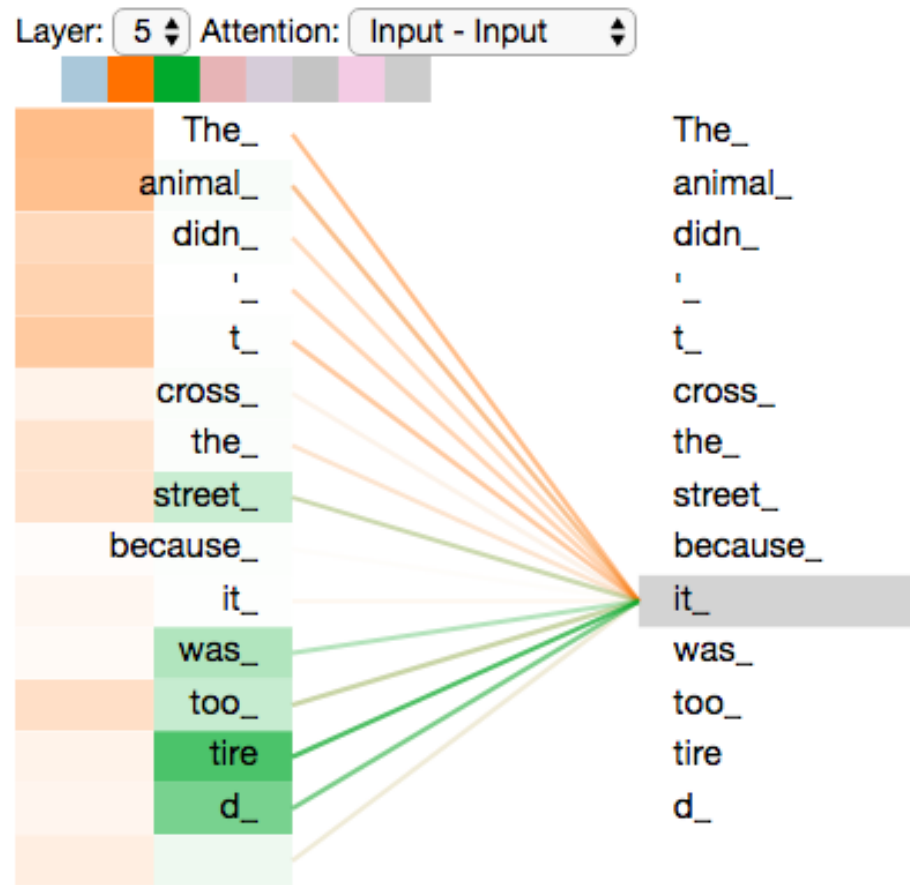


Transformer

2017

Self-Attention

input-input





Transformers  $>$  RNNs (LSTMs)

Attention is all you need!

2018

ELMo

---

2017

Transformer

Attention is all you need  
No need for RNNs

---

2014

Attention

---

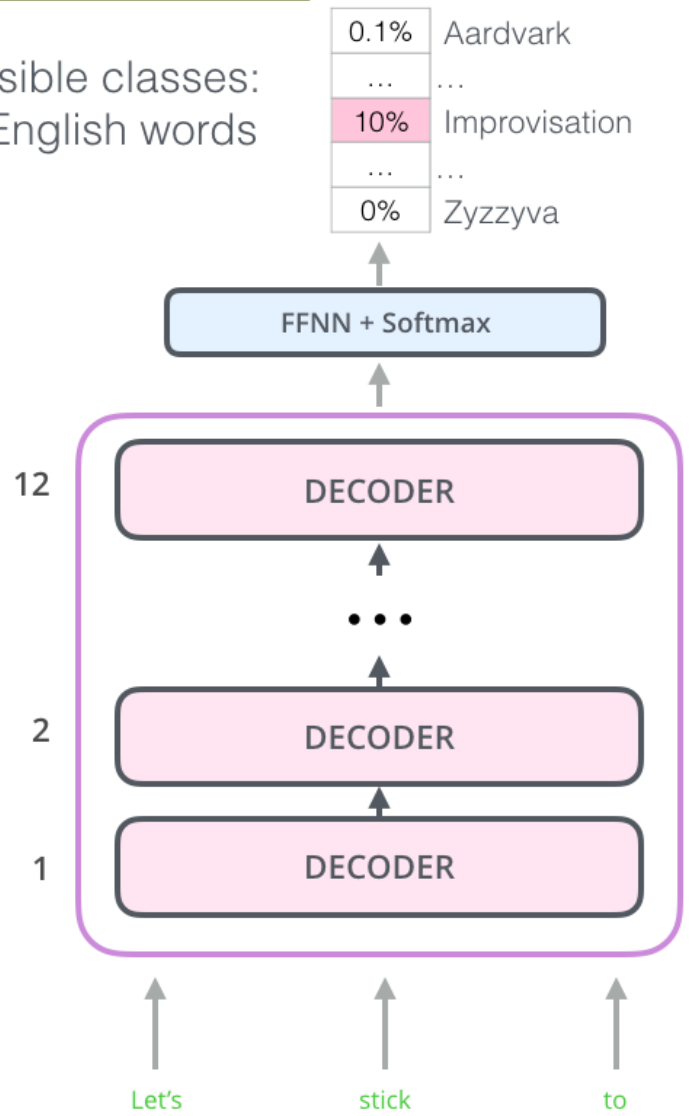
2013

Word2Vec

## Unsupervised LM Training

Possible classes:  
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva



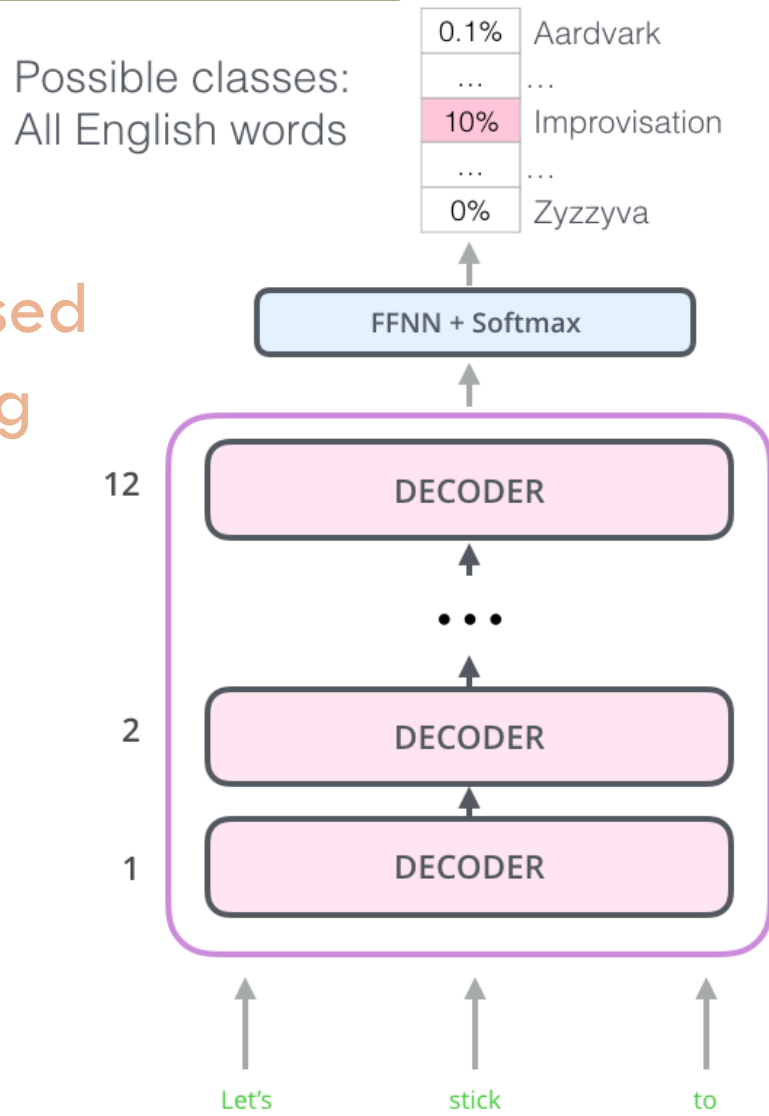
Radford, Alec, et al. "Improving language understanding by generative pre-training." 2018.

<http://jalamar.github.io/illustrated-bert/>

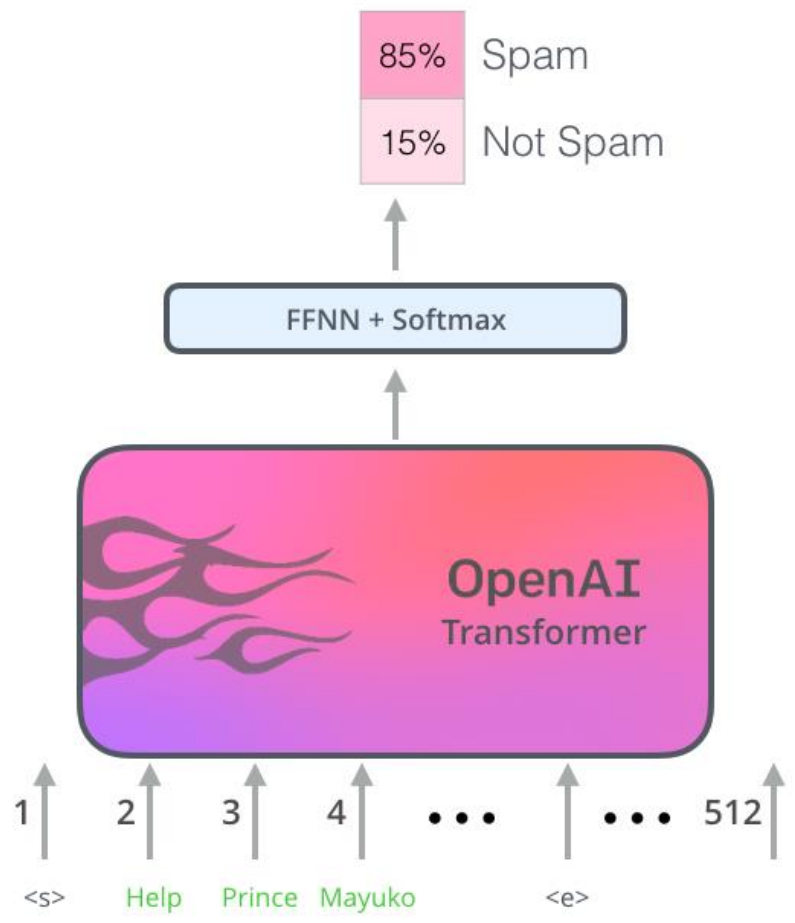
# OpenAI Transformer

2018

Unsupervised  
LM Training

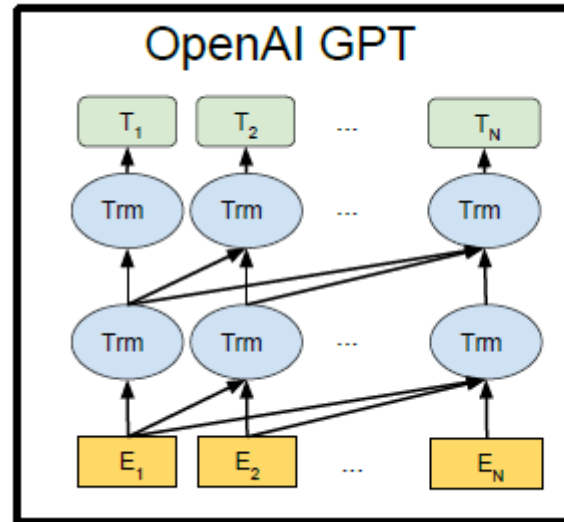


Supervised  
Fine-Tuning

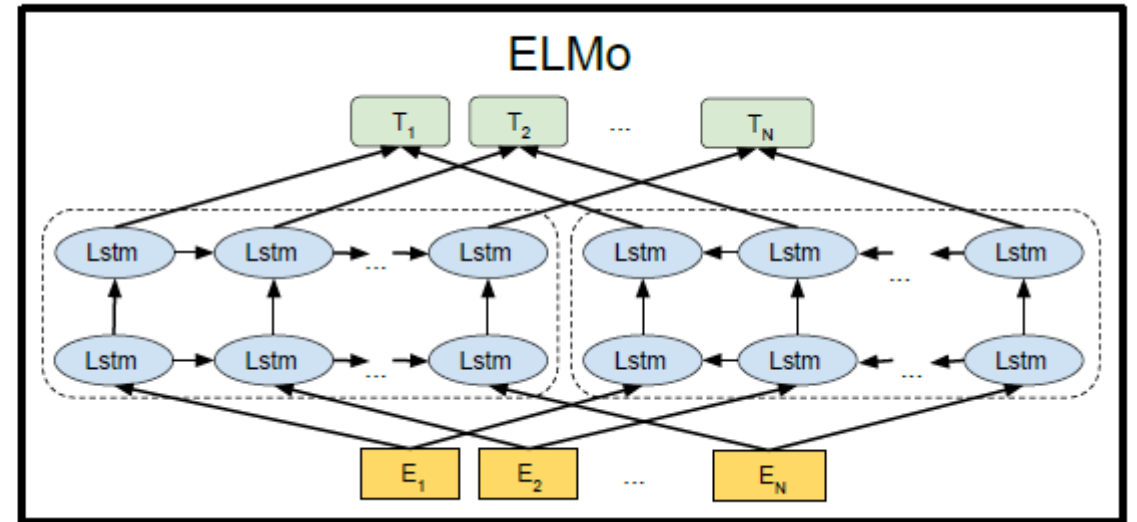


Radford, Alec, et al. "Improving language understanding by generative pre-training." 2018.

<http://jalamar.github.io/illustrated-bert/>



Unidirectional  
Transformer  
(Transformer  
Decoder)



Shallow BiDirectional LSTM

Training a LM using  
both left and right context jointly

Training a LM using  
both left and right context jointly

## Masked Language Model + Transformer Encoders



# BERT

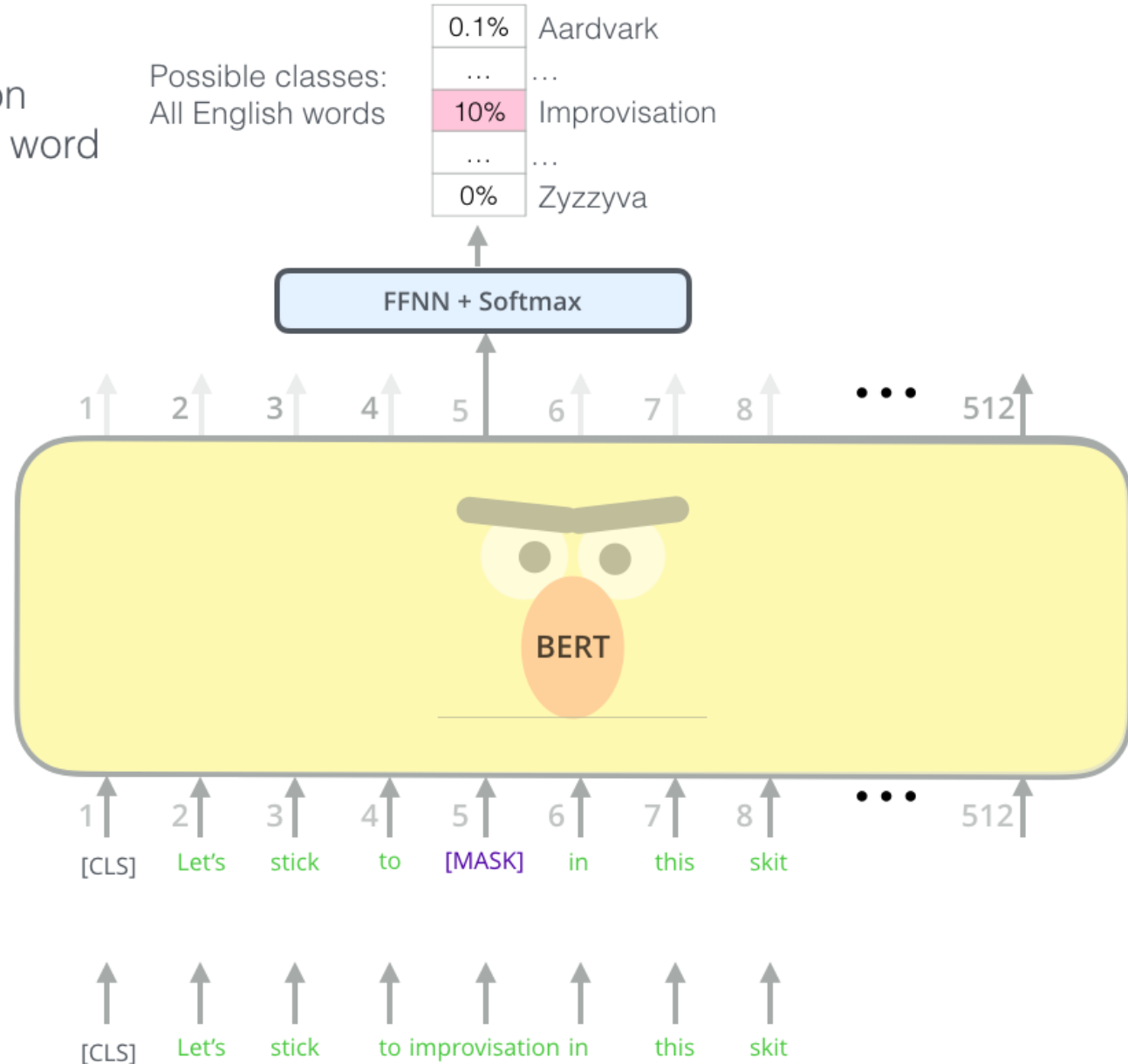
2018

Use the output of the masked word's position to predict the masked word

## Bidirectional Encoder Representations from Transformers

Randomly mask 15% of tokens

Input





1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

### Semi-supervised Learning Step

**Model:**



**Dataset:**



**Objective:**

Predict the masked word  
(language modeling)

2 - **Supervised** training on a specific task with a labeled dataset.

### Supervised Learning Step

**Model:**  
(pre-trained  
in step #1)

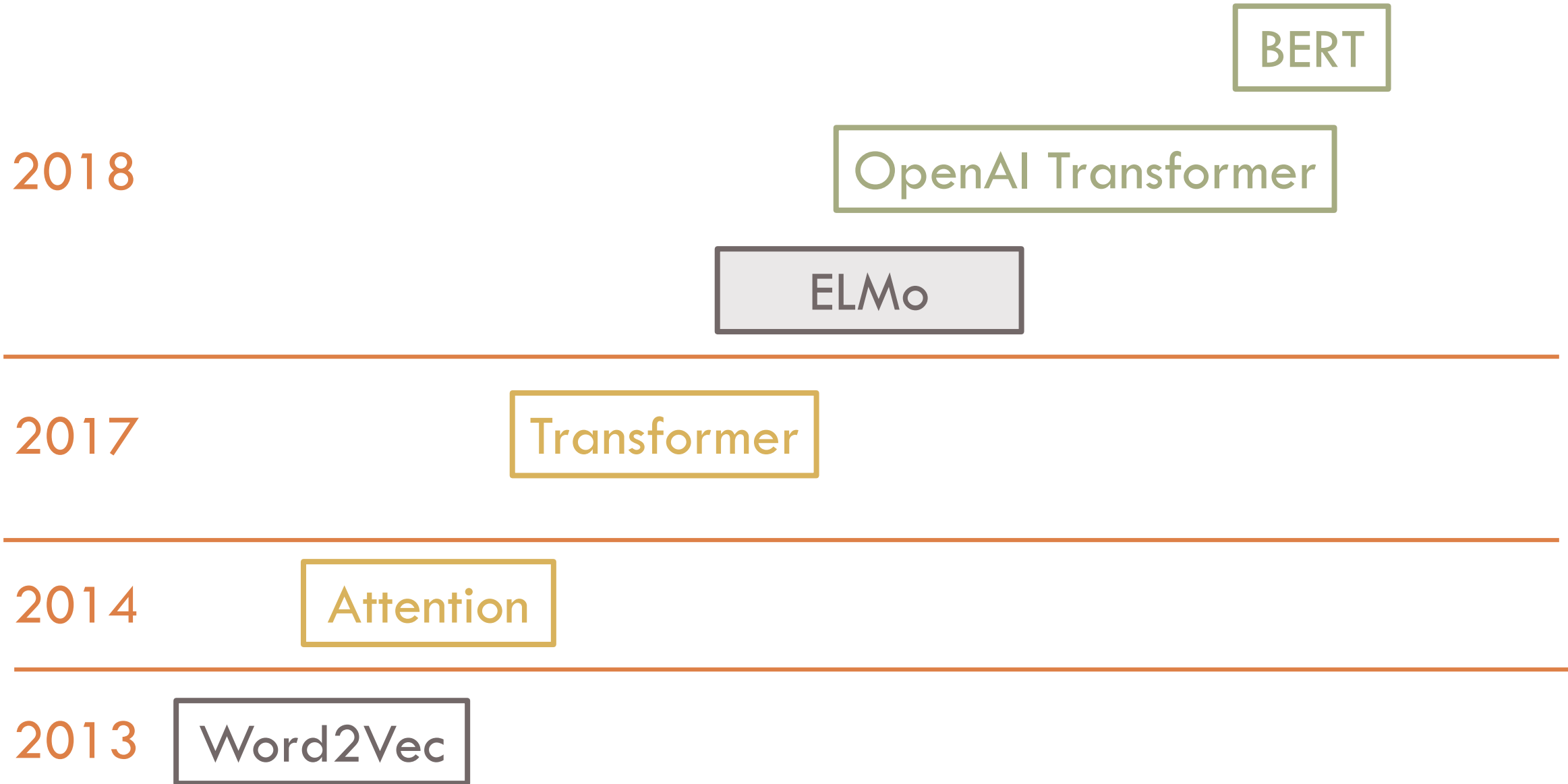


**Classifier**

75% Spam  
25% Not Spam

**Dataset:**

Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam



123

SESAME STREET®



Any Questions?

