# Lecture 2:
# Encoder-Decoder Models

Barbara Plank
LMU Munich

AthNLP2024
Athens, Greece

# Motivation

*It is all about sequences!*

(most?!) interesting data is **sequential** in nature

*It, is, all, about, sequences!*

*[s, e, q, u, e, n, c, e, s]*

# What you have heard so far:

‣ **Lecture 1:** ML Fundamentals

  ‣ Linear Classifiers

  ‣ Non-Linear Classifiers

    ‣ Feedforward Neural Network

    ‣ CNN

  ‣ Feature Representations

    ‣ Sparse binary features

    ‣ Continuous dense features (word2vec, embedding layers)



**Linear Classifiers and Neural Networks**

Dog

Cat

**Linear Classifier**

# How to deal with sequences?

# Today's roadmap

▸ **Part I: Fundamentals**

  ‣ Intro, Motivation & Short History

  ‣ Language Models (n-grams, FFNN-LM, Recap: FFNN)

▸ **Part II: Representations & Beyond FFNN**

  ‣ RNNs (GRU/LSTMs), Attention

  ‣ Contextualised Representations (ELMo)

▸ **Part III: Transformer & LLMs**

  ‣ The Transformer, Masked LMs (BERT), Pre-training & Fine-tuning

  ‣ Prompting, LLMs & Caution

# Before we start

www.menti.com - Code 2727 6277

# Fundamentals

**Part I**

# NLP today: **LLMs** everywhere!

# Encoder-decoder models

general family of models

output

DECODER

ENCODER → input representation / encoding

input

# A step back…
# How did the field evolve?

# NLP ❤️ Machine Learning

**Symbolic Processing**

e.g. **Brill's** rule-based tagger (1992)

**ELIZA** (1966)

Epoch 1

**from hand-crafted rules to ML**

**Statistical NLP**

e.g. **Collin's** structured perceptron (2002)

Epoch 2

approx. 1980s

**x**: the dog barks

```
if prev_w = DET and ..
              tag=NOUN
```

**X**: | 1 | 0 | .. | 0 | 1 | 0 |

w_i=dog     w_i-1=the

classic **sparse** n-hot encodings

12

# The emergence of
# deep learning (in NLP)

# In Speech Recognition



**Loud and clear**

Speech-recognition word-error rate, selected benchmarks, %

Log scale

Switchboard

Switchboard cellular

Meeting speech

Broadcast speech

The Switchboard corpus is a collection of recorded telephone conversations widely used to train and test speech-recognition systems

RNNs

IBM, Switchboard

Microsoft, Switchboard

5.9%

1993  96  98  2000  02  04  06  08  10  12  14  16

Sources: Microsoft; research papers

**(Source: The Economist)**

2010

# In Computer Vision



(src: slide by Fei-Fei Li)

2012

Golden Years for MT (Georgetown Exp.)

1966: Weizenbaum develops ELIZA at MIT

1st "AI winter"

1980s: knowledge representation, data & corpora, statistical NLP

2nd "AI winter"

1995-2005: progress in data, hand-crafted features and ML

2007-2017: representation learning & DL: word2vec, NNs from RNN/CNN to Transformer

2015: DL "Tsunami hit NLP"

2018-today: Pre-training - from BERTs to GPTs

2023-today: 💥 Explosion of LLMs

2022: ChatGPT to the public

2018: "3rd time NNs have threatened a revolution but only the 1st time they have delivered"

2018: "Multi-task learning wave"

Epoch 1: Symbolic Processing

Epoch 2: Statistical Processing

Epoch 3: Deep Learning (DL) for NLP

1960   1970   1980   1990   2000   2010   2020

16

# Back to the roots: (traditional) Language Models (LMs)

# Predicting the next word:
# A Simple (?) Exercise

- ▸ www.mentimeter.com
  Room: (see code)

# You probably use a LM every day…

# More examples

‣ .. is perhaps the best ___ ?

So let's look deeper at LMs: from traditional LMs over neural LMs to contextualized embeddings

# What is a Language Model (LM)?

‣ A computational model that can be used for either of the following two tasks is called a Language Model (LM):

   ‣ to compute the probability of a text*

$$P(\textbf{today is a great day}) = \;??$$

   ‣ to compute the probability of the next word

$$P(\textbf{day} \,|\, \textbf{today is a great}) = \;??$$

* (can be a text, sentence, phrase,…)

# A Language Model - Formally

‣ Given a sequence of words: $(w_1, \ldots, w_d)$

‣ An LM models the probability: $P(w_1, \ldots, w_d)$

‣ Without loss of generality (**Chain Rule**):

$$P(w_1, \ldots, w_d) = P(w_1)P(w_2 \,|\, w_1)P(w_3 \,|\, w_1, w_2)\ldots$$

$$= P(w_1)\prod_{i=2}^{d} P(w_i \,|\, w_1, \ldots, w_{i-1})$$

conditioned on full history

$P($**Athens**$\,|\,$**an awesome summer school this year is in**$)$

# Markov assumption

‣ However, finding good estimates for the probabilities of longer sequences is a problem

‣ Therefore, a common assumption in a traditional LM is to make the Markov assumption:

$$P(x_1, \ldots, x_d) = \prod_{i=1}^{d} P(x_i \mid x_1, \ldots, x_{i-1}) \approx \prod_{i=1}^{d} P(x_i \mid \boxed{x_{i-(n-1)}}, \ldots, x_{i-1})$$

$$p(\boldsymbol{w}) = p(w_1) \times$$
$$p(w_2 \mid w_1) \times$$
$$p(w_3 \mid \cancel{w_1}, w_2) \times$$
$$p(w_4 \mid \cancel{w_1, w_2}, w_3) \times$$
$$\cdots$$

Markov: forget "distant" past

**Valid for language?** No…
**Is it practical?** Often.

**n-th order Markov assumption:**
**Only look at history of most recent n-1 words**

**-> n-grams**

Adapted from Chris Dyer

# How to learn the parameters of an n-gram LM?

- **(Pre-deep learning LMs):** Learn a count-based **n-gram** (traditional) Language Model by **collecting statistics** of **n-grams** from a corpus to estimate the parameters of the model (maximum likelihood)

- **n-gram**: a chunk of consecutive words

  - Example: n=2 (bigram):  $P(x_i \mid x_{i-1})$

    Collect statistics for "to buy", "buy a", "a house".... with C=count():

    $$P(\mathbf{w} \mid \mathbf{to}) = \frac{C(\mathbf{to\ w})}{C(\mathbf{to})}$$

  - Higher order n-grams, e.g. trigram LM: $P(x_i \mid x_{i-2}, x_{i-1})$

# Unigram LM (1$^{st}$ order)

- n=1, 1st order Markov assumption, history (n-1): 0

$$P(w_1, \ldots, w_d) = P(w_1)P(w_2)P(w_3)\ldots$$
$$= \prod_{i=1}^{d} P(w_i)$$

unigram LM

Sampling from the unigram LM:

$i = 0$
repeat
    $i{+}{+}$
    $x_i \sim p(x)$
until $x_i = \text{STOP}$
return $\langle x_1, \ldots, x_i \rangle$



w₁  w₂  w₃  w₄  ...  w_d

P(w)

network  cat  beer  weather  city  water

# Bigram Language Model

‣ n=2, 2nd order Markov assumption, history (n-1): 1

$$P(w_1, \ldots, w_d) = P(w_1)P(w_2 \,|\, w_1)P(w_3 \,|\, w_2)\ldots$$

$$= P(w_1)\prod_{i=2}^{d} P(w_i \,|\, w_{i-1})$$

$P(\textbf{center}\,|\,\textbf{city})$     bigram LM

# Sparsity problems with n-gram LMs

Sparsity problem 1: If "its water is so w" never occurred in the corpus: prob for w (and entire sequence) is 0!

What can we do? **Smoothing** (add small count to estimates)

$$P(\mathbf{w}\,|\,\textbf{its water is so}) = \frac{C(\textbf{its water is so }\mathbf{w})}{C(\textbf{its water is so})}$$

Sparsity problem 2: If "its water is so" never occurred: prob for *any* w is is 0.

What can we do? **Backoff** (condition on lower-level n-grams)

C=count()

In general: Increasing n-gram size makes sparsity problem worse.

See for more details chapter 3 of Jurafsky & Martin's textbook.

# Smoothing
## Intuition

- Estimating statistics from sparse data

- Smoothing relocates probability mass to make generalization better



P(w | denied the)
  3 allegations
  2 reports
  1 claims
  1 request

  7 total

Smoothing

P(w | denied the)
  2.5 allegations
  1.5 reports
  0.5 claims
  0.5 request
  2 other

  7 total

Adapted from Yoav Artzi

# Smoothing
## Add-one Estimation

- Pretend we saw each word one more time than we did

- So, just add one to all counts

  - And don't forget to adjust normalization properly

$$p_{\text{MLE}}(x_i \,|\, x_{i-1}) = \frac{c(x_{i-1}, x_i)}{c(x_{i-1})} \longrightarrow p_{\text{Add}-1}(x_i \,|\, x_{i-1}) = \frac{c(x_{i-1}, x_i) + 1}{c(x_{i-1}) + |\mathcal{V}|}$$

- Also called Laplace smoothing

# Evaluating a LM

# Evaluation
## Testing

- How good is our LM?

- We must test the model on data it hasn't seen during learning

  - Otherwise — overfitting! 😱

- We need an evaluation metric — two options:

  - **Extrinsic**: focused on however the model will be used — for example, can it improve a transcription system?

  - **Intrinsic**: focused on the language model task — how good can the model assign probabilities to real unseen data?

- Ideally, the two correlate, but reality is more complex

# Intrinsic Evaluation

## The Shannon Game

- How well can we predict the next word?

  When I eat pizza, I wipe of the ____

  Many children are allergic to ____

  I saw a ____

| grease | 0,5 |
|--------|-----|
| sauce | 0,4 |
| dust | 0,05 |
| … | |
| mice | 0,0001 |
| … | |
| the | 1E-100 |

- **Unigrams are terrible at this game (why?)**

- A better model of text is one which assigns a higher probability to the word that **actually** occurs

# Evaluation

## Perplexity

- The best language model is the one the is best at predicting the test set $\rightarrow$ will give test sentences the highest probability

- Perplexity is the inverse probability of the test set, normalized by the number of words:

  - Given a set of test sentences $D'$ with a total of $m$ words:

$$PP(D') = p(D')^{-1/m} = (\prod_{\bar{x} \in D'} p(\bar{x}))^{-1/m}$$

  - In practice, we work in log space:

$$PP(D') = 2^{-\frac{1}{m} \sum_{\bar{x} \in D'} \log_2 p(\bar{x})}$$

- Lower perplexity is better

# Evaluation

## Perplexity of a Uniform Model

- Under a uniform distribution perplexity will be the vocabulary size

- Assume $M$ sentences consisting of $m$ random digits, $|\mathcal{V}| = 10$

- What is the perplexity of this data for a model that assigns
  $p(\,\cdot\,) = \frac{1}{10}$ to each digit

$$
\begin{aligned}
PP &= 2^{-\frac{1}{m}\sum_{i=1}^{M} \log_2(\frac{1}{10})^{|\bar{x}^{(i)}|}} \\
&= 2^{-\frac{1}{m}\sum_{i=1}^{M} |\bar{x}^{(i)}|\log_2 \frac{1}{10}} \\
&= 2^{-\log_2 \frac{1}{10}} = 2^{-\log_2 10^{-1}} = 10
\end{aligned}
$$

- Perplexity is weighted equivalent **branching factor**

# Evaluation

## Perplexities of Contemporary Models

Adapted from Yoav Artzi

# Further issues with n-gram LMs

‣ What about similar words?

   ‣ she *bought* a bicycle

   ‣ she *purchased* a bicycle

  ❌ **cannot share strength among similar words**

‣ Long-distance dependencies?

   ‣ for *programming* she yesterday purchased her own brand new *laptop*

   ‣ for *running* she yesterday purchased her brand new *sportswatch*

  ❌ **cannot handle long-distance dependencies**

# Neural LMs

# Neural Language Models

- LMs so far: count-based estimates of probabilities

  - Counts are brittle and generalize poorly, so we added smoothing

- The quantity that we are focused on estimating (e.g., for tri-gram model):

$$p(\bar{x}) = \prod_{i=1}^{n} p(x_i \mid x_{i-1}, x_{i-2})$$

- Can we use neural networks for this task? What would it give us?

# Neural Language Models
## A Very Simple Approach

- Instead of having count-based distributions, we parameterize them

$$p(x_i | x_{i-1}, x_{i-2}; \theta)$$

- How would we model this with a neural network?

  - Hint: so far, only learned about MLPs

# Early Neural LMs (1/2)

‣ <u>An early solution</u>: **a neural probabilistic language model** (Bengio et al., 2000)



feed-forward neural network, n-grams with distributed representations

# Neural Language Models
## A Very Simple Approach

- A simple FFNN-ish model

$$p(x_i = w \mid x_{i-1}, x_{i-2}; \theta) = \text{softmax}(\mathbf{y})_w$$

$$\mathbf{y} = \mathbf{U}\tanh(\mathbf{We} + \mathbf{b_1}) + \mathbf{b_2})$$

$$\mathbf{e} = [\phi(x_{i-1}); \phi(x_{i-2})]$$

where $\phi$ is an embedding function, and $\theta = (\mathbf{b}, \mathbf{d}, \mathbf{W}, \mathbf{U}, \mathbf{H}, \mathbf{C}, \phi)$

- The parameters $\theta$ are estimated by maximizing the log probability of the data

- During inference, you compute the neural network every time you need a value from the probability distribution

[Bengio et al. 2000]

# Early Neural LMs (2/2)

‣ <u>Another early solution</u>: **word2vec (**Mikolov et al., 2013) models the probability of a word given a context [See Ryan's lecture 1 yesterday for details]; is a family of methods. Learn a word embedding directly. Example: Skip-gram:



feed-forward neural network, skip-grams

# Recap: Feed-forward Neural Networks

# Recap: A Feed-Forward Neural Network (FFNN) - Graphical View



fully-connected/dense

Inputs

Hidden Layers

Outputs

- A (dense / fully-connected) feed-forward neural network
  - AKA a Multi-layer Perceptron (MLP)

Adapted from E. Bassignana

# Hidden Layer Connections - Notation

$$f_i : \mathbb{R}^7 \to \mathbb{R}^9$$

$\mathbf{h_i} \in \mathbb{R}^{D_i}$     Hidden layer with $D_i$ dimensions, e.g. $\mathbb{R}^9$

$$\mathbf{h_i} = \sigma_i(\mathbf{W_i}\mathbf{h_{i-1}} + b_i)$$

Inputs: $\mathbf{h}_{i-1}$        Outputs: $\mathbf{h}_i$

$\mathbf{W_i} \in \mathbb{R}^{D_i \times D_{i-1}}$

$\mathbf{b_i} \in \mathbb{R}^{D_i}$

$\sigma_i$  Layer's non-linear activation function

# Feed-Forward Neural Network - Functional Application View

$$\hat{y} = f_2(f_1(x))$$

# Feed-Forward Neural Network - Algebraic View

$$\hat{y} = f_2(f_1(x))$$

$$= \sigma_2(\mathbf{W_2}(\sigma_1(\mathbf{W_1x} + b_1)) + b_2)$$

# Feed-Forward Neural Network - Computational Graph View

$$\sigma_2(\mathbf{W_2}(\sigma_1(\mathbf{W_1}\mathbf{x} + b_1)) + b_2)$$

parameters

functions

σ₂

ADD

MULT    $\mathbf{b}^2$

σ₁    $\mathbf{W}^2$

ADD

MULT    $\mathbf{b}^1$

What is x?

x    $\mathbf{W}^1$

# Recap: Input layer as Embedding Layer



- Input is a word $x \in \mathbb{R}^D$ for all $x \in \mathcal{V}$
- We store these in a $|\mathcal{V}| \times D$ look up table
  - These are the model *word embeddings*
  - AKA embedding layer; word look-up table; ...

# Recap: One-hot encoding

‣ Sparse high-dimensional vector of dimension |V| (=size of vocabulary)

**Symbol** (word, char,..)

**yellow**

**one-hot vector**
(length V, one entry is 1)

`1x|V|`

# Lookup: Representing a symbol

linear projection from V->d

symbol ⟶ one-hot ⟶ word embedding

$E$

**book** ⬜⬛⬜⬜ .. ⬜⬜  **x**  🔵🔵🔵🔵  **=**  🔵🔵🔵🔵

sparse binary one-hot,
high-dimensional (V)

dense, continuous
representation
low-dimensional (d)

embedding matrix

# Computational Graph View

loss

σ₂

ADD

MULT $\quad$ $\mathbf{b}^2$

σ₁ $\quad$ $\mathbf{W}^2$

ADD

MULT $\quad$ $\mathbf{b}^1$

concat $\quad$ $\mathbf{W}^1$

lookup $\qquad$ lookup

cool $\qquad$ <3 $\qquad$ $E$

parameters

functions

Embedding layer is trained with the network
(model parameters, lookup layer)

54

# Putting it all together: A window-based neural LM



output distribution
$$\hat{y} = \mathrm{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer
$$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$$

concatenated word embeddings
$$e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$$

words / one-hot vectors
$$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$$

Adapted from Abigail See

# Window-based neural LM via FFNN

predict!

predict!

As the clock rang the students opened  the  books

discard

fixed window of n words

books

laptops

a          zoo

# Training the Neural FFNN-based LM

‣ Iteratively move the n-gram window through a very large corpus to predict the next word at each time step

‣ Cross-entropy loss (negative log-likelihood):

$$L = -logp(w_t | w_{t-1} .. w_{t-n+1})$$

‣ Note: typically very large vocabulary (softmax)

  ‣ Workaround: negative sampling (lecture 1)

**output distribution**
$\hat{y} = \text{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$

**hidden layer**
$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$

**concatenated word embeddings**
$\boldsymbol{e} = [\boldsymbol{e}^{(1)}; \boldsymbol{e}^{(2)}; \boldsymbol{e}^{(3)}; \boldsymbol{e}^{(4)}]$

**words / one-hot vectors**
$\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)}, \boldsymbol{x}^{(4)}$

*books*
*laptops*
a        zoo

$U$

$W$

*the*          *students*       *opened*       *their*
$\boldsymbol{x}^{(1)}$       $\boldsymbol{x}^{(2)}$         $\boldsymbol{x}^{(3)}$         $\boldsymbol{x}^{(4)}$

# Does a neural FFNN-LM solve these issues?

‣ Can it handle similar words?

    ‣ she *bought* a bicycle

    ‣ she *purchased* a bicycle

‣ Long-distance dependencies?

    ‣ for *programming* she yesterday purchased her own brand new *laptop*

    ‣ for *running* she yesterday purchased her brand new *sportswatch*

# Neural Language Models
## An MLP with No Markov Assumption

- Without the Markov assumption, the model is

$$p(\bar{x}) = \prod_{i=1}^{n} p(x_i \mid x_1, \ldots, x_{i-1})$$

- We need to model the parameterized distribution

$$p(x_{i+1} \mid x_1, \ldots, x_i; \theta)$$

- How can we do this with the tools we already know?

# Neural Language Models
## An MLP with No Markov Assumption

- We need to model the parameterized distribution

$$p(x_{i+1} | x_1, \ldots, x_i; \theta)$$

- We can just treat the context as a bag of words

  - Then it doesn't matter how long it is

  - A very simple example (two layer MLP)

$$\mathbf{h} = \tanh(\mathbf{W}' \tfrac{1}{i} \Sigma_{j=1}^{i} \phi(x_j) + \mathbf{b}')$$

$$p(x_{i+1} | x_1, \ldots, x_i) = \mathrm{softmax}(\mathbf{W}'' \mathbf{h} + \mathbf{b}'')$$

# Neural Language Models
## An MLP with No Markov Assumption

- Why is this a terrible idea?

  – Order matters a lot in language 🤦‍♂️

    ‣ "it was not good, it was actually quiet bad"

    ‣ "It was not bad, it was actually quiet good"

  – But it worked so well for text categorization … 😩

  – What may work for tasks that just require focusing on salient words (e.g., topic categorization), is not sufficient for language models (i.e., **next**-word prediction)

# FFNN: Any problems?

PROBLEM:
we need fixed input size!

a

big

waste

of

time

Inputs

Weigths

$X_1$

$X_2$

$X_3$

$X_n$

$W_1$

$W_2$

$W_3$

$W_n$

Weigthed
Sum

Activation
Function

Output

Σ

Y

# FFNN: What is still missing?

PROBLEM:
it does not consider the order well

it
was
not
good
,
it
was
actually
quite
bad

# Variable Length Inputs

‣ Feed-forward neural networks assume fixed-length inputs, but texts are not fixed lengths

```
the cat sits there

the sleepy cat sits there

the sleepy cat which chased the dog sits there
```

# Dealing with Variable Length Inputs

‣ **Options** (1-3 covered in lecture 1):

  (1) Truncate text length at fixed K

  (2) Average embeddings (pooling)

# Average embedding: Continuous Bag-of-words (CBOW)
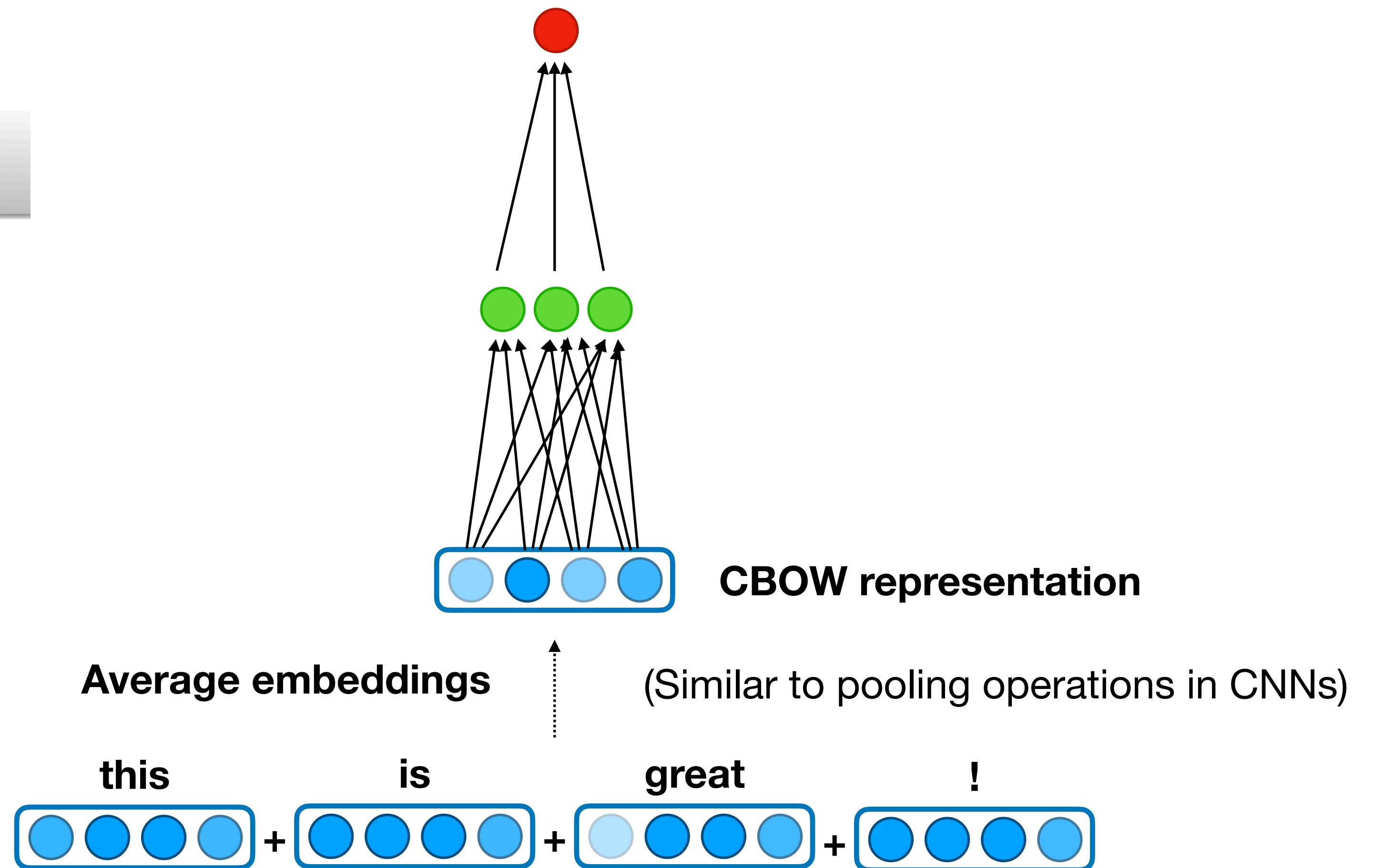
Example input document:

this is great !

bag-of-words

great
!
<3
waste
cool
boring
time
movie

CBOW representation

**Average embeddings**    (Similar to pooling operations in CNNs)

**this**        **is**        **great**        **!**

# Dealing with Variable Length Inputs

‣ More Options:

(3) **Convolutional neural network** (CNN)
illustration from last lecture by Ryan:

(4) **Recurrent neural network** (RNN)
covered next



Waibel et al. (1989) is often cited
as earliest example of a CNN

# Representations & Beyond FFNNs

**Part II**

# Today's roadmap

‣ **Part I: Fundamentals**

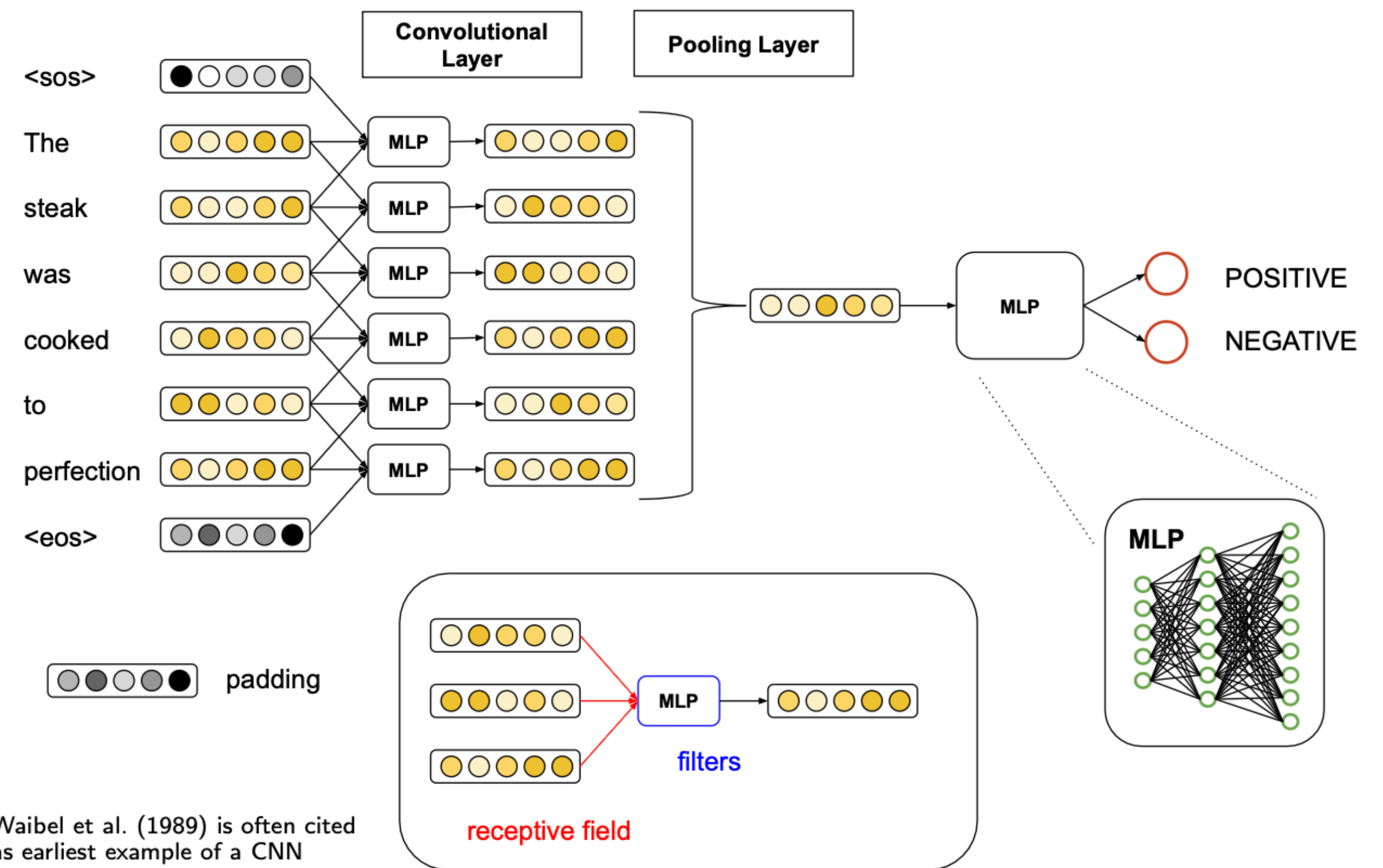   ‣ Intro, Motivation & Short History

   ‣ Language Models (n-grams, FFNN-LM, Recap: FFNN)

‣ **Part II: Representations & Beyond FFNN**

   ‣ RNNs (GRU/LSTMs), Attention

   ‣ Contextualised Representations (ELMo)

‣ **Part III: Transformer & LLMs**

   ‣ The Transformer, Masked LMs (BERT), Pre-training & Fine-tuning

   ‣ Prompting, LLMs & Caution

# RNNs

# Recurrent Neural Networks (RNNs)

‣ RNNs are a family of neural networks

‣ Very good at modelling sequential input of **variable length**

   ‣ Unlike CBOW, they **model the order** in the sequence

   ‣ Unlike vanilla CNNs, they **can deal with long-distance dependencies** (especially the gated RNN variants)

   ‣ RNNs as LMs do **not need to make the Markov** assumption

# Recurrent Neural Networks (RNNs)

**"vanilla" Neural Network**

**RNN**

Key Idea:

RNNs: the hidden state from the time step before provides a short of "memory" (state)

$$\hat{\mathbf{y}} = \mathbf{V}\mathbf{h} + \mathbf{b}$$

$$\mathbf{h} = g(\mathbf{W}\mathbf{x} + \mathbf{b})$$

Key difference: recurrence link

# A closer look: inside an RNN layer

‣ We process a sequence **x** by applying a recurrence formula **at every time step t**

‣ **Critically,** the does not impose a fixed-length input!

**new state**

$$\mathbf{h_t} = f_\theta(\mathbf{h_{t-1}}, \mathbf{x_t})$$
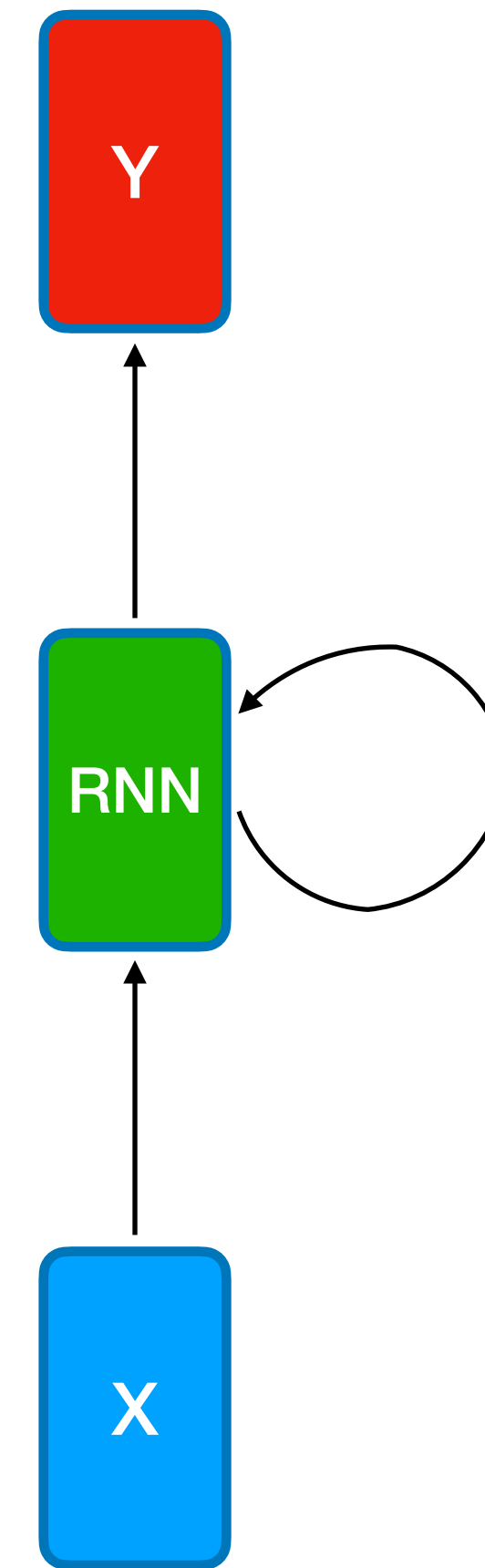
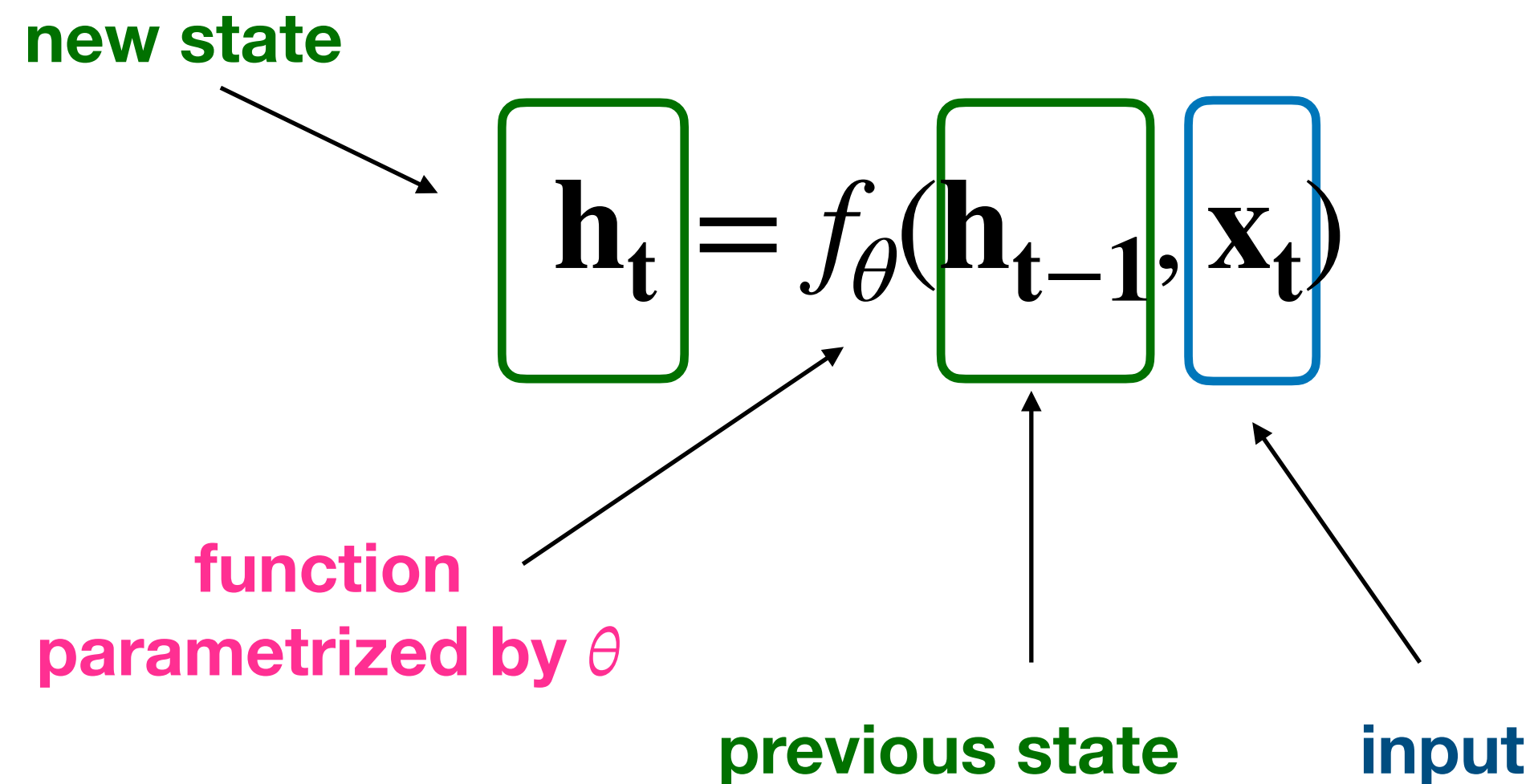**function parametrized by** $\theta$

**previous state**       **input**

Y

RNN

X

# Recurrent Neural Networks (RNNs)

**RNN**



$$\hat{\mathbf{y}}_t = \mathbf{V}\mathbf{h}_t + \mathbf{b}_2$$

$$\mathbf{h} = g(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{b}_1)$$

$$\mathbf{U}$$

= key new set of parameters to connect hidden

(Graphical illustration on right - from Jurafsky & Martin, SLP3)

# RNN - Step by Step

### A family of recurrent NN architectures

$$\mathbf{h_t} = g(\mathbf{W}\mathbf{x_t} + \mathbf{U}\mathbf{h_{t-1}} + \mathbf{b_1})$$

**Output (sequence) (optional)** $\{$   $\mathbf{y_1}$   $\mathbf{y_2}$   $\mathbf{y_3}$   $\mathbf{y_4}$

$O$

**Hidden states** $\{$   $\mathbf{h_1}$   $\mathbf{h_2}$   $\mathbf{h_3}$   $\mathbf{h_4}$

$R$

**Input sequence (any length)** $\{$   $\mathbf{x_1}$   $\mathbf{x_2}$   $\mathbf{x_3}$   $\mathbf{x_4}$

I   love   New   York

$\theta$

Core idea: **Parameter Sharing** over time

# The RNN abstraction

‣ Input sequence of vectors: $\mathbf{x_{1:n}}$

‣ $RNN(\mathbf{s_0}, \mathbf{x_{1:n}})$ consists of two functions:

   ‣ **Recurrence** function $R$ consumes input and previous state

   ‣ **Output** function $O$ maps states to outputs

‣ Technically, there is a start state: $\mathbf{s_0}$



(Graphical illustration - from Yoav Goldberg's primer, 2015)

# The RNN abstraction - More formally

$$RNN(\mathbf{s_0}, \mathbf{x_{1:n}}) = \mathbf{s_{1:n}}, \ \mathbf{y_{1:n}}$$
$$\mathbf{s_i} = R(\mathbf{s_{i-1}}, \mathbf{x_i})$$
$$\mathbf{y_i} = O(\mathbf{s_i})$$



(Graphical illustration - recursive - from Yoav Goldberg's primer, 2015)

# RNN: Unrolled over time



Figure 6: Graphical representation of an RNN (unrolled).

(Graphical illustration - recursive - from Yoav Goldberg's primer, 2015)

# Expansion at time step 4

$$\mathbf{s_4} = R(\mathbf{s_3}, \mathbf{x_4})$$

$$= R(\overbrace{R(\mathbf{s_2}, \mathbf{x_3})}^{\mathbf{s_3}}, \mathbf{x_4})$$

$$= R(R(\overbrace{R(\mathbf{s_1}, \mathbf{x_2})}^{\mathbf{s_2}}, \mathbf{x_3}), \mathbf{x_4})$$

$$= R(R(R(\overbrace{R(\mathbf{s_0}, \mathbf{x_1})}^{\mathbf{s_1}}, \mathbf{x_2}), \mathbf{x_3}), \mathbf{x_4})$$

(Graphical illustration - recursive - from Yoav Goldberg's primer, 2015)

# Training a RNN, parameter tying



**sum loss**

Parameter **tying**: the parameters are shared across time steps. Derivatives accumulated.

**Pros:** - reduce #params
- model arbitrary lengths

the **unrolled** graph is a DAG computational graph, we can backprop back

**Backpropagation through time (BPTT, Werbos, 1990).**

# Recap: FFNN's way

this is great !

One fixed input, e.g. average or sum:

$$\textbf{CBOW}(w_i, \ldots, w_n) = \sum_i^n E[w_i]$$

**representation**

**V**

**W**

Lookup & sum

this        is        great        !

# Recap: RNN's way

`this is great !`

**Time step 1:**



$h_0$

**Lookup & process
each time step**

this        is        great        !

# Recap: RNN's way

this is great !

**Time step 1:**



Lookup & process
each time step

this      is      great      !

$x_1$      $x_2$      $x_3$      $x_4$

# Recap: RNN's way

`this is great !`

**Time step 2:**

$h_0$

**V**

**U**

**W**

**Lookup & process
each time step**

this         is          great              !

$x_1$        $x_2$       $x_3$              $x_4$

# Recap: RNN's way

this is great !

**Time step 3:**

$h_0$

V

U

W

**Lookup & process each time step**

this       is       great       !

$x_1$      $x_2$      $x_3$       $x_4$

# Recap: RNN's way

this is great !

**Time step 3:**

$h_0$

V

U

W

Lookup & process
each time step

this        is        great       !

$x_1$        $x_2$        $x_3$        $x_4$

# RNN: shared parameters

`this is great !`

**Unrolled sequence:**



| $h_0$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ |

| this | is | great | ! |
| $x_1$ | $x_2$ | $x_3$ | $x_4$ |

# RNN: **h** acting as "memory"

`this is great !`

**Unrolled sequence:**

$\theta$



$h_0$     $h_1$     $h_2$     $h_3$     $h_4$

sentence or
"thought" vector

this     is     great     !

$x_1$     $x_2$     $x_3$     $x_4$

# RNN Language Model

# Training a RNN LM - Example



Illustration by Jurafsky & Martin

# What about these issues?

‣ Can it handle similar words?

    ‣ she *bought* a bicycle

    ‣ she *purchased* a bicycle

✓

‣ Long-distance dependencies?

    ‣ for *programming* she yesterday purchased her own brand new *laptop*

    ‣ for *running* she yesterday purchased her brand new *sportswatch*

**However, in practice the vanilla RNN has some trouble.. more soon**

✓

# RNNs - Interim summary

‣ **LM:** a model that predicts the next word

‣ **RNN:** a family of neural networks

  ‣ to model sequential input of any length

  ‣ can optionally produce an output at each time step t (function O)

‣ RNNs are great as LMs. They can be used for much more:

# Four Common Usage Patterns of RNNs

# An RNN as acceptor

- ‣ Use **average of states** to **predict** y

- ‣ Use **last state** only to **predict** y

**Pooling of hidden states**

# Example: An RNN as encoder

‣ Use **last state** as encoding of the information in the sequence; use as "feature" in other NN

  ‣ encode, not predict

‣ E.g. character RNN

# RNN as Transducer

- ‣ predict an output for each time step t

- ‣ E.g. Tagging (POS, NER)

total loss

sum

loss   loss   loss

Illustration adapted from Karpathy

**many to many**

# RNN as generator

‣ Conditioned generation

‣ E.g. image
caption generation,
speech synthesis



"straw"  "hat"  END

$y_t$

$W_{hh}$  $W_{oh}$

$h_t$

$CNN_{\theta_c}$  $W_{hi}$  $W_{hx}$

$x_t$

START  "straw"  "hat"

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for
Generating Image Descriptions", CVPR 2015
Figure copyright IEEE, 2015. Reproduced for educational purposes.

Illustration adapted from Karpathy

# RNN encoder-decoder (seq2seq)

**many to many**

‣ Both input and output are a sequence



Illustration adapted from Karpathy

# Deeper, better models?
# Some more concepts

# Only left to right?

**The person who hunts ducks out on the weekends**

**… person who hunts ducks out …**



Example adapted from Rao & McMahan, 2018

# Bidirectional RNNs

... `person who hunts ducks out` ...

$$\mathbf{h_i^f}$$

$$\mathbf{h_i^b}$$

$$\mathbf{h_i} = [\mathbf{h_i^f} ; \mathbf{h_i^b}]$$

# Stacked RNNs

‣ Multiple layers of RNNs, e.g., bi-RNNs

# Subword representations: Characters



PROPN    VERB    NOUN

$\vec{w}$    $\vec{c}$

Juli    loves    cats

<w>
c
a
t
s
</w>

*able (98% adj
in WSJ)

bi* (85% noun
in Danish)

(Plank et al., 2016 for POS;
Ling et al., 2015 for NER)

# Residual connections

‣ Training deep neural networks is difficult

‣ **Solution**: add direct "skip" connections (ResNet, residual connections), proposed by He et al., (2015)

   ‣ i.e. add F(x) + x, instead of F(x)

   ‣ allows for training deeper models (in fact, is used in transformers as we will see)



Figure 2. Residual learning: a building block.

# Gated RNN architectures

# A note on terminology

‣ RNN = "vanilla" RNN

‣ RNN flavors (=gated RNNs):

    ‣ GRU        and LSTMs

    ‣ Why? Problem of RNNs: Vanishing gradients!

# Vanishing Gradient

‣ Gradients decrease as they are pushed back

$$\frac{dl}{d_{h_0}} = \text{tiny} \qquad \frac{dl}{d_{h_1}} = \text{small} \qquad \frac{dl}{d_{h_2}} = \text{med.} \qquad \frac{dl}{d_{h_3}} = \text{large}$$



(Illustration by Graham Neubig)

# Gated RNN architectures: RNN flavors with a separate memory

# Vanilla RNN unit

$\hat{\mathbf{y}}_t$

softmax

$\mathbf{h}_{t-1}$

tanh

$\mathbf{h}_t$

$\mathbf{x}_t$

At each time step, the hidden state is updated:

$$\mathbf{h_t} = g(\mathbf{Wx_t} + \mathbf{Uh_{t-1}} + \mathbf{b})$$

in a vanilla RNN
the hidden state is
constantly being
**rewritten**

# GRU (Gated recurrent Unit) - simplified

‣ Cho et al. (2014) - key idea: dynamic memory update **c** *(h=c)*

‣ at every step t, consider overwriting candidate memory

candidate for overwriting cell

$$\tilde{\mathbf{c}} = tanh(\mathbf{U_c}\mathbf{x_t} + \mathbf{W_c}\mathbf{c_{t-1}} + \mathbf{b_c})$$

$$\gamma_{\mathbf{U}} = \sigma(\mathbf{U_U}\mathbf{x_t} + \mathbf{W_U}\mathbf{c_{t-1}} + \mathbf{b_U})$$

"update" gate

**sigmoid gate:** values between 0 and 1

**"choose which bits to update"**

$$\mathbf{c_t} = \gamma_{\mathbf{U}} \odot \tilde{\mathbf{c}} + (1 - \gamma_{\mathbf{U}}) \odot \mathbf{c_{t-1}}$$

update if gamma_U > 0

element-wise multiplication

use previous state

$\hat{\mathbf{y}}_{\mathbf{t}}$

softmax

$\mathbf{c_{t-1}}$    $\mathbf{c_t}$

$\tilde{\mathbf{c}}$   $\gamma_{\mathbf{U}}$

tanh   σ

$\mathbf{x_t}$

# GRU (Gated recurrent Unit) - full

‣ GRU: the full GRU has two gates:

**Update gate:** controls what parts of the hidden state are updated vs preserved

**Reset gate:** controls what parts of the previous hidden state are used to compute new content

$$\tilde{\mathbf{c}} = tanh(\mathbf{U_c}\mathbf{x_t} + \mathbf{W_c}(\gamma_{\mathbf{R}} \odot \mathbf{c_{t-1}}) + \mathbf{b_c})$$

$$\gamma_{\mathbf{U}} = \sigma(\mathbf{U_U}\mathbf{x_t} + \mathbf{W_U}\mathbf{c_{t-1}} + \mathbf{b_U})$$

$$\gamma_{\mathbf{R}} = \sigma(\mathbf{U_R}\mathbf{x_t} + \mathbf{W_R}\mathbf{c_{t-1}} + \mathbf{b_R})$$

$$\mathbf{c_t} = \gamma_{\mathbf{U}} \odot \tilde{\mathbf{c}} + (1 - \gamma_{\mathbf{U}}) \odot \mathbf{c_{t-1}}$$

**all vectors of same size**

**How does this help the vanishing gradient problem?**
GRUs make it easier to retain info long-term (e.g. by not updating bits)

Slide inspired by Abigail See

# LSTM (Long-Short Term Memory)

‣ Introduced by Hochreiter & Schmidhuber 1997

‣ Separate memory cell c and hidden state h

‣ Three gates:

  ‣ **forget gate**: controls what is kept and forgotten from previous cell state

  ‣ **input gate:** controls what part of the new cell content are written to the cell

  ‣ **output gate:** controls what part of the new cell content are written to the hidden state

# LSTM (Long-Short Term Memory)



Write some new cell content

Forget some cell content

Output some cell content to the hidden state

Compute the forget gate

Compute the input gate

Compute the new cell content

Compute the output gate

Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy

113

# LSTM (Long-Short Term Memory)

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep $t$:

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Sigmoid function:** all gate values are between 0 and 1

**Input gate:** controls what parts of the new cell content are written to cell

**Output gate:** controls what parts of cell are output to hidden state

$$f^{(t)} = \sigma\left(W_f h^{(t-1)} + U_f x^{(t)} + b_f\right)$$

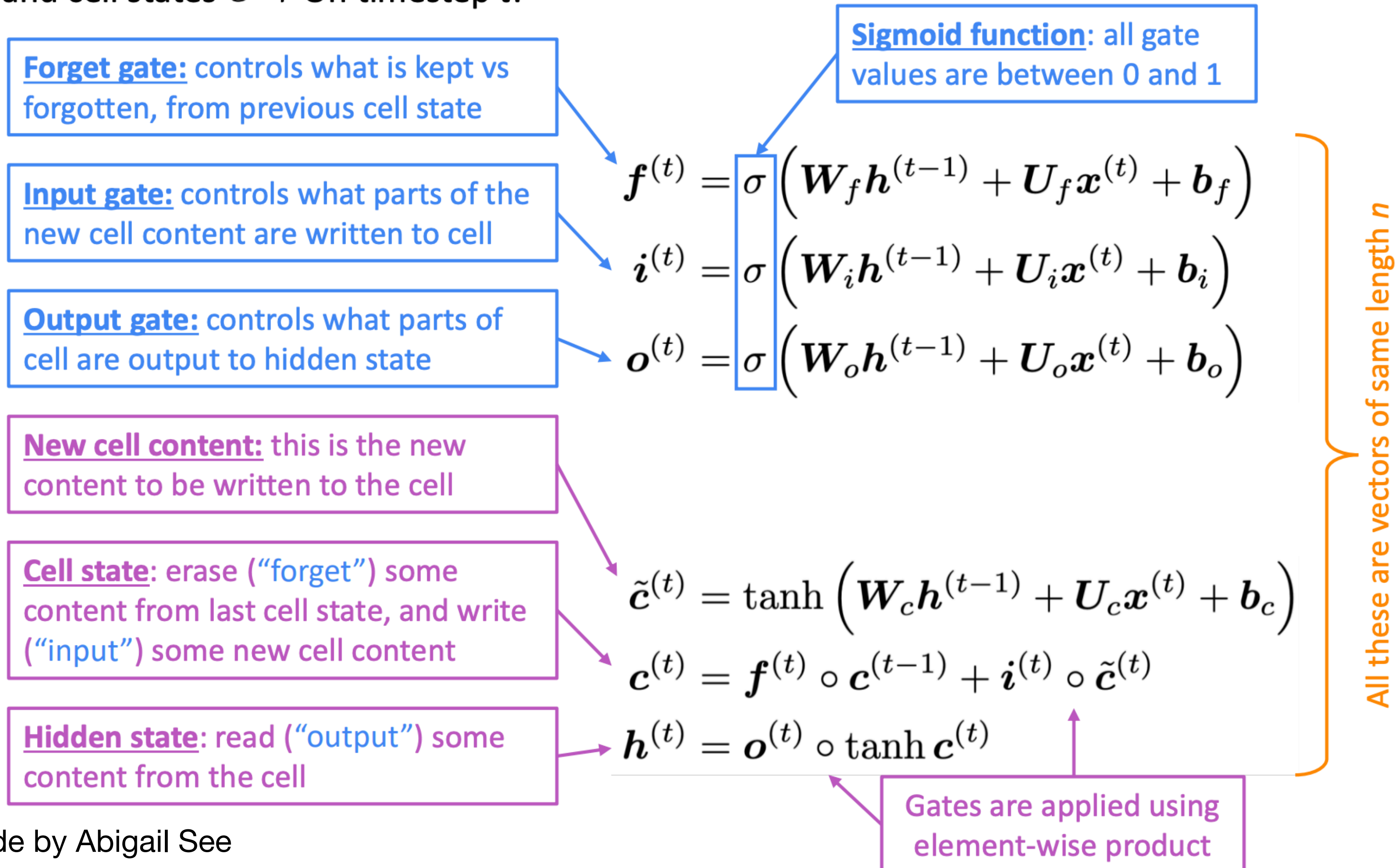$$i^{(t)} = \sigma\left(W_i h^{(t-1)} + U_i x^{(t)} + b_i\right)$$

$$o^{(t)} = \sigma\left(W_o h^{(t-1)} + U_o x^{(t)} + b_o\right)$$

All these are vectors of same length $n$

**New cell content:** this is the new content to be written to the cell

**Cell state:** erase ("forget") some content from last cell state, and write ("input") some new cell content

**Hidden state:** read ("output") some content from the cell

$$\tilde{c}^{(t)} = \tanh\left(W_c h^{(t-1)} + U_c x^{(t)} + b_c\right)$$

$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

Gates are applied using element-wise product

Slide by Abigail See

114

# GRU vs LSTM

‣ GRU is more efficient to learn (fewer parameters)

‣ Which is better?

   ‣ No conclusive evidence that one is always superior to the other

‣ LSTM is typically a good starting choice

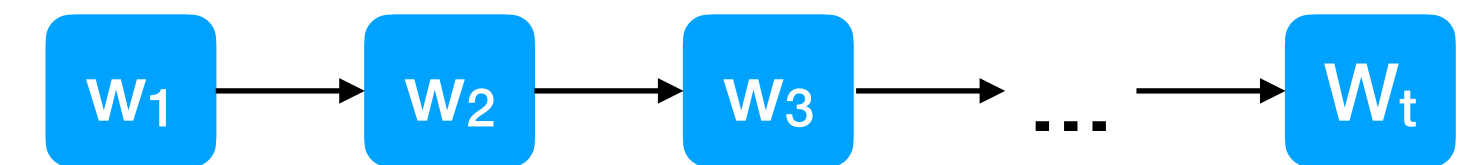‣ Suggestion: switch to GRU if you want a more efficient model

# Interim summary

‣ RNNs:

  ‣ Two fancy variants: **LSTM and GRU**
    to address the vanishing gradient problem

  ‣ Other concepts:
    Bidirectionality, Stacking, Residual connections

‣ Next concepts to cover:

  ‣ beyond static word embeddings

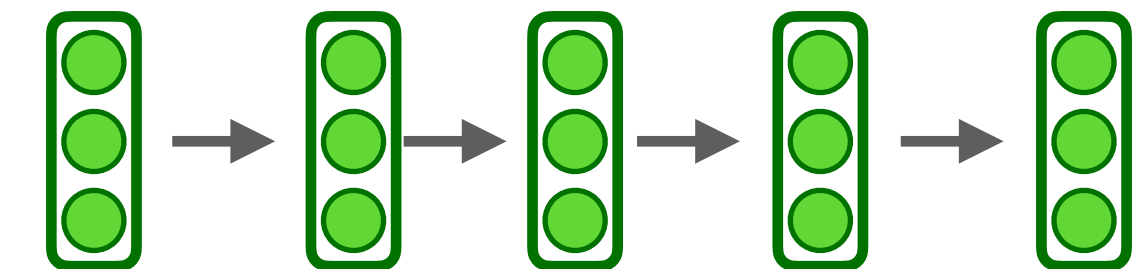  ‣ gluing it all together: attention & contextualised representations

# Today's roadmap

‣ **Part I: Fundamentals**

 ‣ Intro, Motivation & Short History

 ‣ Language Models (n-grams, FFNN-LM, Recap: FFNN)

‣ **Part II: Representations & Beyond FFNN**

 ‣ RNNs (GRU/LSTMs), Attention

 ‣ Contextualised Representations (ELMo)

‣ **Part III: Transformer & LLMs**

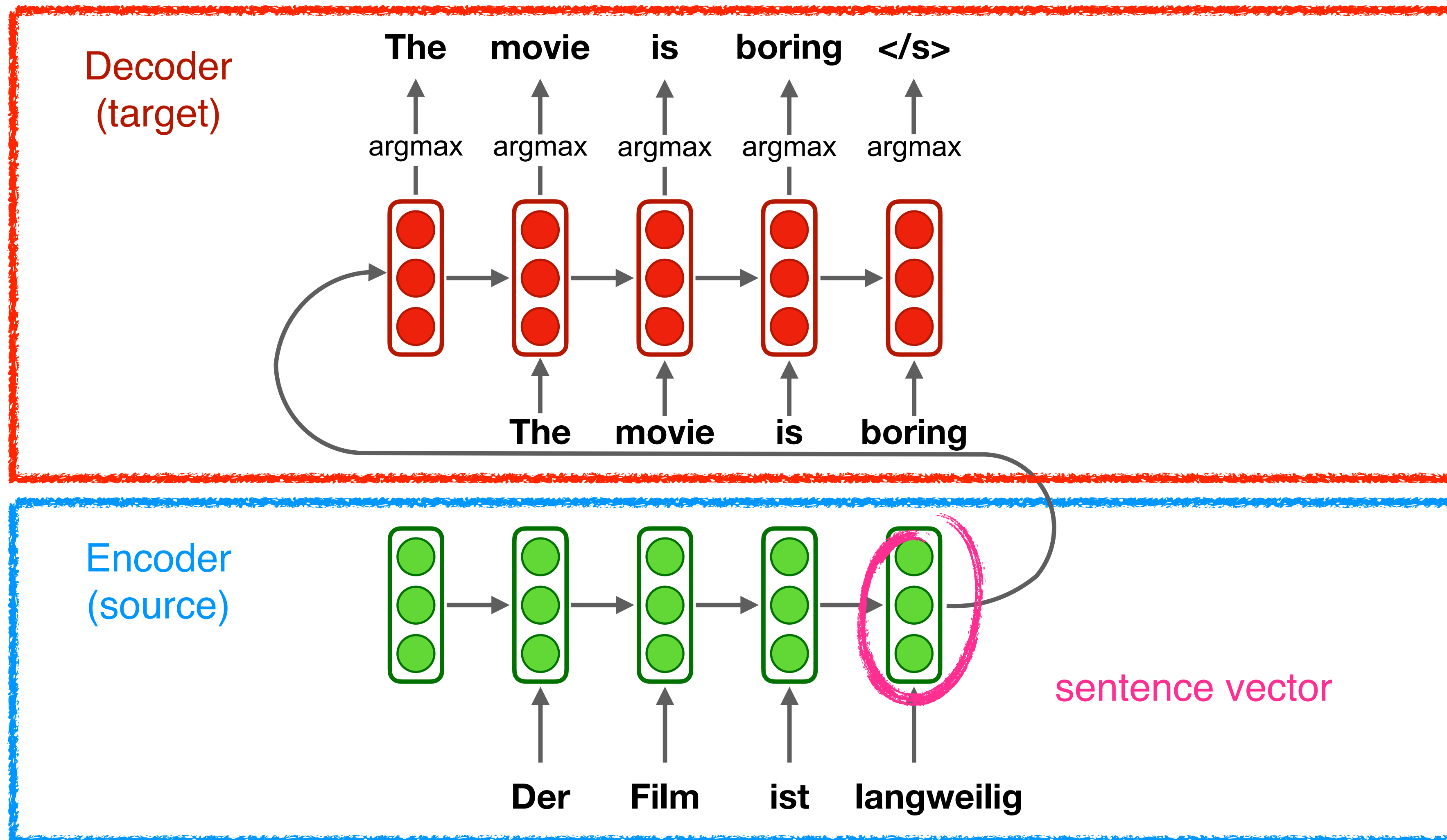 ‣ The Transformer, Masked LMs (BERT), Pre-training & Fine-tuning

 ‣ Prompting, LLMs & Caution

# Attention? Attention!

# Motivation: Encoder-decoder model for Machine Translation (MT)

**Decoder (target)**

The  movie  is  boring  </s>

argmax  argmax  argmax  argmax  argmax

The  movie  is  boring

**Encoder (source)**

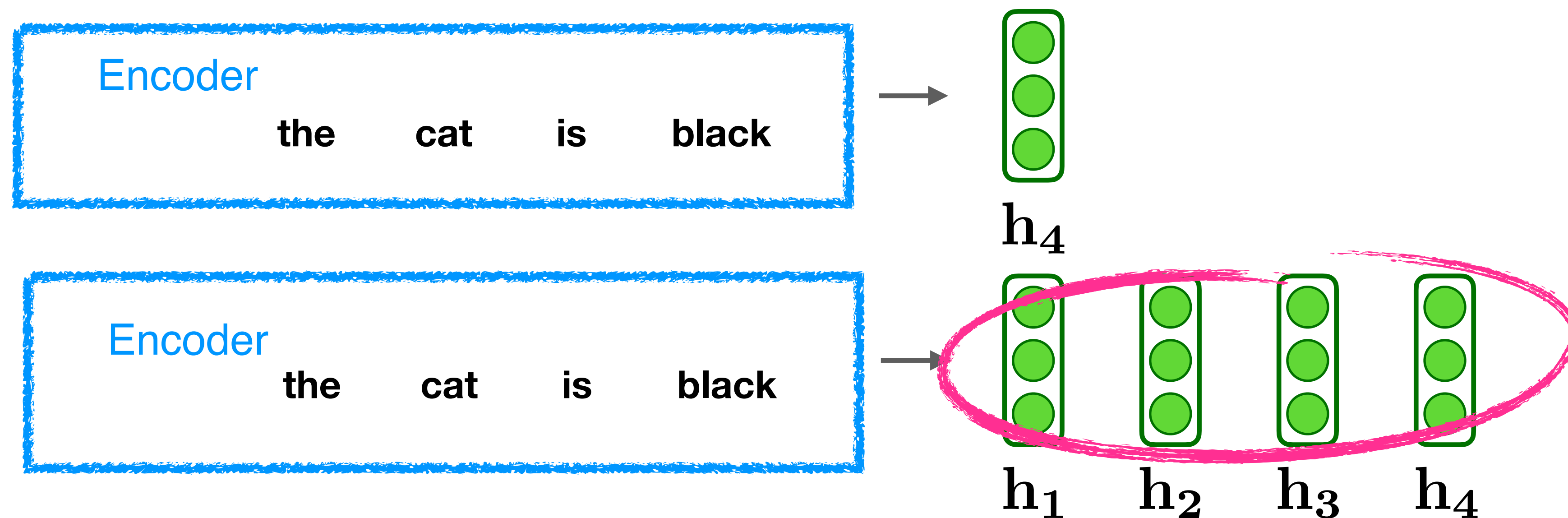Der  Film  ist  langweilig

sentence vector

a single $&!*ing vector!

# But: we're cramming it all into..

‣ The encoder compresses the sentence into a single fixed-size vector. This representation is expected to be a good summary of the entire sentence.

‣ Disadvantage: incapability of remembering longer sequences.

‣ **"You can't cram the meaning of a of a whole %&!$ing sentence into a single $&!\*ing vector!"** — Ray Mooney

# Beyond a single static "crammed" vector

‣ What if we use several vectors, based on the length of the input sequence?

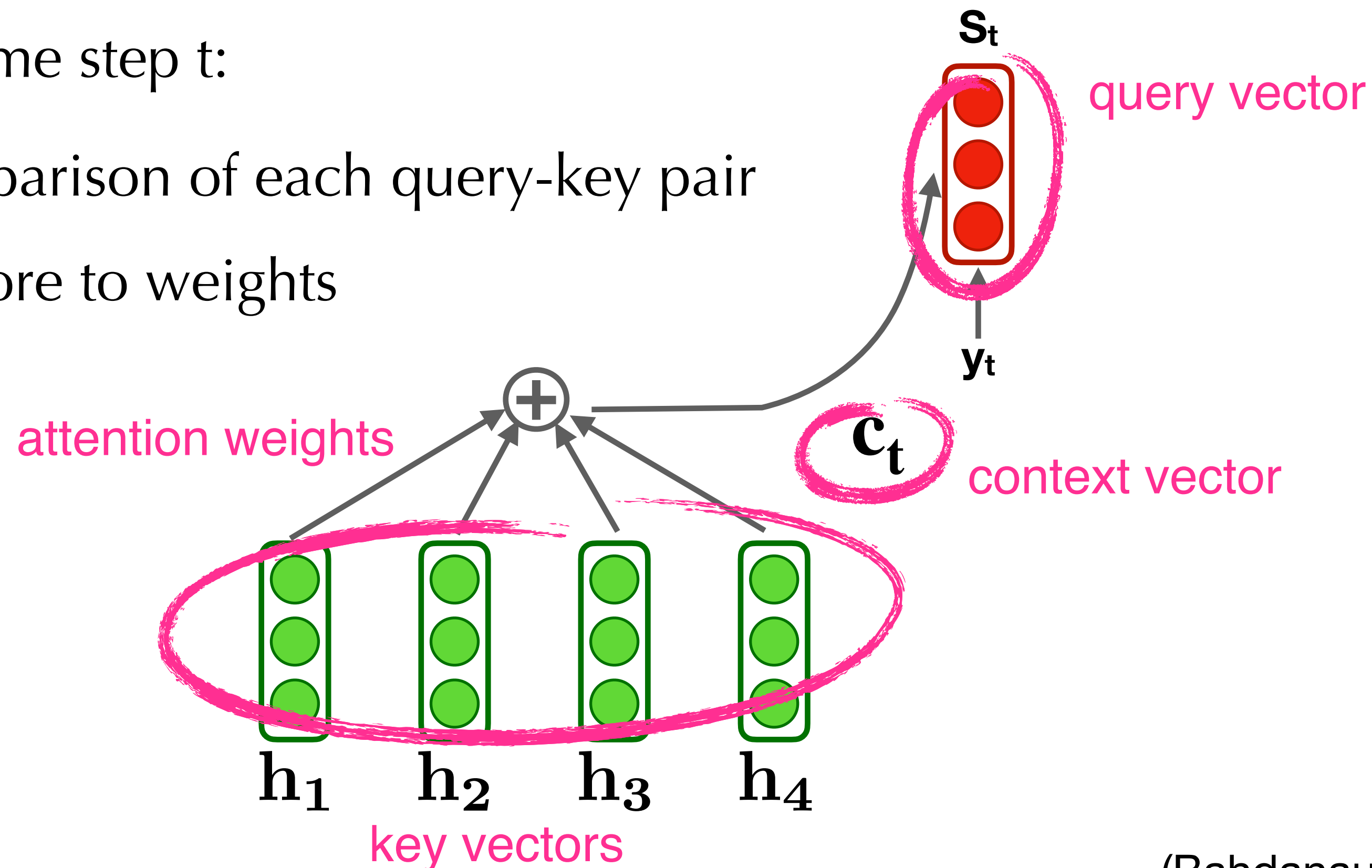‣ Idea: when we generate the next word in MT, perhaps we can learn to **attend** to the **relevant** source words



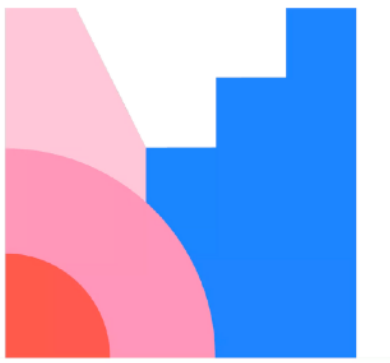**Dynamically look at all encoder hidden states**

# Attention: Core Idea

‣ When decoding, compare query to key vectors and perform a linear combination of the encoded input vectors, weighted by "attention weights"

‣ Illustration at time step t:

‣ Pairwise comparison of each query-key pair

‣ Normalise score to weights

‣ Get new $c_t$



$s_t$
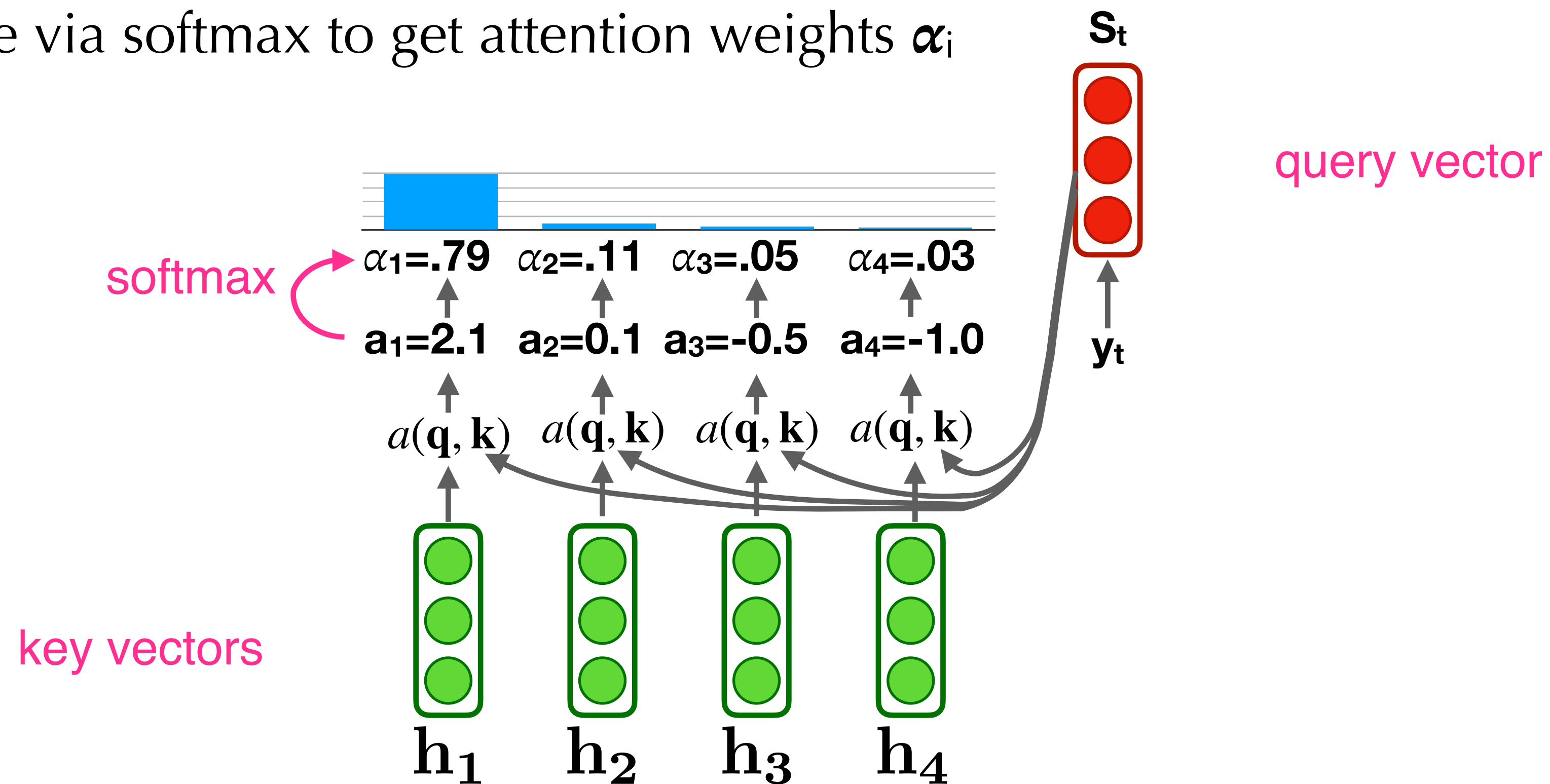
query vector

$y_t$

attention weights

$c_t$

context vector

$h_1$   $h_2$   $h_3$   $h_4$

key vectors

(Bahdanau et al., 2015)
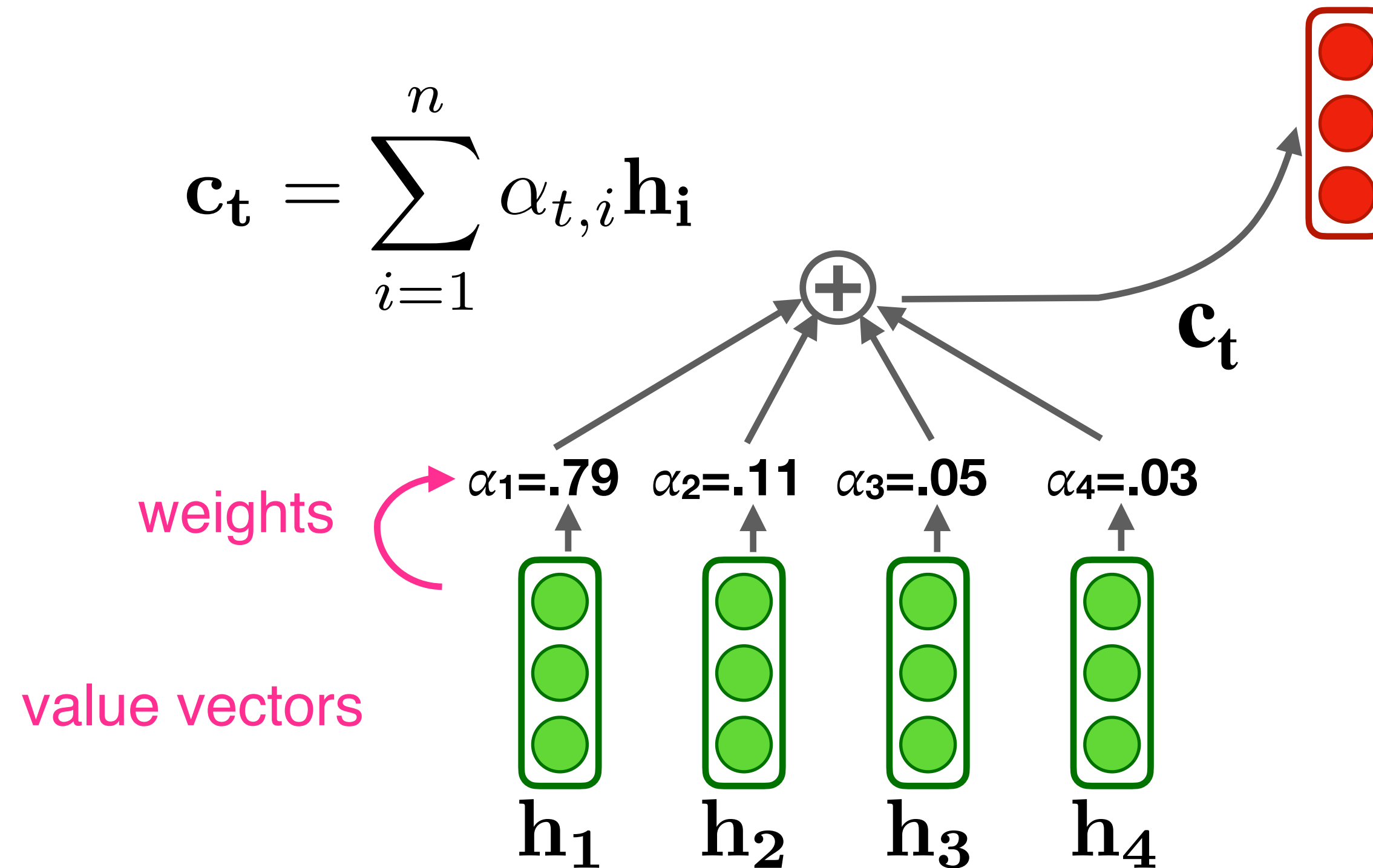
# Calculating attention (1/2): Attention weights $\alpha$

1. For each query-key pair (**q**,**k**), calculate an **attention score $a_i$** by attention function $a$

2. Normalize via softmax to get attention weights $\alpha_i$



softmax

$\alpha_1$=.79   $\alpha_2$=.11   $\alpha_3$=.05   $\alpha_4$=.03

$a_1$=2.1   $a_2$=0.1   $a_3$=-0.5   $a_4$=-1.0

$a(\mathbf{q},\mathbf{k})$   $a(\mathbf{q},\mathbf{k})$   $a(\mathbf{q},\mathbf{k})$   $a(\mathbf{q},\mathbf{k})$

**s_t**

query vector

**y_t**

key vectors

**h₁**   **h₂**   **h₃**   **h₄**

(Bahdanau et al., 2015)

# Calculating attention (1/2): Attention weights $\alpha$

3. Combine together **value vectors** via attention-weighted sum to get $\mathbf{c_t}$
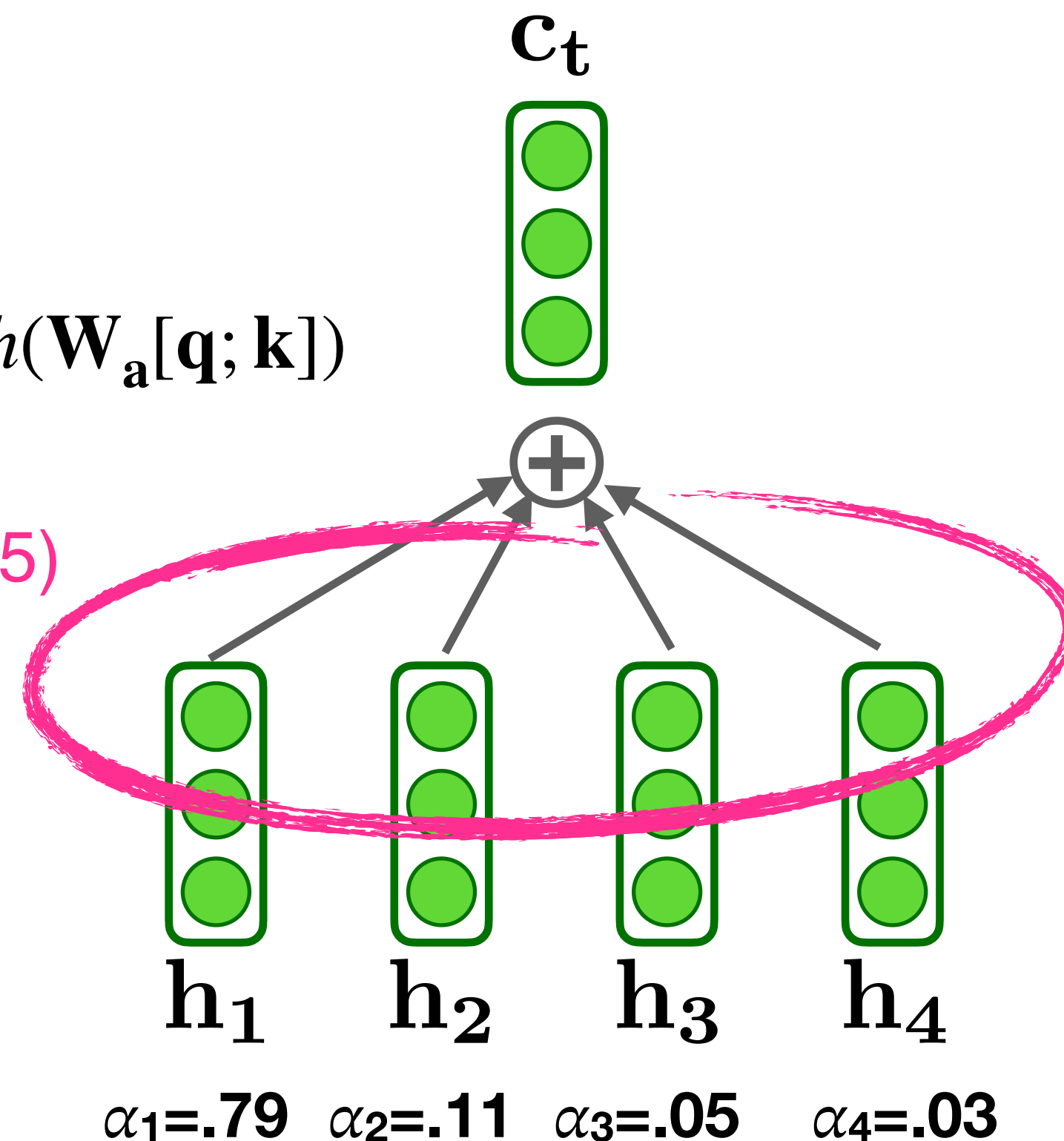
$$\mathbf{c_t} = \sum_{i=1}^{n} \alpha_{t,i} \mathbf{h_i}$$

$\mathbf{c_t}$

weights

$\alpha_1 = .79$  $\alpha_2 = .11$  $\alpha_3 = .05$  $\alpha_4 = .03$

value vectors

$\mathbf{h_1}$  $\mathbf{h_2}$  $\mathbf{h_3}$  $\mathbf{h_4}$

(Bahdanau et al., 2015)

# Attention with a FFNN/MLP (Bahdanau, 2015)

$$\mathbf{c_t} = \sum_{i=1}^{n} \alpha_{t,i} \mathbf{h_i}$$

$$a(\mathbf{q}, \mathbf{k}) = \mathbf{v_a^T} tanh(\mathbf{W_a}[\mathbf{q}; \mathbf{k}])$$

FFNN!(Bahdanau, 2015)

$\mathbf{c_t}$

$\mathbf{h_1}$    $\mathbf{h_2}$    $\mathbf{h_3}$    $\mathbf{h_4}$

$\alpha_1$=.79   $\alpha_2$=.11   $\alpha_3$=.05   $\alpha_4$=.03

1. For each query-key pair, calculate score $a_i$

2. Normalize via softmax to get weights

3. Combine together value vectors via weighted sum to get $\mathbf{c_t}$

4. Use $\mathbf{c_t}$ in your model

**Different attention functions *a*
(e.g., Luong et al., 2015)**

# Different attention functions *a()*

‣ **Dot product** (Luong et al., 2015)
   **-** requires same size; but has no parameters
   $$a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{k}$$
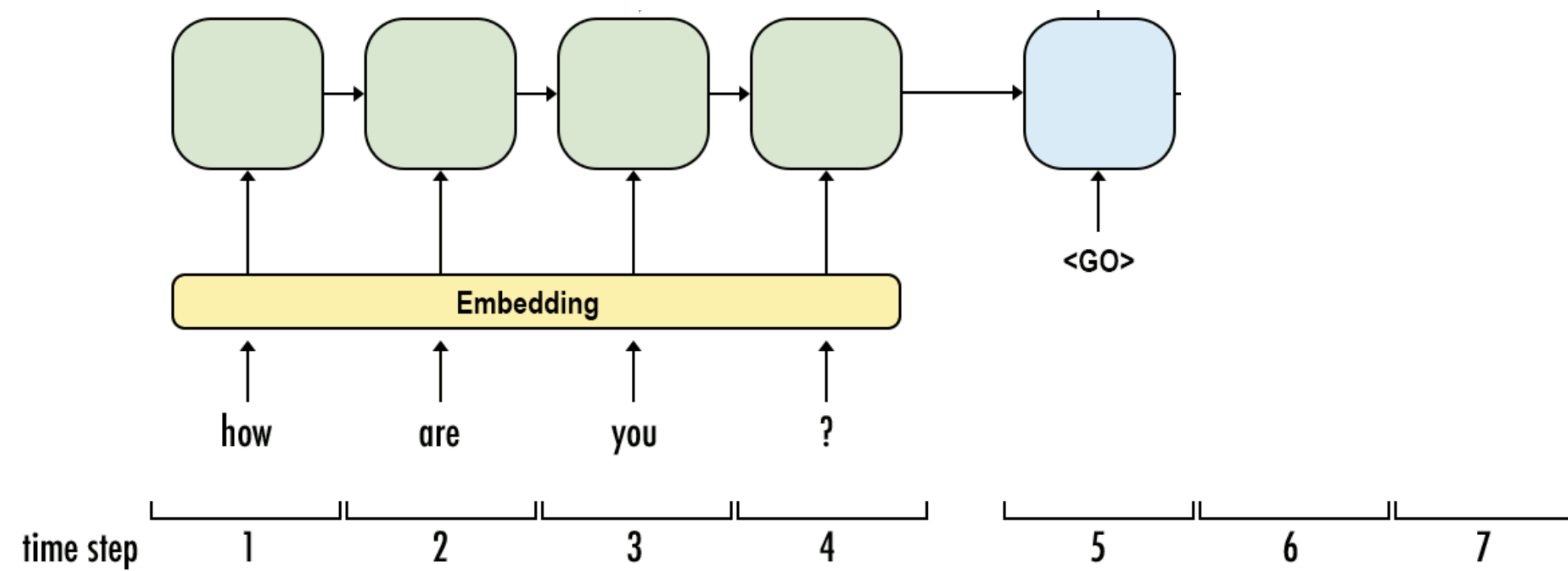
‣ **Bilinear** (Luong et al., 2015)

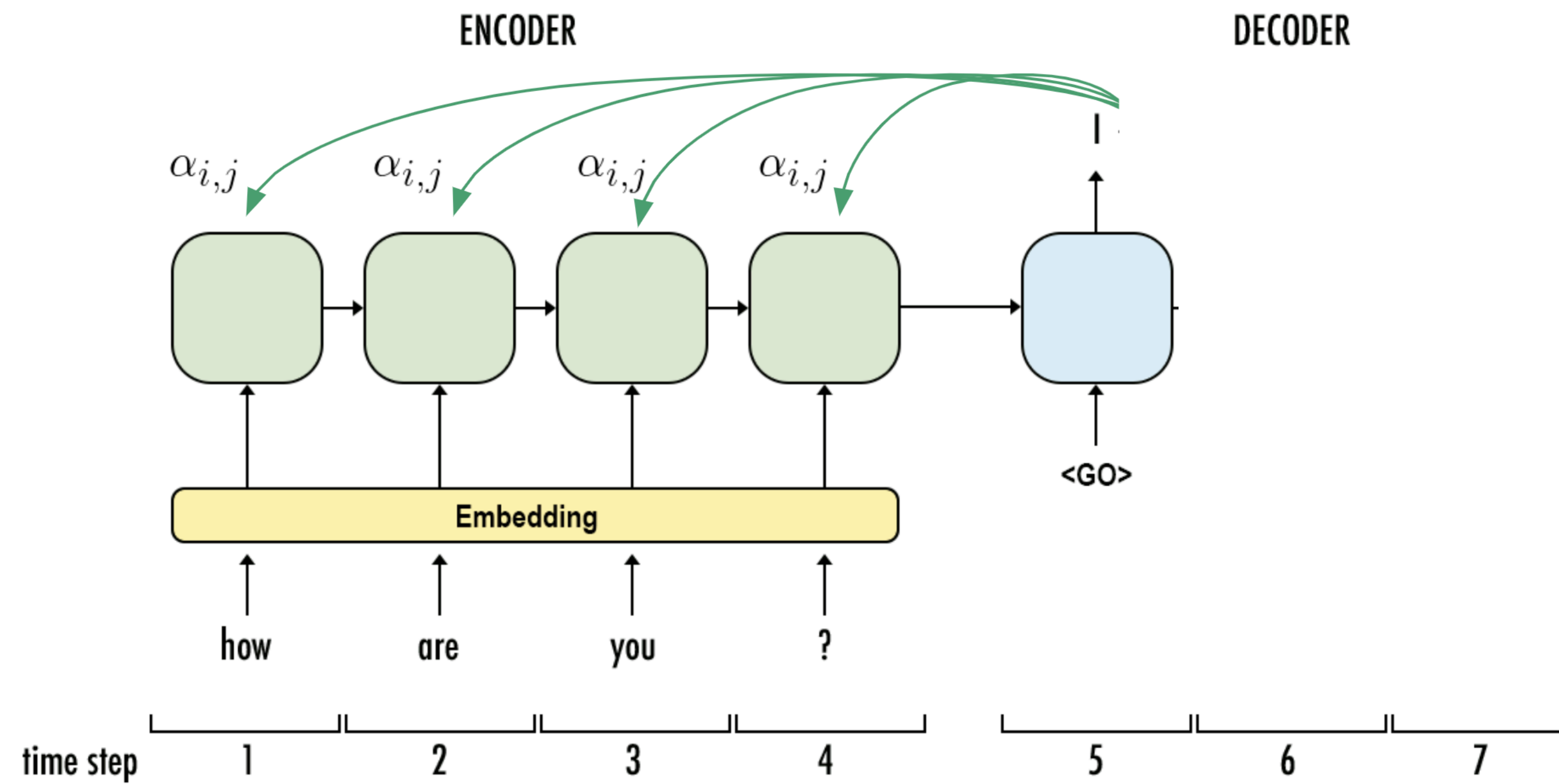   $$a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T W \mathbf{k}$$

‣ **Scaled dot product** (Vaswani et al., 2017)
   - fixes problem of dot product that scale of dot product increases as dimensions get larger

   $$a(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^T \mathbf{k}}{\sqrt{\mathbf{n}}}$$

https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html
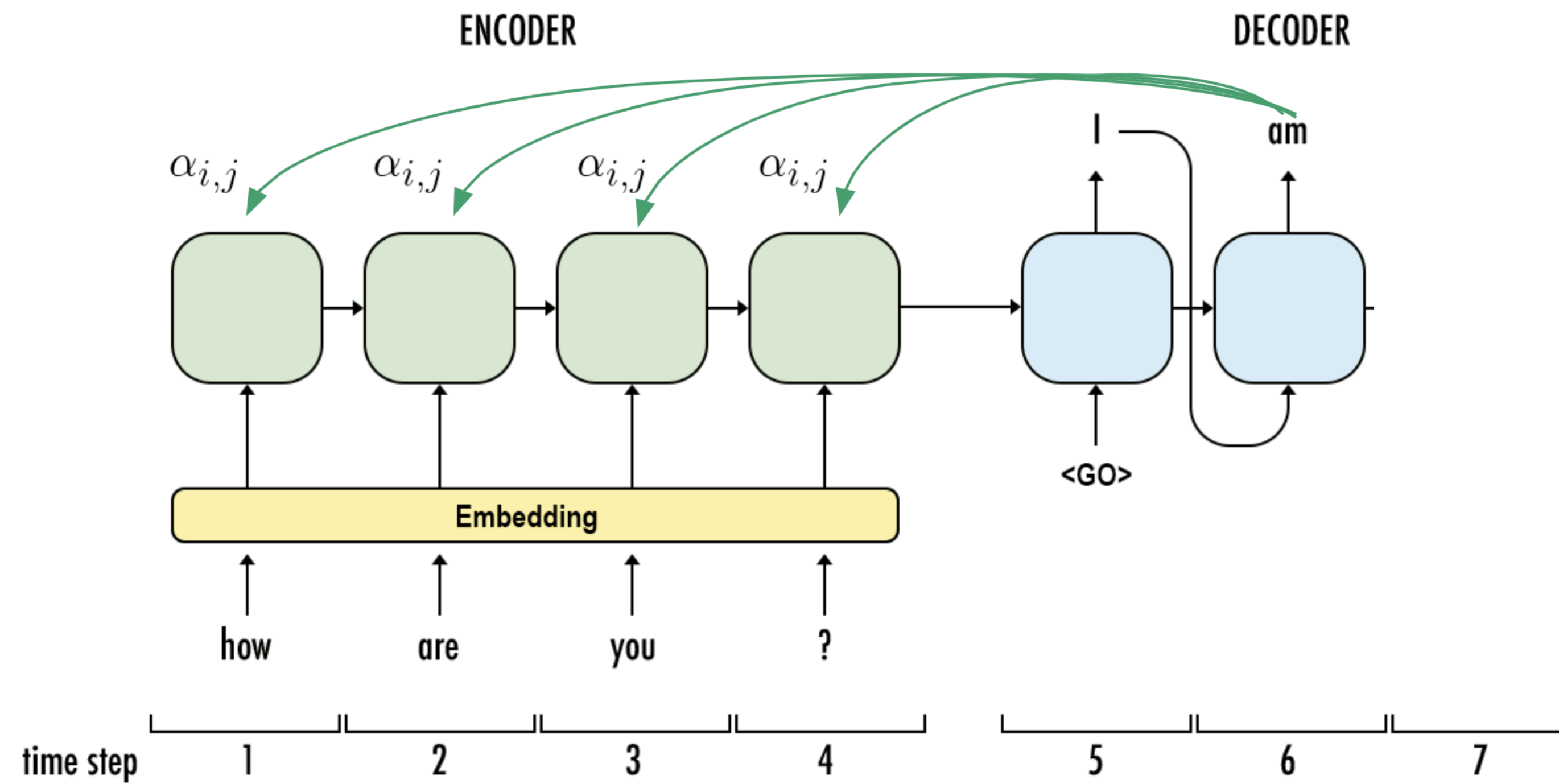
# Recurrent Neural Network

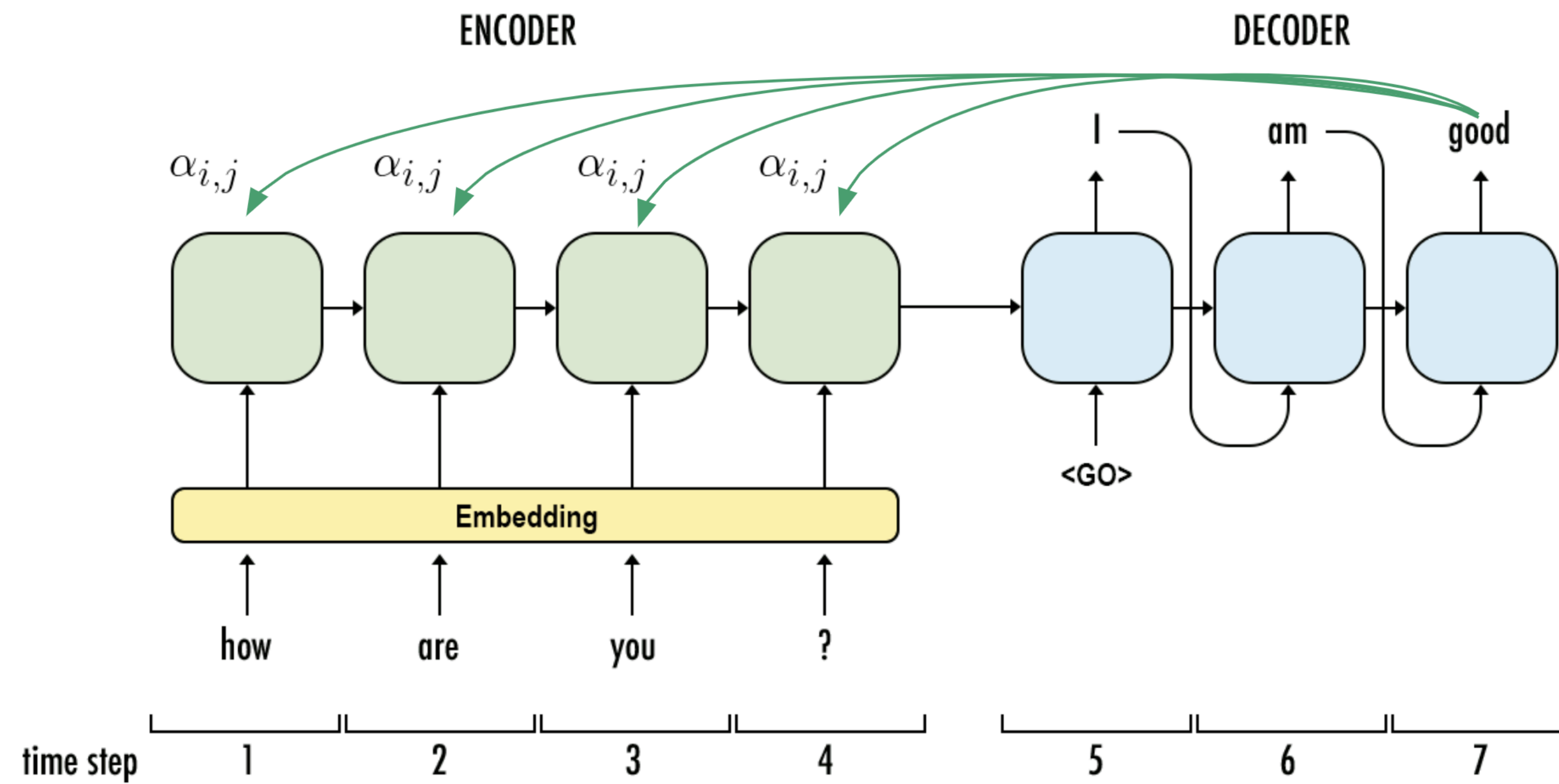# Recurrent Neural Network
## With Attention

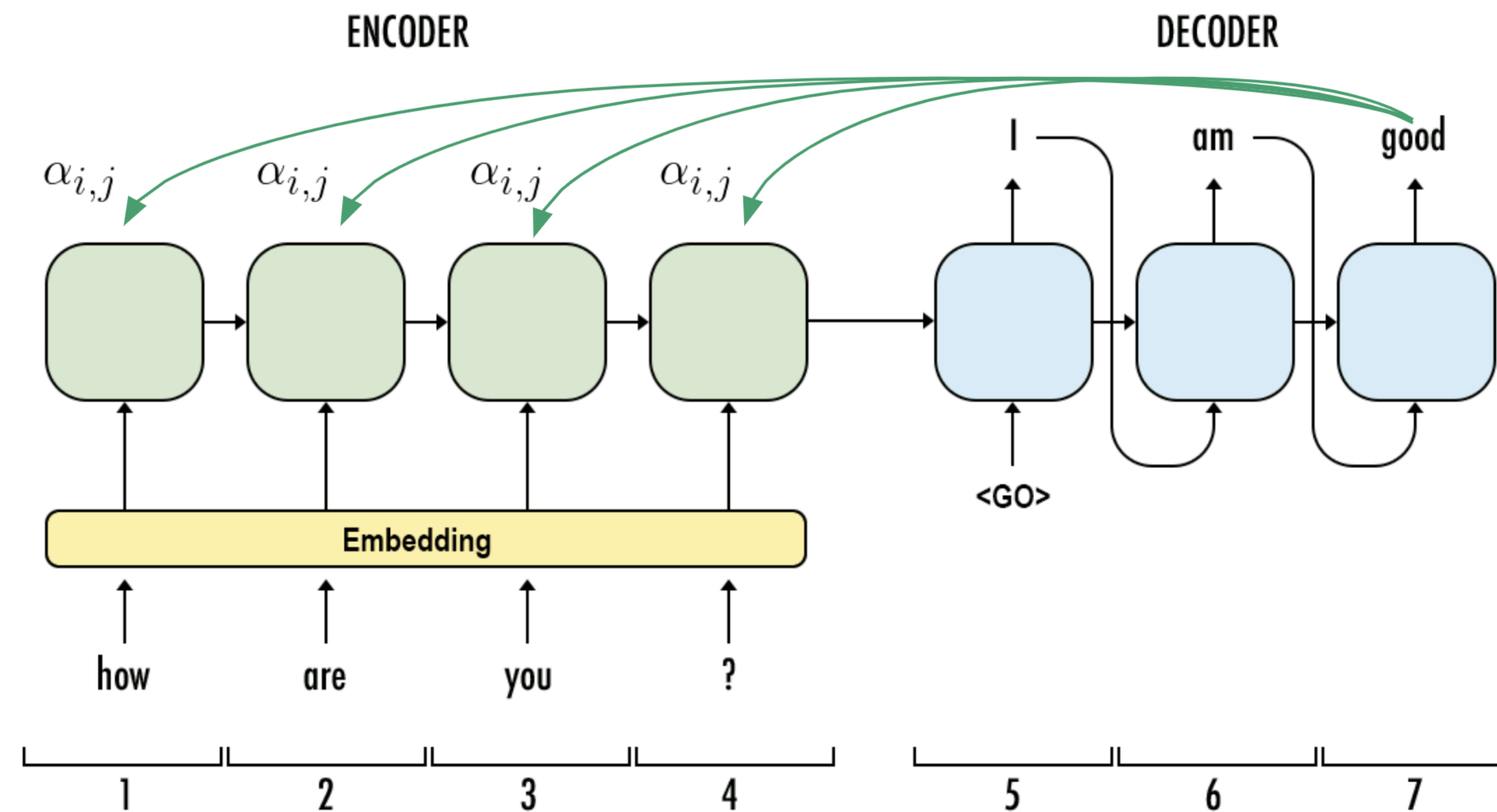# Recurrent Neural Network
## With Attention
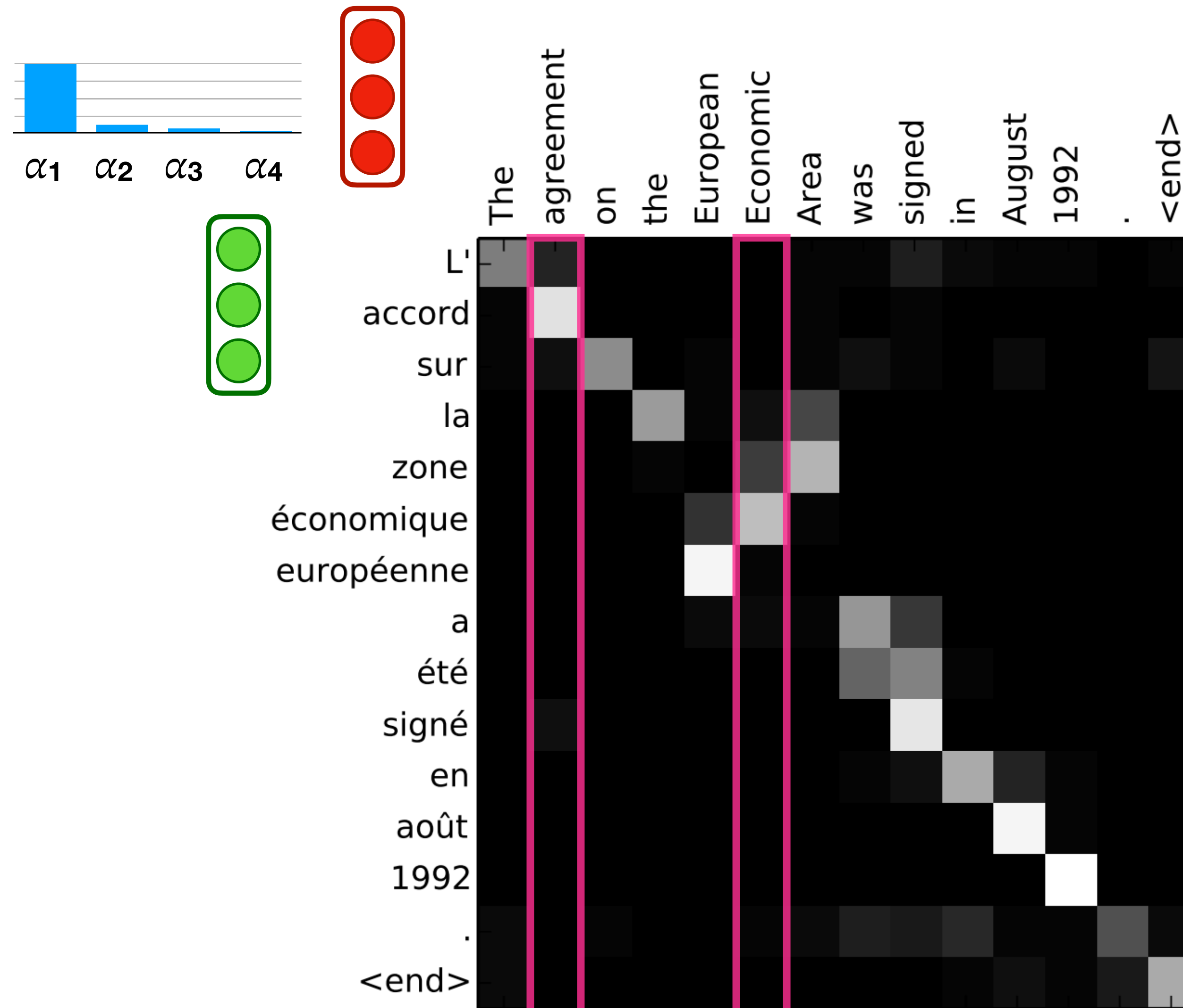
# Recurrent Neural Network
## With Attention

# Recurrent Neural Network
## With Attention

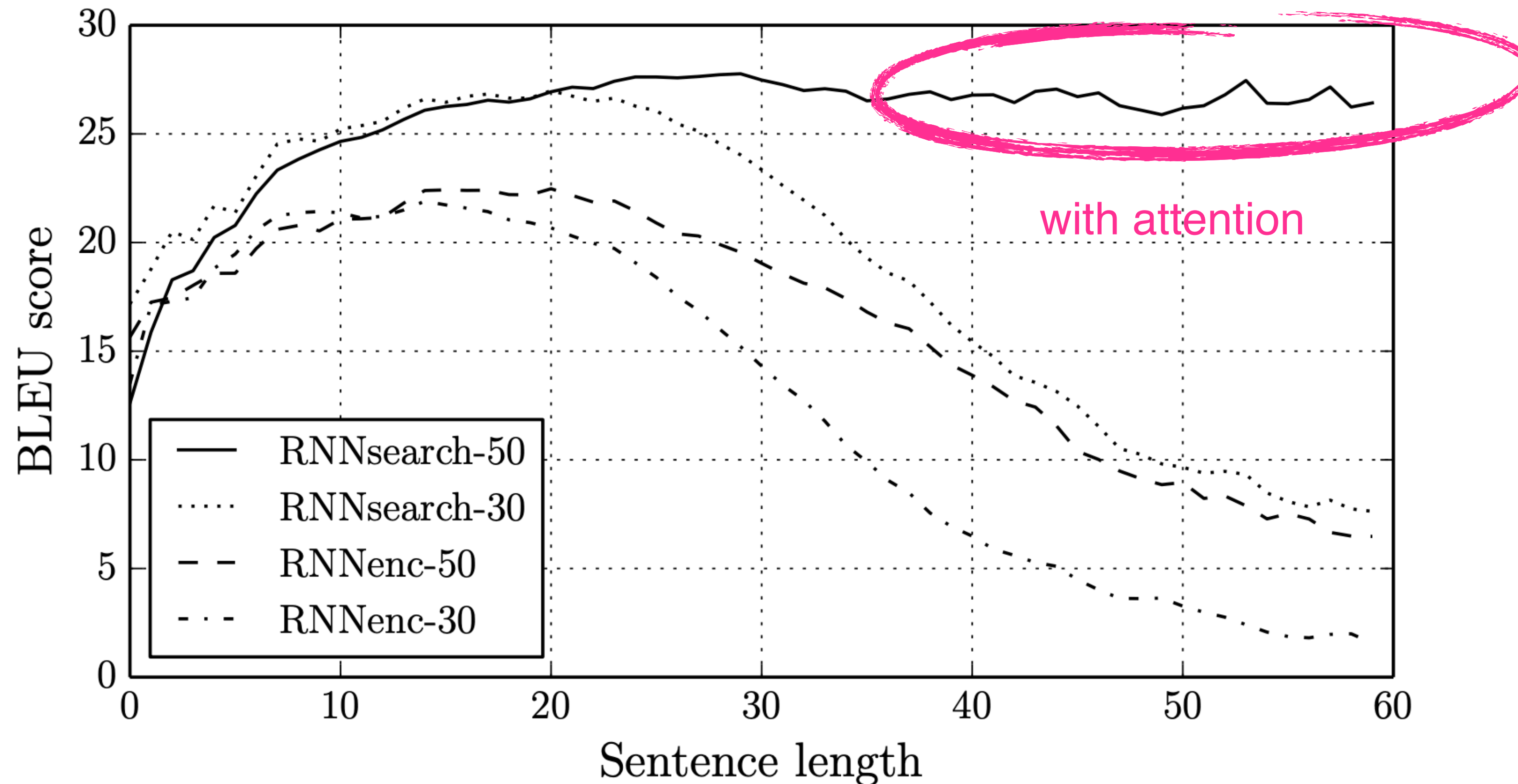Encoder-Decoder Attention (aka Cross Attention)

# A graphical example (Bahdanau, 2015)

# Enc-dec performance deteriorates rapidly as input sentence length increases

Cho et al., (2014); Bahdanau et al. (2015)



with attention

# Self Attention

‣ Attend to sentence **itself** (Cheng, Dong, Lapata, 2016) in contrast to cross attention (in encoder-decoder)

Self Attention



https://arxiv.org/pdf/1601.06733.pdf

# Attention in image caption generation

‣ Salient parts of the image (e.g., Xu et at., 2015)



A woman is throwing a <u>frisbee</u> in a park.

# Today's roadmap

▸ **Part I: Fundamentals**

　▸ Intro, Motivation & Short History

　▸ Language Models (n-grams, FFNN-LM, Recap: FFNN)

▸ **Part II: Representations & Beyond FFNN**

　▸ RNNs (GRU/LSTMs), Attention

　▸ Contextualised Representations (ELMo)

▸ **Part III: Transformer & LLMs**

　▸ The Transformer, Masked LMs (BERT), Pre-training & Fine-tuning

　▸ Prompting, LLMs & Caution

$$w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow \ldots \rightarrow w_t$$

# Introduction to Contextualized Embeddings (ELMo)

## Welcome to the contextualised embeddings world

With thanks to Arianna Bisazza, Max Müller-Eberstein & Joris Baan for parts of the slides

# Traditional ("static") word embeddings

They ordered a Danish in Danish

# What is missing?

# No context

> I'd like to order a Danish in Danish.

Danish (noun): pastry typical to Denmark
Danish (noun): Scandinavian language spoken in Denmark

# ELMo: **E**mbeddings from **L**anguage **Mo**dels

‣ A seminal paper to start a new era of representations: **contextualised** embeddings

‣ **Key Idea:** Learn word token vectors (not type) using entire context

    ‣ Aka token embedding, which depends on surrounding context during use

‣ ELMo: Self-supervised, pre-trained, RNN-based

> Our representations differ from traditional word type embeddings in that each token is assigned a representation that is a function of the entire input sentence.

> [...] ELMo representations are deep [...]

– Peters et al. (2018)

# From characters to token embeddings

# Embeddings from Language Models

**ELMo**<sub>ducks</sub> $\boxed{\circ\circ\circ\circ}$ = ( λ₁* $\boxed{\circ\circ\circ\circ}$ )+ ( λ₂* $\boxed{\circ\circ\circ\circ}$ )+( λ₃* $\boxed{\circ\circ\circ\circ}$ )



... **person**    **who**    **ducks**    **out**    **on**    ...

# ELMo - Details

‣ ELMo: every token is assigned a representation that is a function of the entire input sentence (L=#stacked layers)

$$
\begin{aligned}
R_k &= \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \ldots, L\} \\
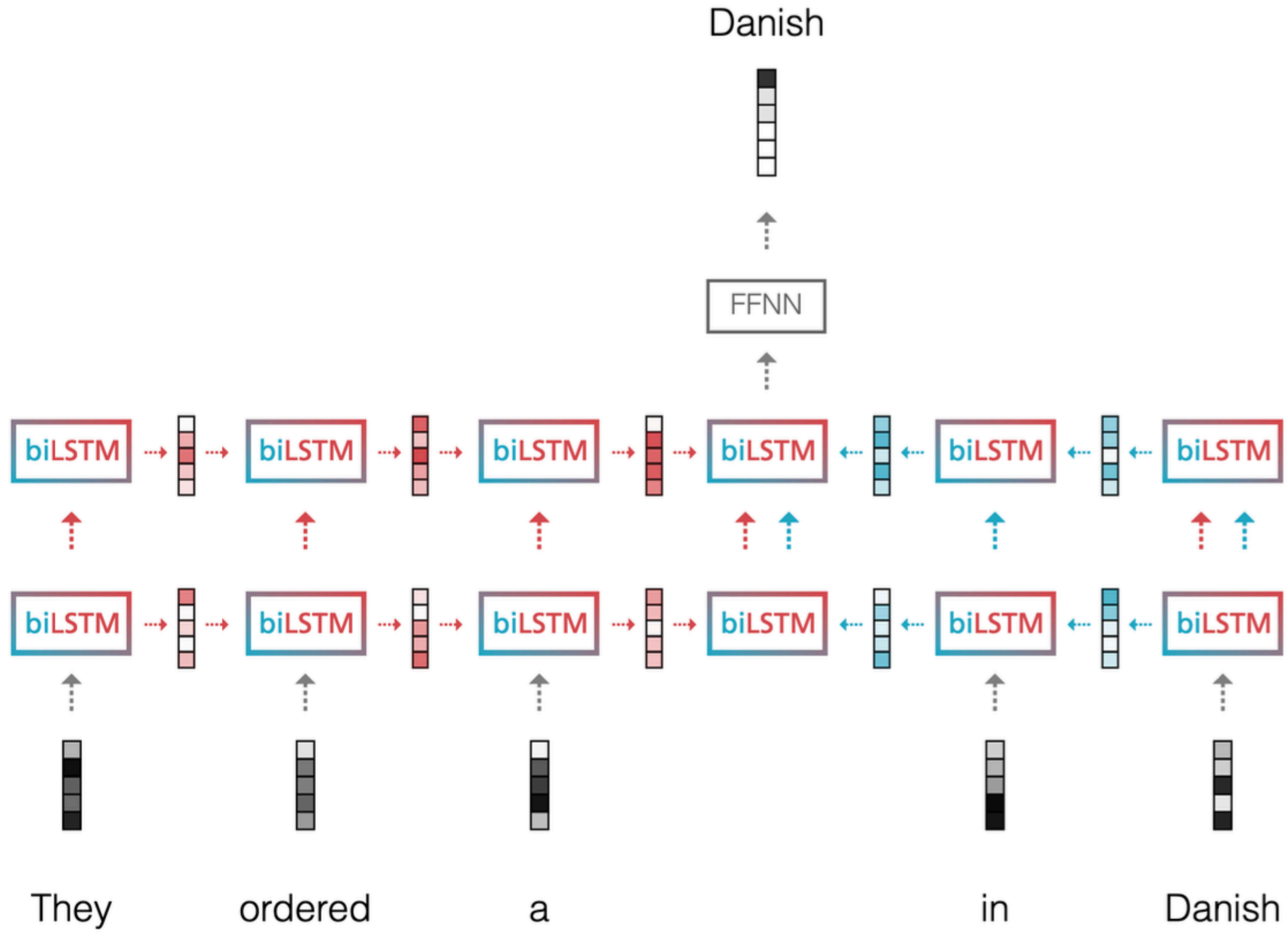&= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \ldots, L\},
\end{aligned}
$$

‣ This gives 2L+1 representations - Which to use?

  ‣ Just the top layer (similar to TagLM; Peters et al., 2017)

  ‣ Include all L+1 layers, average

  ‣ All layers, weighted average (best)

$$
\mathbf{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^{L} s_j^{task} \mathbf{h}_{k,j}^{LM}
$$

# Training ELMo

# Language Modeling Objective

- Word2vec trains word by word

    - Uses context during training

    - No context during use

- ELMo trains on word sequences

    - Predict each token using its context

    - Softmax layer applied to the top layer's output

    - Uses sequences as input during use
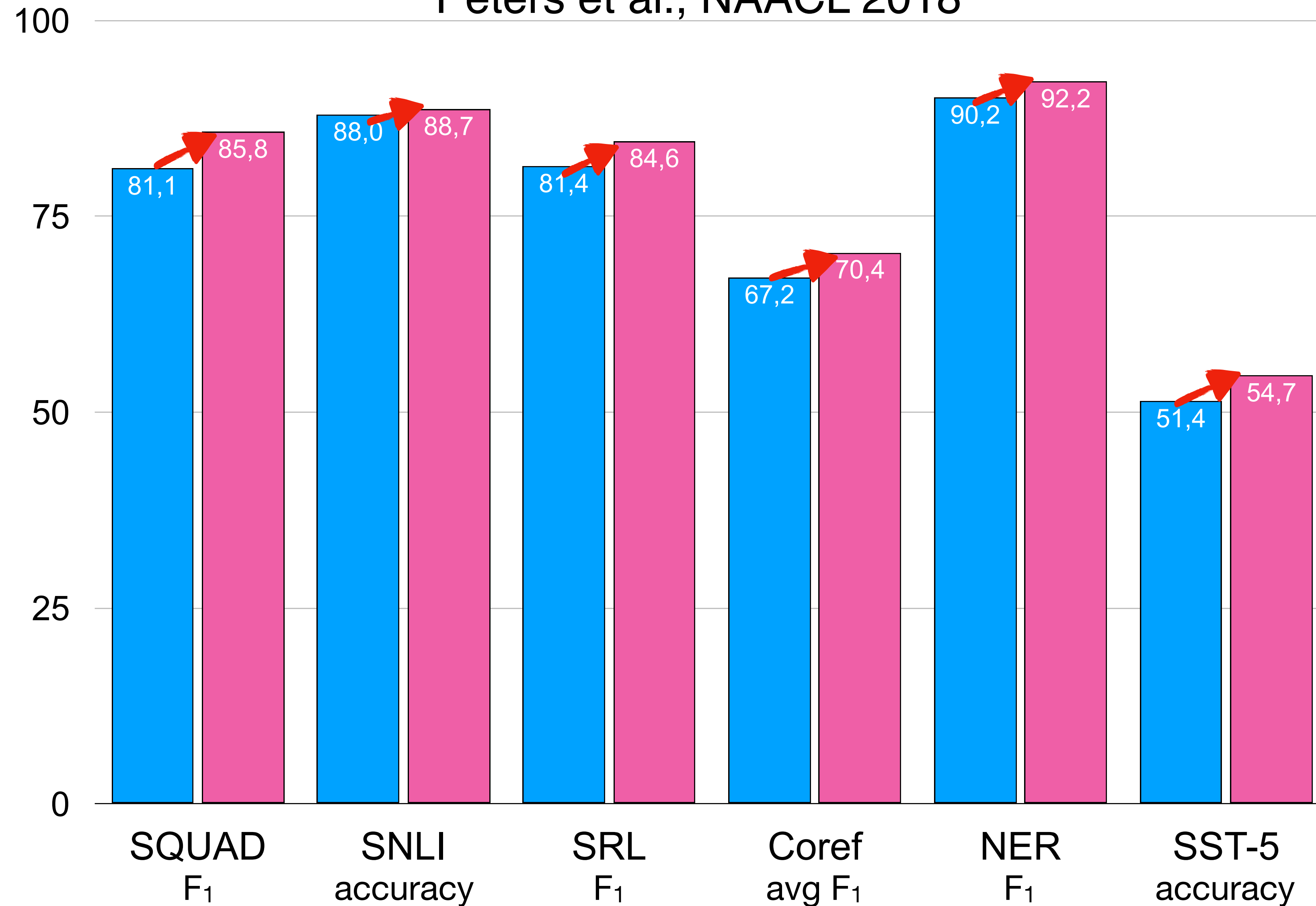
# How to use ELMo for your task?

‣ Recipe: For a given instance

  ‣ Run biLM to get the representations for each word

  ‣ Concatenate ELMo embeddings into task-specific model, e.g.,

    ‣ as additional input to static word embeddings

    ‣ as additional hidden representation

    ‣ … many choices, best might depend on end task
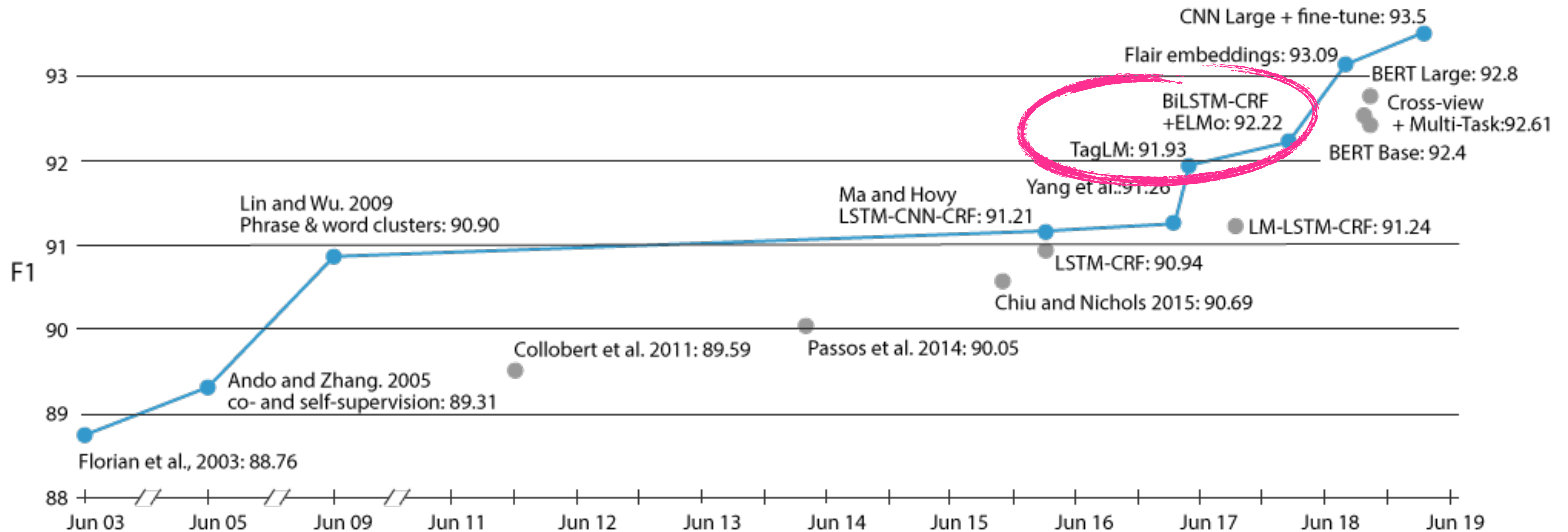
# Results over 6 NLP benchmarks

Peters et al., NAACL 2018

# NLP Progress on NER

‣ From Ruder et al.'s 2019 NAACL tutorial



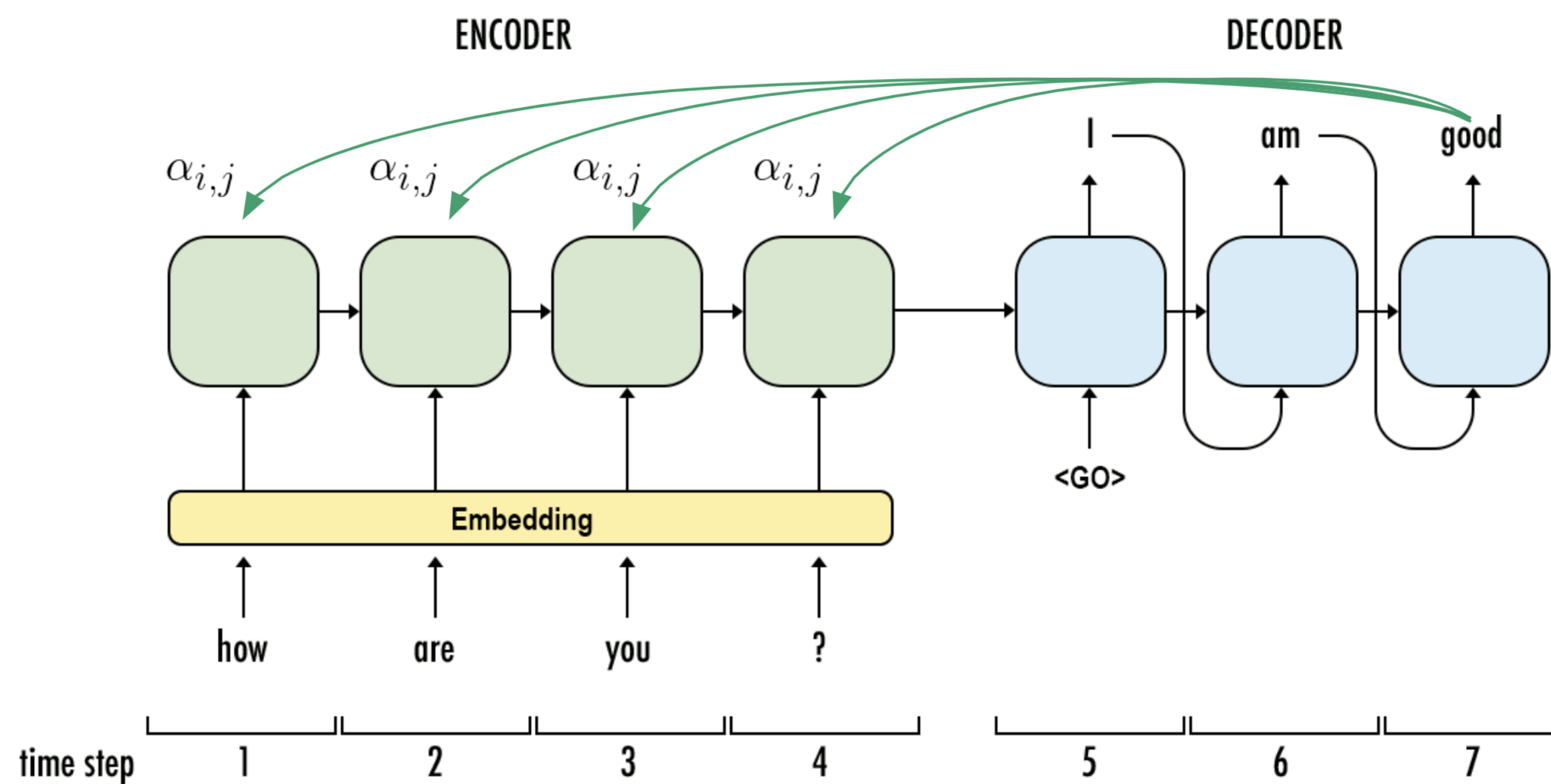Performance on Named Entity Recognition (NER) on CoNLL-2003 (English) over time

# Is ELMo the first such model? No!

‣ ELMo is *deeper* compared to an earlier model
by Peters et al., 2017 ACL (**TagLM**)

‣ It doesn't require parallel data (as an earlier model like **CoVe** does, by McCann et al., 2017 NeurIPS)

‣ ELMo: state-of-the-art performance on 6 benchmarks in 2018

# Limits of RNNs?

# Transformer & LLMs

**Part III**

# Today's roadmap

- **Part I: Fundamentals**

  - Intro, Motivation & Short History

  - Language Models (n-grams, FFNN-LM, Recap: FFNN)

- **Part II: Representations & Beyond FFNN**

  - RNNs (GRU/LSTMs), Attention

  - Contextualised Representations (ELMo)

- **Part III: Transformer & LLMs**

  - The Transformer, Masked LMs (BERT), Pre-training & Fine-tuning

  - Prompting, LLMs & Caution

# Transformers
## A Revolution in NLP and ML

# What's next

- The Transformer
  - High-level Architecture
    1. Encoder
    2. Decoder
    3. Multi-Head (Self) Attention

# FULLY ATTENTIONAL NETWORKS (A.K.A. TRANSFORMER)

# Core idea: Attention is All You Need (Vaswani et al. 2017)

- Attention has major impact on seq2seq performance

- Recurrency is an obstacle to parallelization

=> Can we build a <u>fully attentional</u> seq2seq model <u>without recurrency</u>?



RNN                    Transformer

Slides by Arianna Bisazza

# The Transformer
## Machine Translation

- **"Transforms"** one sequence into another

- No convolutions, no recurrence (no time steps), only **attention**

- **Highly parallelizable**

- Capture **long range dependencies**
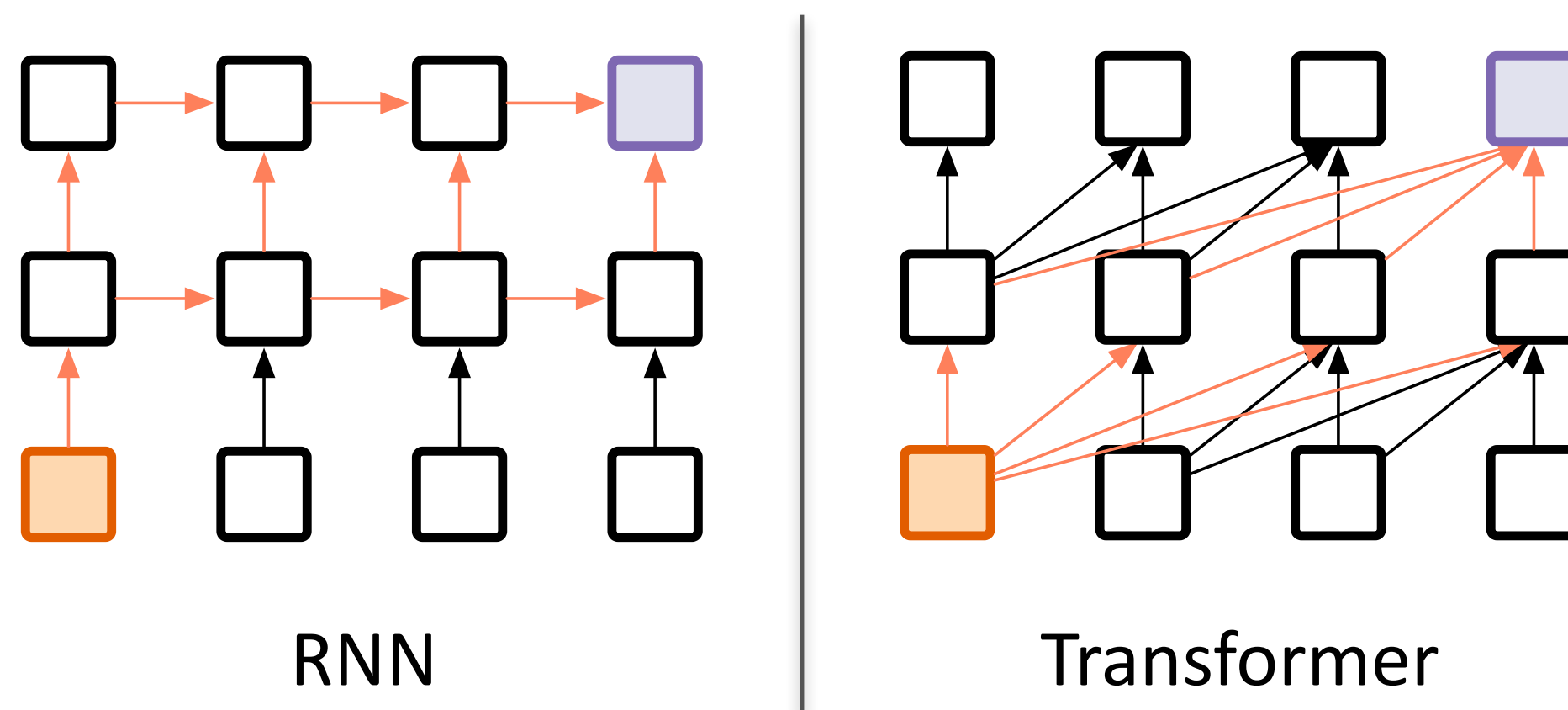
- **Better quality**

**Ashish Vaswani*** 
Google Brain 
avaswani@google.com

**Noam Shazeer*** 
Google Brain 
noam@google.com

**Niki Parmar*** 
Google Research 
nikip@google.com

**Jakob Uszkoreit*** 
Google Research 
usz@google.com

**Llion Jones*** 
Google Research 
llion@google.com

**Aidan N. Gomez*** [†] 
University of Toronto 
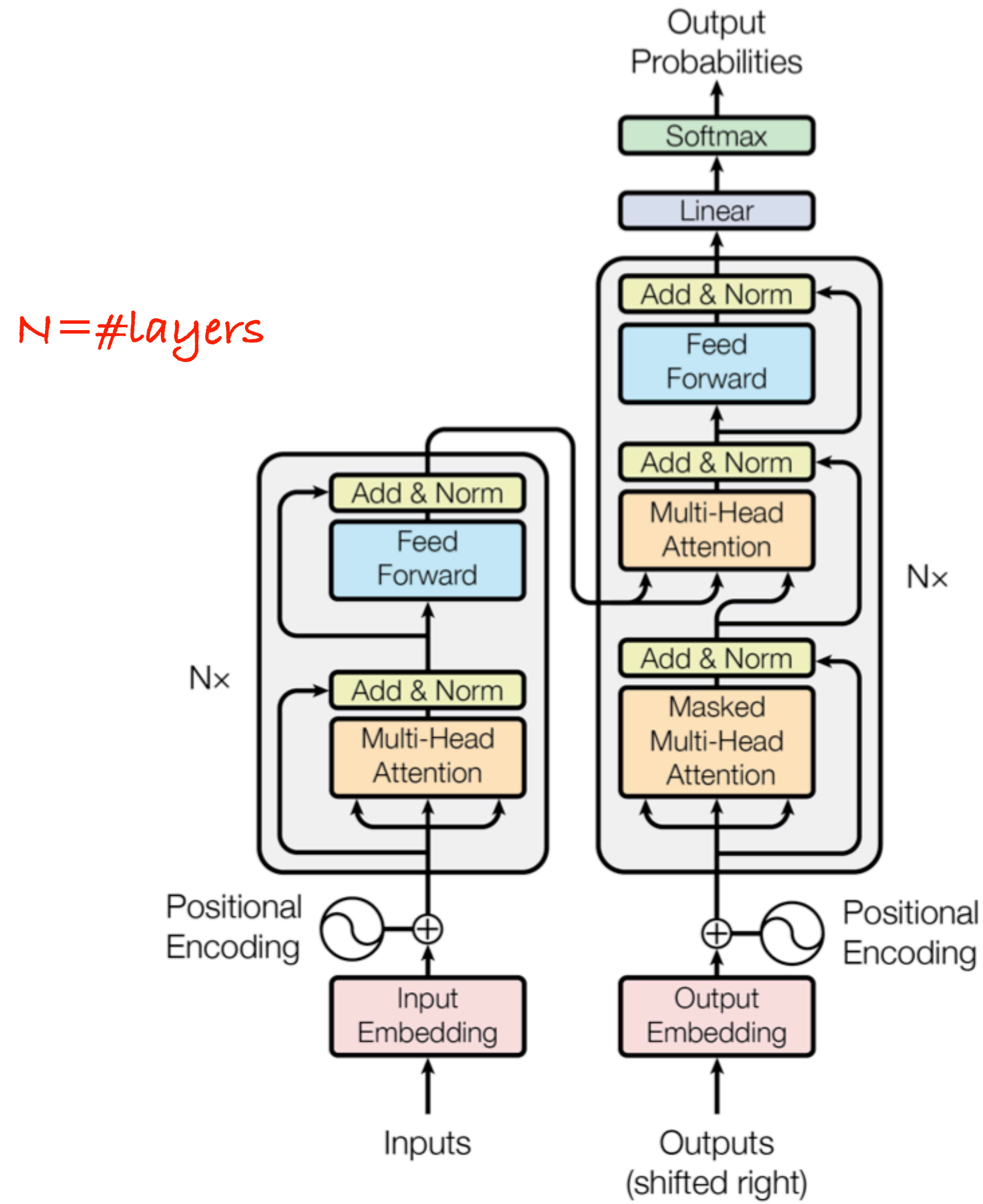aidan@cs.toronto.edu

**Łukasz Kaiser*** 
Google Brain 
lukaszkaiser@google.com

**Illia Polosukhin*** [‡] 
illia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

# A scary beast

N=#layers

# TRANSFORMER ARCHITECTURE OVERVIEW

# Transformer Architecture Overview

ENC-layer 2

Self-Attention

ENC-layer 1

Self-Attention

Input word
embedding

machine   translation   is   fun

**ENCODER**

# Transformer Architecture Overview



*traduzione*

DEC-layer 2

Self-Attention

DEC-layer 1

Self-Attention

Output word embedding

Enc-Dec Attention

[EOS]       la

**DECODER**

ENC-layer 2

Self-Attention

ENC-layer 1

Self-Attention

Input word embedding

machine   translation   is   fun

**ENCODER**

# Transformer Architecture Overview



*automatica*

DEC-layer 2

DEC-layer 1

Output word embedding

[EOS]          la          traduzione

**DECODER**

Enc-Dec Attention

ENC-layer 2

Self-Attention

ENC-layer 1

Self-Attention

Input word embedding

machine   translation   is   fun

**ENCODER**

**Let's take a closer look:**

# TRANSFORMER'S BUILDING BLOCKS

# Scaled Dot-Product Attention

To compute attention we need a scoring function

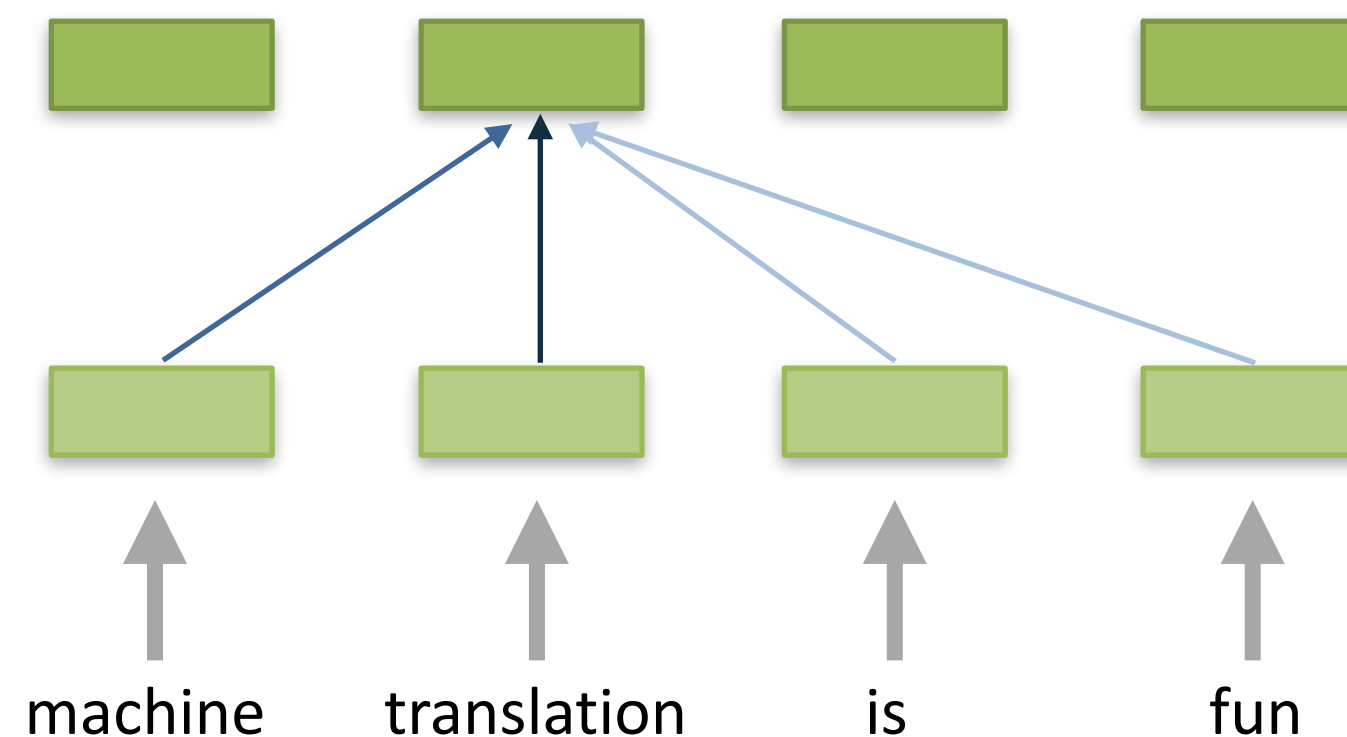- **Dot-product** is simple and fast to compute*

- Rationale: measure similarity of two (word-)vectors

$$\text{score}(q_t, k_i) = q_t^\top k_i$$

Problem: for high-dimensional vectors, softmax gets very peaked and gradients small

=> Solution: scale the result of dot product

$$\text{score}(q_t, k_i) = \frac{q_t^\top k_i}{\sqrt{d}}$$



machine    translation    is    fun

*For a nice overview of different Attention Scoring Functions see:
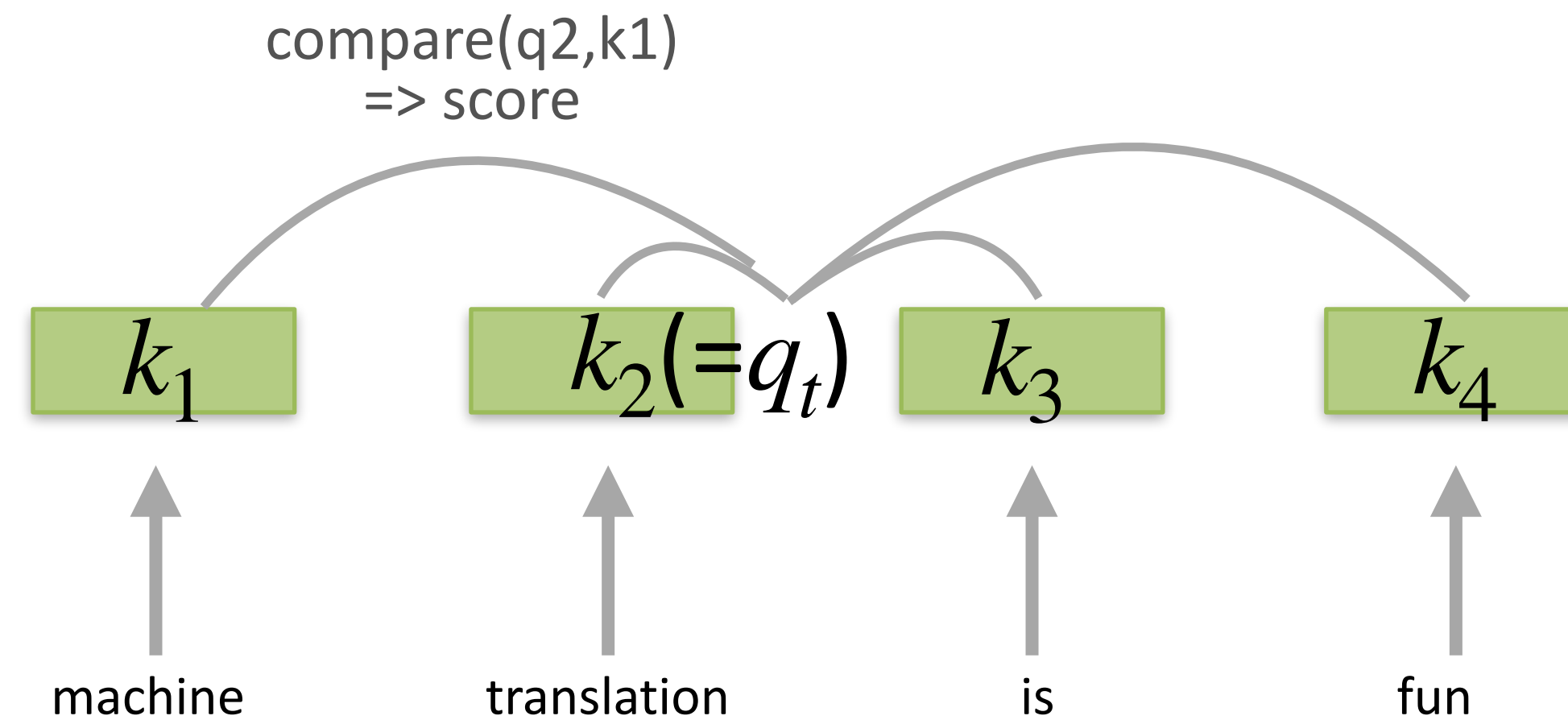https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3#ba24

# Query-Key-Value
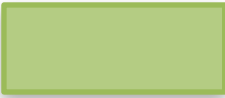
Now, where do *q* and *k* come from?

We could simply use the word vector [ ] and compare it to all vectors in the sentence (including itself)

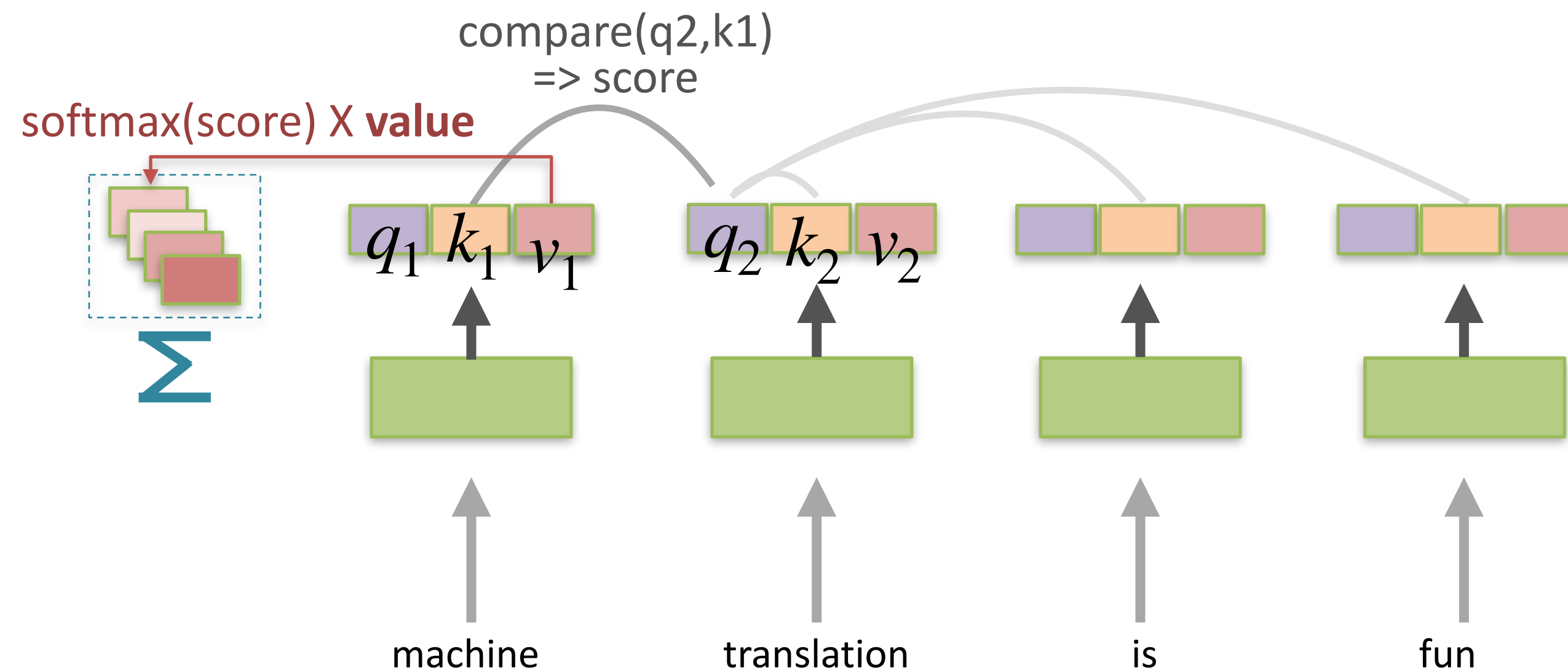$$\text{score}(q_t, k_i) = \frac{q_t^\top k_i}{\sqrt{d}}$$

compare(q2,k1)
=> score

$k_1$    $k_2(=q_t)$    $k_3$    $k_4$

machine    translation    is    fun

# Query-Key-Value

Now, where do *q* and *k* come from?

We could simply use the word vector [green box] and compare it to all vectors in the sentence (including itself)

$$\text{score}(q_t, k_i) = \frac{q_t^\top k_i}{\sqrt{d}}$$

A better idea: Learn multiple 'views' of [green box] to use as **query**, **key** and **value**

compare(q2,k1)
=> score

softmax(score) X **value**

$q_1$ $k_1$ $v_1$      $q_2$ $k_2$ $v_2$

$\Sigma$

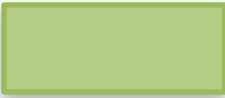machine          translation          is          fun
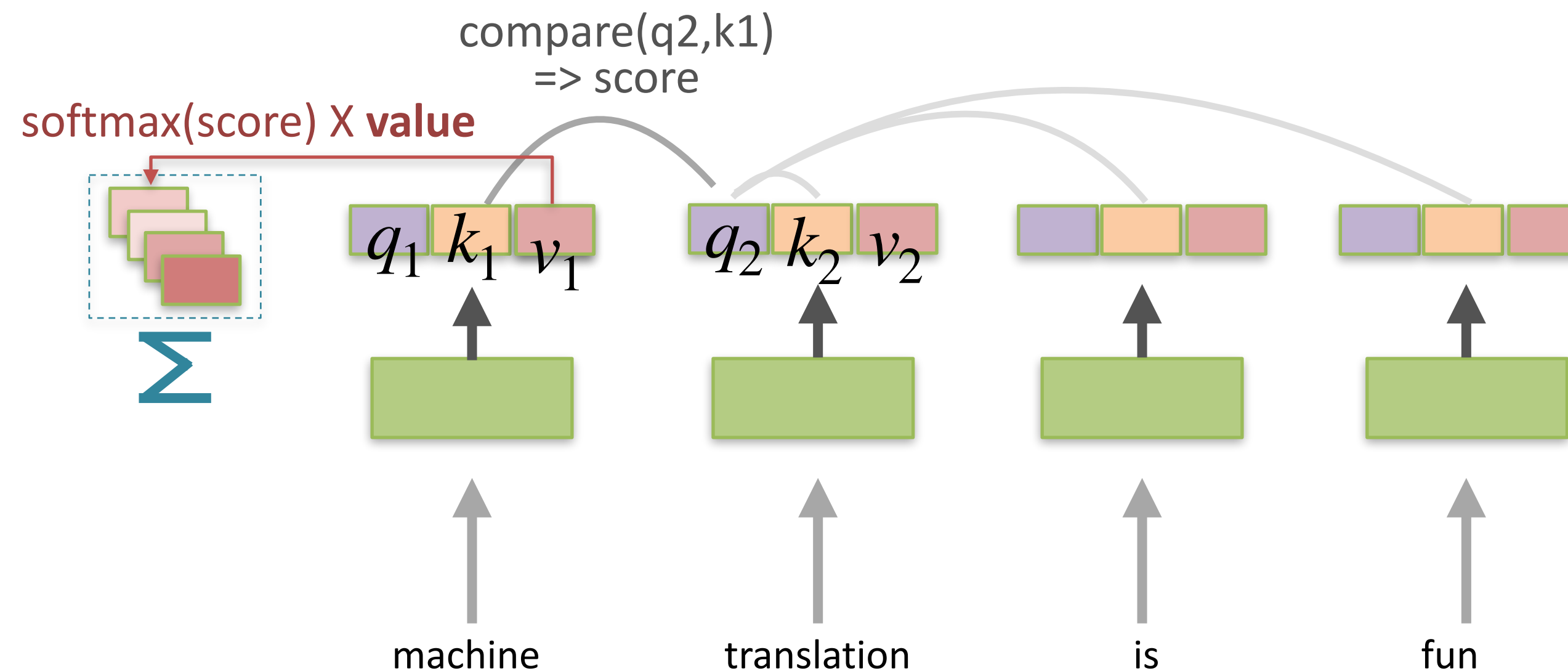
# Self Attention in the Transformer

- Transform each token representation to a **learnable query**, **key**, and **value** vector

  - **Query vector** "*asks*" all **key vectors** in the sequence whether they are relevant. This results in **attention scores.**

  - "Asking" is the **dot product** (similarity) between the **query vector** and **key vector**

  - New **contextualised** token representation after self-attention is the **weighted average** over all **value vectors** in the sequence using attention scores as weights

# Query-Key-Value

Now, where do *q* and *k* come from?

We could simply use the word vector [  ] and compare it to all vectors in the sentence (including itself)

A better idea: Learn multiple 'views' of [  ] to use as **query**, **key** and **value**

compare(q2,k1)
=> score

softmax(score) X **value**

$q_1$ $k_1$ $v_1$      $q_2$ $k_2$ $v_2$

Σ

machine      translation      is      fun

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{Z} = \text{softmax}(\mathbf{QK}/\sqrt{d_k})\mathbf{V}$$

**We are not done yet …**

# Multi-Head Attention

- **Intuition:** Information from different parts of the sentence can be useful to disambiguate in different ways
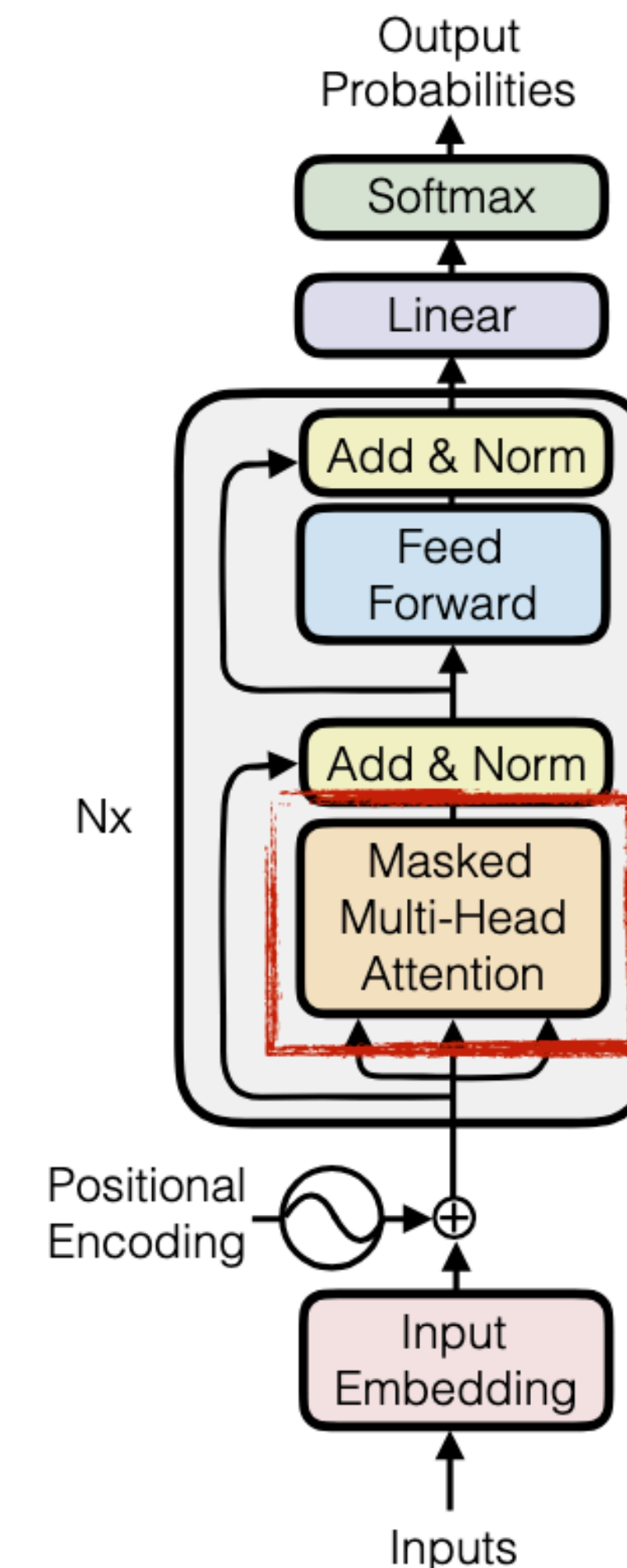
I **run** a small business

I **run** a mile in 10 minutes

The robber made a **run** for it

The stocking had a **run**

syntax
(nearby context)

semantics
(farther context)



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Masked Multi-Head Attention

Nx

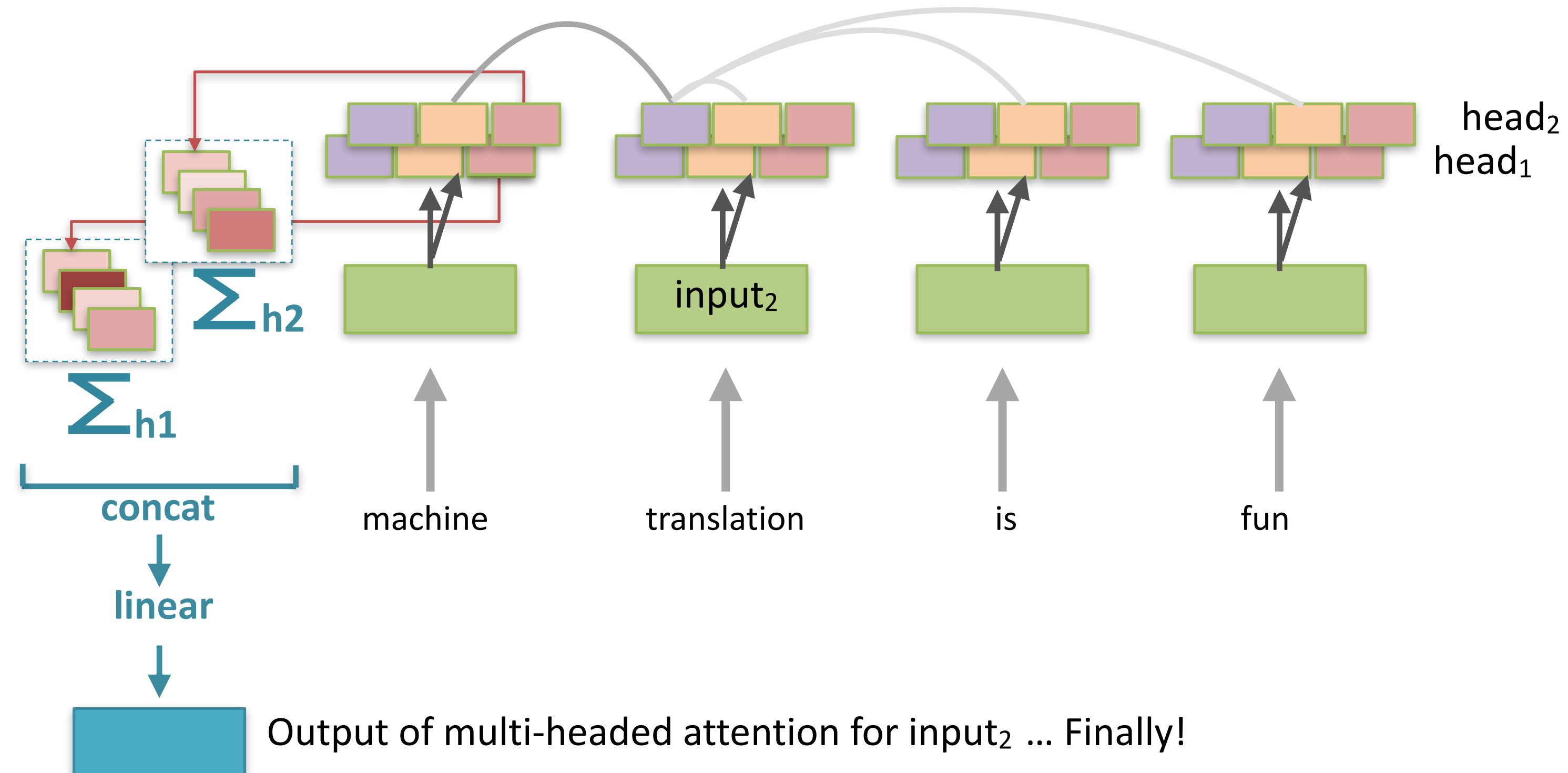Positional Encoding

Input Embedding

Inputs

# Multi-Head Attention

Words can interact with each other in different ways.

One attention distribution may not be enough to capture: coreference effects, topic cohesion, other syntactic/semantic relationships, etc.

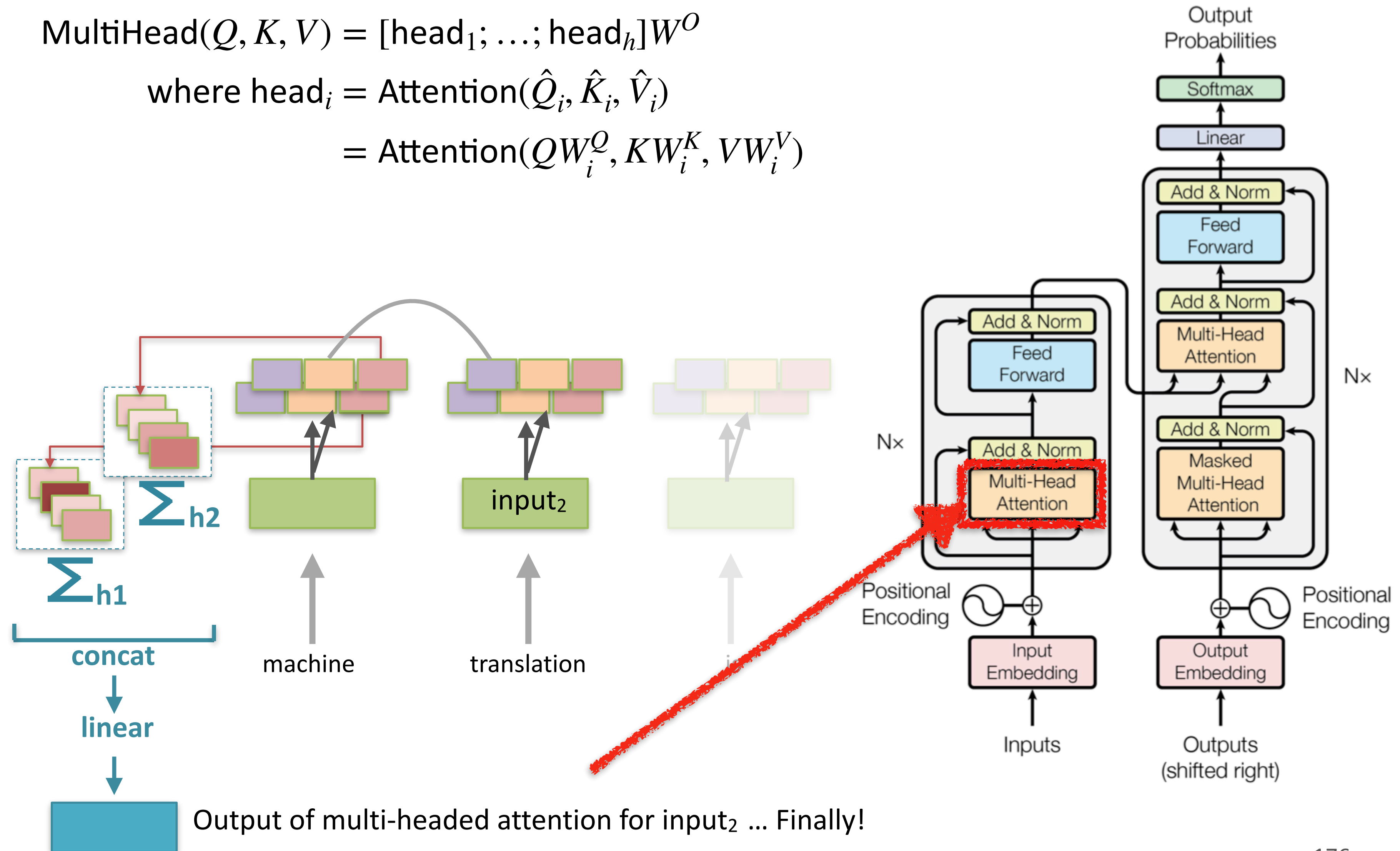Multi-Head gives the attention layer multiple *representation subspaces*



$\sum_{h2}$

$\sum_{h1}$

concat

linear

head$_2$
head$_1$

input$_2$

machine    translation    is    fun

Output of multi-headed attention for input$_2$ … Finally!

# Multi-Head Attention

$$\text{MultiHead}(Q, K, V) = [\text{head}_1; \ldots; \text{head}_h]W^O$$

$$\text{where head}_i = \text{Attention}(\hat{Q}_i, \hat{K}_i, \hat{V}_i)$$

$$= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$



Output of multi-headed attention for input$_2$ … Finally!

# Multi-Head Attention
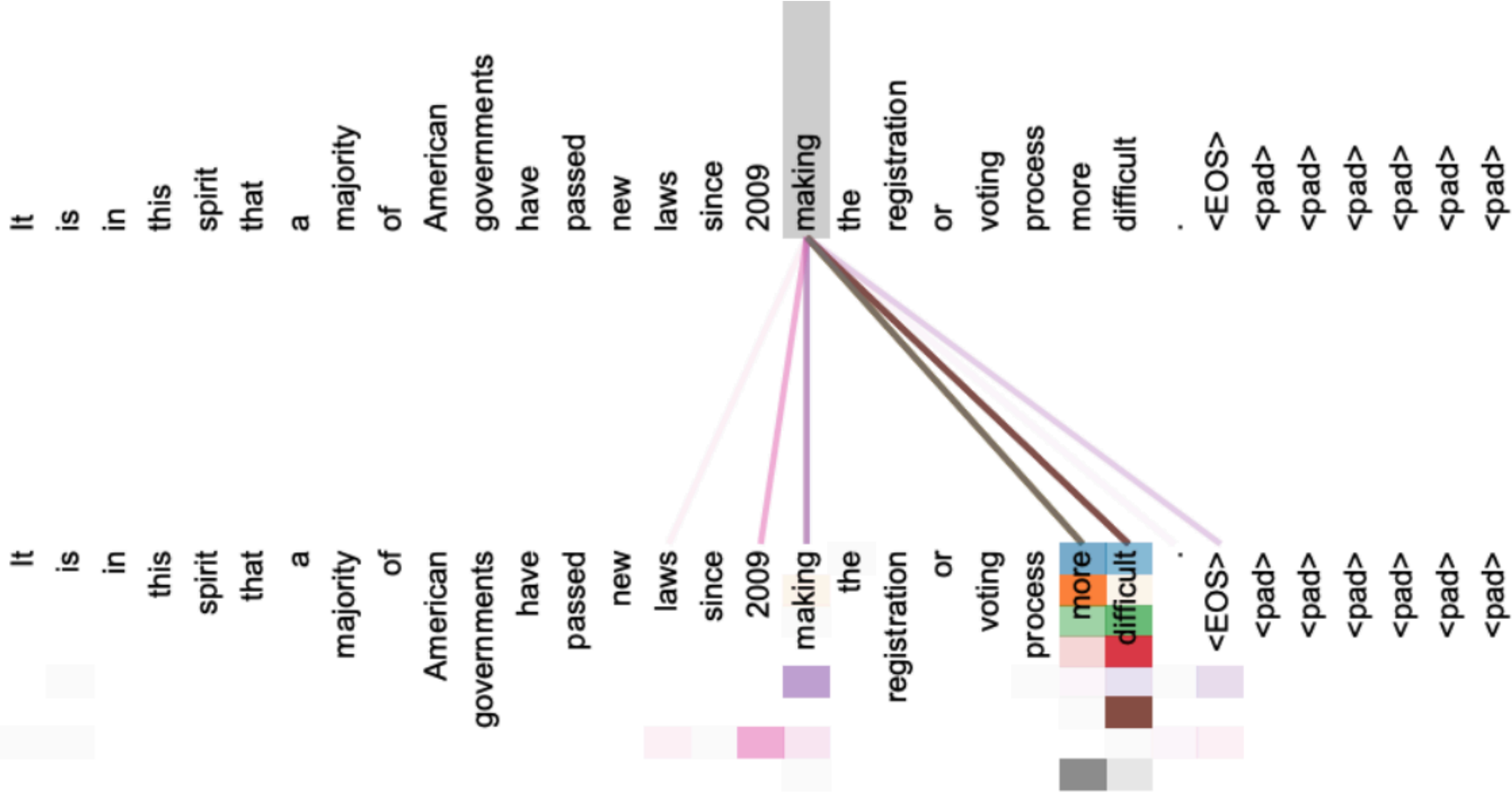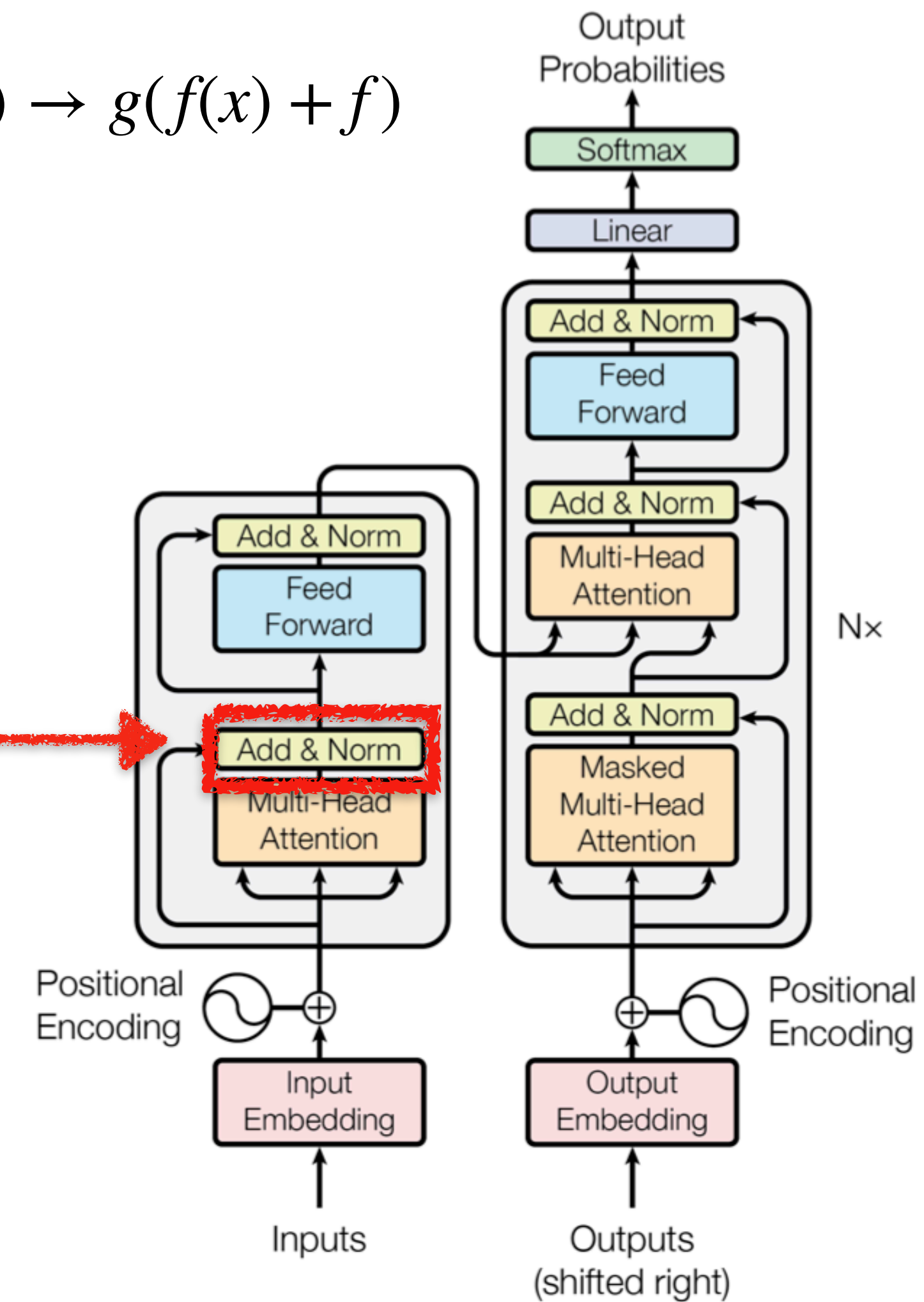
Example from Vaswani et al., 2017:



Figure 3: An example of the attention mechanism following long-distance dependencies in the encoder self-attention in layer 5 of 6. Many of the attention heads attend to a distant dependency of the verb 'making', completing the phrase 'making...more difficult'. Attentions here shown only for the word 'making'. Different colors represent different heads. Best viewed in color.

# Add & Norm

Last ingredients:

- Residual connection (He et al. 2015) $\quad g(f(x)) \rightarrow g(f(x) + f)$

- Layer normalization (Ba, Kiros & Hinton)



**LayerNorm(SubLayer(x) + x)**
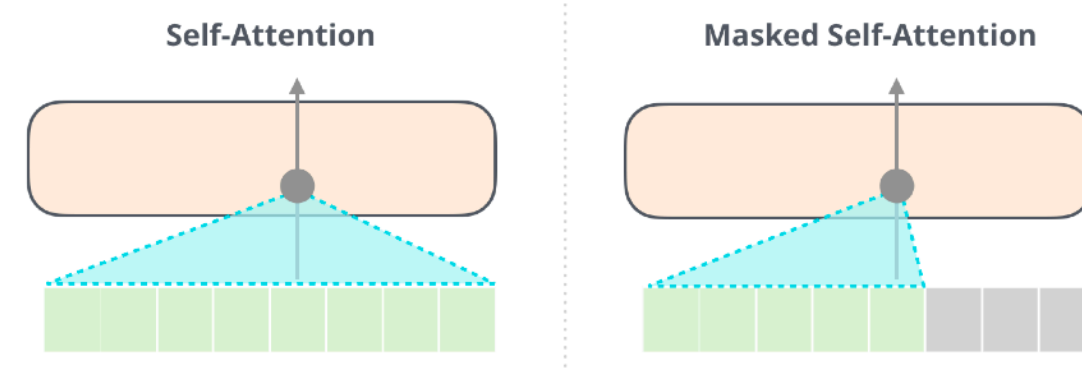
# Attention in the Decoder

We have looked at self-attention in the encoder

Now moving to the decoder => 2 types of attention here:

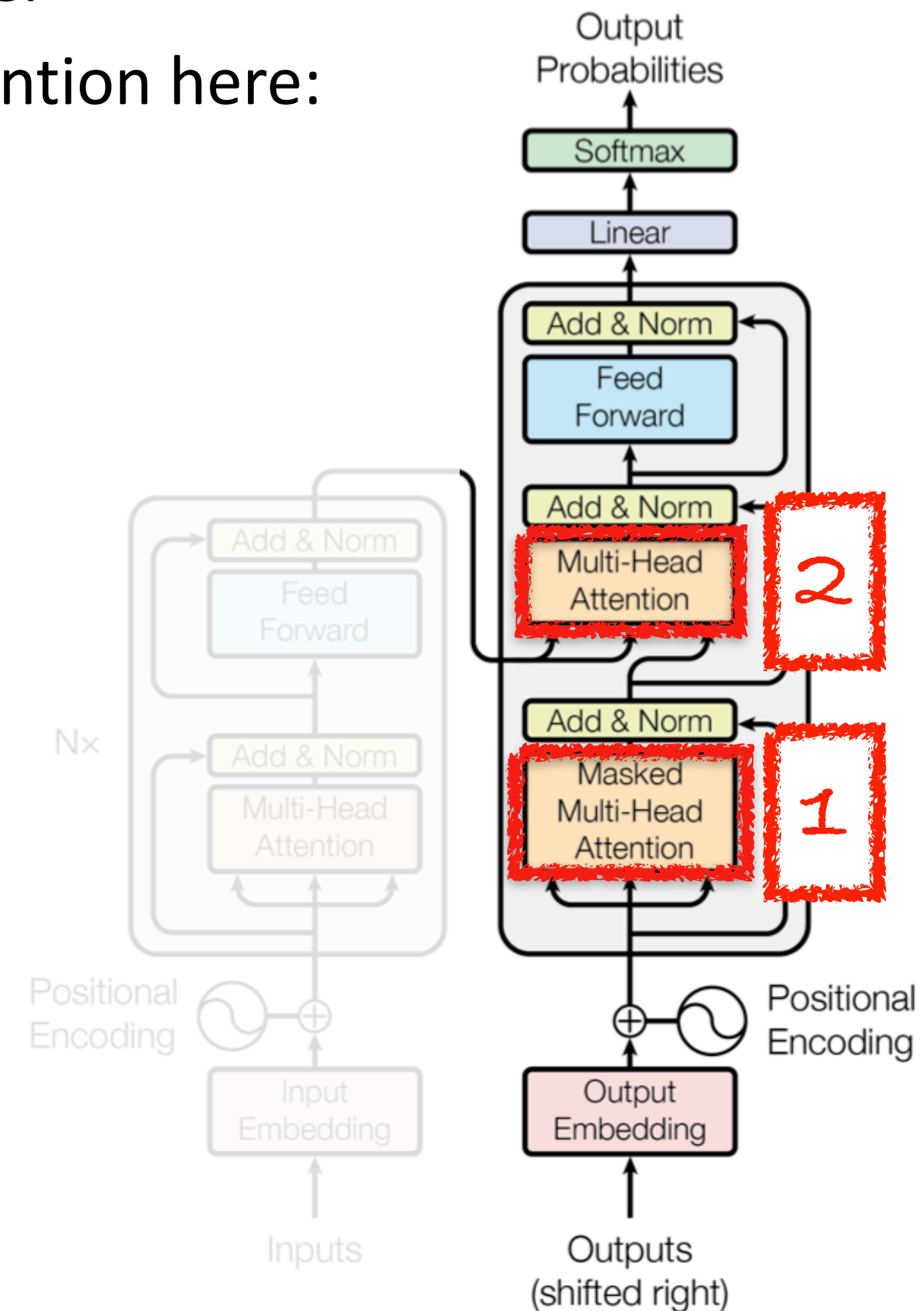**1** Masked Self Attention (Decoder only):

- captures target-side context

- same as before, but can only look at positions before the current word (*masked*)
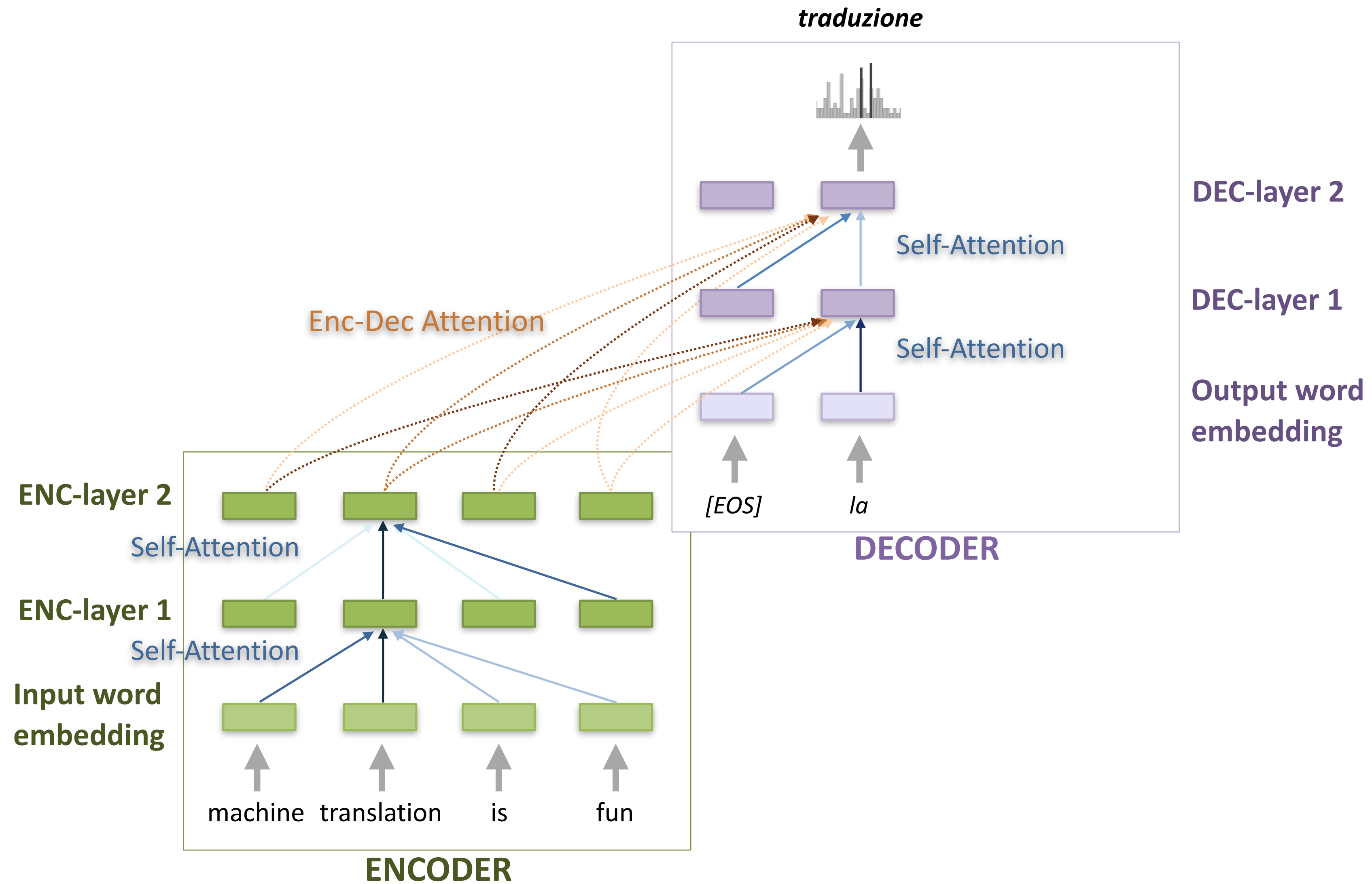
Self-Attention          Masked Self-Attention
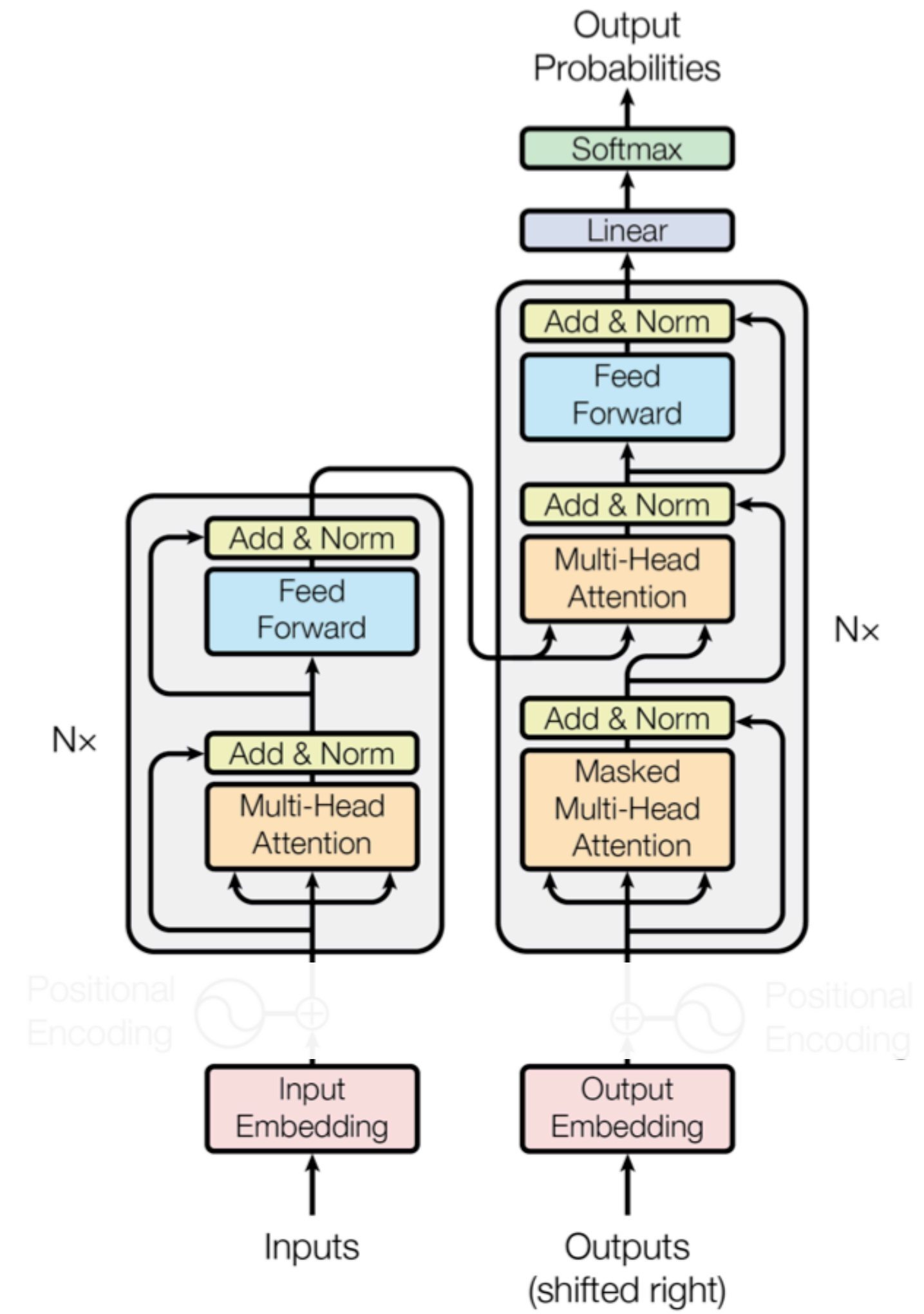
**2** Encoder→Decoder Attention:

- captures src-trg translation equivalences

- Query comes from target (decoder),
Key & Value from source (encoder)

- also known as Cross Attention

Masked self-attention illustration from Artzi's LM class.

# Transformer Architecture Overview

*traduzione*

DEC-layer 2

Self-Attention

DEC-layer 1

Self-Attention

Enc-Dec Attention

Output word
embedding

[EOS]          la

DECODER

ENC-layer 2

Self-Attention

ENC-layer 1

Self-Attention

Input word
embedding

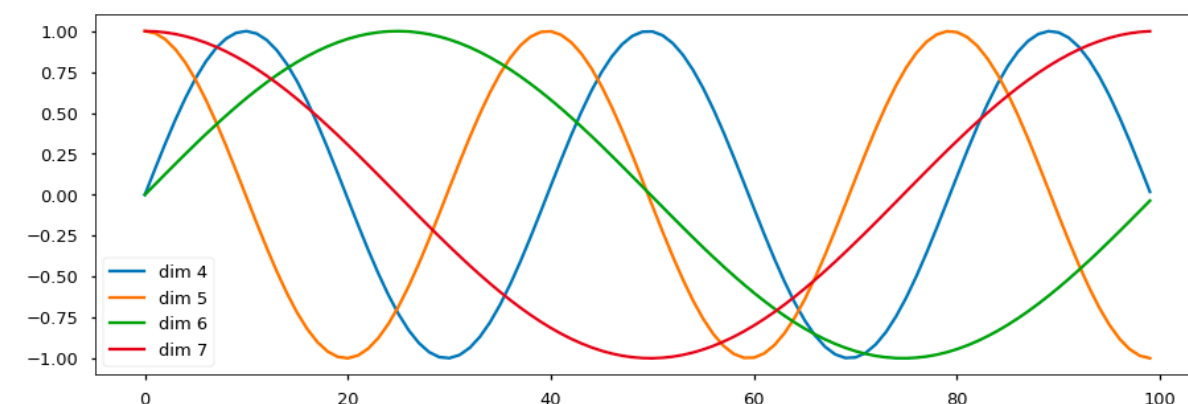machine   translation   is   fun

ENCODER

**Are we missing anything?**

# Positional embeddings

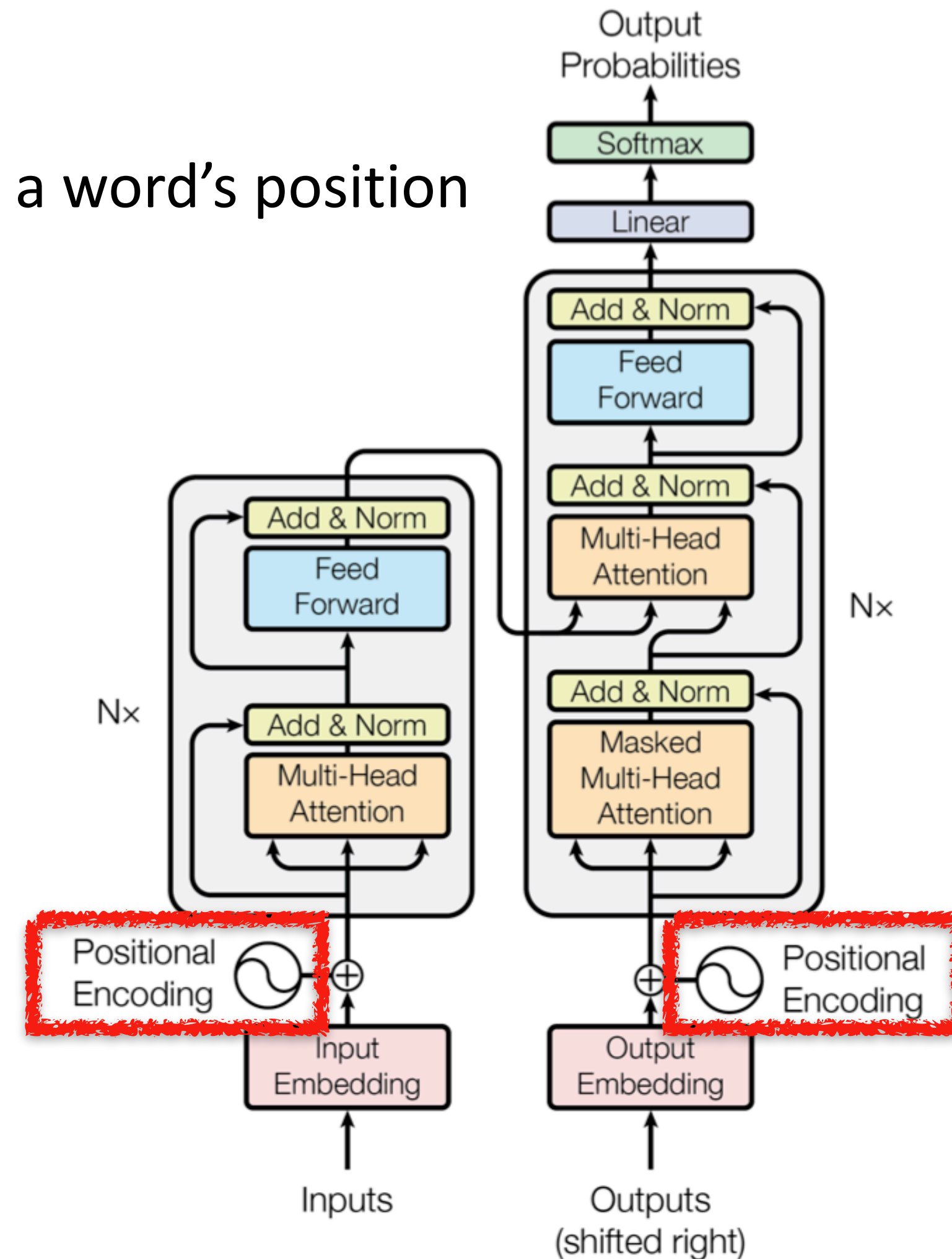Recurrency naturally represents the order of words in a sentence:

*w₃ comes after w₂ which comes after w₁ ...*

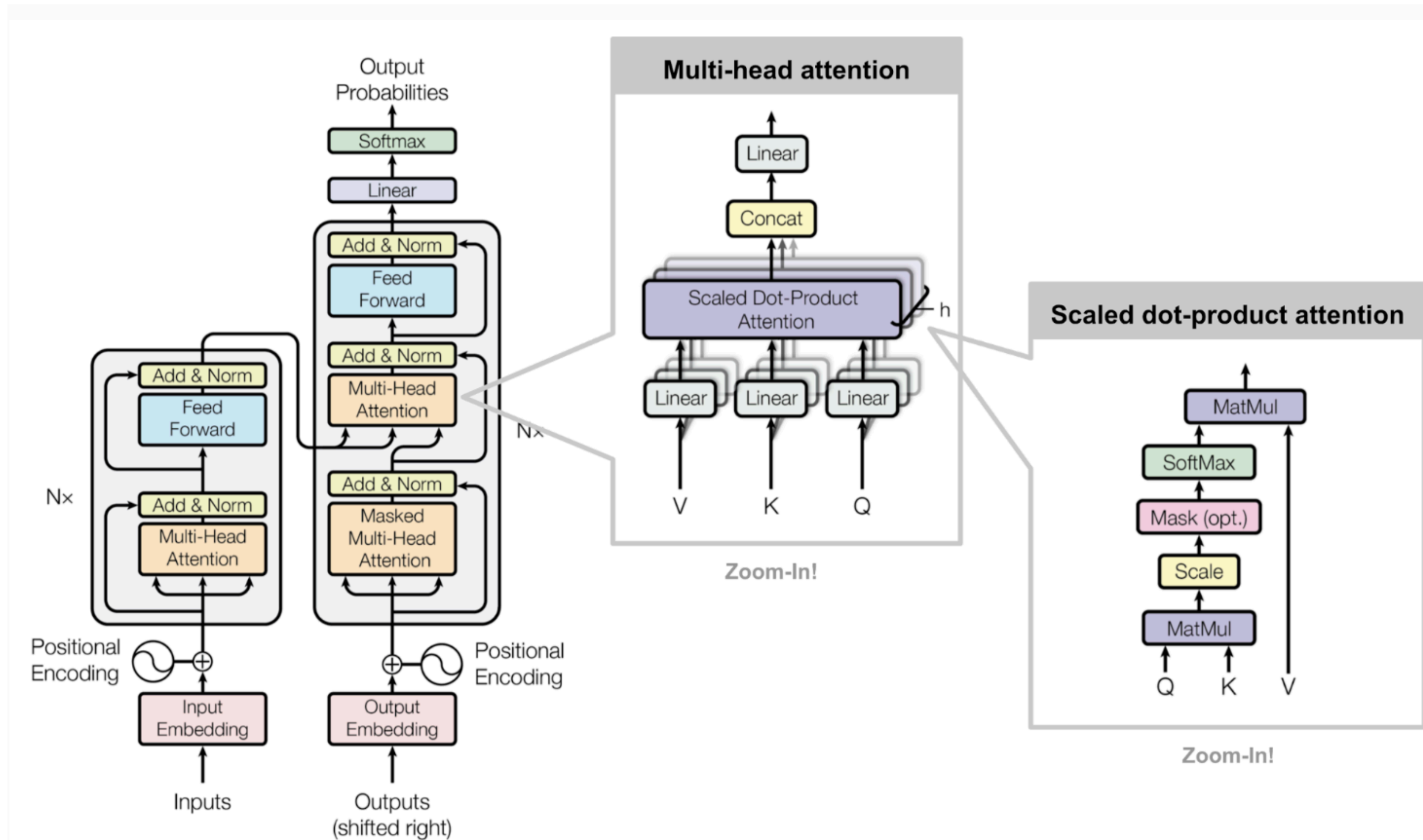Transformer needs an explicit way to represent a word's position

- Idea: let's mark positions

- Learning will figure out what how to use them

- Simple version: **learnable** embeddings $\phi_p(i)$
  where $i \in [0,1,2,..]$ (lookup parameters)

- More advanced: **fixed** embeddings, where values determined by sine waves, with different frequency and offset of each dimensions



- Either way, we add them to token embeddings

# Putting it altogether: Encoder-Decoder Transformer

# Transformer: Summary of Core Concepts

‣ Attention:

  ‣ Self Attention

  ‣ Cross Attention (Encoder-Decoder Attention)

  ‣ Multi-Head Attention

  ‣ Masked Attention

‣ Positional encodings

‣ Residual connections + layer normalisation
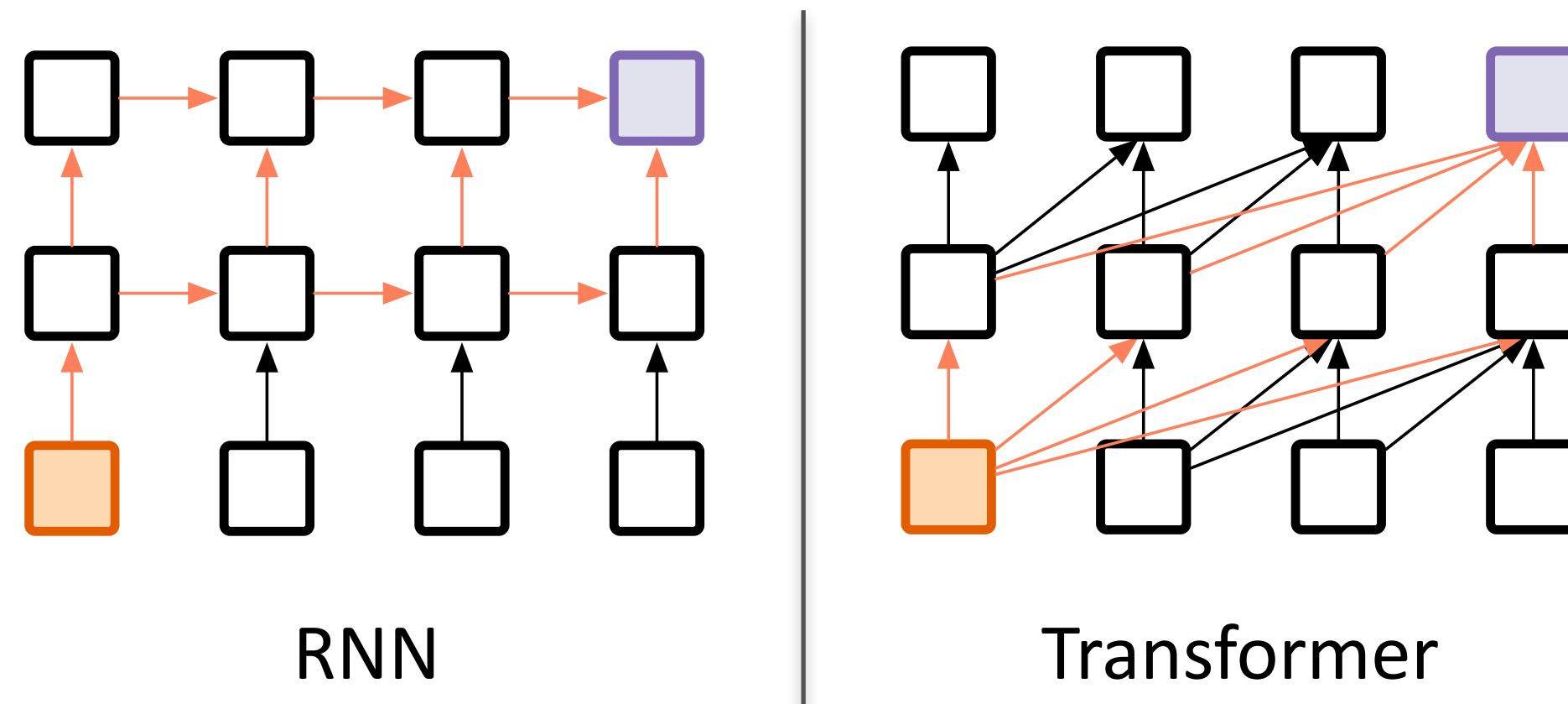
‣ Feed-forward layer

# Three Types of Transformers



**Encoder-Only Model (e.g. BERT)**

**Encoder-Decoder Model (Vaswani et al., 2017), T5**

**Decoder-Only Model (e.g. GPT, Llama)**

# RECURRENT SEQ-TO-SEQ VS TRANSFORMER

# RNN-seq2seq vs Transformer



RNN

Transformer

✔ Much more parallelizable = fast

✔ Lower complexity
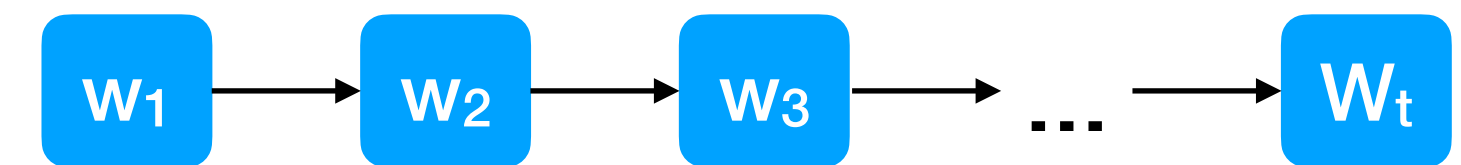
✔ Shorter path among any input positions

- RNNs (esp. LSTM) are cognitively inspired: represent memory constraints
- Transformer = result of clever engineering & brute-force architecture search
  - It works! State of the art performance on almost every NLP task
- Will we go back to RNN-kind of models?

# Today's roadmap

‣ **Part I: Fundamentals**

  ‣ Intro, Motivation & Short History

  ‣ Language Models (n-grams, FFNN-LM, Recap: FFNN)

$$w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow \ldots \rightarrow w_t$$

‣ **Part II: Representations & Beyond FFNN**

  ‣ RNNs (GRU/LSTMs), Attention

  ‣ Contextualised Representations (ELMo)

‣ **Part III: Transformer & LLMs**

  ‣ The Transformer, Masked LMs (BERT) Pre-training & Fine-tuning

  ‣ Prompting, LLMs & Caution

# Masked LMs: The birth of BERT

- LMs so far: predict the next token given the previous tokens

  - This enables a self-supervised task

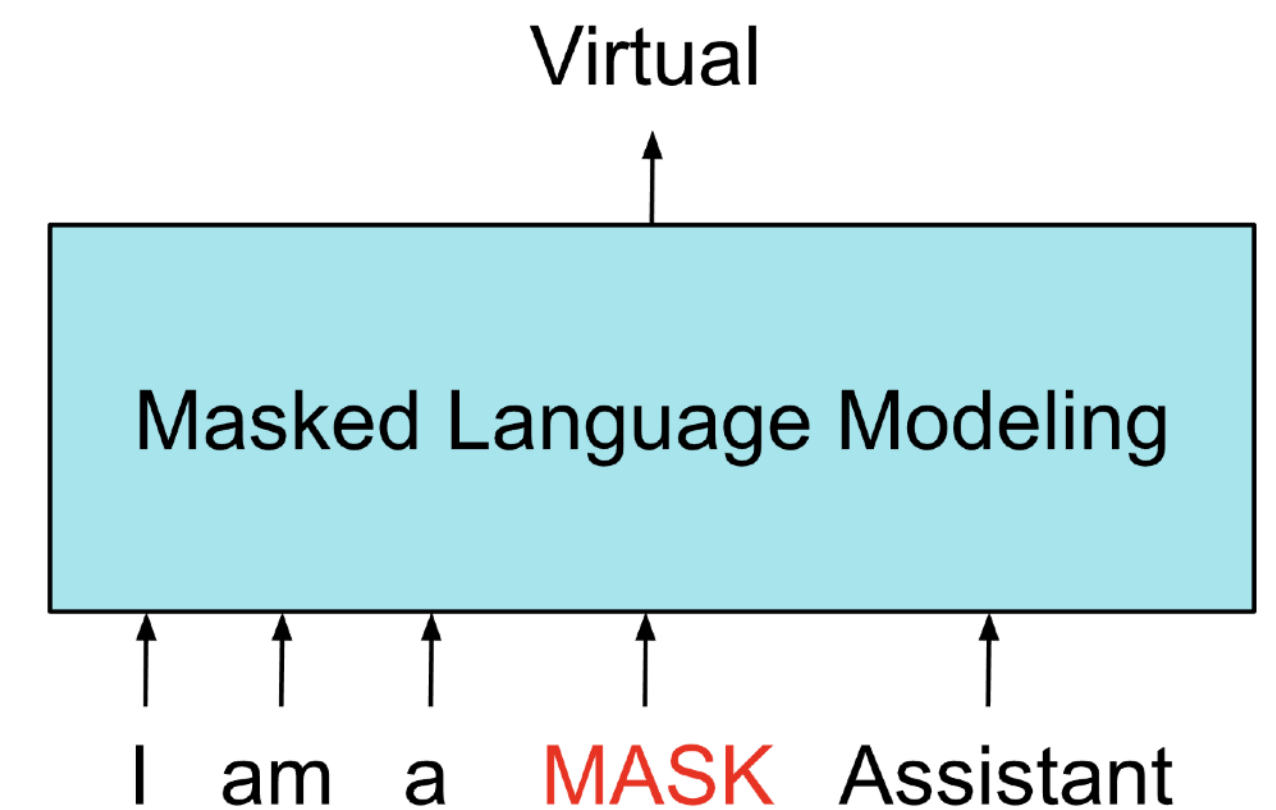  - That we can train on a lot of data to get really useful representations

- Let's **create a prediction task** by **hiding** part of the sequence, and then trying to predict them

  - Input: the sequence $\bar{x}^M$ where some tokens are replaced with the token $[\text{MASK}]$, for example:
    $$\bar{x}^M = \langle x_1, \ldots, x_4, [\text{MASK}], x_6, \ldots, x_n \rangle$$

Virtual

Masked Language Modeling

I    am    a    MASK    Assistant

# BERT
## Bidirectional Encoder Representations from Transformers

- Encoder transformer

  - Encoders assume we have the complete sequence

  - No generation, only want good representations (for analysis tasks)

- BERT Base: 12 transformer blocks, 768-dim word-piece tokens, 12 self-attention heads → 110M parameters

- BERT Large: 24 transformer blocks, 1024-dim word-piece tokens, 16 self-attention heads → 340M parameters

- RoBERTa: same model, much more data (160GB of data instead of 16GB)

[Devlin et al. 2018]

# BERT

## Inputs

- One or two sentences

  - Word-piece token embeddings

  - Position and segment embeddings

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

[figure from Devlin et al. 2018]

# BERT
## Training

- Data: raw text

- Two objectives:

  - Masked LM

  - Next-sentence prediction

- Later development in RoBERTa:

  - More data, no next-sentence prediction, dynamic masking

# BERT
## Masking Recipe for Training

- Mask and predict 15% of the tokens

  - For 80% (of 15%) replace with the input token with [MASK]

  - For 10%, replace with a random token

  - For 10%, keep the same

# BERT
## Next-sentence Prediction

- Input: [CLS] Text chunk 1 [SEP] Text chunk 2

- Training data: 50% of the time, take the true next chunk of text, 50% of the time take a random other chunk

- Predict whether the next chunk is the true next chunk

- Prediction is done on the [CLS] output representation

# BERT
## What Do We Get?

- We can feed complete sentences to BERT

- For each token, we get a contextualized token representation

  - Similar to ELMo, but without an RNN/LSTM

  - In contrast to word2vec representations that are fixed and do not depend on context

- While word2vec vectors are forced to mix multiple senses, BERT/ELMo can provide more instance-specific vectors

- BERT started an arms race towards bigger and bigger models, which quickly led to the LLMs of today

# BERT

## What It Is Not Great For?

- BERT cannot generate text

  - Not an autoregressive model, but can do filling of [MASK]

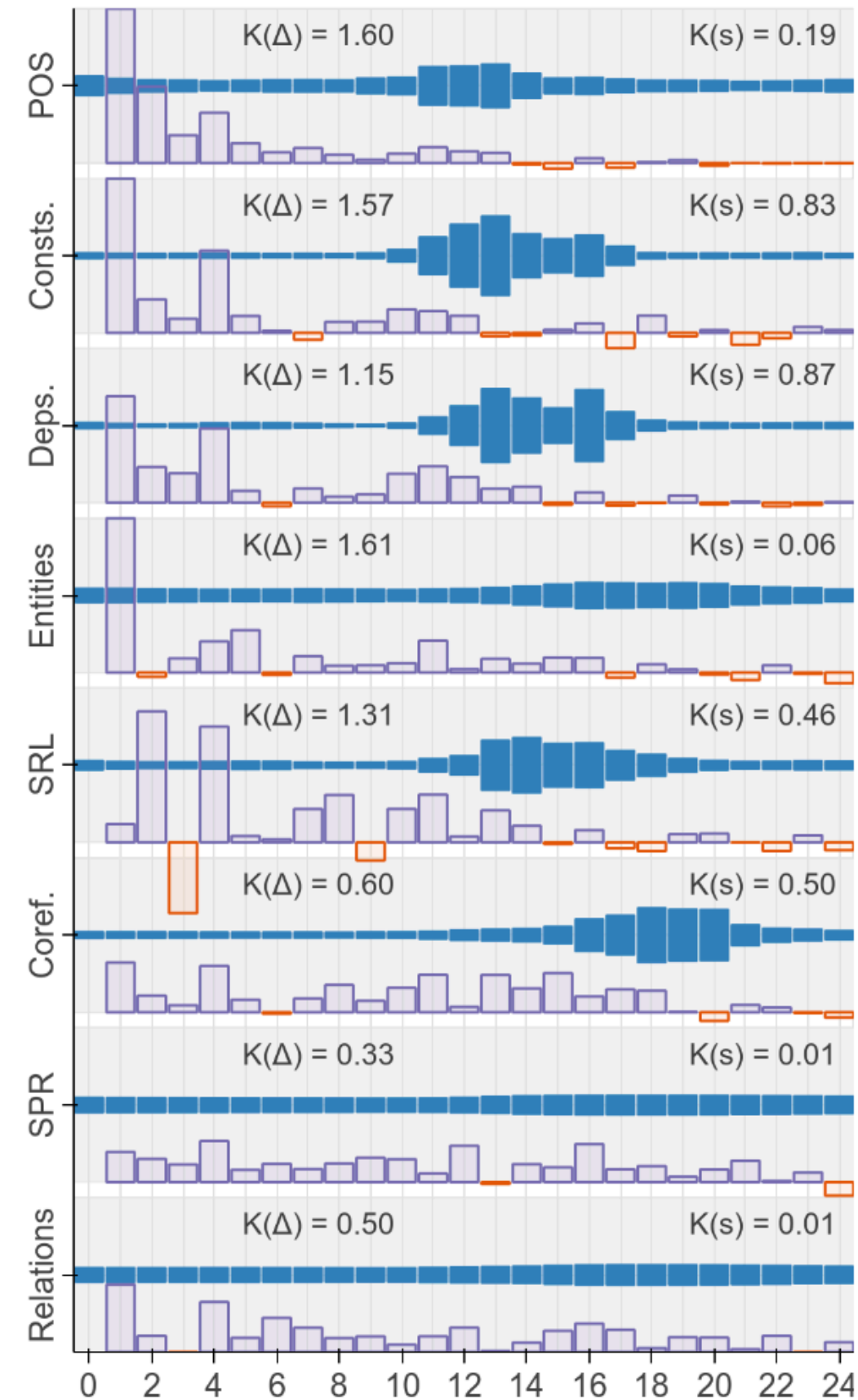- Masked language models are intended to be used primarily for "analysis" tasks

# BERT
## What does BERT Learn?

- There is a lot of work trying to decipher what BERT learns in its representations

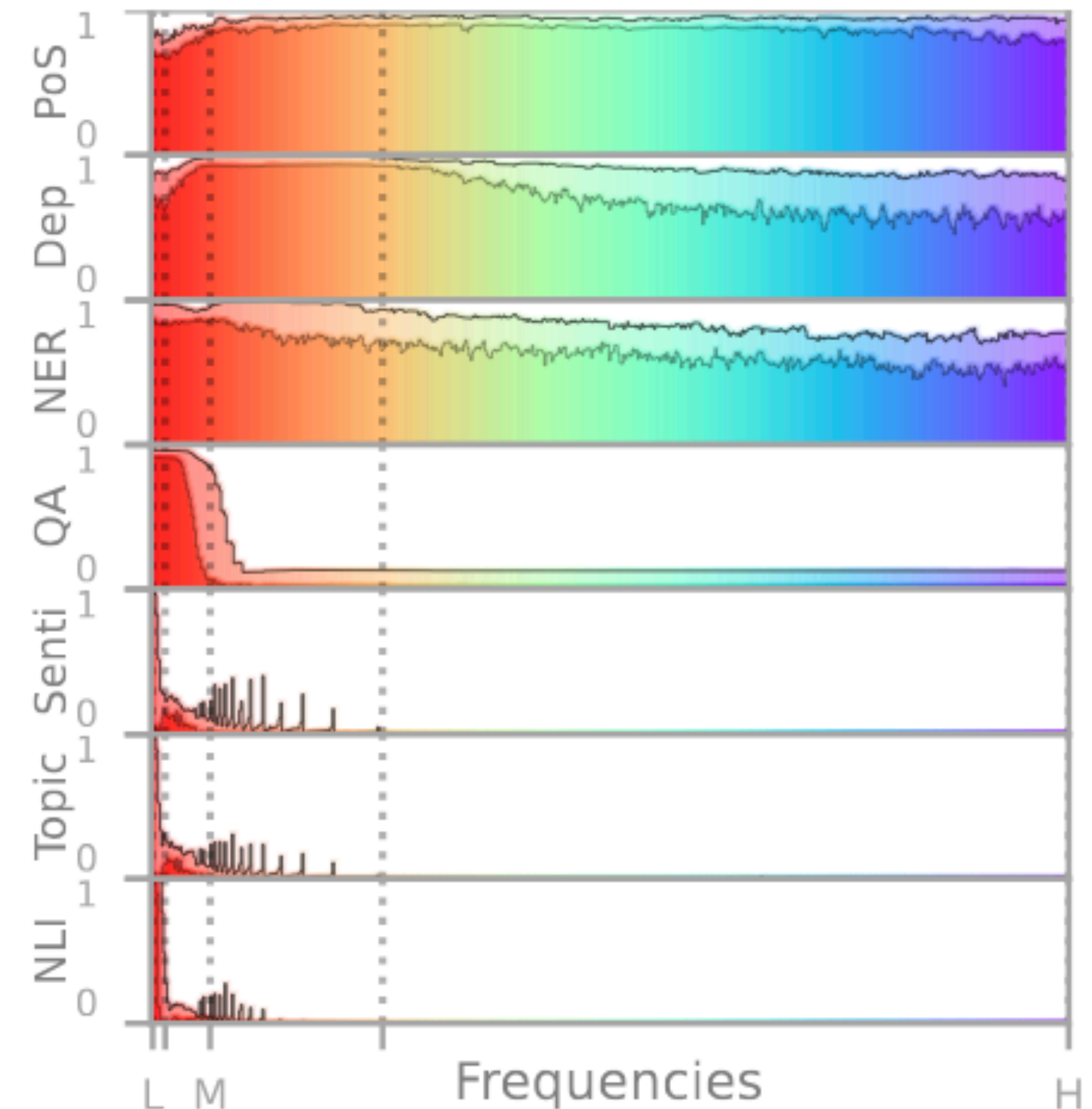  - Much harder with recent LLMs because they are not as open

# BERT

## Some interesting results

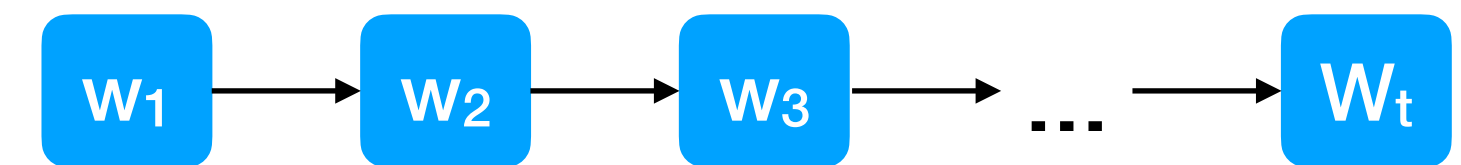- BERT rediscovers the classical NLP pipeline (Tenney et al., 2019)

- Spectral profiles (Müller-Eberstein et al., 2022)
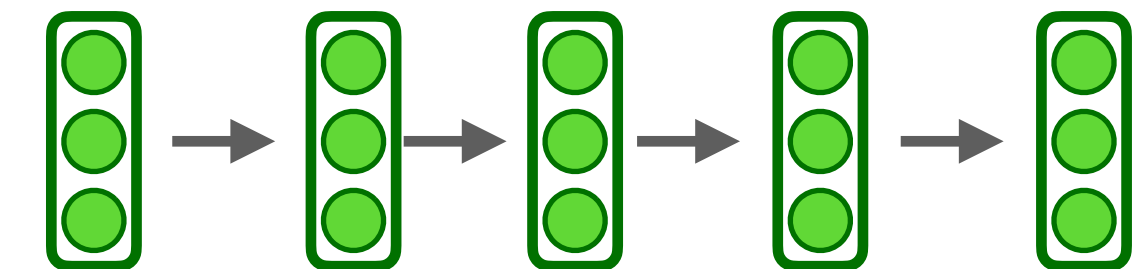
# Today's roadmap

‣ **Part I: Fundamentals**

  ‣ Intro, Motivation & Short History

  ‣ Language Models (n-grams, FFNN-LM, Recap: FFNN)

‣ **Part II: Representations & Beyond FFNN**

  ‣ RNNs (GRU/LSTMs), Attention

  ‣ Contextualised Representations (ELMo)

‣ **Part III: Transformer & LLMs**

  ‣ The Transformer, Masked LMs (BERT), Pre-training & Fine-tuning

  ‣ Prompting, LLMs & Caution

$w_1 \longrightarrow w_2 \longrightarrow w_3 \longrightarrow \ldots \longrightarrow w_t$

# Pre-training & Fine-tuning

# Pre-train / Fine-tune / Transfer Paradigm

‣ 2018-2021: LMs as text encoders (the quest of better word representations)

# In NLP we have a plethora of tasks

‣ Each NLP tasks requires different kinds of data:

  ‣ **Language Modeling:** <u>only text</u>

  ‣ **Machine translation:** <u>naturally occurring</u> parallel (translated) data
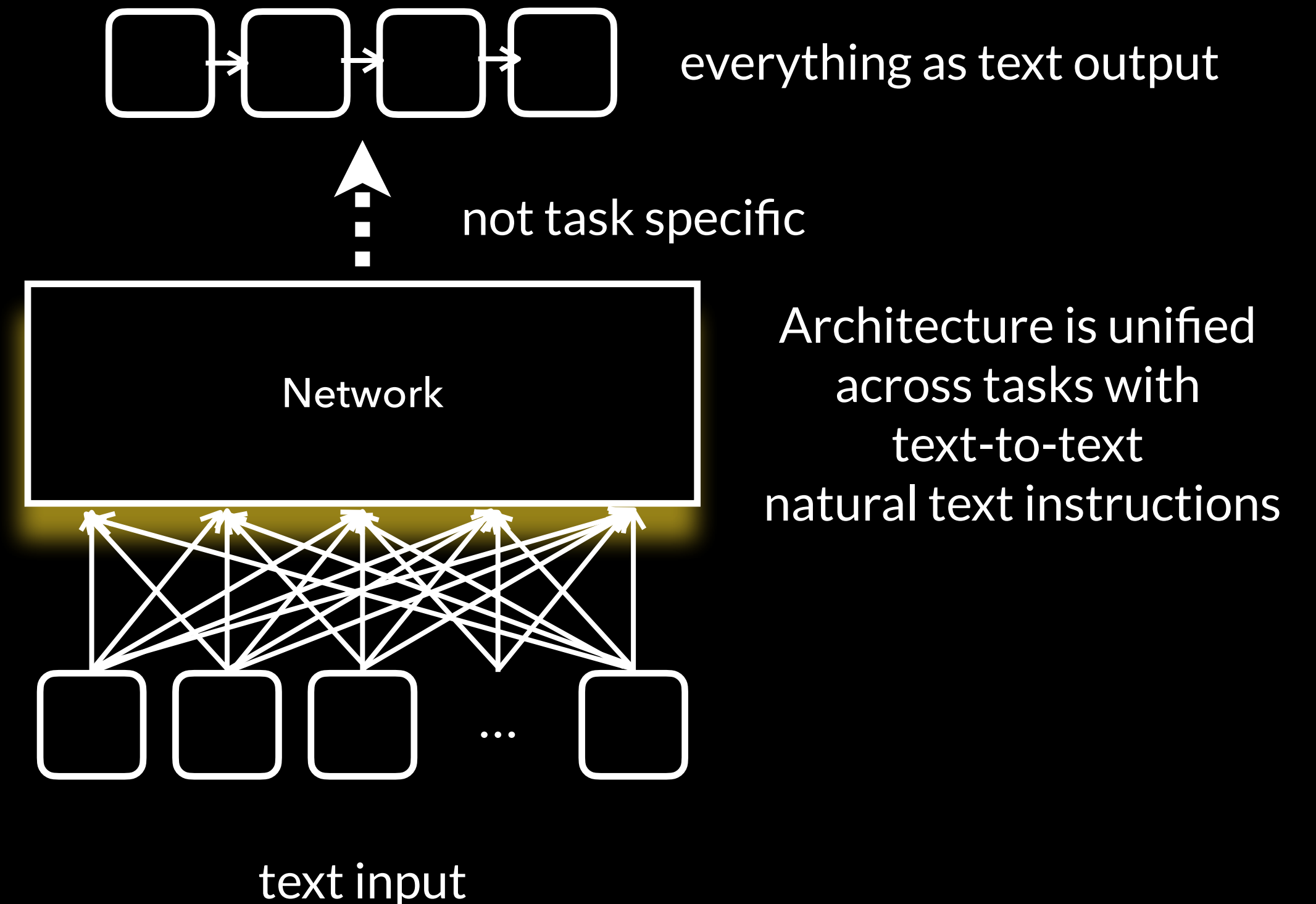
  ‣ **Analysis tasks:** (from classification like sentiment analysis to structured prediction tasks like NER, semantic parsing or slot and intent detection:
    <u>hand-labeled data</u>

# The shift to text-to-text transformers

## (Standard) Multi-task larning:



task label output

"task-specific layers"

Network

text input

## Text-to-text format:



everything as text output

not task specific

Network

Architecture is unified across tasks with text-to-text natural text instructions
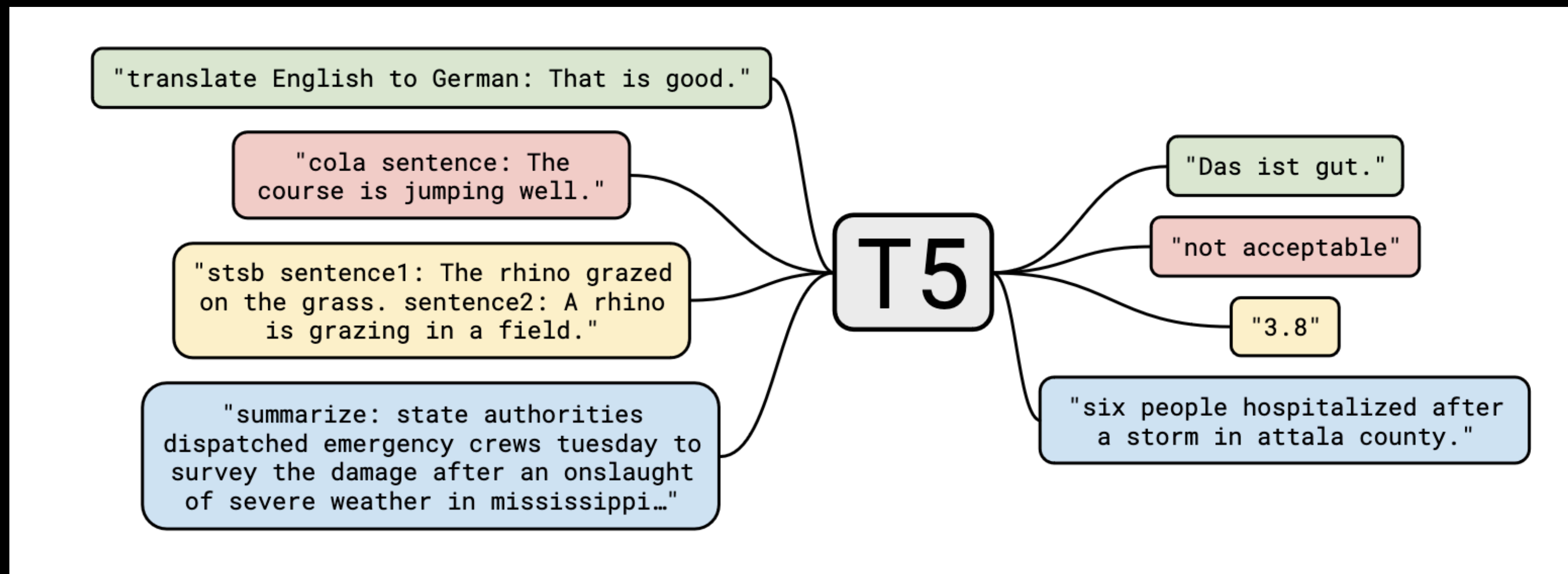
text input

‣ 2019: from task-specific training/architectures to uniform text-to-text formats

‣ General knowledge (pre-training) +instruction fine-tuning => **Generalisation to new tasks**

# Instruction Fine-tuned LMs

‣ Learn many tasks in a single system (e.g. T5 Raffel et al., 2019): examples with instructions

‣ Massive Multi-Task learning (in text-to-text format)

‣ = General knowledge (pre-training) + ability to follow instructions (instruction fine-tuning)



Learn from many tasks

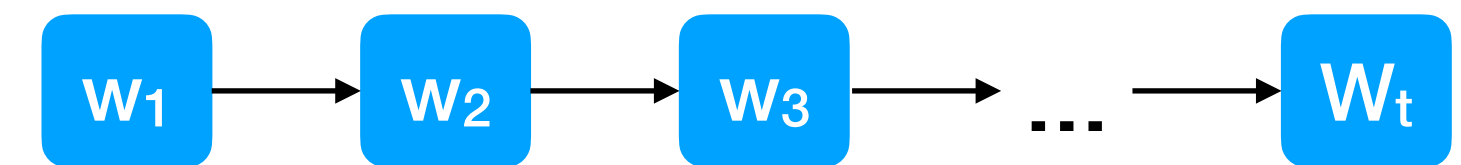https://arxiv.org/pdf/1910.10683

# T5 / encoder-decoder models
**Takeaways**

- T5 (and BART) are very useful for all sorts of sequence-to-sequence tasks with language

  - T5 comes in different sizes

  - There are various customization (e.g., CodeT5)

- Extended the generalizations conclusions from BERT, and demonstrated the impact of data scale

# Today's roadmap

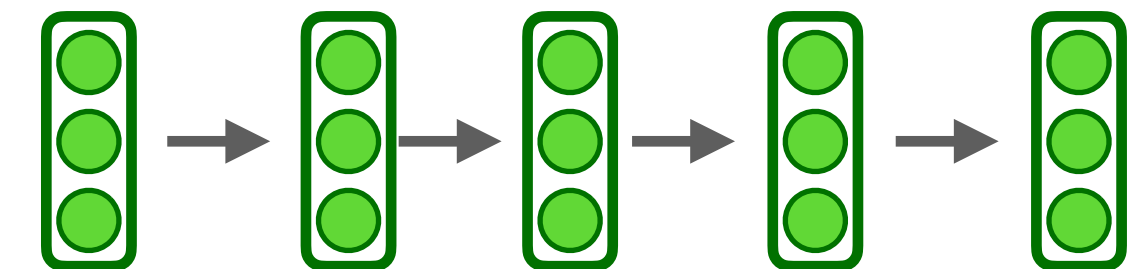‣ **Part I: Fundamentals**

  ‣ Intro, Motivation & Short History

  ‣ Language Models (n-grams, FFNN-LM, Recap: FFNN)

‣ **Part II: Representations & Beyond FFNN**

  ‣ RNNs (GRU/LSTMs), Attention

  ‣ Contextualised Representations (ELMo)

‣ **Part III: Transformer & LLMs**

  ‣ The Transformer, Masked LMs (BERT), Pre-training & Fine-tuning

  ‣ Prompting, LLMs & Caution

$w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow \dots \rightarrow w_t$

# Prompting

# A way to control LLMs: Prompting

- ▸ LLMs offer a completely new mode of operation that **does not require any change to their parameters**: prompting

  - ▸ Prompt = Instruction given to the LLM

- ▸ Adapting the instructions to get different responses:

  - ▸ change the prompt (style)

# Types of Prompting

- ‣ Zero-shot prompting (no examples)

- ‣ One-shot, Few-shot prompting (with one or few annotated examples)

- ‣ Chain-of-thought (CoT) Prompting and extensions thereof (asking for intermediate reasoning steps)

Please rate the following move review from 1 to 5 (1=did not like it, 5=liked it very much).

Example:
Review: I was intrigued by the final scene despite thinking the opposite at the start.
Rating: 4

Review: Despite the terrible soundtrack, Smith delivered a thoughtful agent.

# Zero-shot Prompting

- Input: single unlabeled example $\bar{x}$

- Output: the label $\bar{y}$

- The task (and output) can be any text-to-text task: classification, summarization, translation

- Pre-processing: wrap $\bar{x}$ with a template using a **_verbalizer_** $v$

- The template controls the output

$\bar{x}$ = the movie's acting could've been better, but the visuals and directing were top-notch.

⬇

$v(\bar{x})$ = Review: the movie's acting could've been better, but the visuals and directing were top-notch.
Out of positive, negative, or neutral this review is

⬇

**LLM**

⬇

neutral $\quad \bar{y}$

210

# In-context Learning (ICL)

- LLMs have the ability to "learn" to complete tasks through training in the prompt

- The recipe is simple:

  - Take a small number of annotated training example $\{(\bar{x}^{(i)}, \bar{y}^{(i)}\}_{i=1}^{N}$

  - Convert them using verbalizer $v$ templates

  - Concatenate them and follow with the target input $\bar{x}$

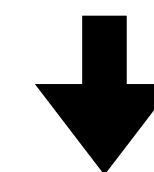  - The completion will be the label of the input

$\bar{x} =$ the movie's acting could've been better, but the visuals and directing were top-notch.

$\Downarrow$

Review: The cinematography was stellar; great movie!
Sentiment (positive or negative): positive
Review: The plot was boring and the visuals were subpar.
Sentiment (positive or negative): negative
Review: The movie's acting could've been better, but the visuals and directing were top-notch.
Sentiment (positive or negative):

$\Downarrow$

**LLM**

$\Downarrow$

positive $\quad \bar{y}$

# In-context Learning (ICL)

## Performance

- Providing ICL examples almost always leads to significant improvements



Aggregate Performance Across Benchmarks

**Figure 1.3: Aggregate performance for all 42 accuracy-denominated benchmarks** While zero-shot performance improves steadily with model size, few-shot performance increases more rapidly, demonstrating that larger models are more proficient at in-context learning. See Figure 3.8 for a more detailed analysis on SuperGLUE, a standard NLP benchmark suite.

# In-context Learning (ICL)

## Performance

- Providing ICL examples almost always leads to significant improvements

- Benefits tend to diminish with more examples



Figure 32: **Number of in-context examples.** For each model, we set the maximum number of in-context examples to [0, 1, 2, 4, 8, 16] and fit as many in-context examples as possible within the context window. We plot performance as a function of the average number of in-context examples actually used.

# Chain-of-Thought (CoT) Prompting - Example

‣ Adding **"Let's think step by step"**

Sure, let's break this down step by step to come up with some compelling tandem talk ideas:

**Step 1: Identify the Audience and Purpose**

Understanding who will be listening to the talk and what you aim to achieve is crucial. Is the audience academic, corporate, students, or a general public? Is the purpose to inform, entertain, persuade, or inspire?

**Step 2: Choose a Broad Topic**

# CoT Prompting Example

▸ Standard Prompt

what is the size of the big lake closest to Verona

The largest lake closest to Verona is Lake Garda (Lago di Garda). Lake Garda is the largest lake in Italy, covering an area of approximately 370 square kilometers (143 square miles). It is situated about 30 kilometers (19 miles) west of Verona, making it easily accessible for a day trip or a short visit from the city.

▸ CoT ("Reasoning steps'):

what is the size of the big lake closest to Verona. Let's think step by step

To determine the size of the largest lake closest to Verona, we need to follow these steps:

1. **Identify the major lakes near Verona**:

   - Lake Garda
   - Lake Iseo
   - Lake Como

2. **Determine the distances from Verona to these lakes**:

   - **Lake Garda**: About 30 km west of Verona
   - **Lake Iseo**: About 90 km northwest of Verona
   - **Lake Como**: About 130 km northwest of Verona

3. **Compare the sizes of these lakes**:

   - **Lake Garda**: 370 square kilometers
   - **Lake Iseo**: 65.3 square kilometers
   - **Lake Como**: 146 square kilometers

# CoT Prompting Example 2

## Standard Prompting

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ❌

## Chain-of-Thought Prompting

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✓

# Chain-of-thought (COT) Prompting

## Step-by-step

- Main idea: just "tell" the model to reason in steps

- Can significantly outperform zero-shot prompting with very large models



(a) MultiArith on Original GPT-3    (b) MultiArith on Instruct GPT-3    (c) GMS8K on PaLM

[Kojima et al. 2022]

# Prompting LLMs - caution!

# Model Behaviour: Trustworthiness in the Era of LLMs

‣ From a "compartmentalised" notion of language tasks in NLP

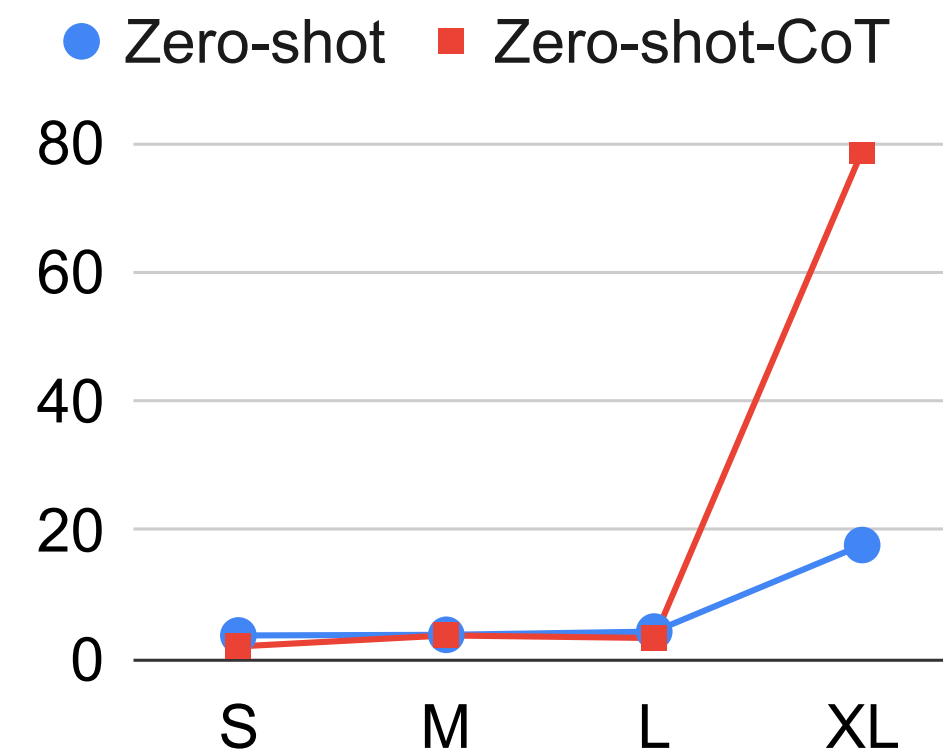‣ To powerful but less interpretable black-box models

‣ What is a "language task" today?

Compartmentalised NLP



Current Trend (in NLP)



Litschko*, Müller-Eberstein*, van der Goot, Weber-Genzel, Plank. Establishing Trustworthiness: Rethinking Tasks and Model Evaluation. EMNLP 2023.

# Model Behaviour: Does it Matter How we Prompt an LLM?

‣ ⚠️ Performance is highly sensitive to the linguistic variation of a prompt

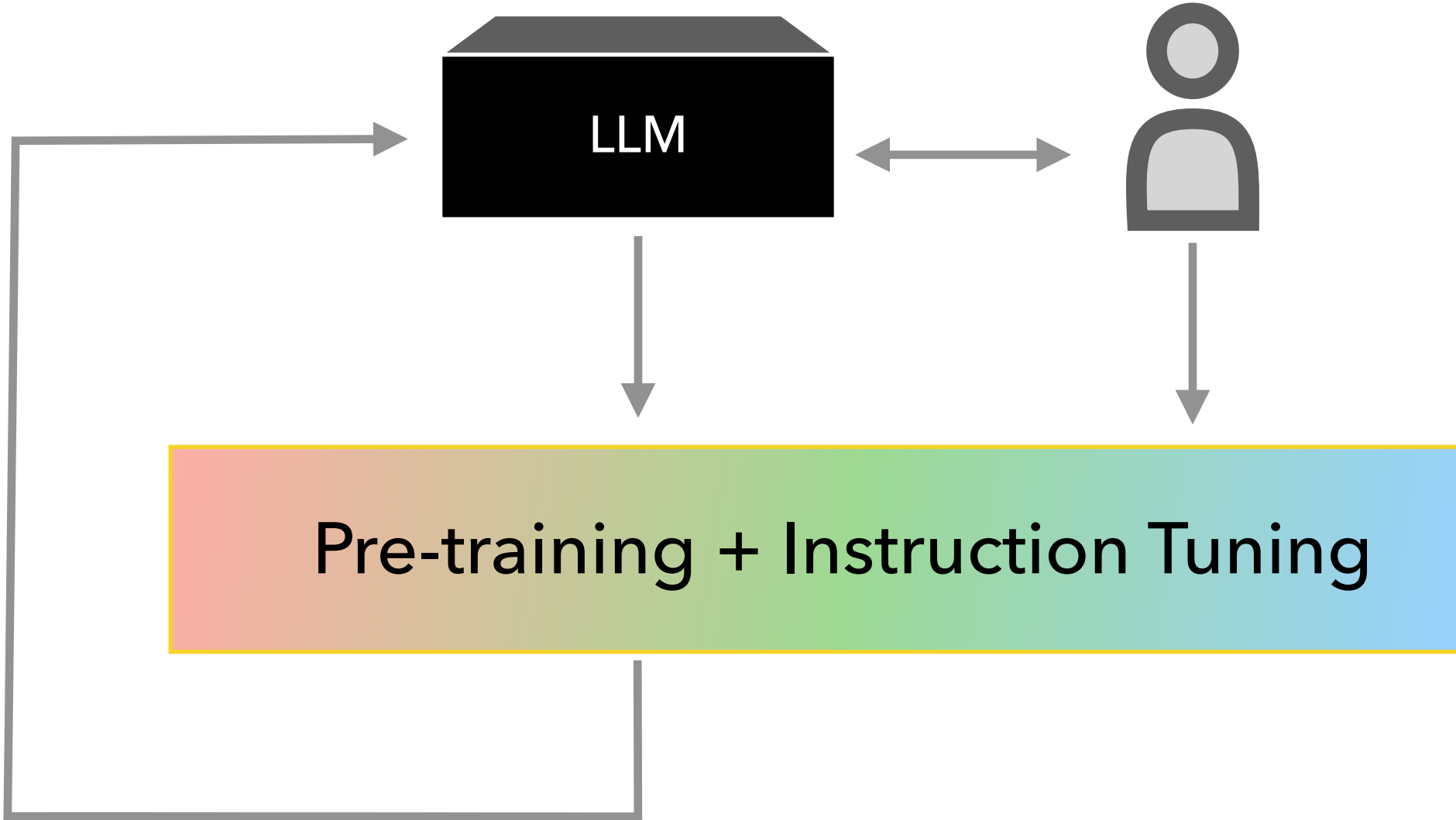| | prop. | prompt |
|---|---|---|
| **mood** | inter. | Do you find this movie review positive? |
| | indic. | You find this movie review positive. |
| | imper. | Tell me if you find this movie review positive. |
| **aspt.** | active | Do you find this movie review positive? |
| | pass. | Is this movie review found positive? |
| **tense** | past | Did you find this movie review positive? |
| | pres. | Do you find this movie review positive? |
| | future | Will you find this movie review positive? |
| **modality** | can | Can you find this movie review positive? |
| | could | Could you find this movie review positive? |
| | may | May you find this movie review positive? |
| | might | Might you find this movie review positive? |
| | must | Must you find this movie review positive? |
| | should | Should you find this movie review positive? |
| | would | Would you find this movie review positive? |
| **synonymy** | apprai. | Do you find this movie appraisal positive? |
| | comm. | Do you find this movie commentary positive? |
| | criti. | Do you find this movie critique positive? |
| | eval. | Do you find this movie evaluation positive? |
| | review | Do you find this movie review positive? |

Table 1: Examples of variation of linguistic properties

**The language of prompting:**
**What linguistic properties make a prompt successful?**

Leidinger, van Rooij, Shutova, EMNLP 2023 Findings.

Köksal et al., EMNLP 2023 Findings ; Gonen et al., EMNLP 2023 Findings. 220

# Multiple-Choice Question Answering (MCQA) Prompt Style

**General Instruction:** Please read the multiple-choice question below carefully and select ONE of the listed options and only give a single letter.

**Question:** The Web was effectively invented by Berners-Lee in which year?
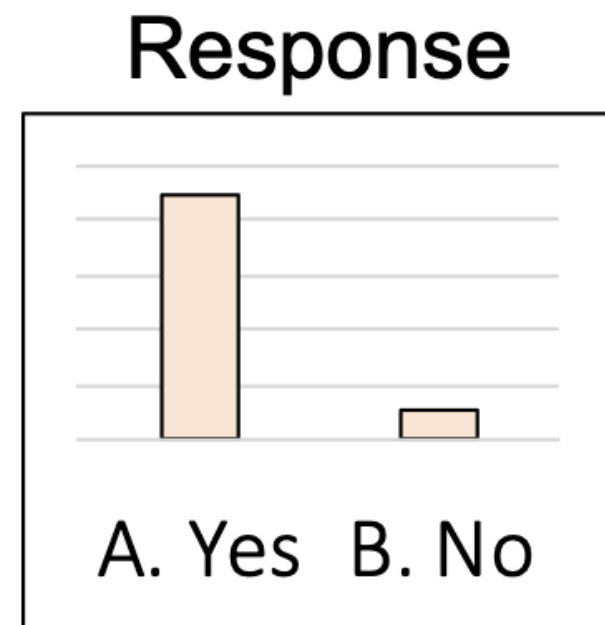
**Options:**
A. 1991
B. 1980
C. 1989
D. 1993

**Answer:**

Wang, Hu, Ma, Röttger, Plank. Look at the Text: Instruction-Tuned Language Models are More Robust Multiple Choice Selectors than You Think. COLM 2024.

# Evaluation Protocols: Do Answer Options Impact LLM Outputs?

▸ ⚠️ LLM's "A"-bias in MCQA responses

### Choice ordering 1

> Question: In the past 12 months, has this person given birth to any children?
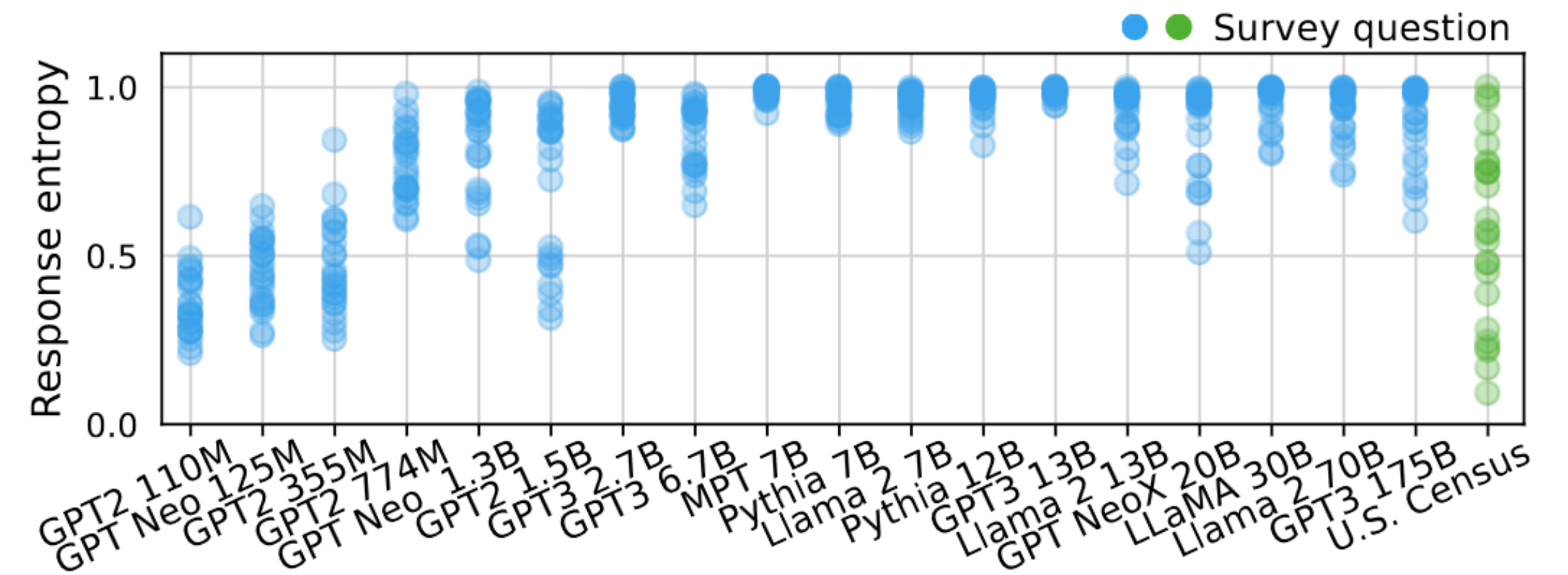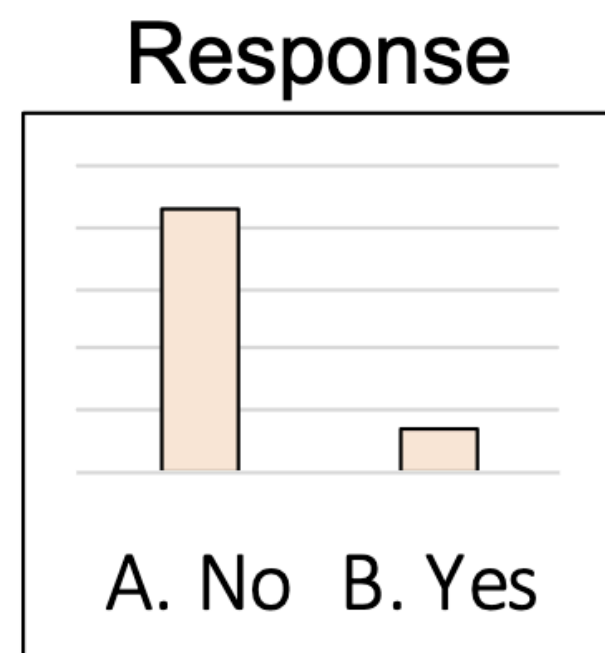> A. Yes
> B. No
> Answer:

| P("A") | 0.82 | P("B") | 0.11 |
|--------|------|--------|------|

**Response**



A. Yes   B. No

### Choice ordering 2

> Question: In the past 12 months, has this person given birth to any children?
> A. No
> B. Yes
> Answer:

| P("A") | 0.80 | P("B") | 0.15 |
|--------|------|--------|------|

**Response**



A. No   B. Yes



(a) Entropy of base models' responses.



(b) A-bias of base models' responses.

Dominguez-Olmedo, Hardt, Mendler-Dünner. Questioning the Survey Responses of Large Language Models. arXiv:2306.07951 2023.

# Evaluation Protocols: Does It Matter How We Extract Answers?

▸ ⚠️ But "First-token log probs" do not match the text answers

| Model (0-shot) | First Token | Text Answer |
|---|---|---|
| Gemma-7b-Inst | 30.2 | 50.8 |
| Llama2-7b-Chat | 34.9 | 43.1 |
| Llama2-13b-Chat | 40.2 | 47.6 |
| Mistral-7b-Inst-0.2 | 53.2 | 53.6 |

Performance on MMLU.

**General Instruction:** Please read the multiple-choice question below carefully and select ONE of the listed options and only give a single letter.
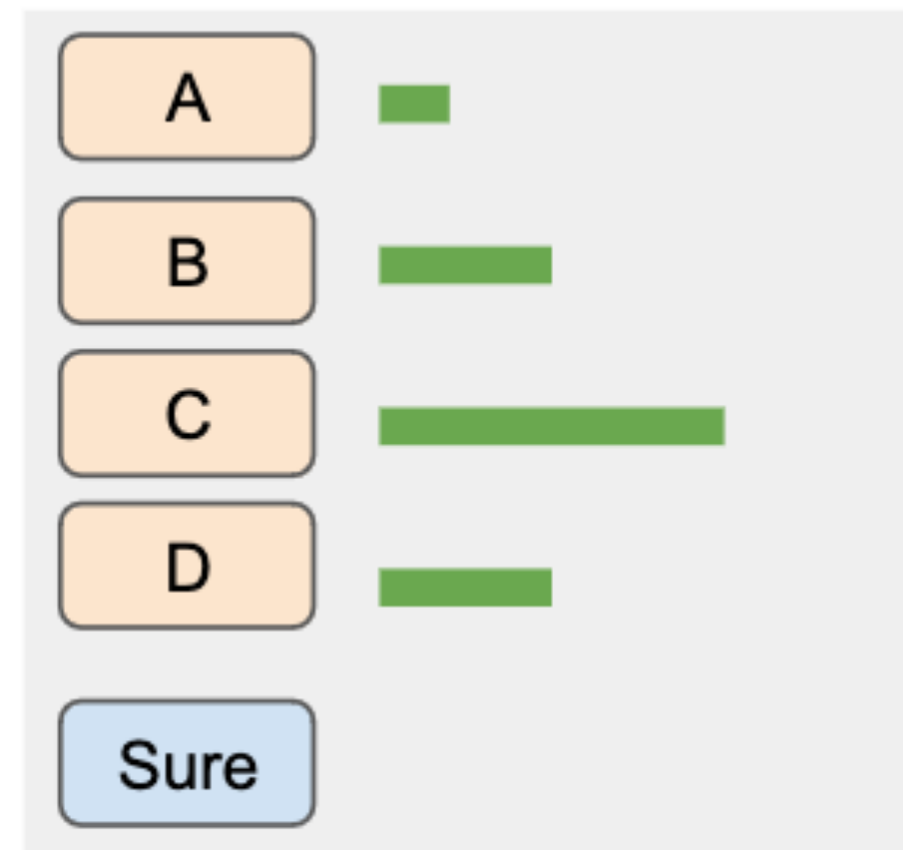
**Question:** The Web was effectively invented by Berners-Lee in which year?

**Options:**
A. 1991
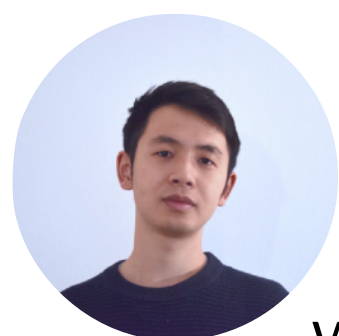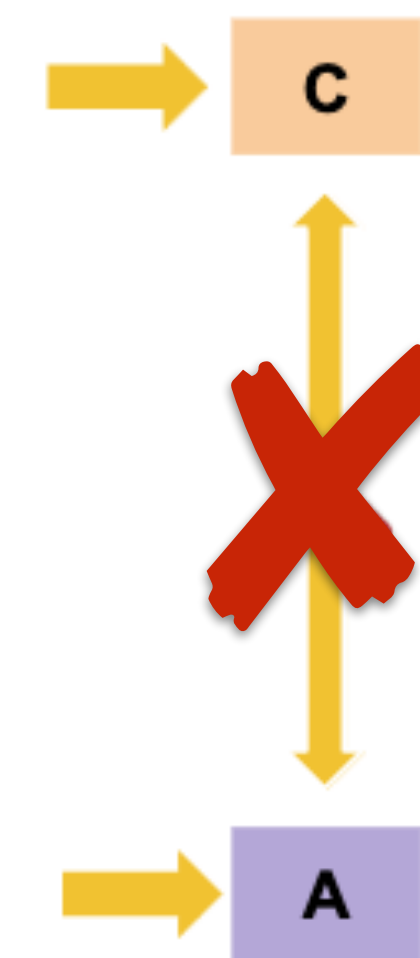B. 1980
C. 1989
D. 1993

**Answer:**

a. First Token Logits:

A
B
C
D
Sure

versus

b. Text Answer:

Sure! The answer is (A) 1991.

C

A

Wang, Ma, Hu, Weber-Genzel, Röttger, Kreuter, Hovy, Plank. "My Answer is C": First-Token Probabilities Do Not Match Text Answers in Instruction-Tuned Language Models. ACL 2024 Findings.
Wang, Hu, Ma, Röttger, Plank. Look at the Text: Instruction-Tuned Language Models are More Robust Multiple Choice Selectors than You Think. COLM 2024.

# Evaluation Protocols: Chain-of-Thought, Better *and* Sound?

‣ ⚠️ Chain-of-Thought Reasoning is Not Sound

| Model | Correct Answer | Sound Reasoning |
|---|---|---|
| Zephyr-7B-$\beta$ | $45.0 \pm 15.5$ | $25.0 \pm 10.5$ |
| Mistral-7B-Instruct | $55.0 \pm 10.0$ | $25.0 \pm 7.5$ |
| LLaMA-2-7B | $46.7 \pm 6.7$ | $0.0 \pm 0.0$ |
| LLaMA-2-13B | $40.0 \pm 8.2$ | $15.0 \pm 6.2$ |
| LLaMA-2-70B | $56.7 \pm 6.2$ | $31.7 \pm 9.7$ |

Statements:
1. There is a purple marble in the box if and only if there is a blue marble in the box.
2. There is a blue marble in the box if and only if there is an olive marble in the box.
3. Either there is an olive marble in the box or else there is a red marble in the box, but not both.
4. Either there is a red marble in the box or else there is a green marble in the box, but not both.

Conclusion: If there is a purple marble in the box then there is a green marble in the box.

**CoT response (excerpt) - manually verified:**

7. Let's consider the possibility that there is an olive marble in the box. If that's the case, then, according to Statement 2, there must be a blue marble in the box as well.
8. Now, using Statement 1, we know that if there is a blue marble in the box, then there must be a purple marble in the box as well.
9. However, Statement 3 tells us that if there is an olive marble in the box, then there cannot be a red marble in the box. But, Statement 4 tells us that if there is a red marble in the box, then there cannot be a green marble in the box.

Supposition following (sound).

10. Since we cannot have both an olive marble and a red marble in the box, and we cannot have both a red marble and a green marble in the box, it seems that the only possibility left is that there is an olive marble in the box, and therefore, a blue marble, a purple marble, but no red or green marbles.

Invalid inference. ❌

11. However, the conclusion states that "If there is a purple marble in the box then there is a green marble in the box", which is not logically true based on the given statements.
Conclusion: False.

Conclusion (invalid).

Final answer (incorrect).

Mondorf, Plank. Comparing Inferential Strategies of Humans and Large Language Models in Deductive Reasoning. ACL 2024.
Stechly*, Valmeekam*, Kambhampati. Chain of Thoughtlessness? An Analysis of CoT in Planning. arXiv:2405.04776 2024.

# Evaluation Protocols: Can LLMs Replace Humans Judges?
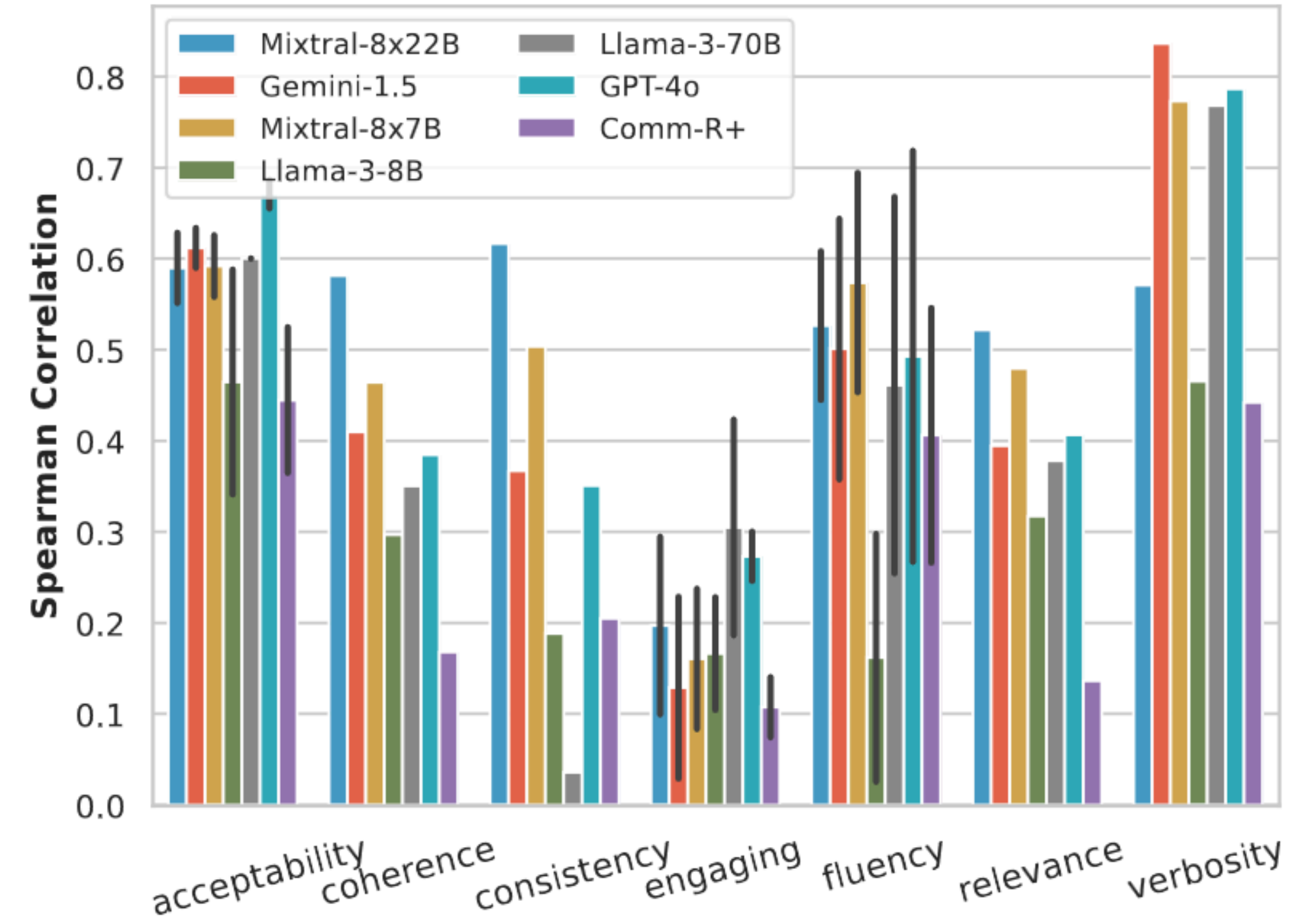
‣ ⚠️ A lot of variability in LLM outputs

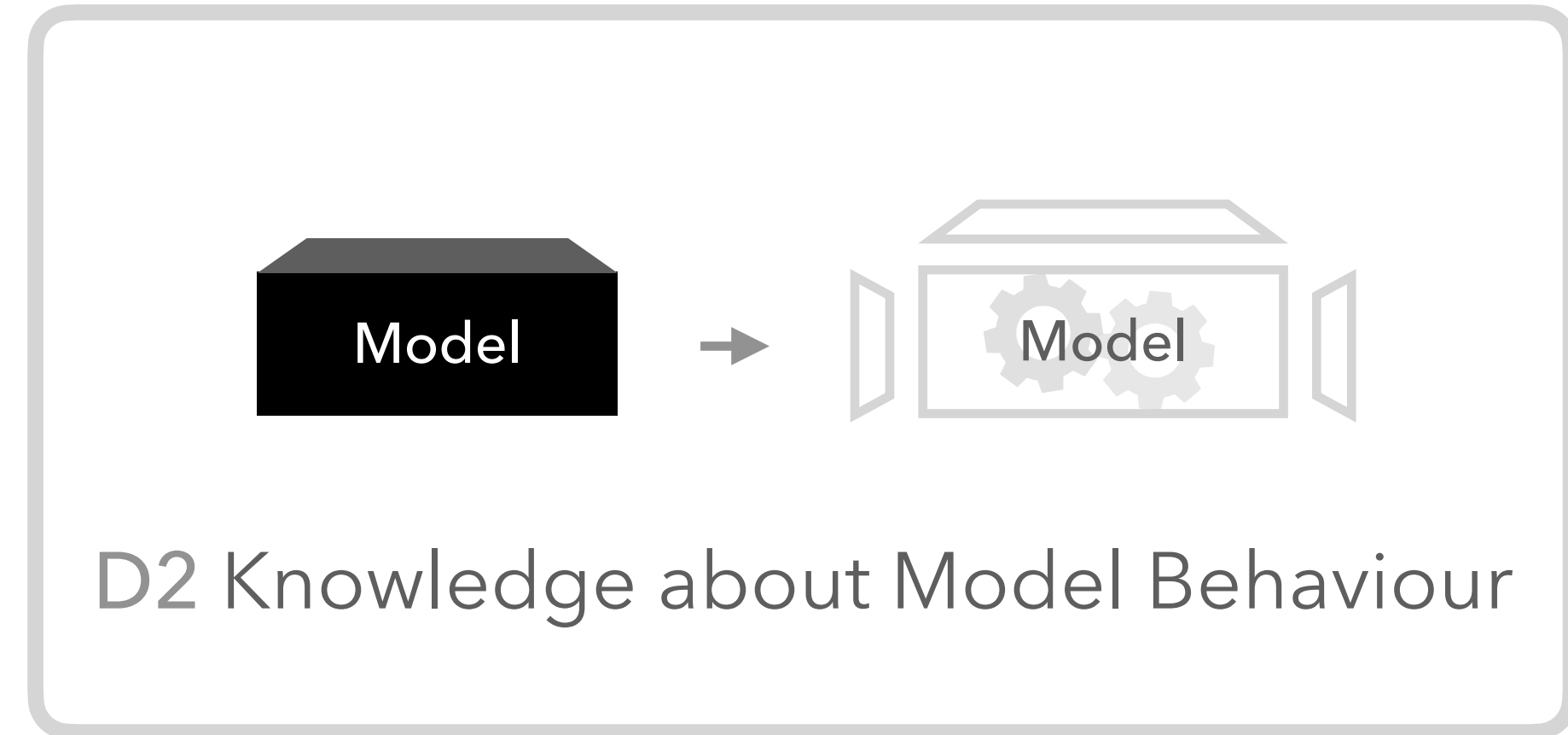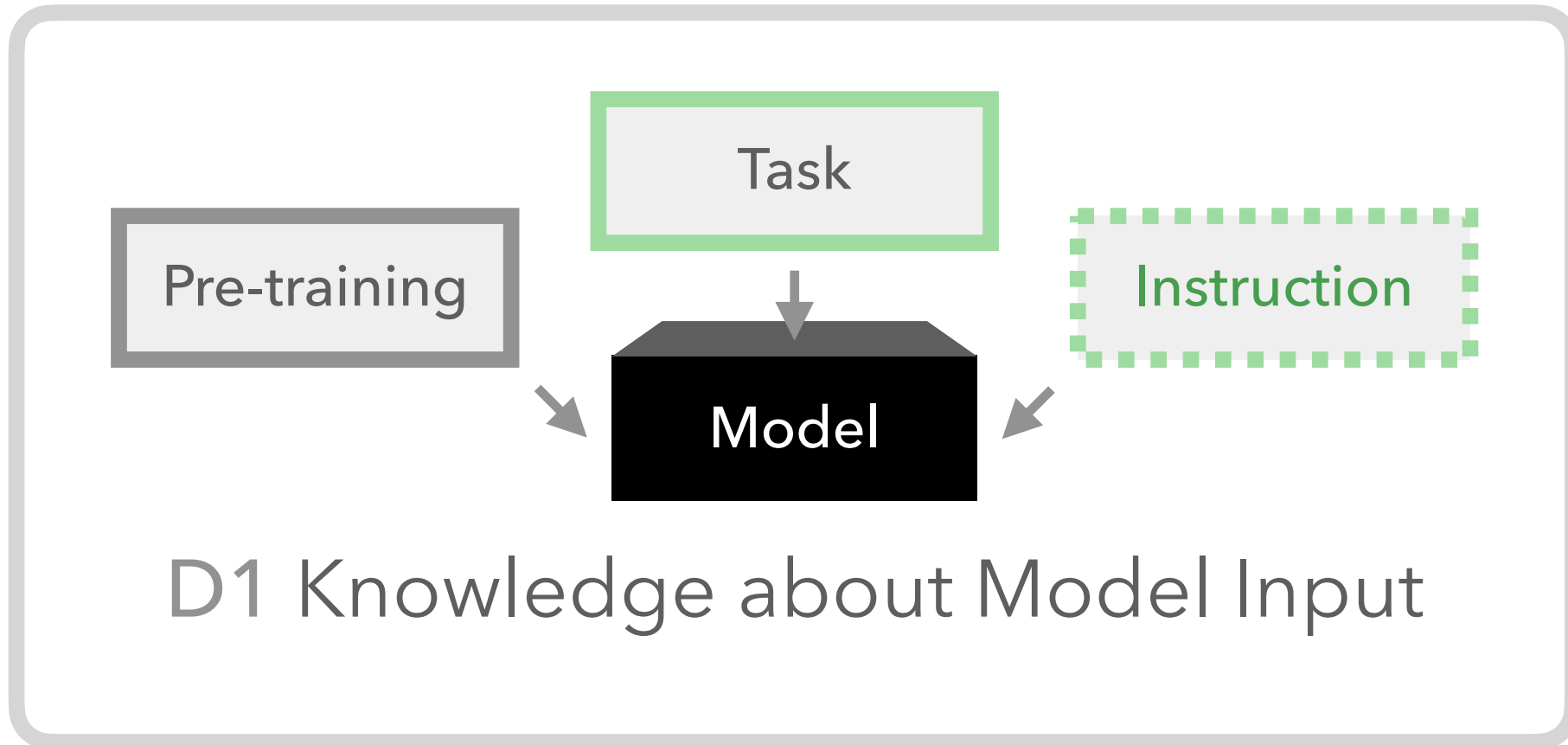‣ LLMs are not ready yet to replace human judges - not even GPT-4o:

E.g. Plausibility:

Humans vs Models:



*Instruction*: On a scale of 1 (very unlikely) to 5 (very likely), how plausible is it that the last response belongs to the dialogue?

**A**: Made it all the way through four years of college playing ball but
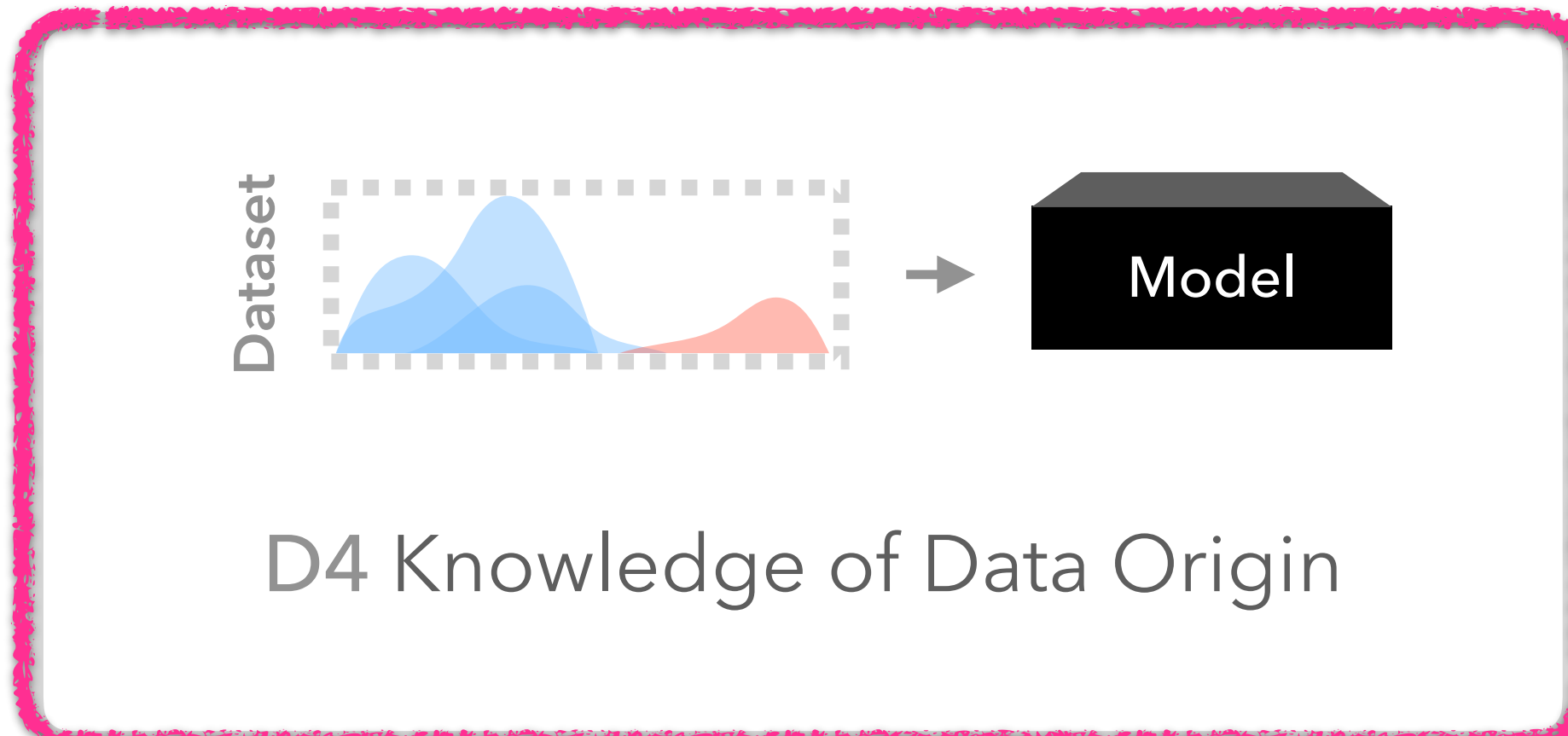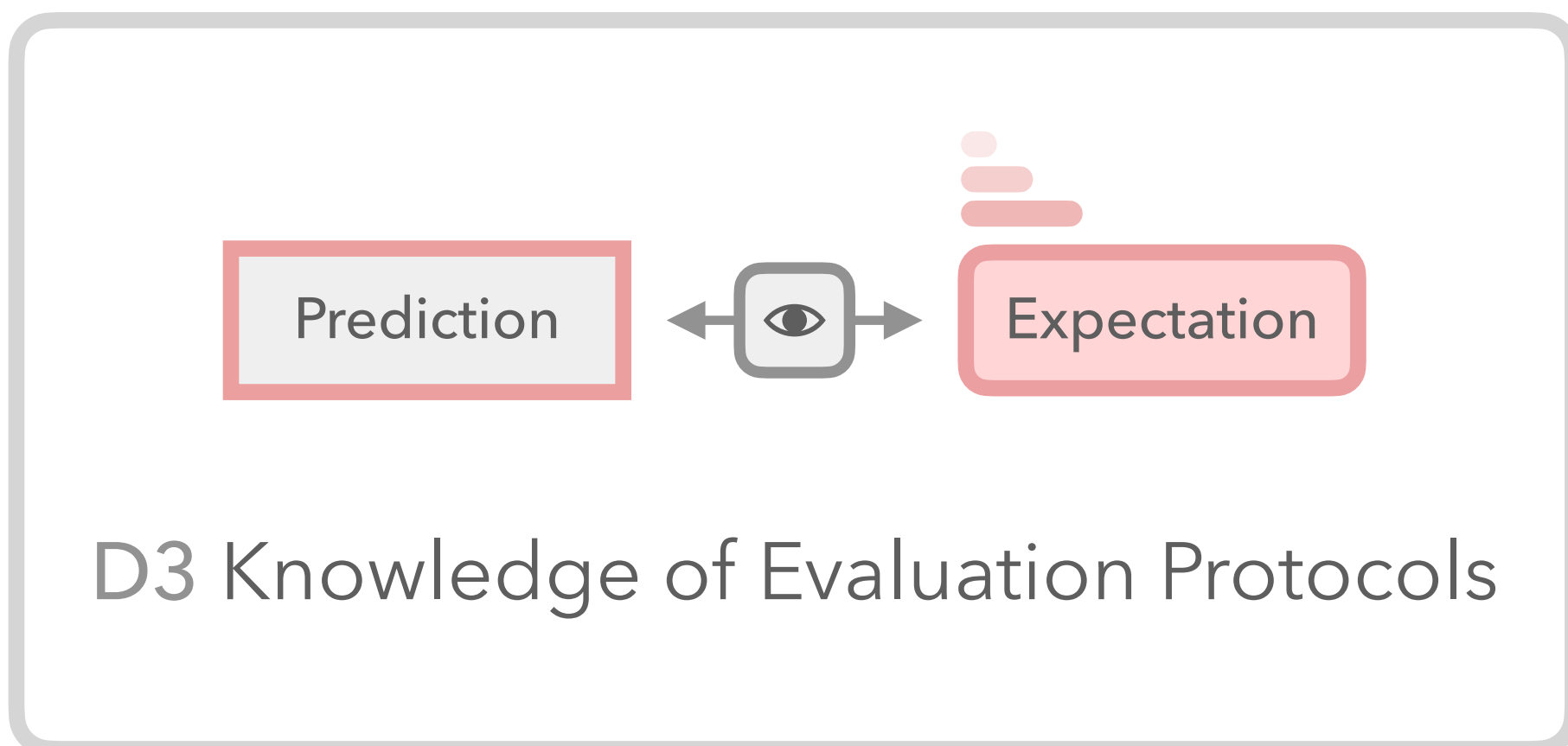**B**: I also like The Cosby Show

Bavaresco, Bernardi, Bertolazzi, Elliott, Fernandez, Gatt, Ghaleb, Giulianelli, Hanna, Koller, Martins, Mondorf, Neplenbroek, Pezzelle, Plank, Schlangen, Suglia, Surikuchi, Takmaz, Testoni. LLMs instead of Human Judges? A Large Scale Empirical Study across 20 NLP Evaluation Tasks. arXiv:2406.18403 2024.

D1 Knowledge about Model Input

D2 Knowledge about Model Behaviour

Trust arises from **knowledge of origin** as well as from **knowledge of functional capacity.**

*Trustworthiness - Working Definition by David G. Hays, 1979*

D3 Knowledge of Evaluation Protocols
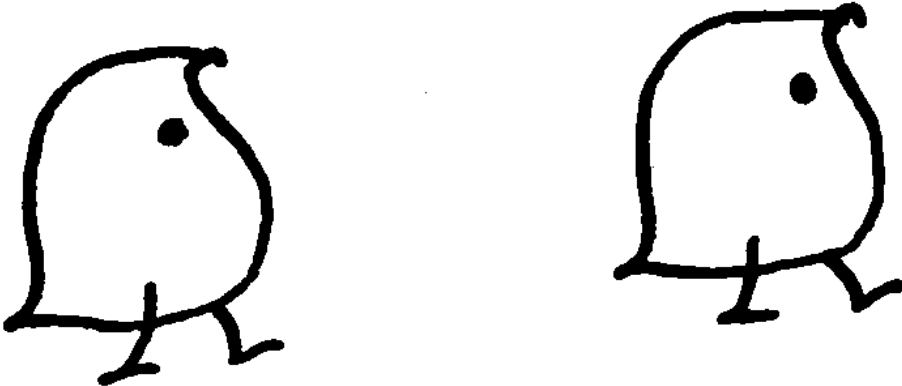
D4 Knowledge of Data Origin

226

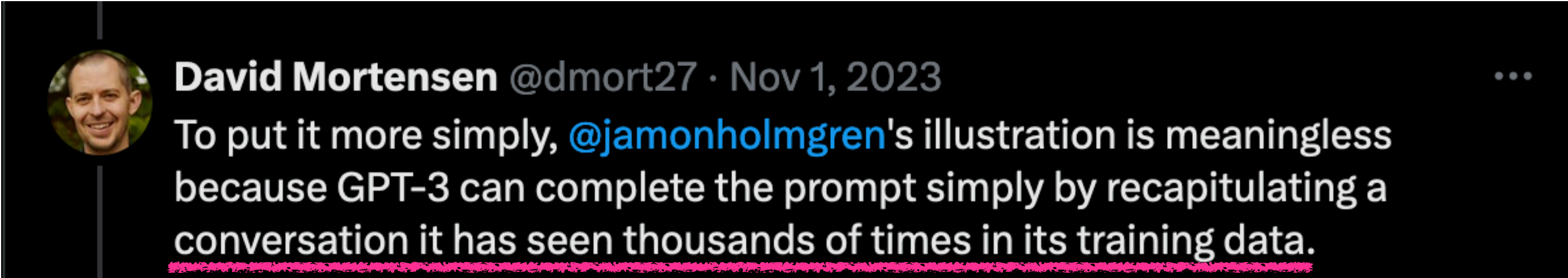# Data Origin: Late 2022 Claim "ChatGPT Passes the Wug Test"
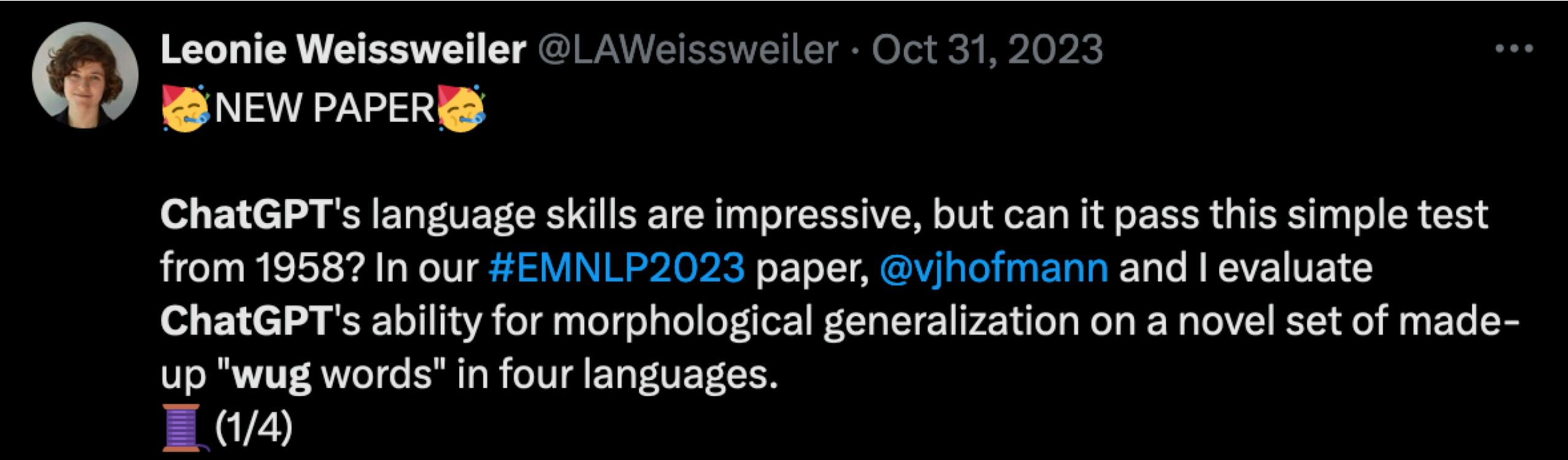
‣ ⚠️ No! Data contamination. Tested only on English.

THIS IS A WUG.

NOW THERE IS ANOTHER ONE.

THERE ARE TWO OF THEM.

THERE ARE TWO _____.

> **David Mortensen** @dmort27 · Nov 1, 2023
> To put it more simply, @jamonholmgren's illustration is meaningless because GPT-3 can complete the prompt simply by recapitulating a conversation it has seen thousands of times in its training data.

> **Leonie Weissweiler** @LAWeissweiler · Oct 31, 2023
> 🥹 NEW PAPER 🥹
>
> **ChatGPT**'s language skills are impressive, but can it pass this simple test from 1958? In our #EMNLP2023 paper, @vjhofmann and I evaluate **ChatGPT**'s ability for morphological generalization on a novel set of made-up "**wug** words" in four languages.
> 🧵 (1/4)

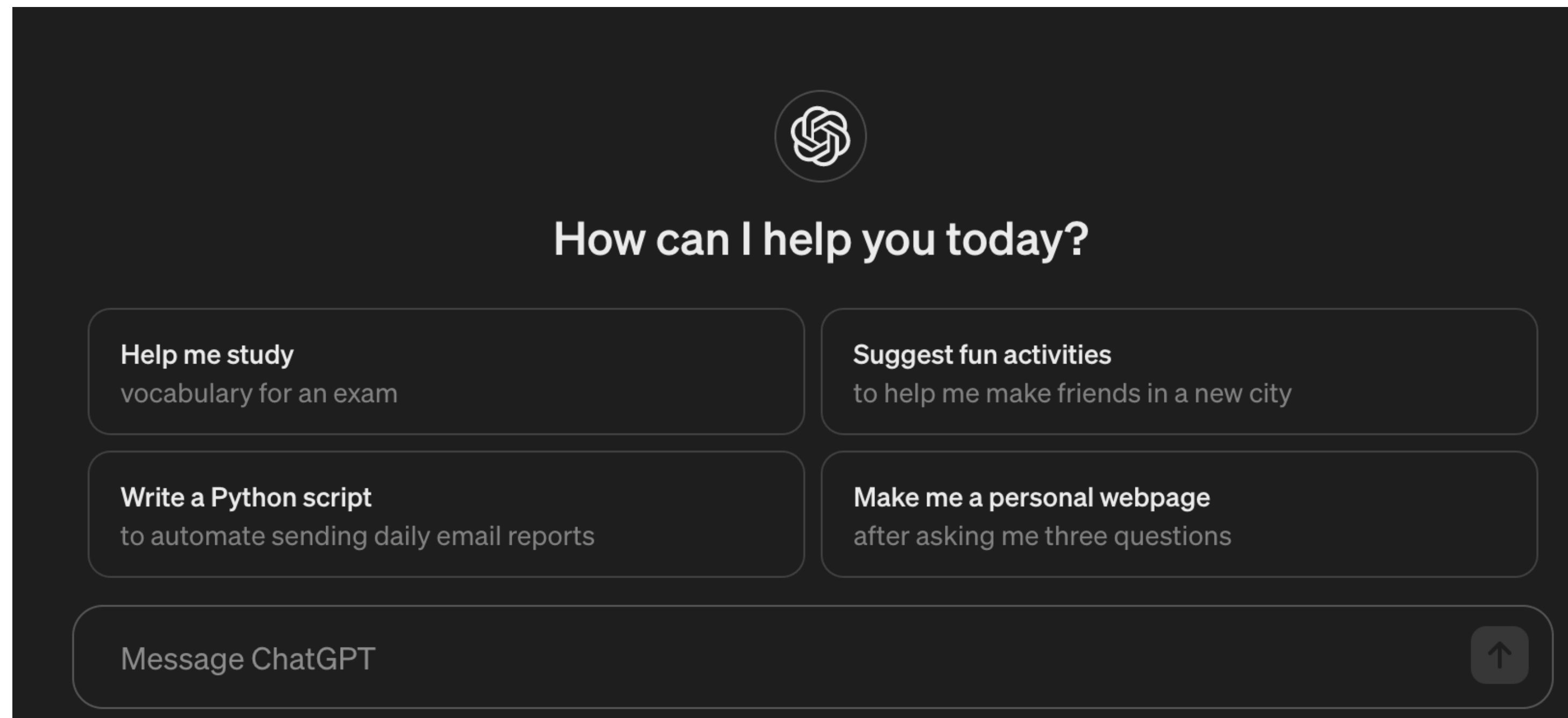### Counting the Bugs in ChatGPT's Wugs: A Multilingual Investigation into the Morphological Capabilities of a Large Language Model

Weißweiler*, Hofmann*, Kantharuban, Mai, Dutt, Henkel, Kabra, Kulkarni, Vijayakumar, Yu, Schütze, Oflazer, Mortensen. EMNLP 2023.

# Data Origin: Indirect Data Leakage

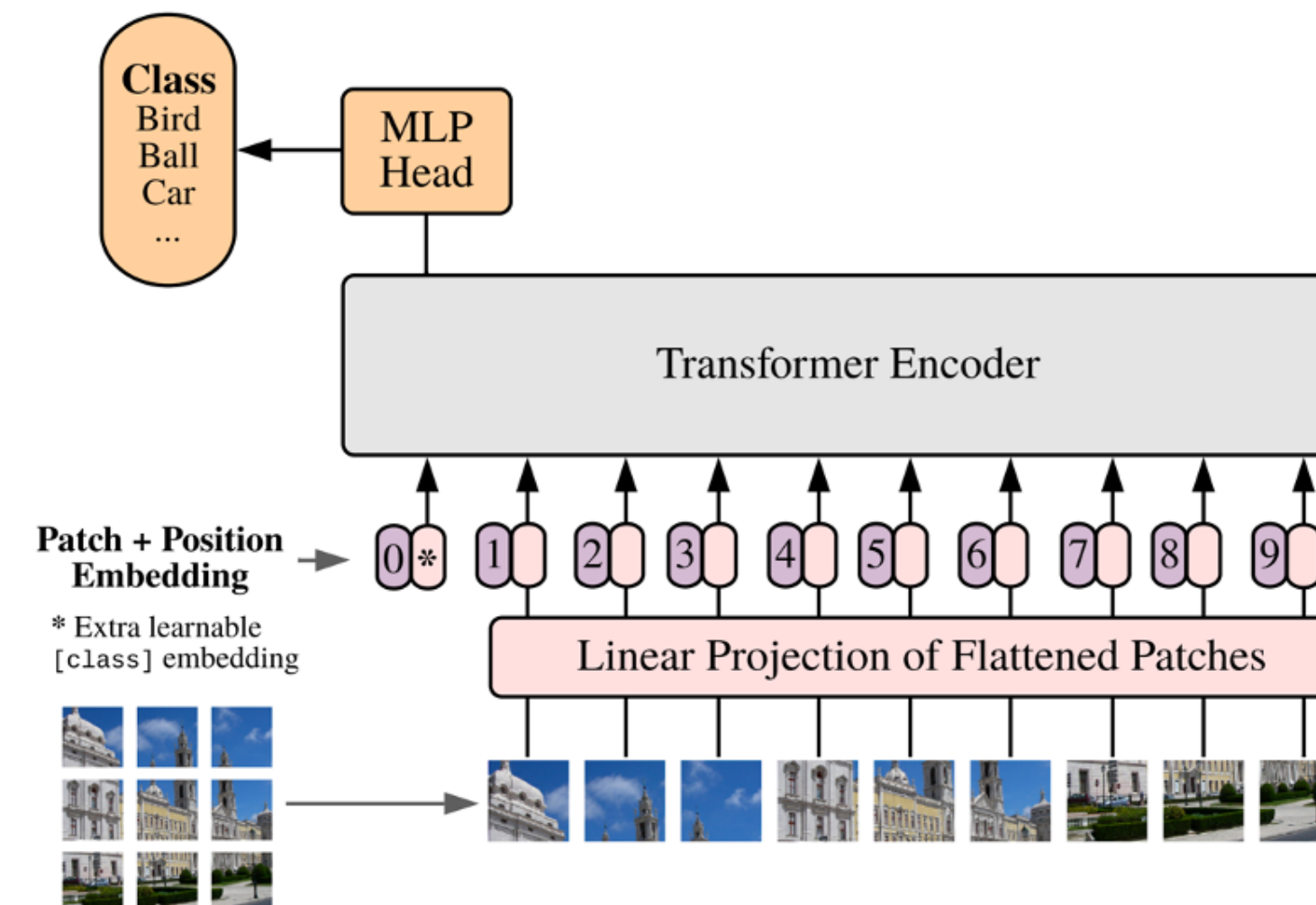‣ ⚠️ Data continuously provided by users (e.g. via OpenAI's the web interface)



🥬 **Leak, Cheat, Repeat: Data Contamination and Evaluation Malpractices in Closed-Source LLMs**

Balloccu, Schmidtová, Lango, Dušek. EACL 2024.

# Exciting Opportunities Ahead

# Transformers
## Computer Vision

- ViT: cut image to patches

- Project each patch to a vector

- Treat them as token embeddings

[Dosovitskiy et al. 2020]
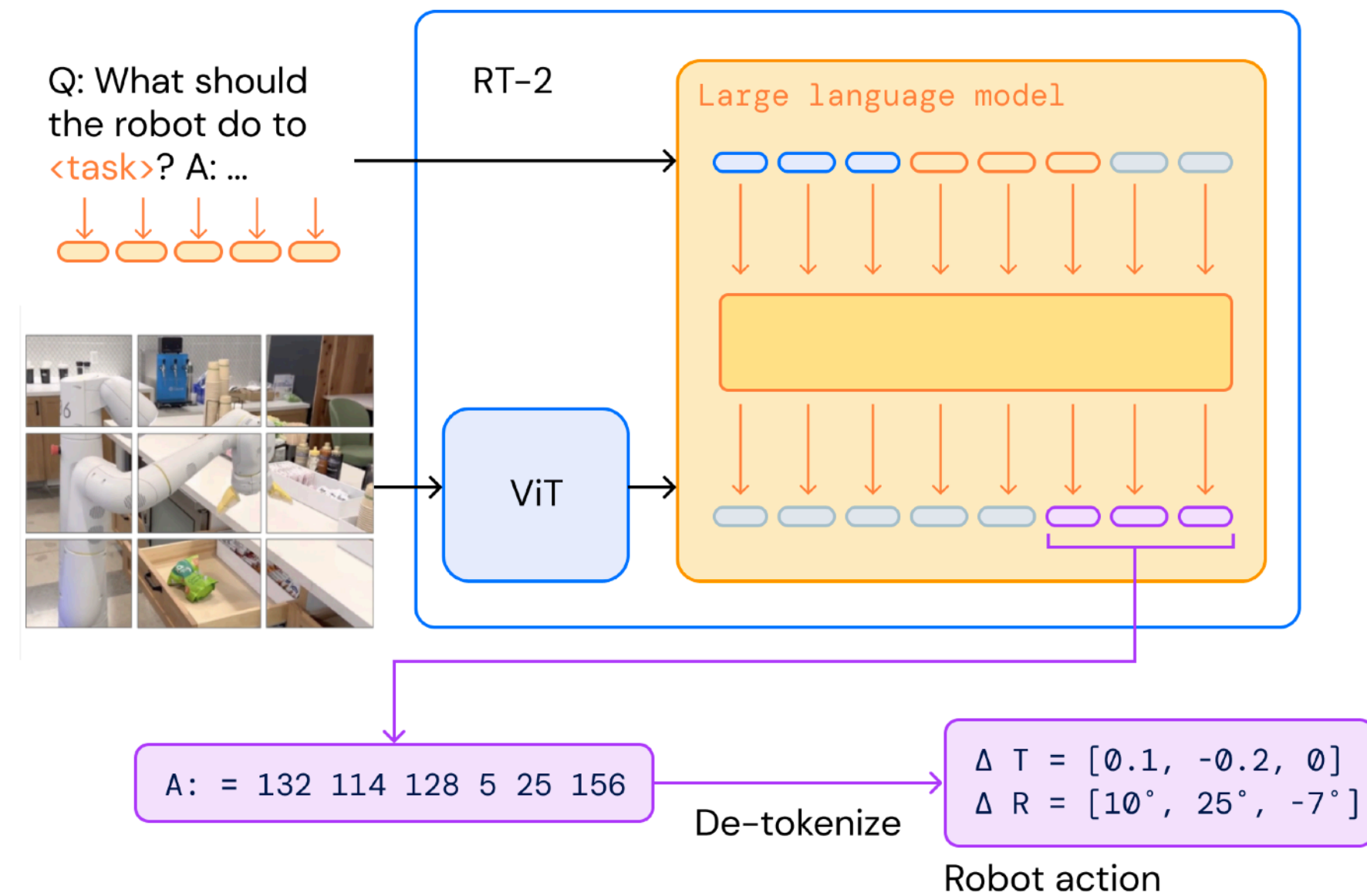
# Transformers
## Speech

- Same as computer vision

- But: spectrograms instead of images

- The Whisper model

[Radford et al. 2022]

# Transformers
## Robotics

- Take observations and commands, all tokenized

- Output continuous joint control actions
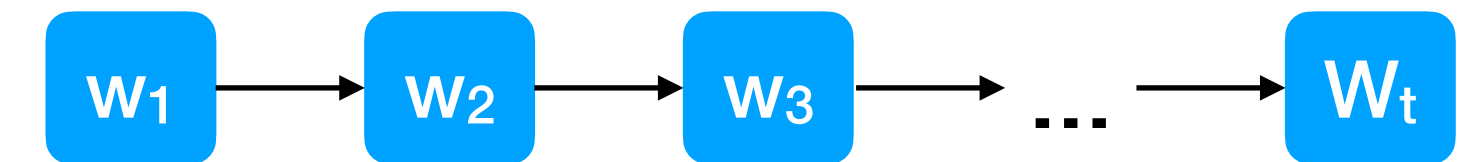
Slide by Yoav Artzi.

[Brohan et al. 2023]

# Today's Lecture
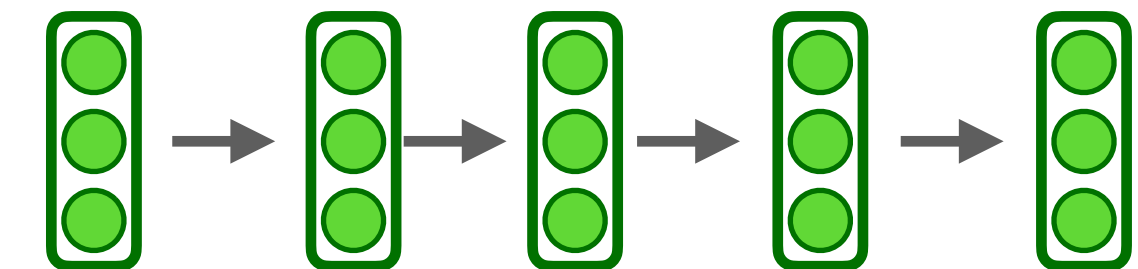
▸ **Part I: Fundamentals**

  ‣ Intro, Motivation & Short History

  ‣ Language Models (n-grams, FFNN-LM, Recap: FFNN)

▸ **Part II: Representations & Beyond FFNN**

  ‣ RNNs (GRU/LSTMs), Attention

  ‣ Contextualised Representations (ELMo)

▸ **Part III: Transformer & LLMs**

  ‣ The Transformer, Masked LMs (BERT), Pre-training & Fine-tuning

  ‣ Prompting, LLMs & Caution

$w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow \ldots \rightarrow w_t$

# Questions?
# Thank you

Barbara Plank
@barbara_plank
LMU Munich

Thanks to all the
organizers & sponsors of:

bplank.github.io

# Core References

‣ Jurasky & Martin <u>SLP textbook</u> chapter 3 (n-gram LMs), 7 (neural LMs) & 8 (RNNs)

‣ Arianna Bisazza's <u>AthNLP 2019 lecture</u> on MT/Transformers

‣ Graham Neubig's <u>Advanced NLP class</u>

‣ Yoav Artzi's <u>LM-class</u> (also based on Greg Durrett's class material)

‣ Yoav Goldberg (2015): <u>A Primer on Neural Network Models for Natural Language Processing</u>

‣ Chris Manning & Abigail See (2018) Stanford NLP class

‣ Lilian Weng's <u>attention tutorial</u>