



# Intermediate Deep Learning for Engineers

---

Spring 2025, Deep Learning for Engineers  
March 27, 2025, 6th Session

Amir Barati Farimani  
*Associate Professor of Mechanical Engineering and Bio-Engineering*  
*Carnegie Mellon University*

# Sequence Modeling

$$y^4 = P(\mathbf{x}^{(5)} | \text{the students opened their})$$

## A RNN Language Model

output distribution

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{U}\mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden states

$$\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

$\mathbf{h}^{(0)}$  is the initial hidden state

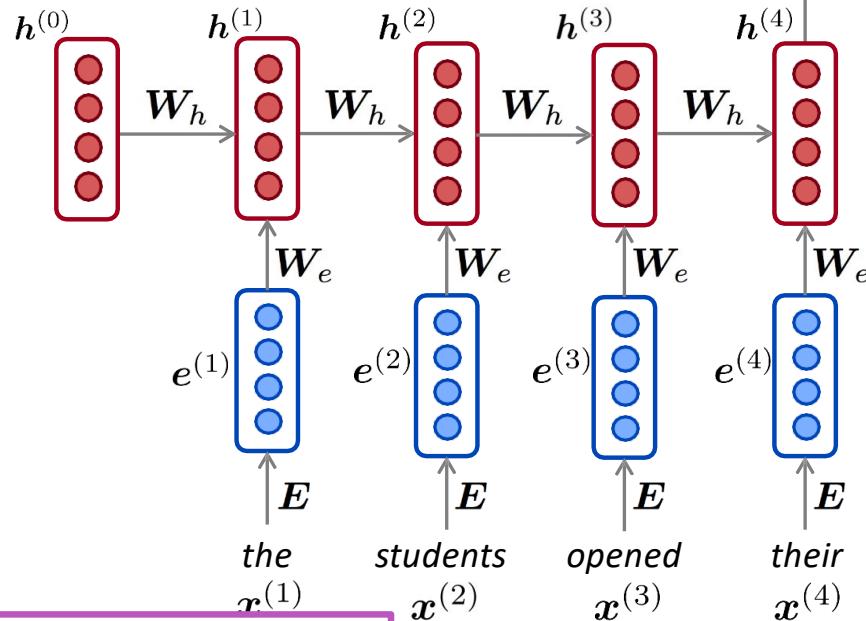
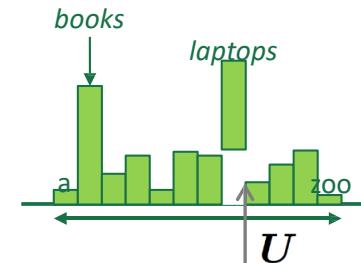
word embeddings

$$\mathbf{e}^{(t)} = \mathbf{E} \mathbf{x}^{(t)}$$

words / one-hot vectors

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$

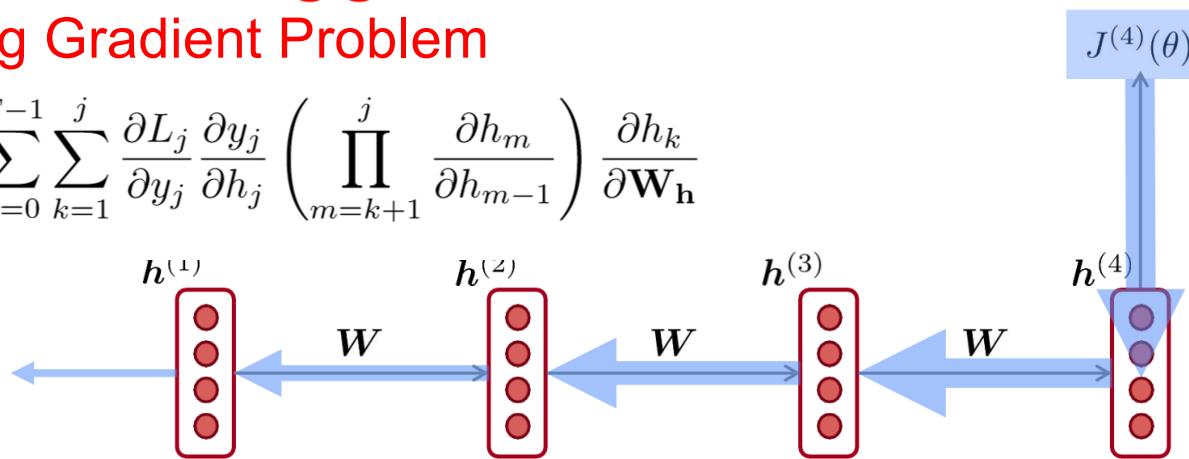
*Note: this input sequence could be much longer, but this slide doesn't have space!*



# Sequence Modeling

## Vanishing Gradient Problem

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{j=0}^{T-1} \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left( \prod_{m=k+1}^j \frac{\partial h_m}{\partial h_{m-1}} \right) \frac{\partial h_k}{\partial \mathbf{W}_h}$$



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \frac{\partial h^{(4)}}{\partial h^{(3)}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

What happens if these are small?

Vanishing gradient problem:  
When these are small, the  
gradient signal gets smaller  
and smaller as it  
backpropagates further

# Sequence Modeling

## Challenges with RNNs

- Long range dependencies
- Gradient vanishing and explosion
- Large # of training steps
- Recurrence prevents parallel computation
- Slow computation for long sequences
- Vanishing or exploding gradients
- Difficulty in accessing information from long time ago



# Attention

## Attention in Computer Vision

- 2014: Attention used to highlight important parts of an image that contribute to a desired output
- Attention in NLP
- 2015: Aligned machine translation
- 2017: Language modeling with **Transformer networks**



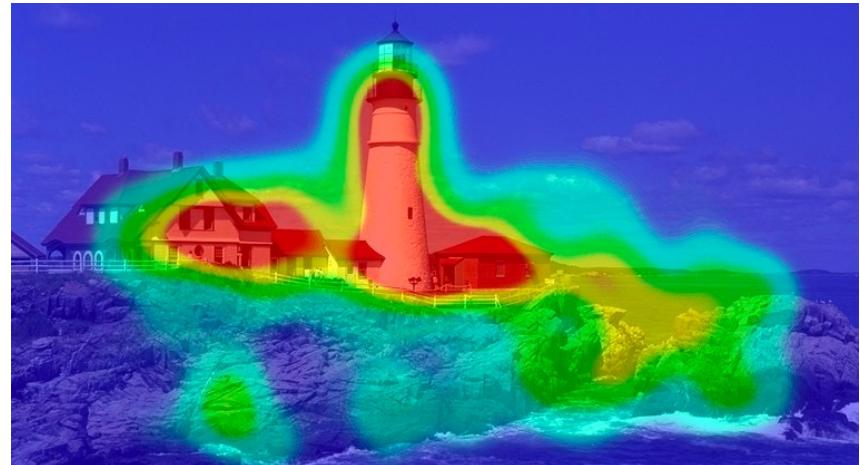
# Why Attention

Attention is a recent and important component to the success of modern neural networks

We want neural nets that **automatically weigh** relevance of the input and **use these weights** to perform a task

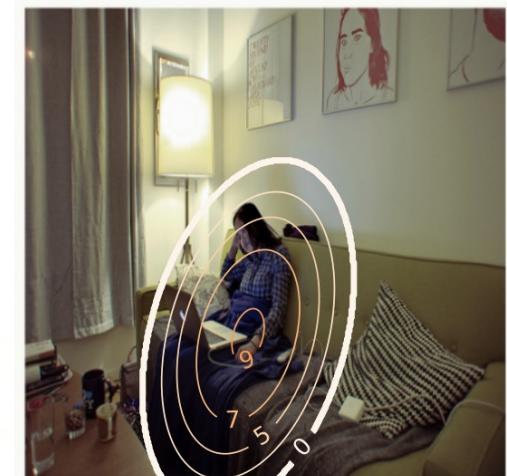
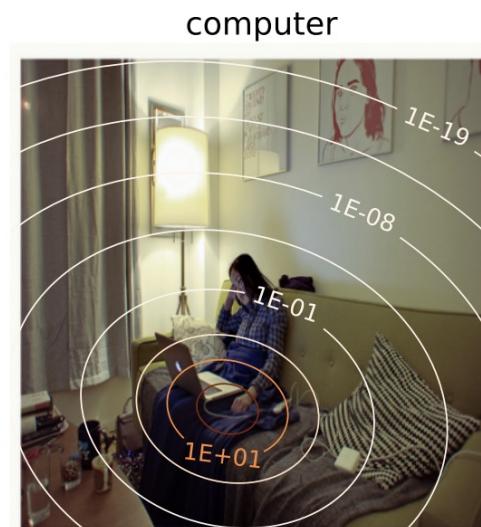
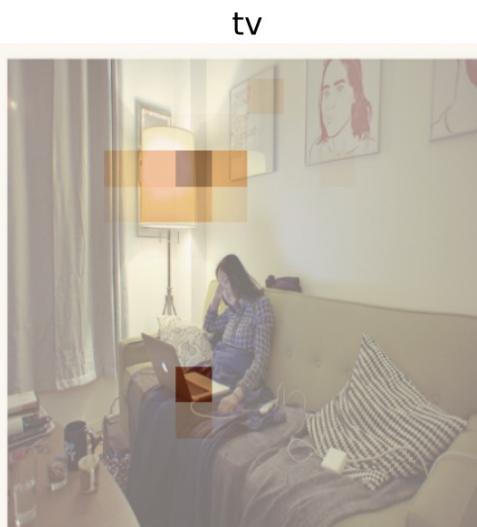
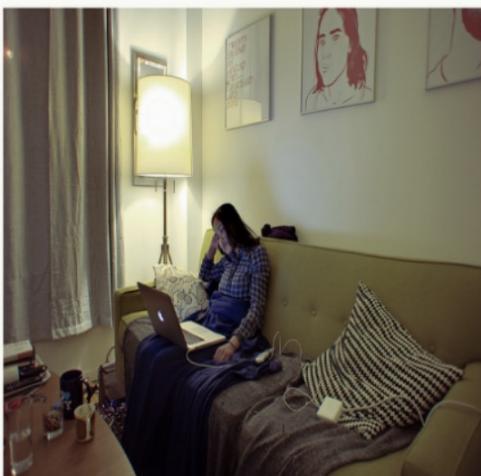
## Main advantages:

- performance gain
- none or few parameters
- fast (easy to parallelize)
- drop-in implementation
- tool for "interpreting" predictions



# Attention Example

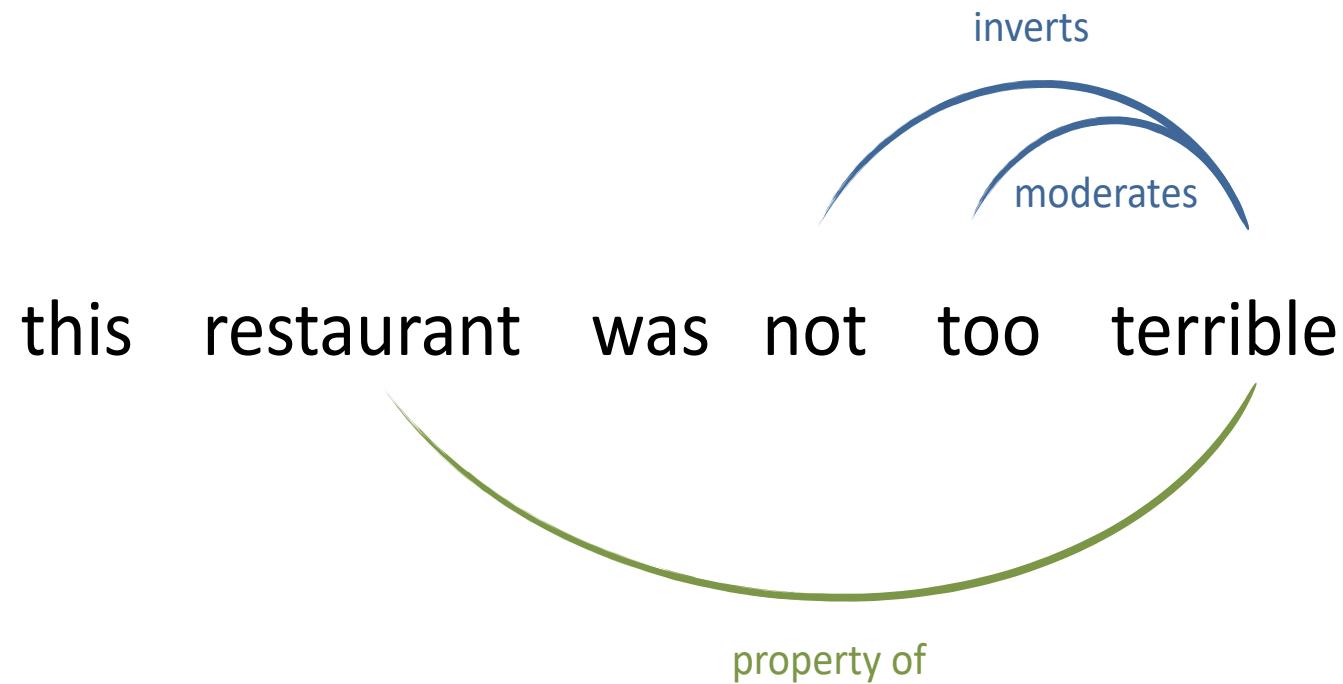
What is the woman looking at?



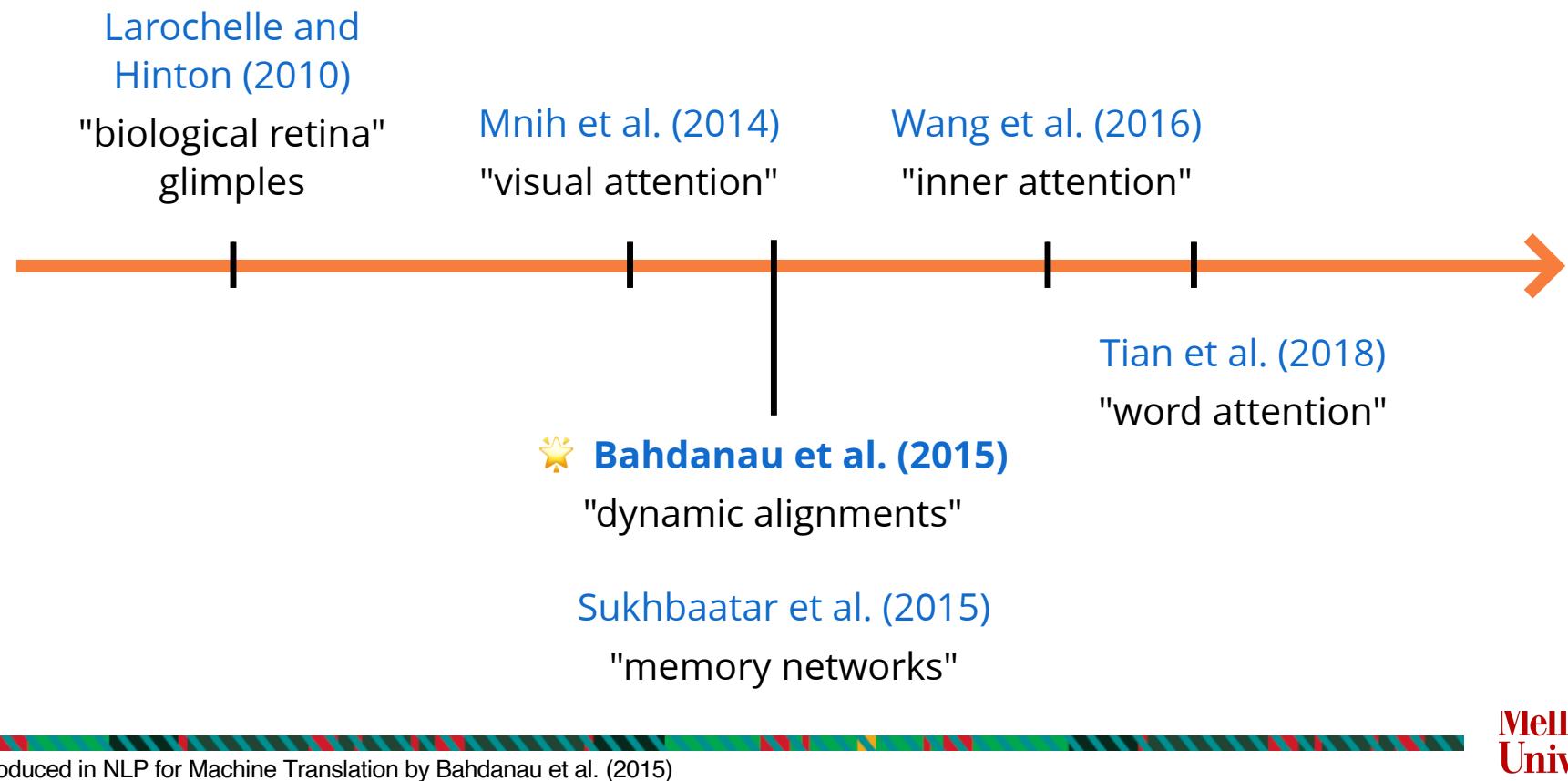
(Martins et al., 2020)

Carnegie  
Mellon  
University

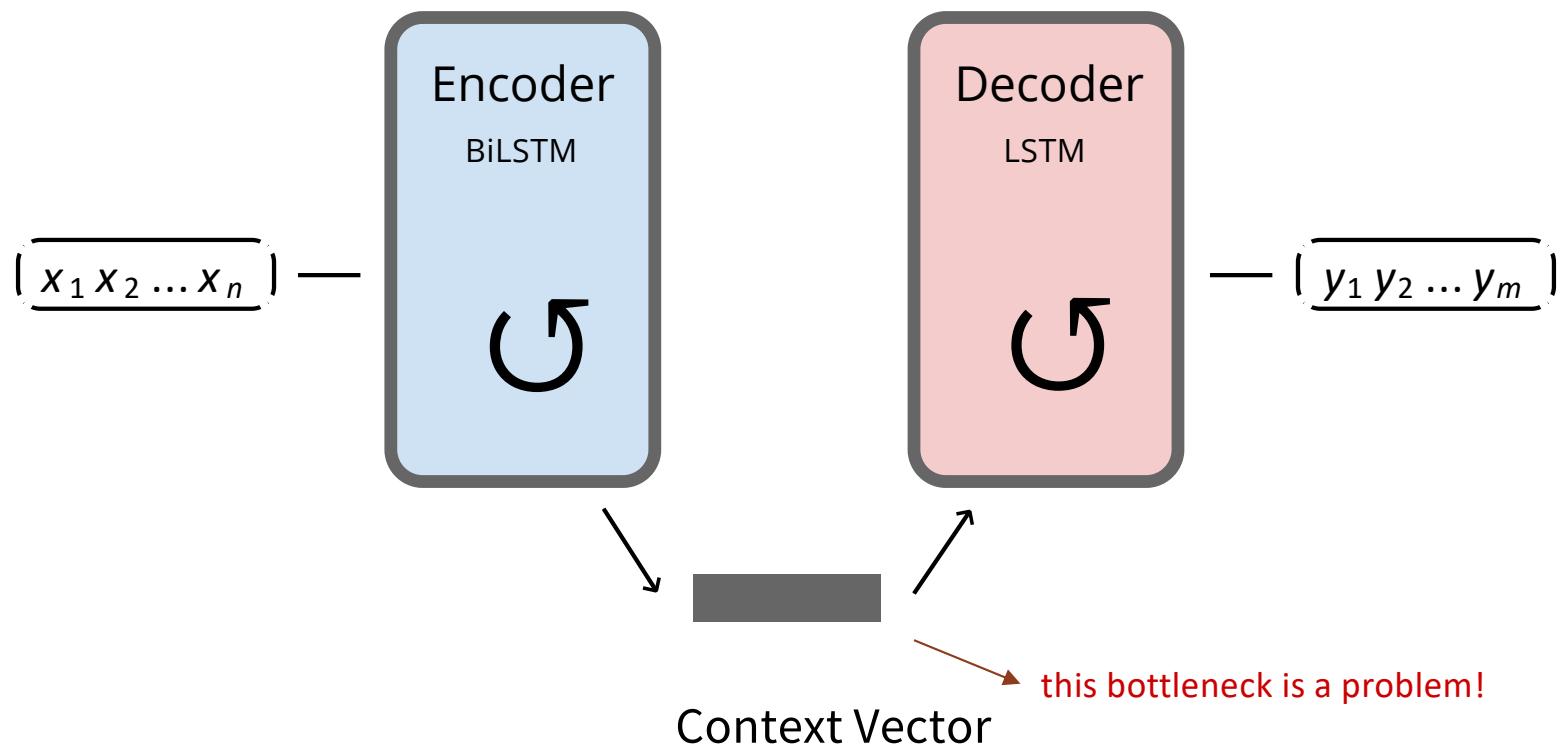
# Attention Example



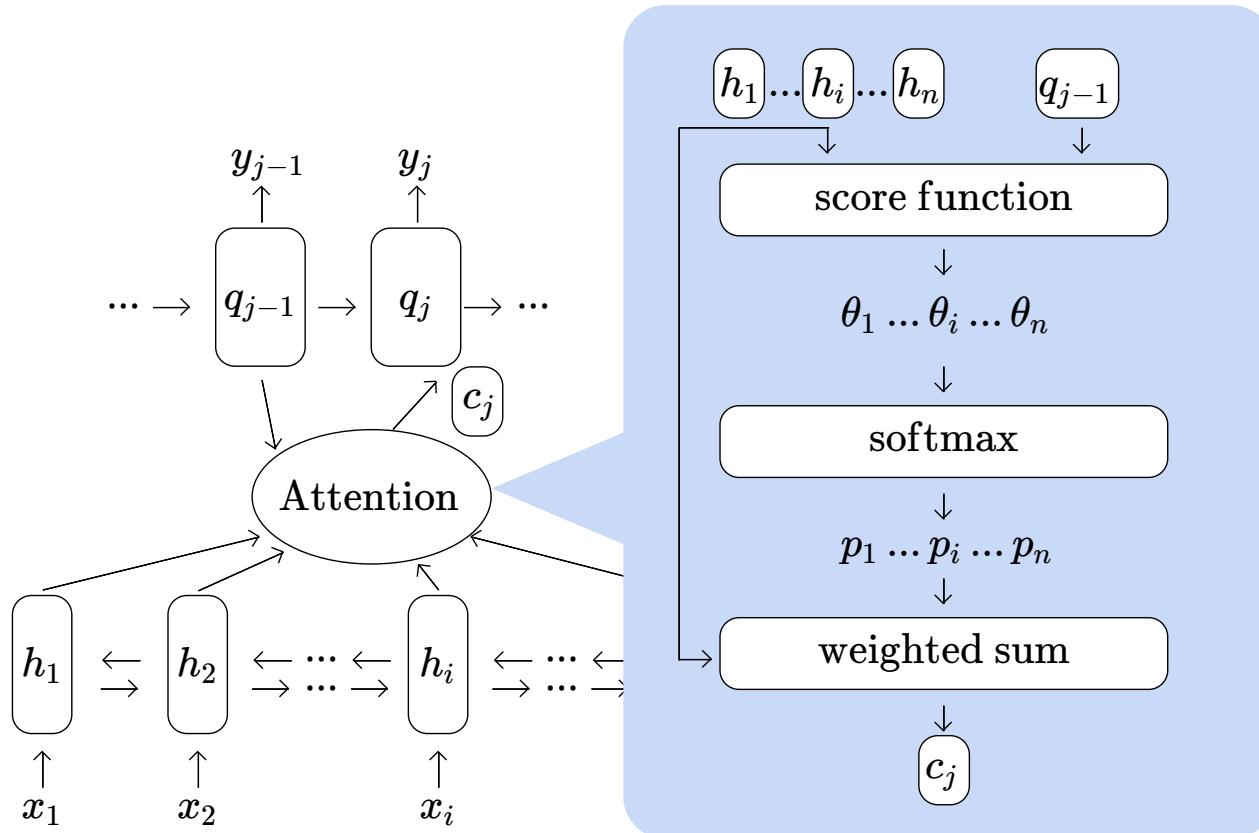
# Attention History



# RNN-based seq2seq



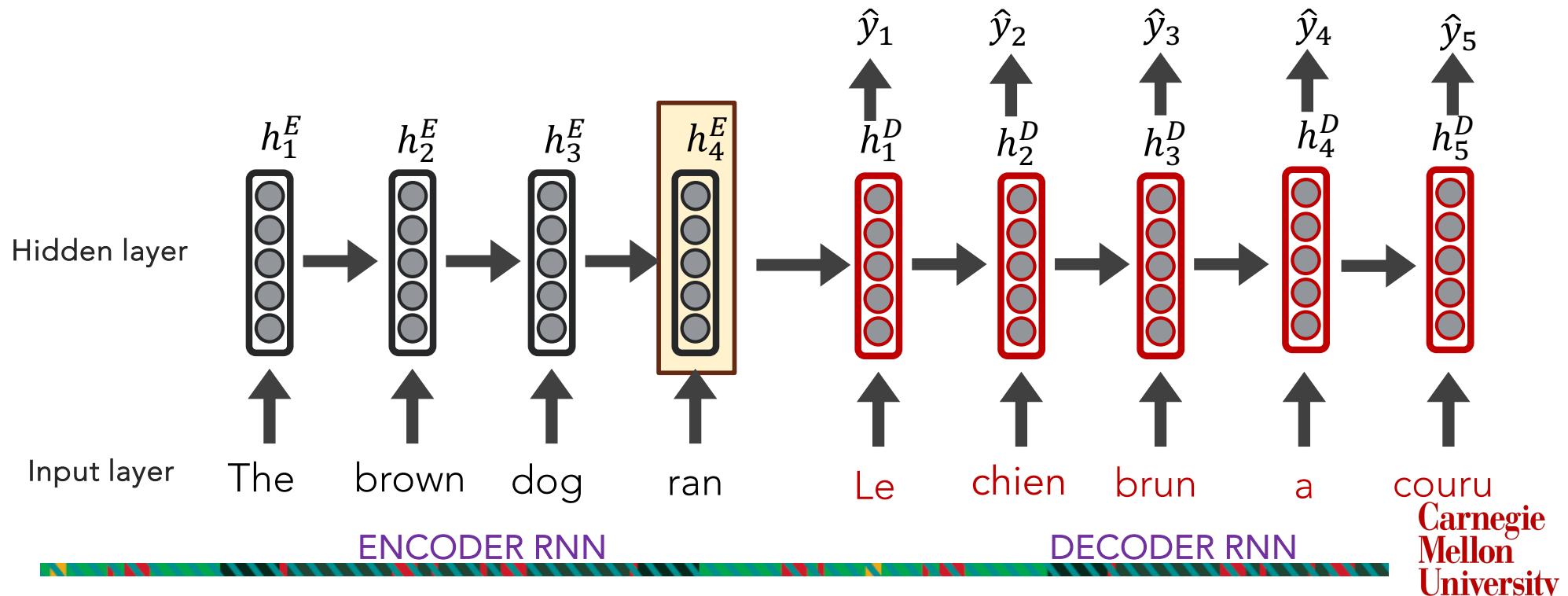
# Attention in Seq2Seq (2015)



Bahdanau et al. (2015)

# Sequence-to-Sequence (seq2seq)

It's crazy that the entire "meaning" of the 1<sup>st</sup> sequence is expected to be packed into this one embedding, and that the encoder then never interacts w/ the decoder again. Hands free.



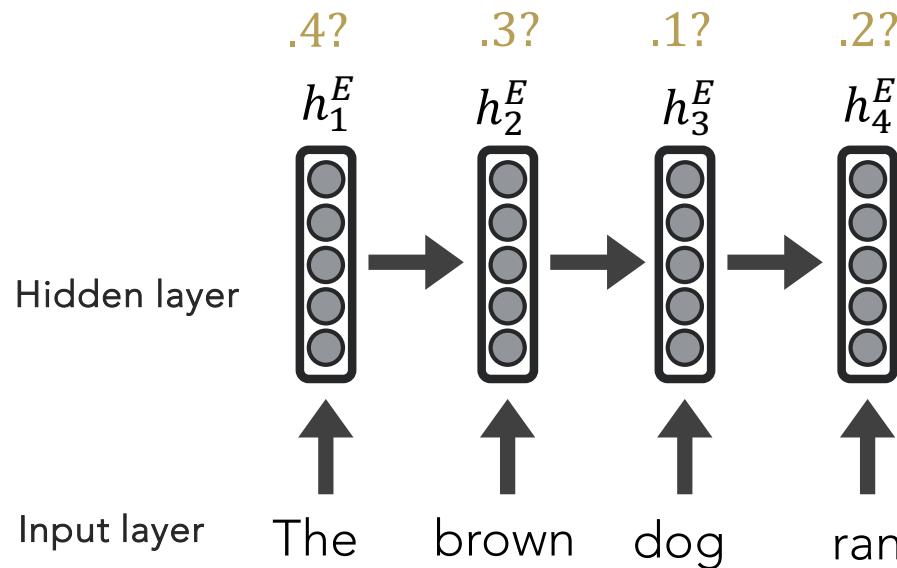
# Seq2Seq models

Instead, what if the decoder, at each step, pays attention to a *distribution* of all of the encoder's hidden states?



# Seq2Seq models

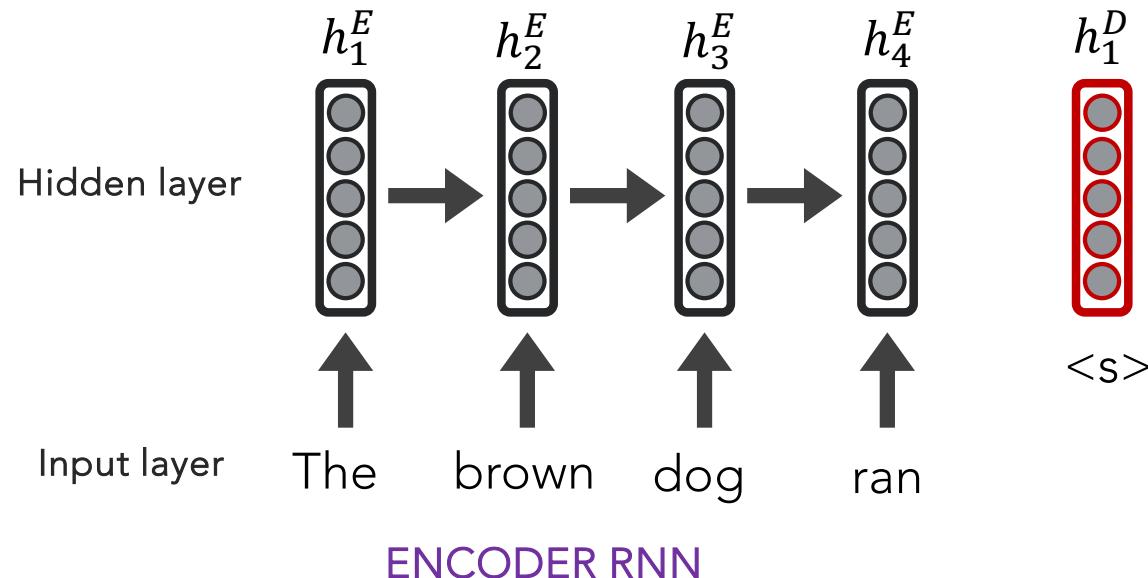
**Q: How do we determine how much to pay attention to each of the encoder's hidden layers?**



# Seq2Seq + Attention models

**Q:** How do we determine how much to pay attention to each of the encoder's hidden layers?

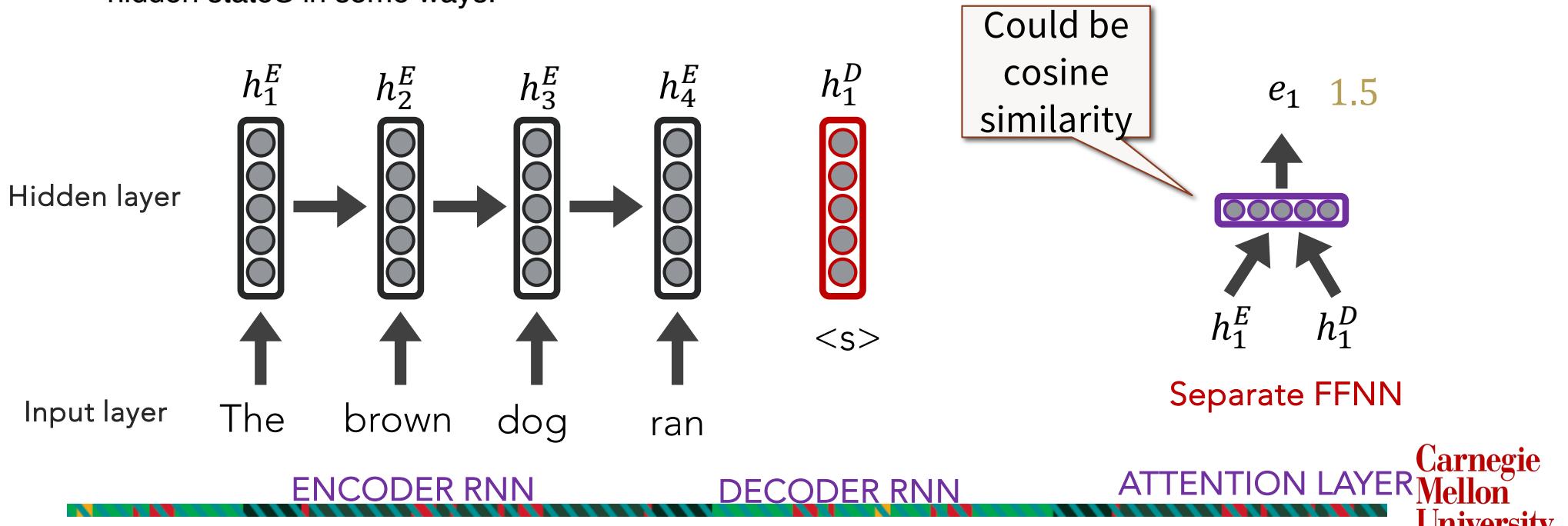
**A:** Let's base it on our decoder's previous hidden state (our latest representation of meaning) and all of the encoder's hidden layers!



# Seq2Seq + Attention

Q: How do we determine how much to pay attention to each of the encoder's hidden layers?

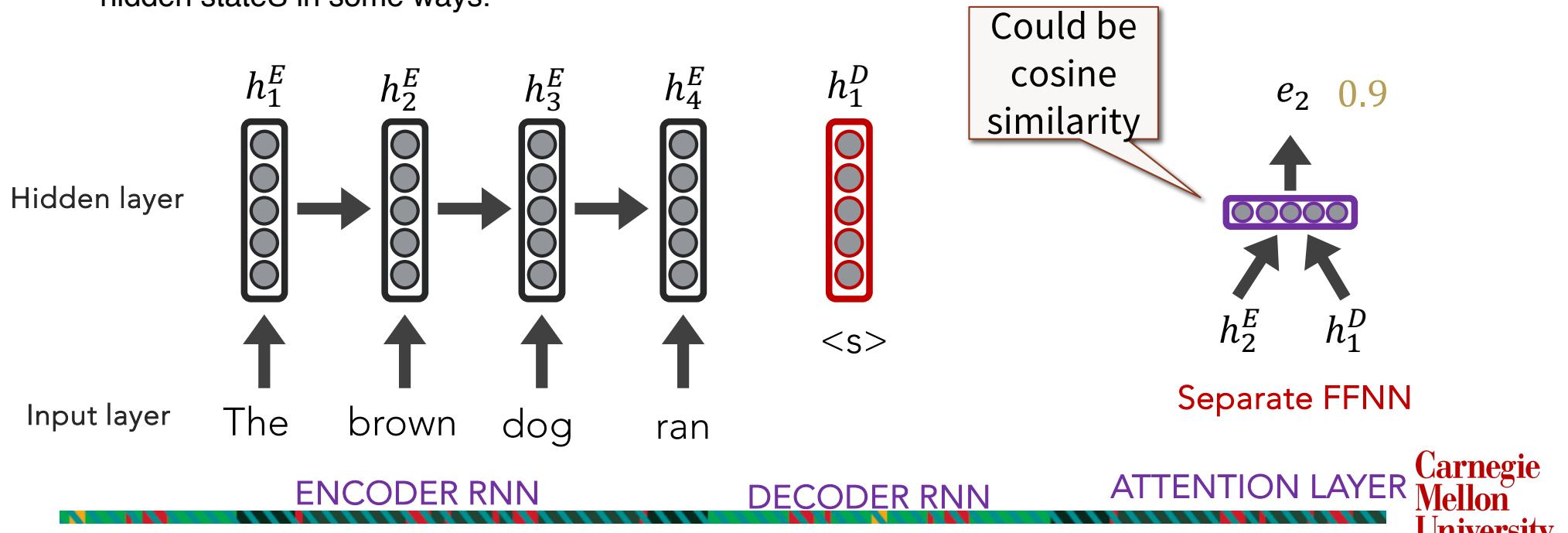
A: Let's base it on our decoder's previous hidden state (our latest representation of meaning) and all of the encoder's hidden layers! We want to measure similarity between decoder hidden state and encoder hidden states in some ways.



# Seq2Seq + Attention

Q: How do we determine how much to pay attention to each of the encoder's hidden layers?

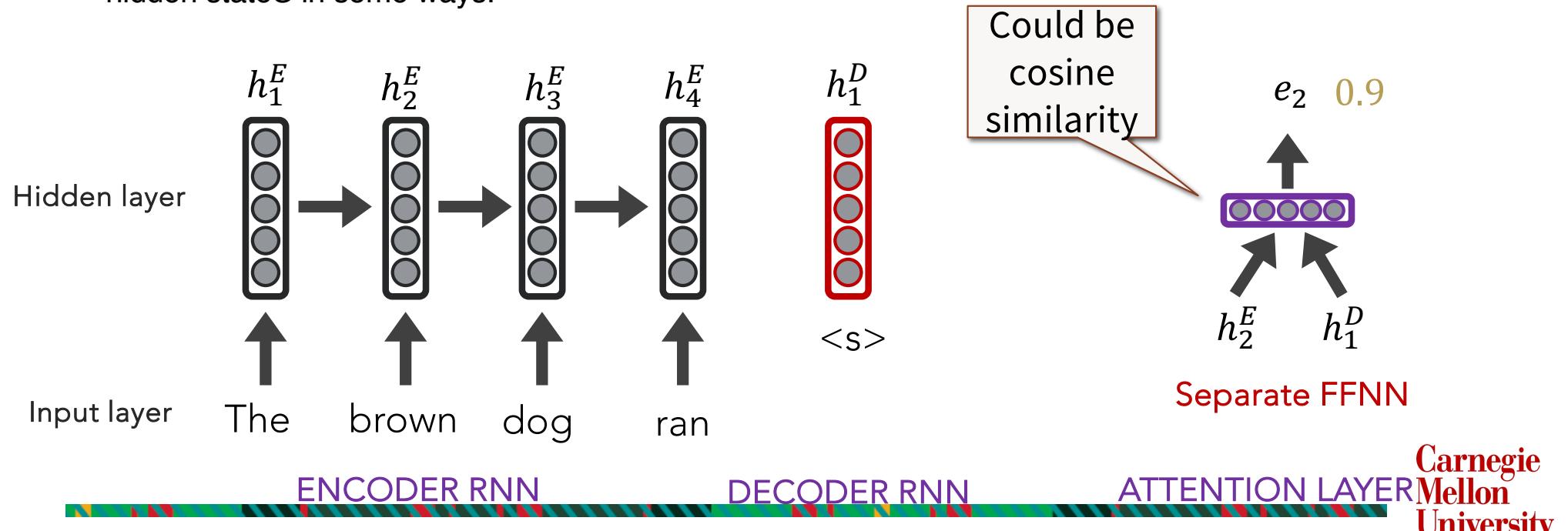
A: Let's base it on our decoder's previous hidden state (our latest representation of meaning) and all of the encoder's hidden layers! We want to measure similarity between decoder hidden state and encoder hidden states in some ways.



# Seq2Seq + Attention

Q: How do we determine how much to pay attention to each of the encoder's hidden layers?

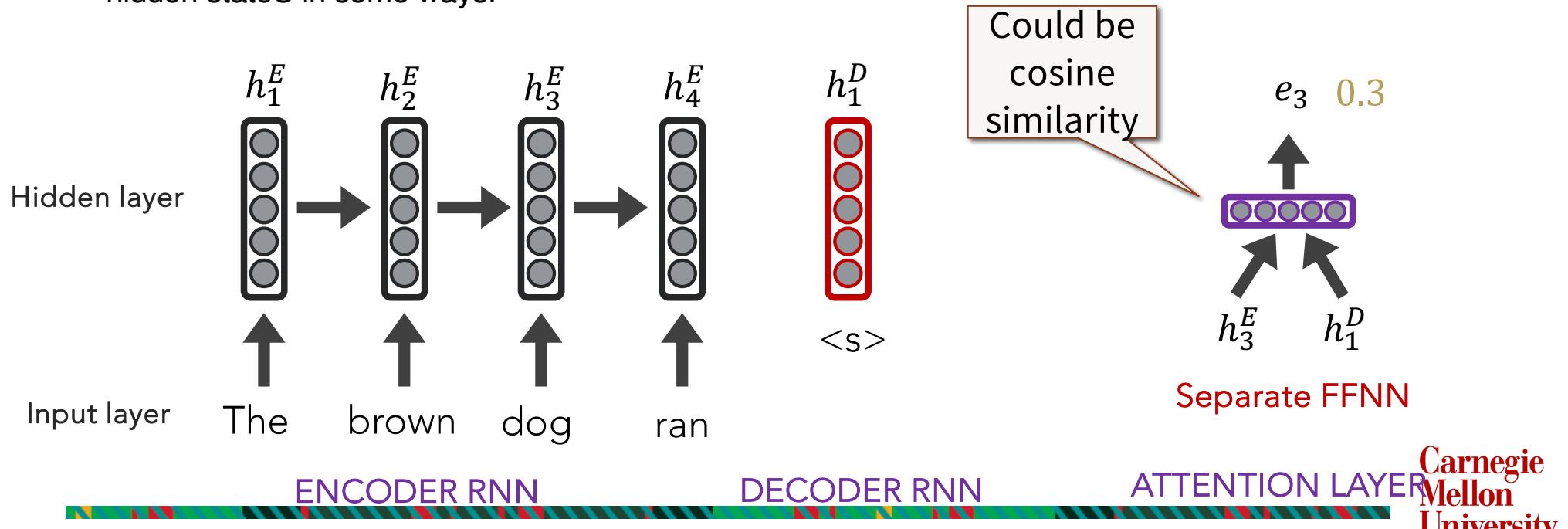
A: Let's base it on our decoder's previous hidden state (our latest representation of meaning) and all of the encoder's hidden layers! We want to measure similarity between decoder hidden state and encoder hidden states in some ways.



# Seq2Seq + Attention

Q: How do we determine how much to pay attention to each of the encoder's hidden layers?

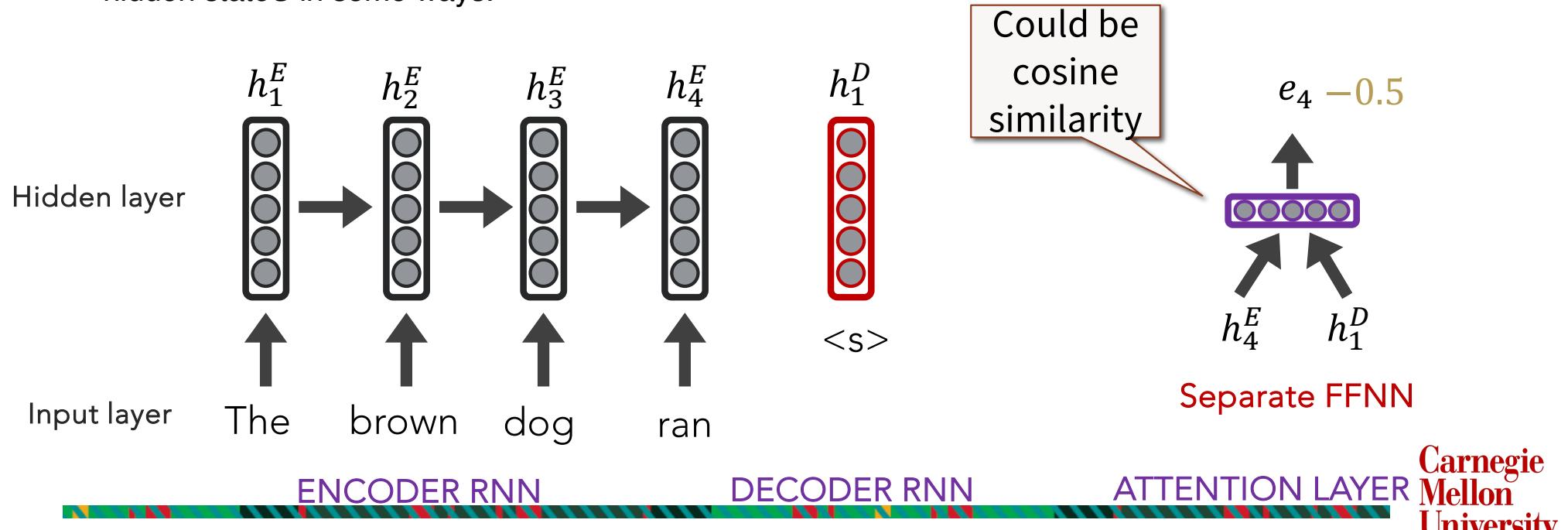
A: Let's base it on our decoder's previous hidden state (our latest representation of meaning) and all of the encoder's hidden layers! We want to measure similarity between decoder hidden state and encoder hidden states in some ways.



# Seq2Seq + Attention

Q: How do we determine how much to pay attention to each of the encoder's hidden layers?

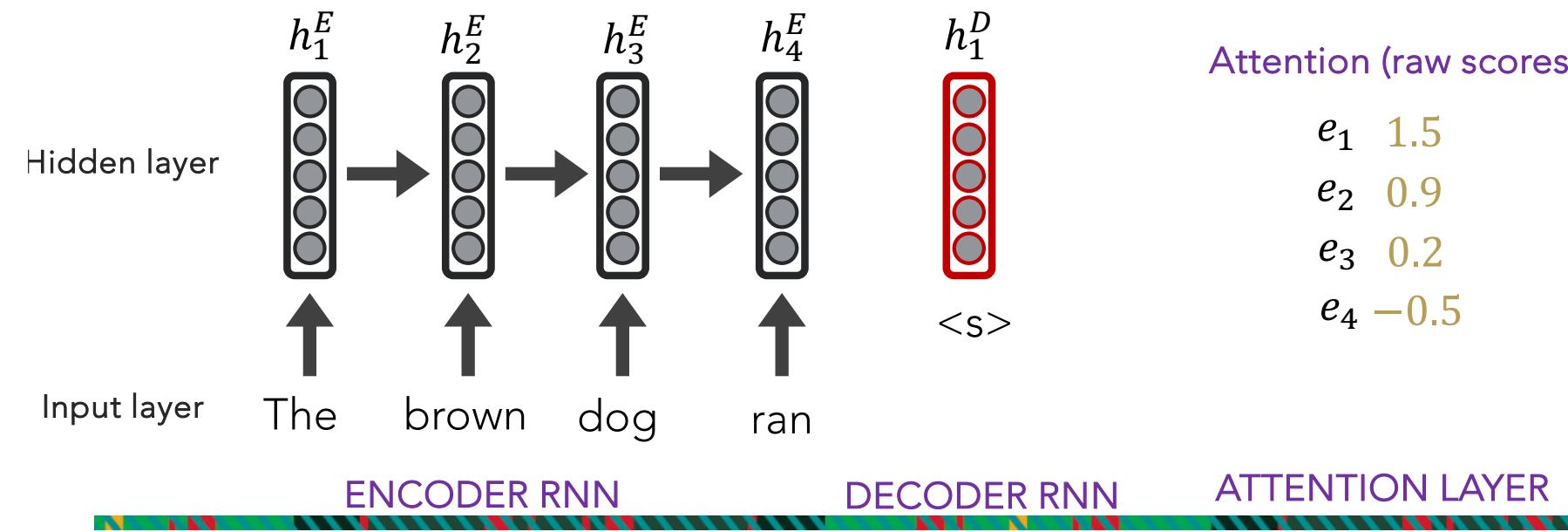
A: Let's base it on our decoder's previous hidden state (our latest representation of meaning) and all of the encoder's hidden layers! We want to measure similarity between decoder hidden state and encoder hidden states in some ways.



# Seq2Seq + Attention

**Q:** How do we determine how much to pay attention to each of the encoder's hidden layers?

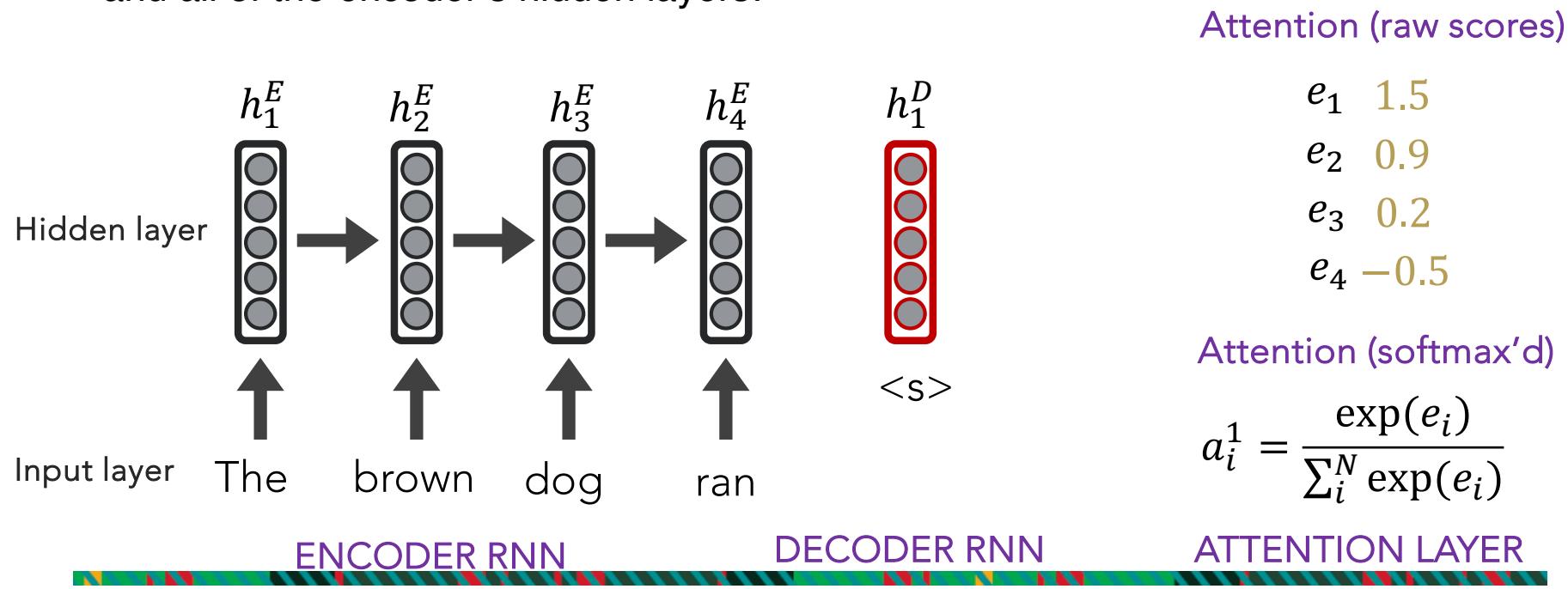
**A:** Let's base it on our decoder's previous hidden state (our latest representation of meaning) and all of the encoder's hidden layers!



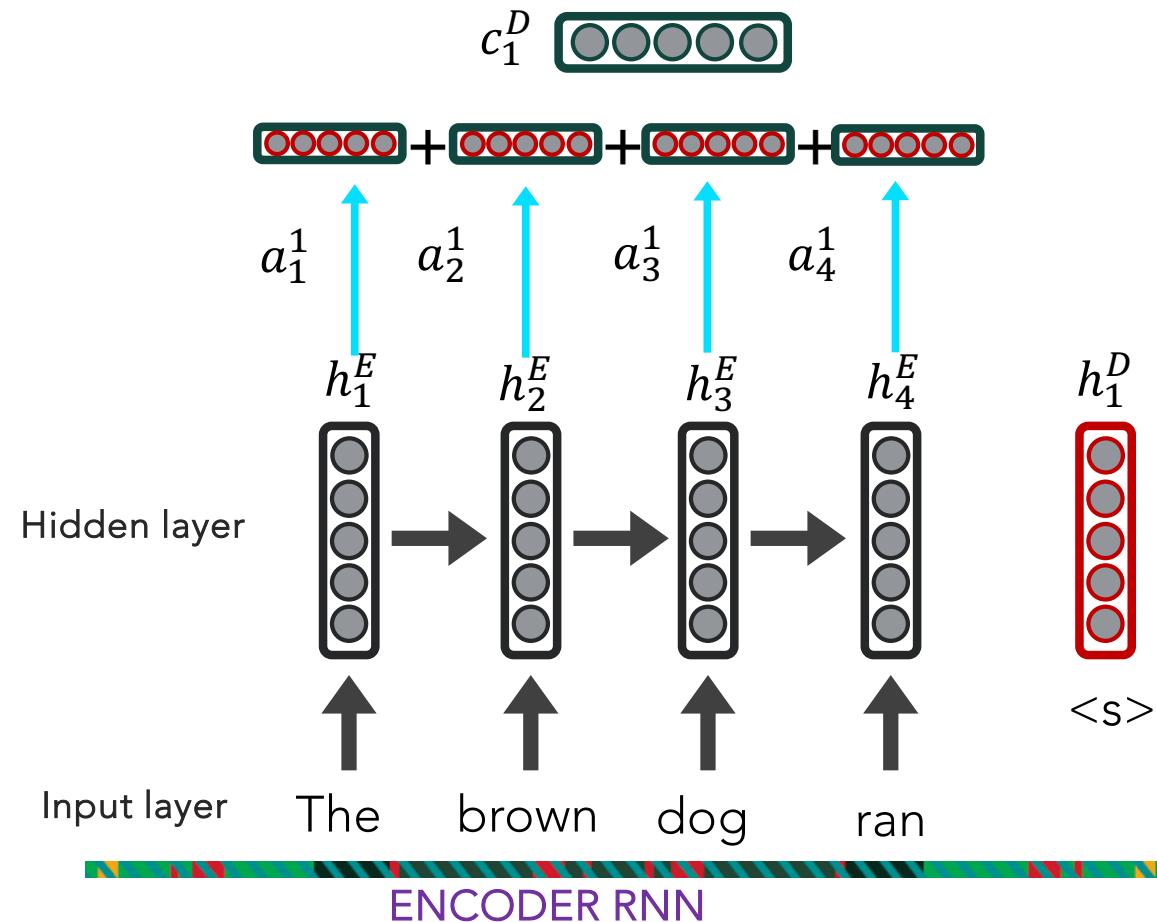
# Seq2Seq + Attention

**Q:** How do we determine how much to pay attention to each of the encoder's hidden layers?

**A:** Let's base it on our decoder's previous hidden state (our latest representation of meaning) and all of the encoder's hidden layers!



# Seq2Seq + Attention



We multiply each hidden layer by its  $a_i^1$  attention weights and then add the resulting vectors to create a context vector  $c_1^D$

Attention (softmax'd)

$$a_1^1 = 0.51$$

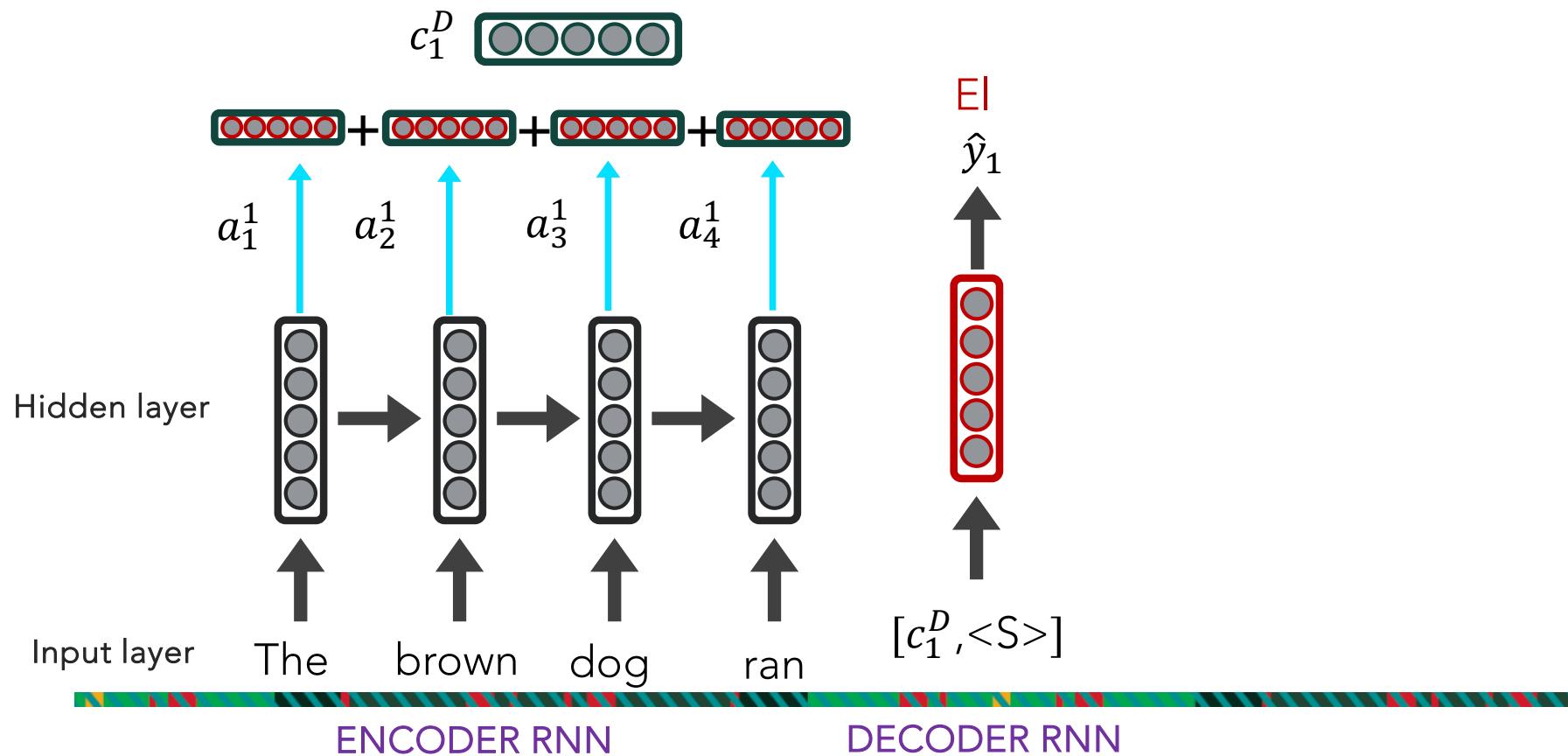
$$a_2^1 = 0.28$$

$$a_3^1 = 0.14$$

$$a_4^1 = 0.07$$

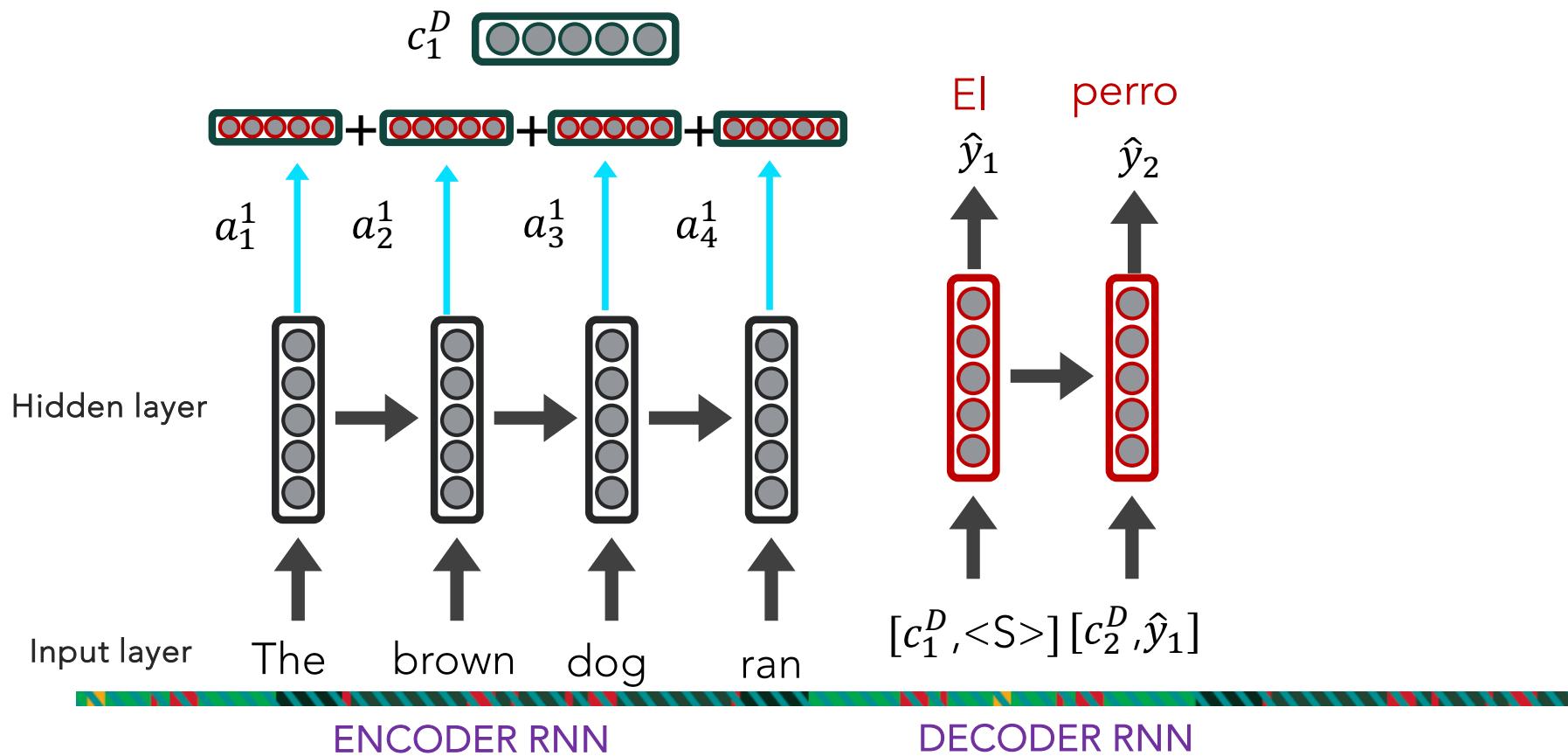
# Seq2Seq + Attention

NOTE: each attention weight  $a_i^j$  is based on the decoder's current hidden state, too.



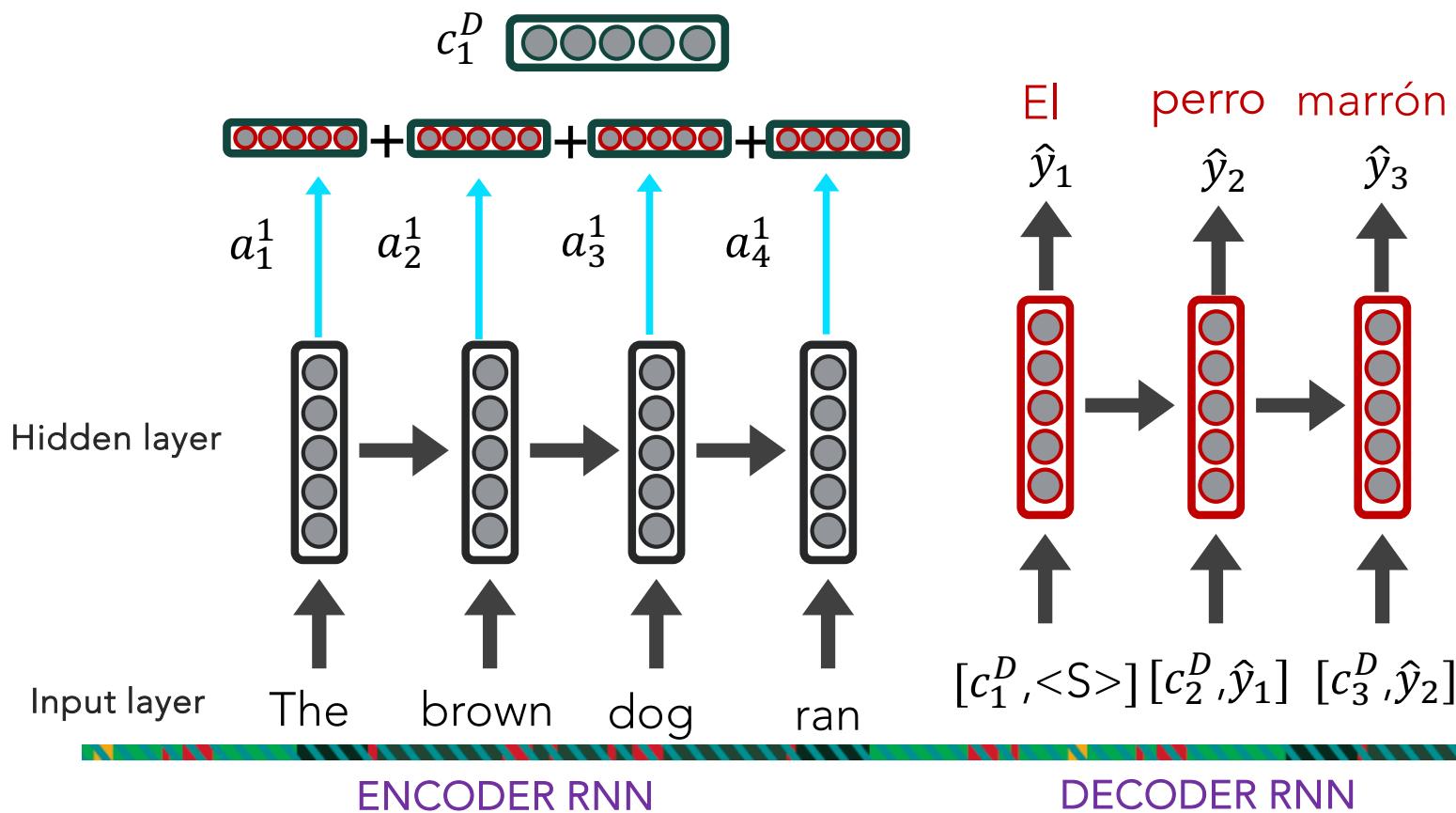
# Seq2Seq + Attention

NOTE: each attention weight  $a_i^j$  is based on the decoder's current hidden state, too.



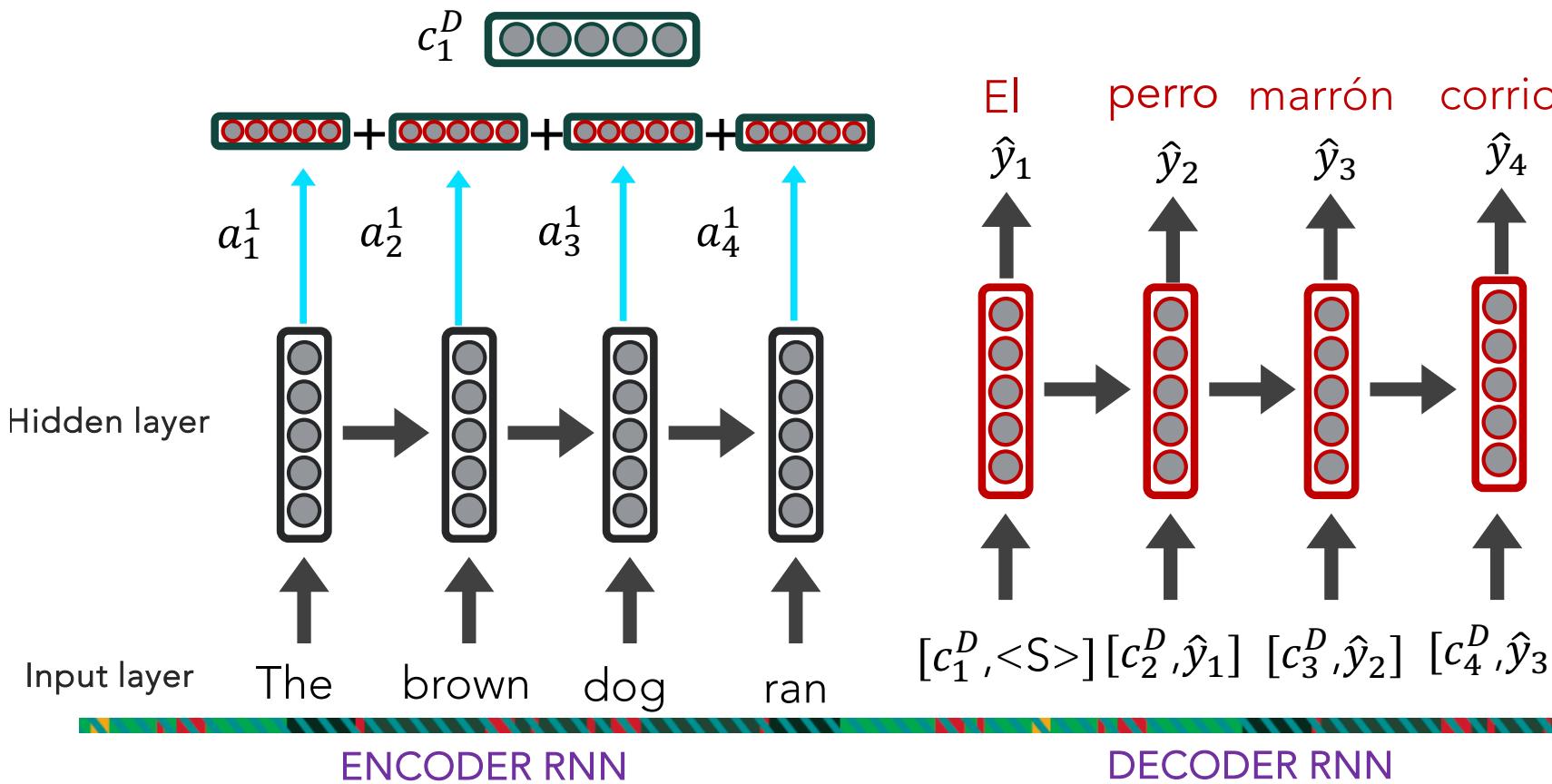
# Seq2Seq + Attention

NOTE: each attention weight  $a_i^j$  is based on the decoder's current hidden state, too.



# Seq2Seq + Attention

NOTE: each attention weight  $a_i^j$  is based on the decoder's current hidden state, too.



# Seq2Seq + Attention

## Attention:

- greatly improves seq2seq results
- allows us to visualize the contribution each word gave during each step of the decoder

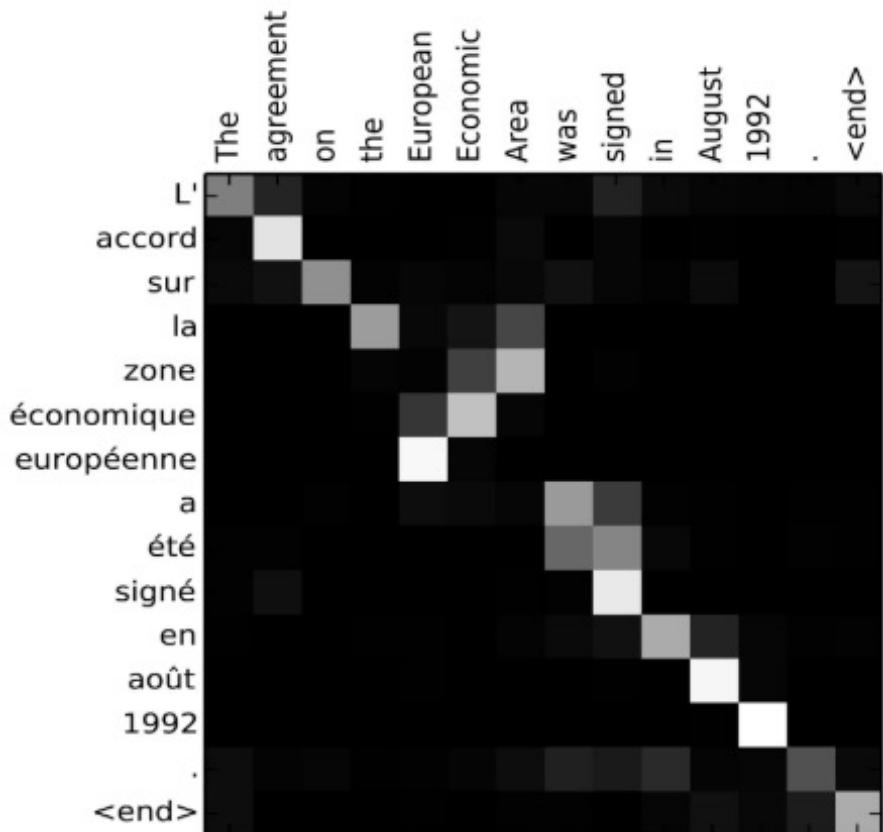


Image source: Fig 3 in [Bahdanau et al., 2015](#)

But still the problem of  
Sequential Computing Exist  
Because Backprop in time  
exists

---

Image source: Fig 3 in [Bahdanau et al., 2015](#)

# Transformers

## Attention Is All You Need

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

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

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

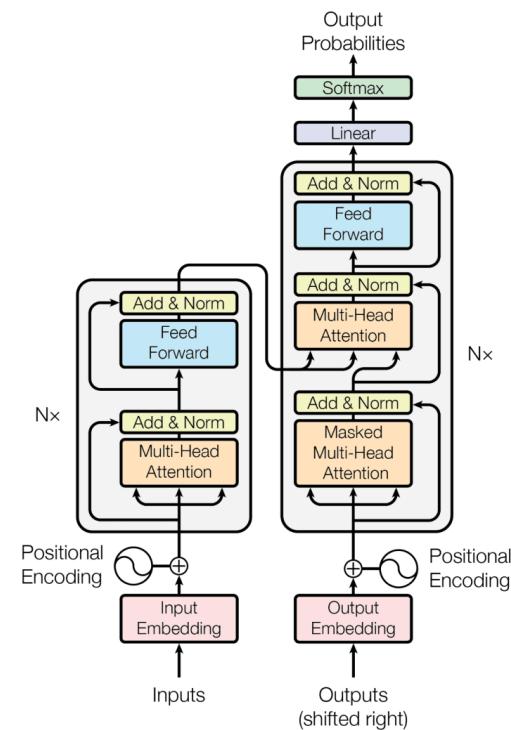
**Jakob Uszkoreit\***  
Google Research  
usz@google.com

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

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

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

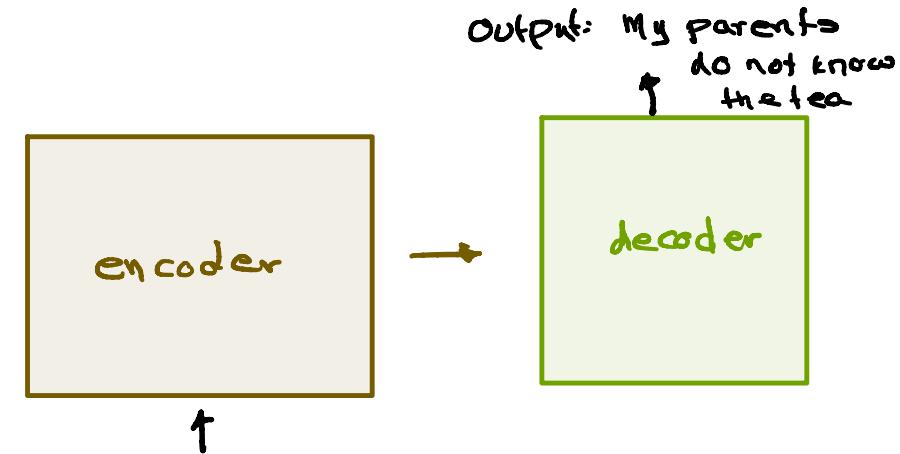
**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com



# Transformers

The **Transformer** is a model that uses attention to boost the speed with which seq2seq with attention models can be trained. The biggest benefit, however, comes from how The Transformer lends itself to **parallelization**. We will break it apart and look at how it functions.

In its heart it contains an encoding component, a decoding component, and connections between them.



Input: Mis padres  
no conocen a los  
maestros

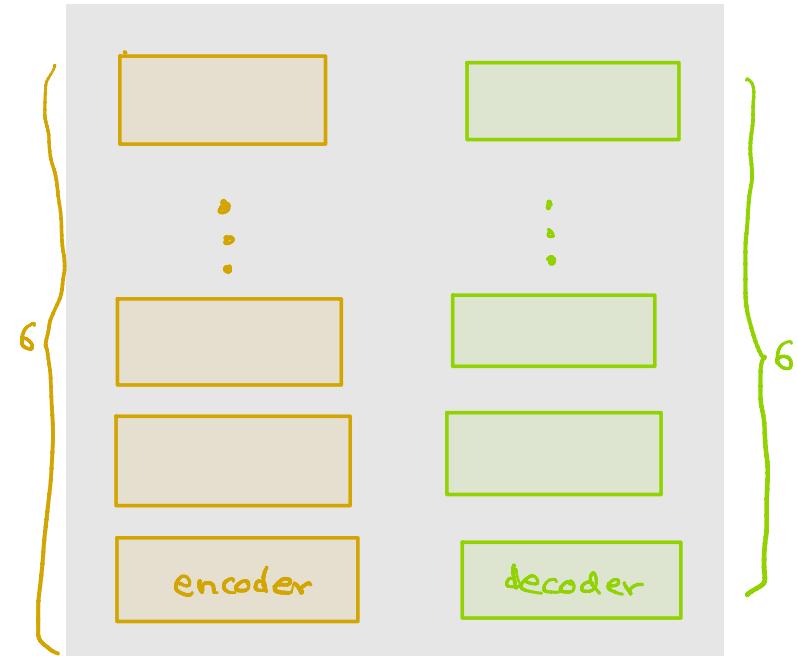
# Transformers

my parents do  
not know the  
teachers

The encoding is a stack of encoders.

The original paper stacks six of them on top of each other – there's nothing magical about the number six, one can definitely experiment with other arrangements).

The decoding is a stack of decoders of the same number.



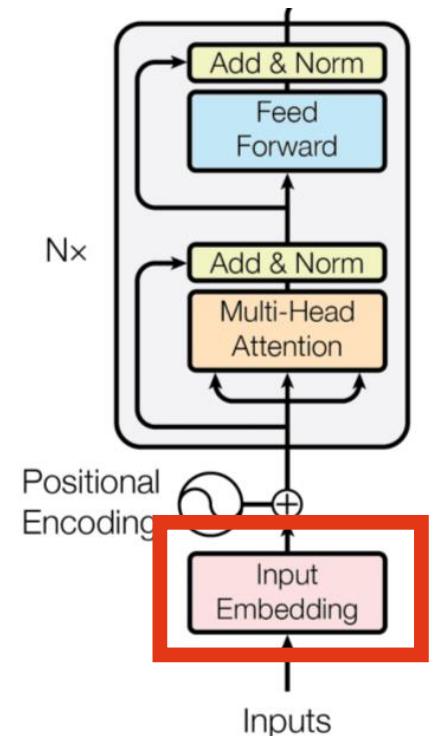
mis padres  
no conocen a  
los maestros

# Transformers (Embedding)

Input embedding of words based on dictionary:

Original sentence (tokens)	YOUR	CAT	IS	A	LOVELY	CAT
Input IDs (position in the vocabulary)	105	6587	5475	3578	65	6587
Embedding (vector of size 512)	952.207 5450.840 1853.448 ... 1.658 2671.529	171.411 3276.350 9192.819 ... 3633.421 8390.473	621.659 1304.051 0.565 ... 7679.805 4506.025	776.562 5567.288 58.942 ... 2716.194 5119.949	6422.693 6315.080 9358.778 ... 2141.081 735.147	171.411 3276.350 9192.819 ... 3633.421 8390.473

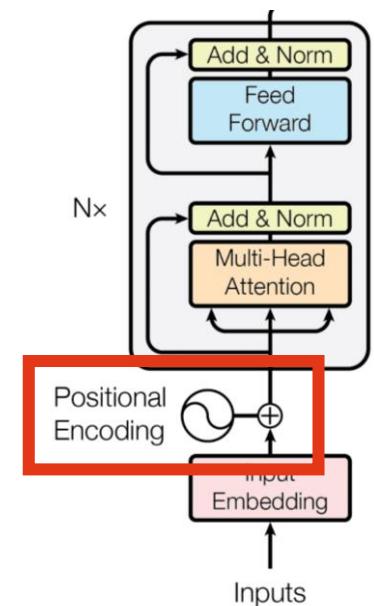
We define  $d_{model} = 512$ , which represents the size of the embedding vector of each word



# Transformers (Pos Encoding)

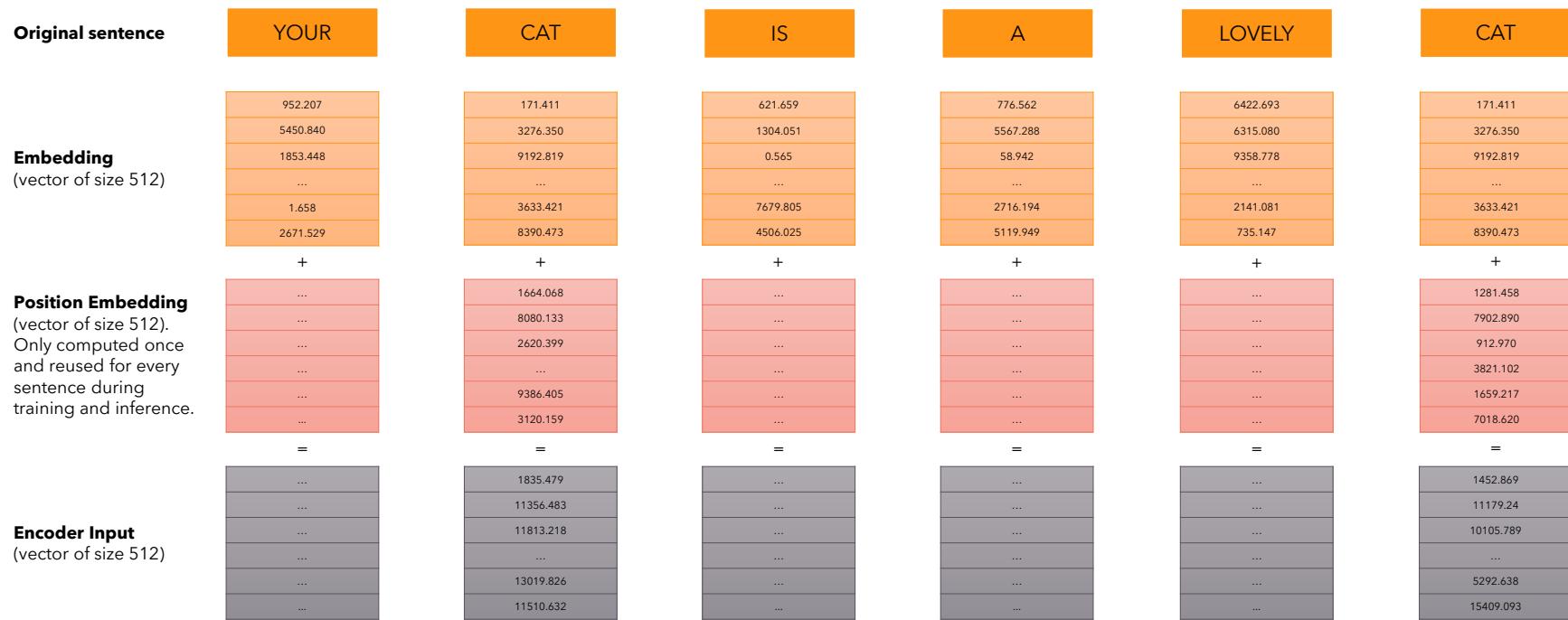
## Positional Encoding

- We want each word to carry some information about its position in the sentence.
- We want the model to treat words that appear close to each other as “close” and words that are distant as “distant”.
- We want the positional encoding to represent a pattern that can be learned by the model.



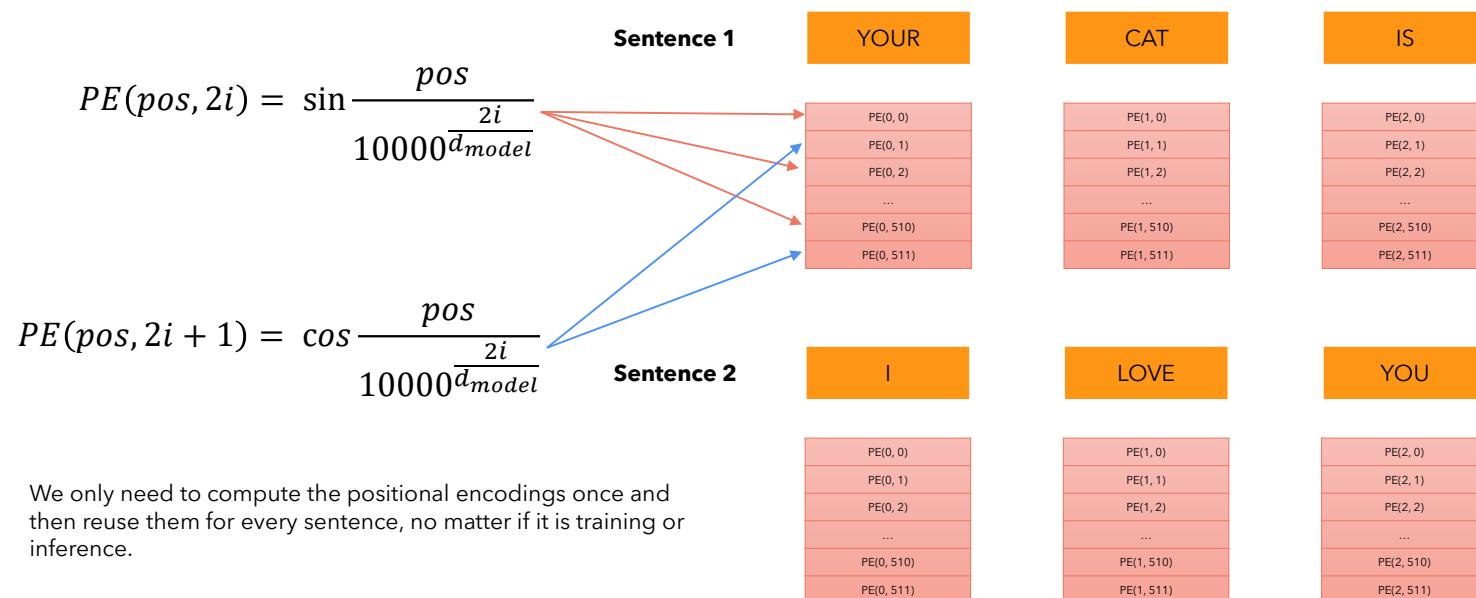
# Transformers (Embed+Pos)

Adding the positional encoding to the embedding



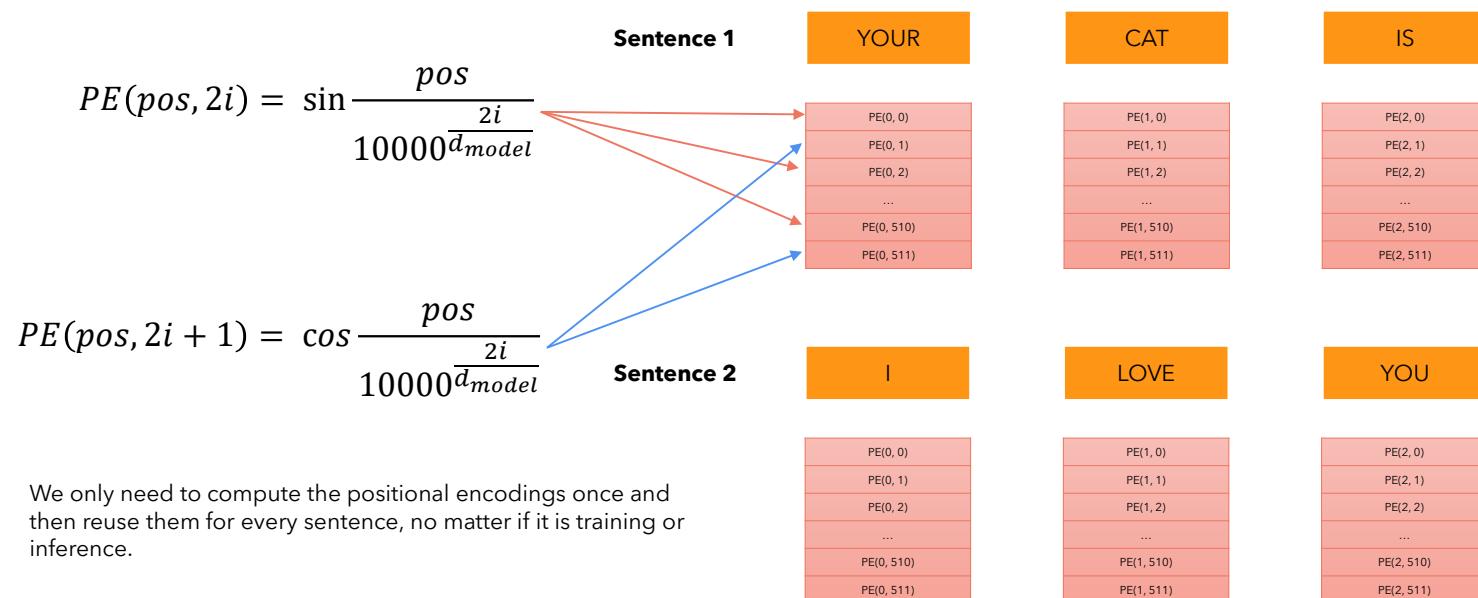
# Transformers (Positional Encoding)

What is positional encoding and how is it implemented



# Transformers (Positional Encoding)

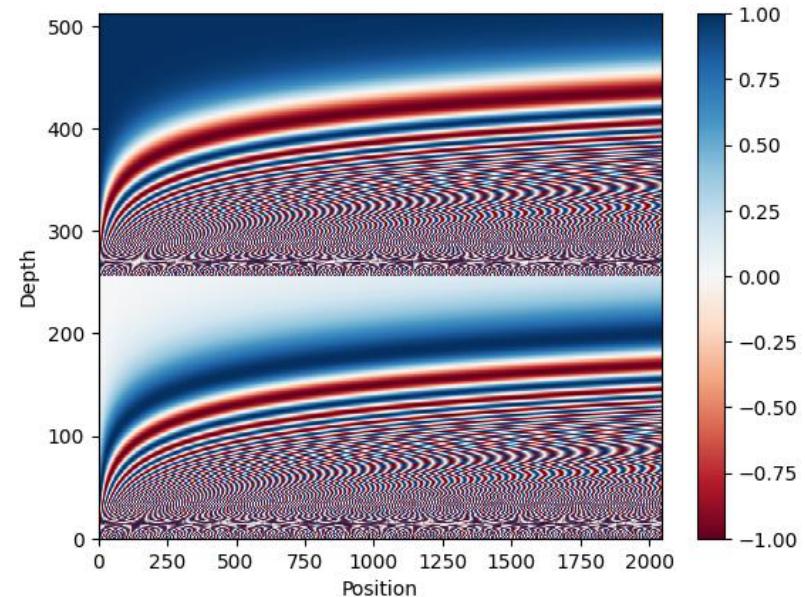
What is positional encoding and how is it implemented



# Transformers (Positional Encoding)

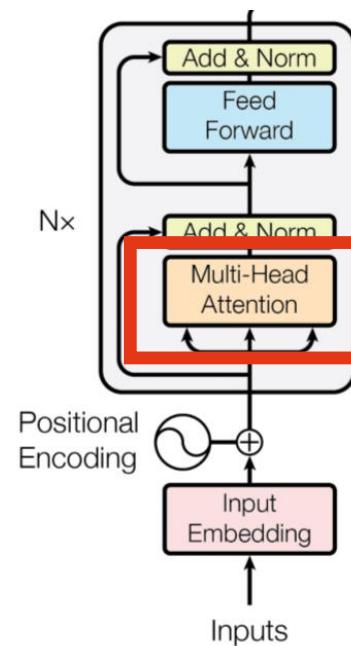
What is positional encoding and how is it implemented

Trigonometric functions like **cos** and **sin** naturally represent a pattern that the model can recognize as continuous, so relative positions are easier to see for the model. By watching the plot of these functions, we can also see a regular pattern, so we can hypothesize that the model will see it too.



# Transformers (Multi-head Attention)

First we need to learn Self attention

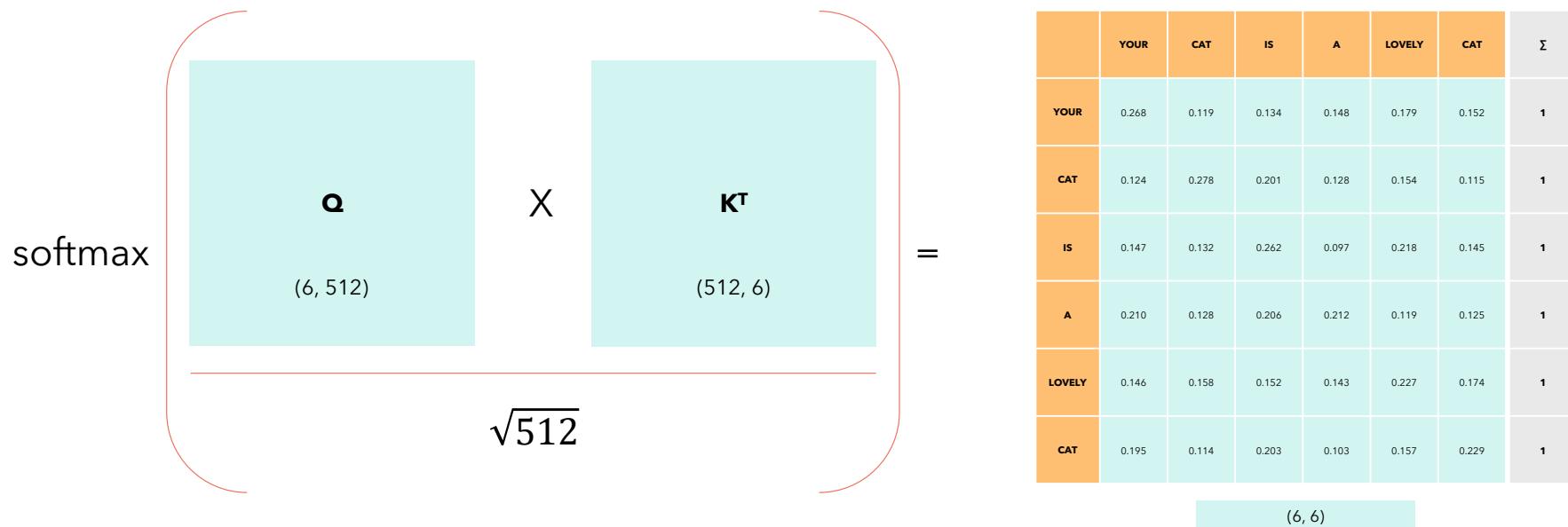


# Self Attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

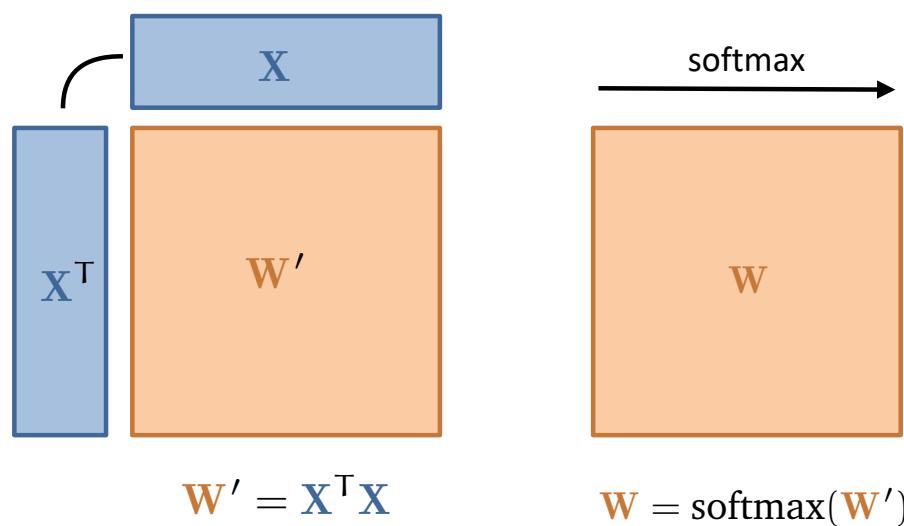
Self-Attention allows the model to relate words to each other.

In this simple case we consider the sequence length  $\text{seq} = 6$  and  $\mathbf{d}_{\text{model}} = \mathbf{d}_k = 512$ . The matrices  $\mathbf{Q}$ ,  $\mathbf{K}$  and  $\mathbf{V}$  are just the input sentence.



# Self Attention

VECTORIZED



$$w'_{ij} = \mathbf{x}_i^T \mathbf{x}_j$$

$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$

# Self Attention

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

(6, 6)

$$\mathbf{V} \times = \text{Attention}$$

(6, 512) (6, 512)

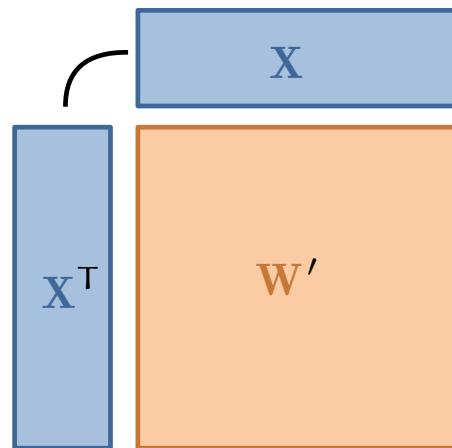
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Each row in this matrix captures not only the meaning (given by the embedding) or the position in the sentence (represented by the positional encodings) but also each word's interaction with other words.



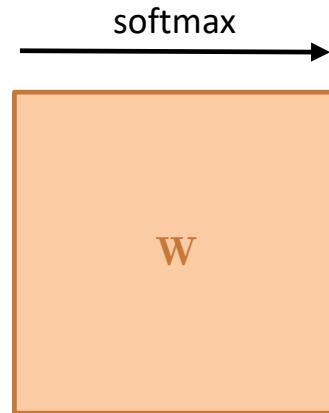
# Self Attention

VECTORIZED

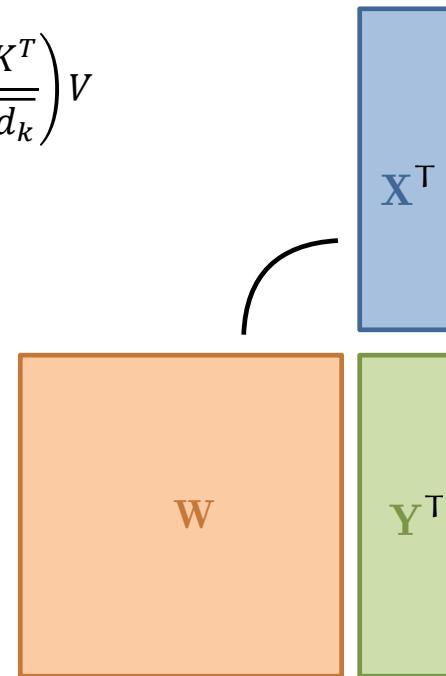


$$\mathbf{W}' = \mathbf{X}^T \mathbf{X}$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



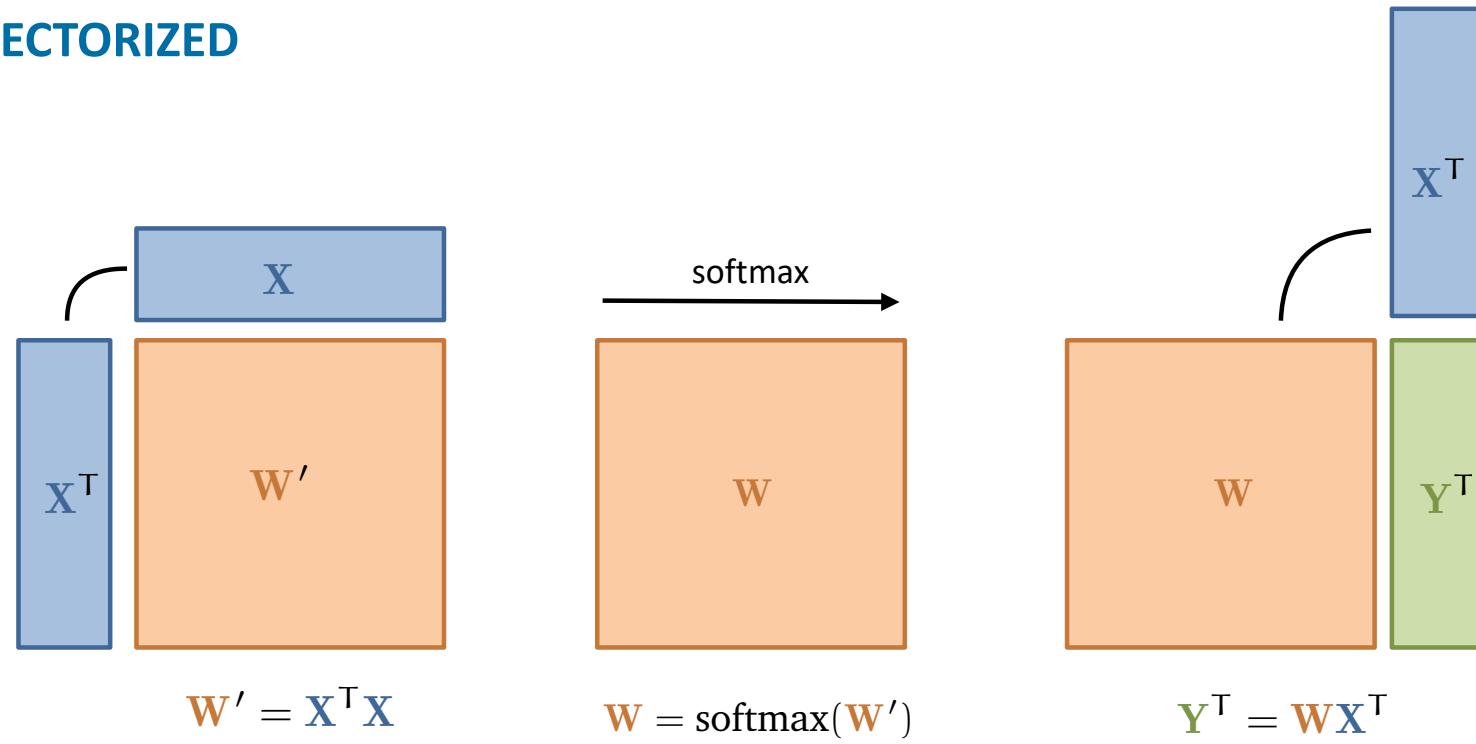
$$\mathbf{W} = \text{softmax}(\mathbf{W}')$$



$$\mathbf{Y}^T = \mathbf{W} \mathbf{X}^T$$

# Self Attention

VECTORIZED

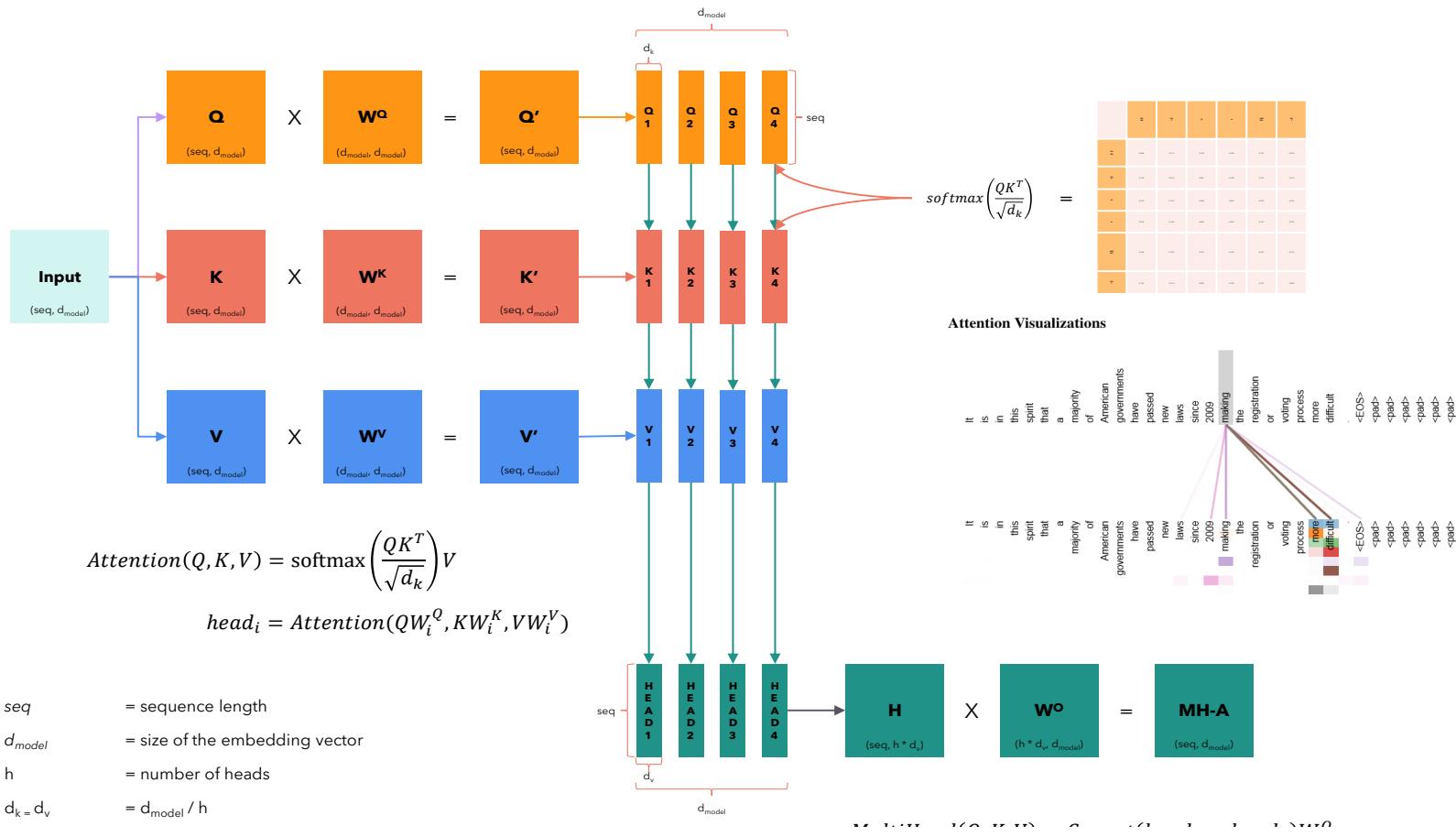


# Self Attention

- Self-Attention is permutation invariant.
- Self-Attention requires no parameters. Up to now the interaction between words has been driven by their embedding and the positional encodings. This will change later.
- We expect values along the diagonal to be the highest.
- If we don't want some positions to interact, we can always set their values to  $-\infty$  before applying the *softmax* in this matrix and the model will not learn those interactions. We will use this in the decoder.

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

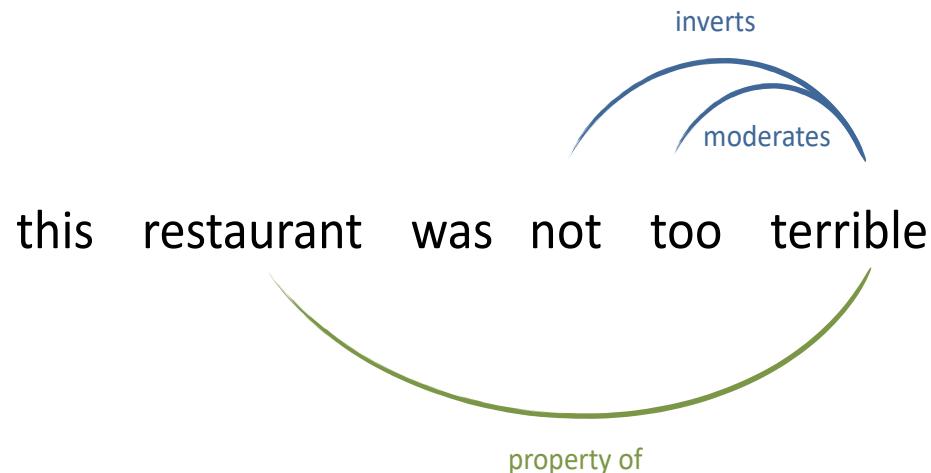
# Multi-head Attention



# Attention Learns functional relation via Multihead

In many sentences, there are different relations to model. Here, the word meaning of the word “terrible” is inverted by “not” and moderated by “too”. Its relation to the word restaurant is completely different: it describes a property of the restaurant.

The idea behind multi-head self-attention is that multiple relations are best captured by different self-attention operations.



# Attention as a Soft Dictionary

- **Attention** is a *soft* dictionary
- key, query and value are vectors
- every key matches the query *to some extent* as determined by their dot-product
- a *mixture* of all values is returned with softmax-normalized dot products as mixture weights
- **Self-attention** Attention with keys, queries and values from the same set.

key	value
a	1
b	2
c	3

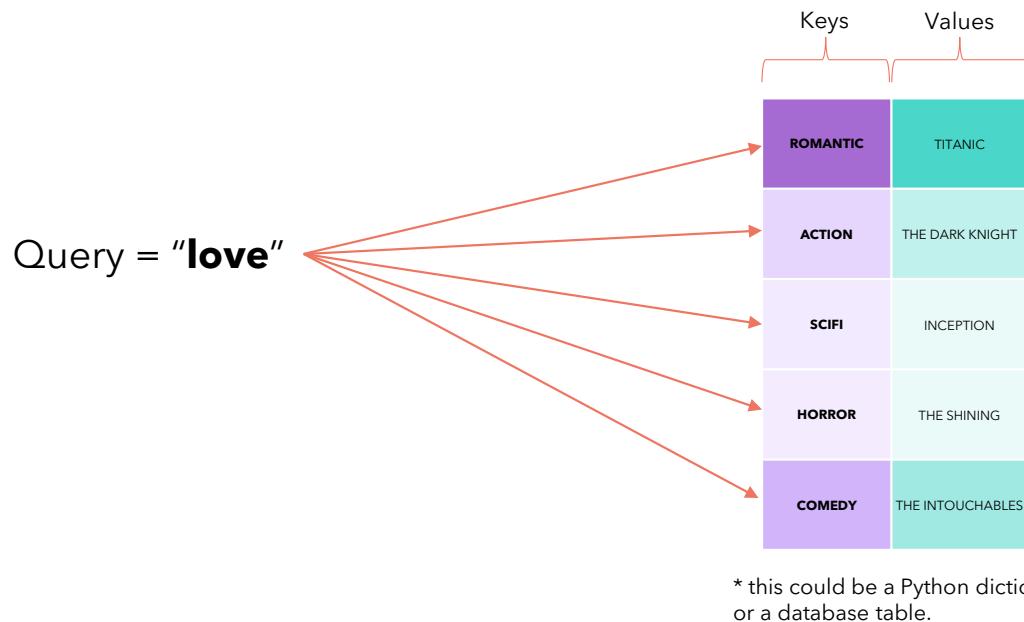
d = {'a' : 1, 'b' : 2, 'c' : 3}

d['b'] = 3

↑ key      ↑ value

↖ query

# K, Q, V interpretation



# Add&Norm

Batch of 3 items

ITEM 1

ITEM 2

ITEM 3

50.147
3314.825
...
...
8463.361
8.021

1242.223
688.123
...
...
434.944
149.442

9.370
4606.674
...
...
944.705
21189.444

$$\mu_1$$

$$\sigma_1^2$$

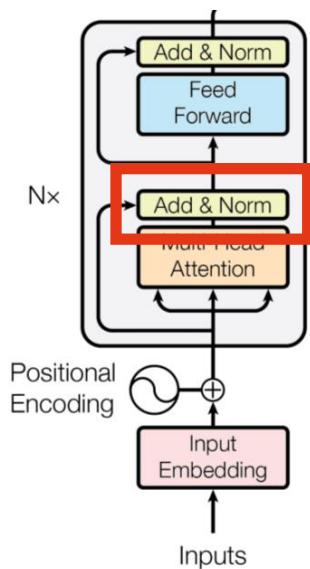
$$\mu_2$$

$$\sigma_2^2$$

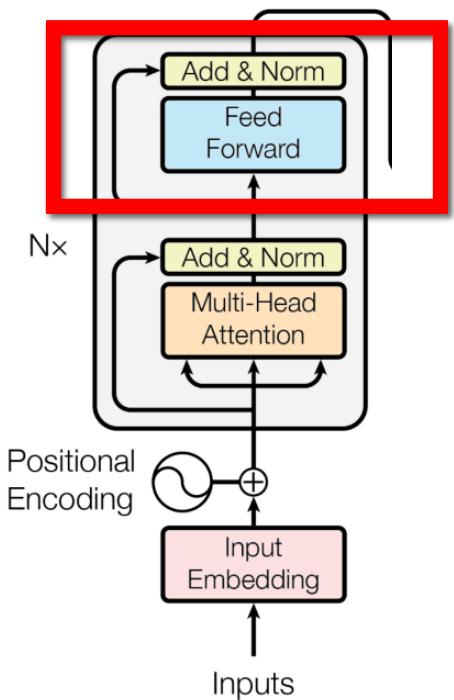
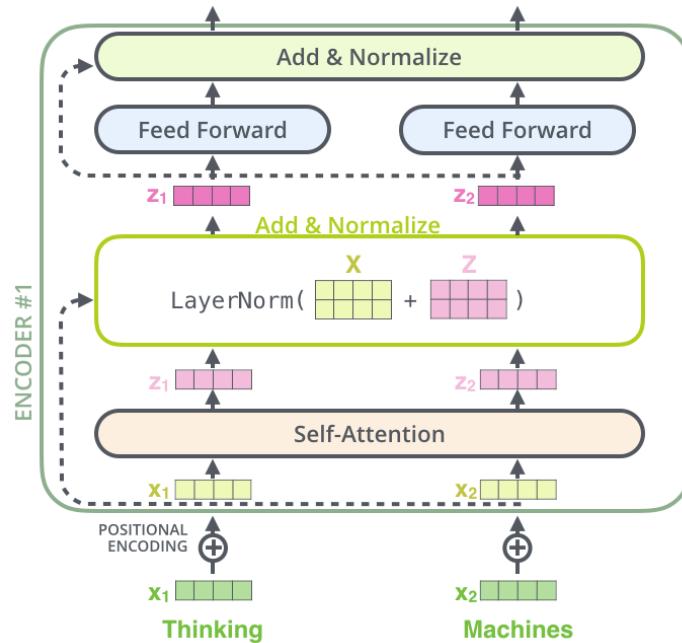
$$\mu_3$$

$$\sigma_3^2$$

$$\hat{x}_j = \frac{x_j - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$



# Feed Forward NN



# Encoders on top of each other

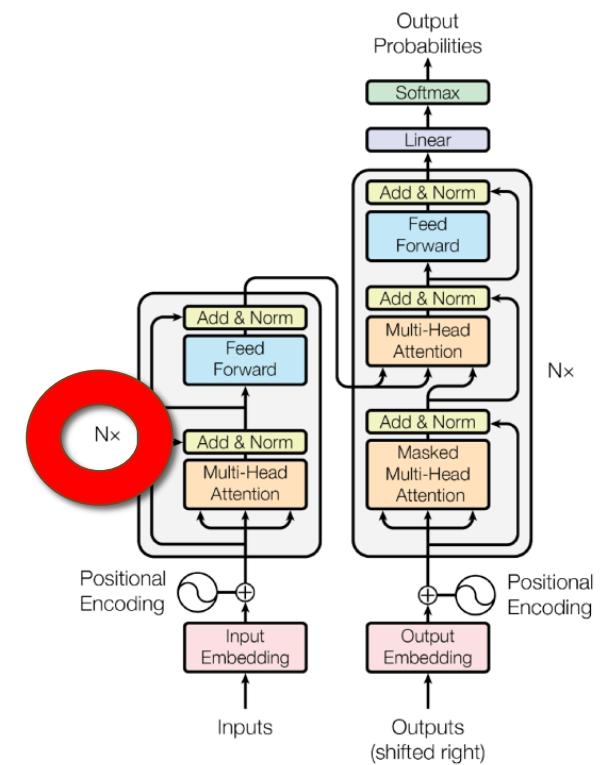
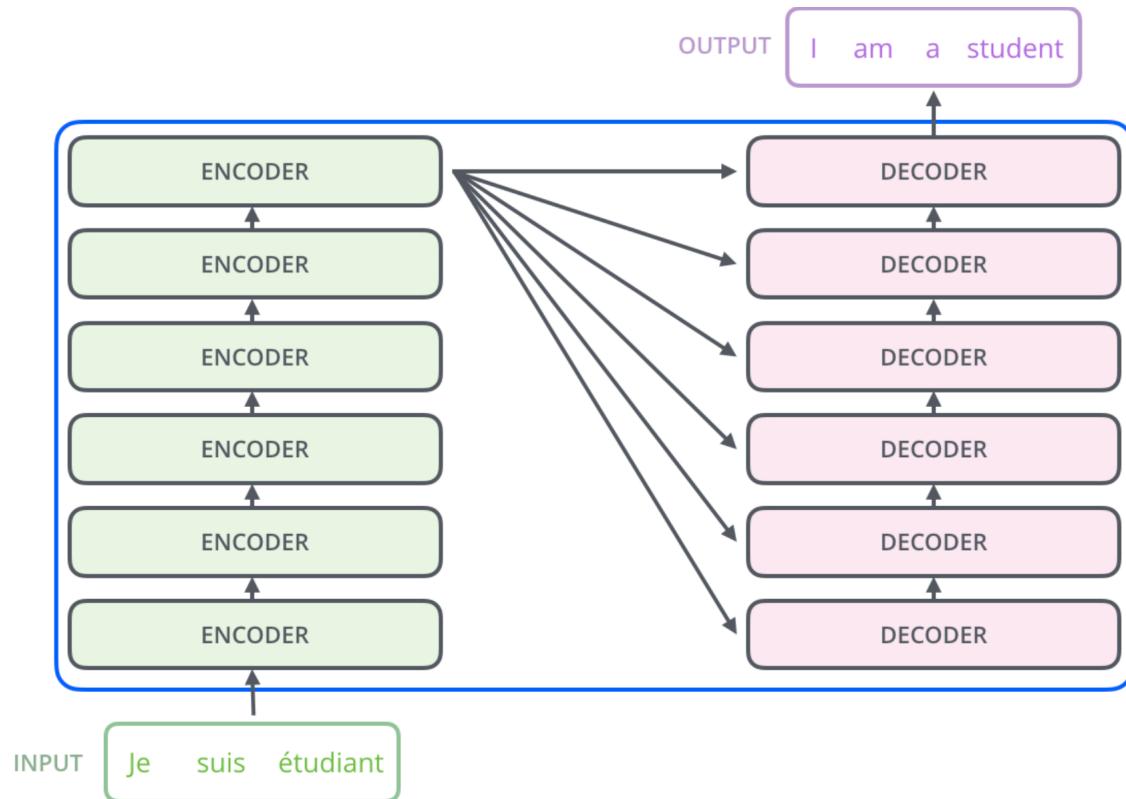
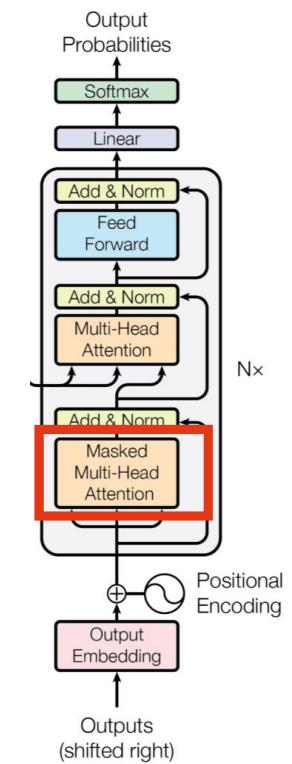


Figure 1: The Transformer - model architecture.

# Decoder

## What is Masked Multi-Head Attention?

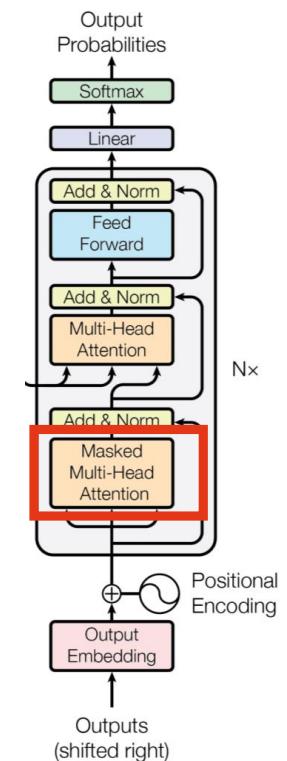
Our goal is to make the model causal: it means the output at a certain position can only depend on the words on the previous positions. The model **must not** be able to see future words.



# Decoder

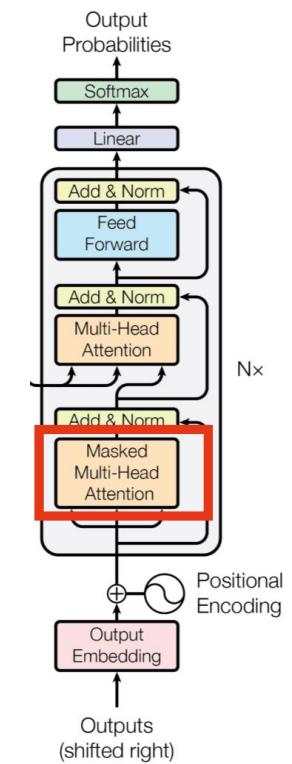
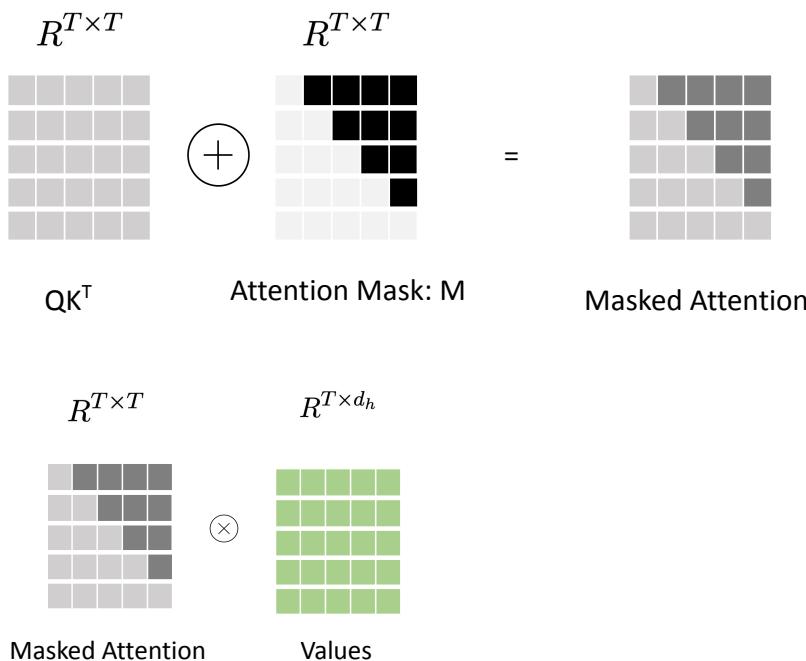
What is Masked Multi-Head Attention?

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	-	0.134	0.148	0.179	0.152
CAT	0.124	0.278	-	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

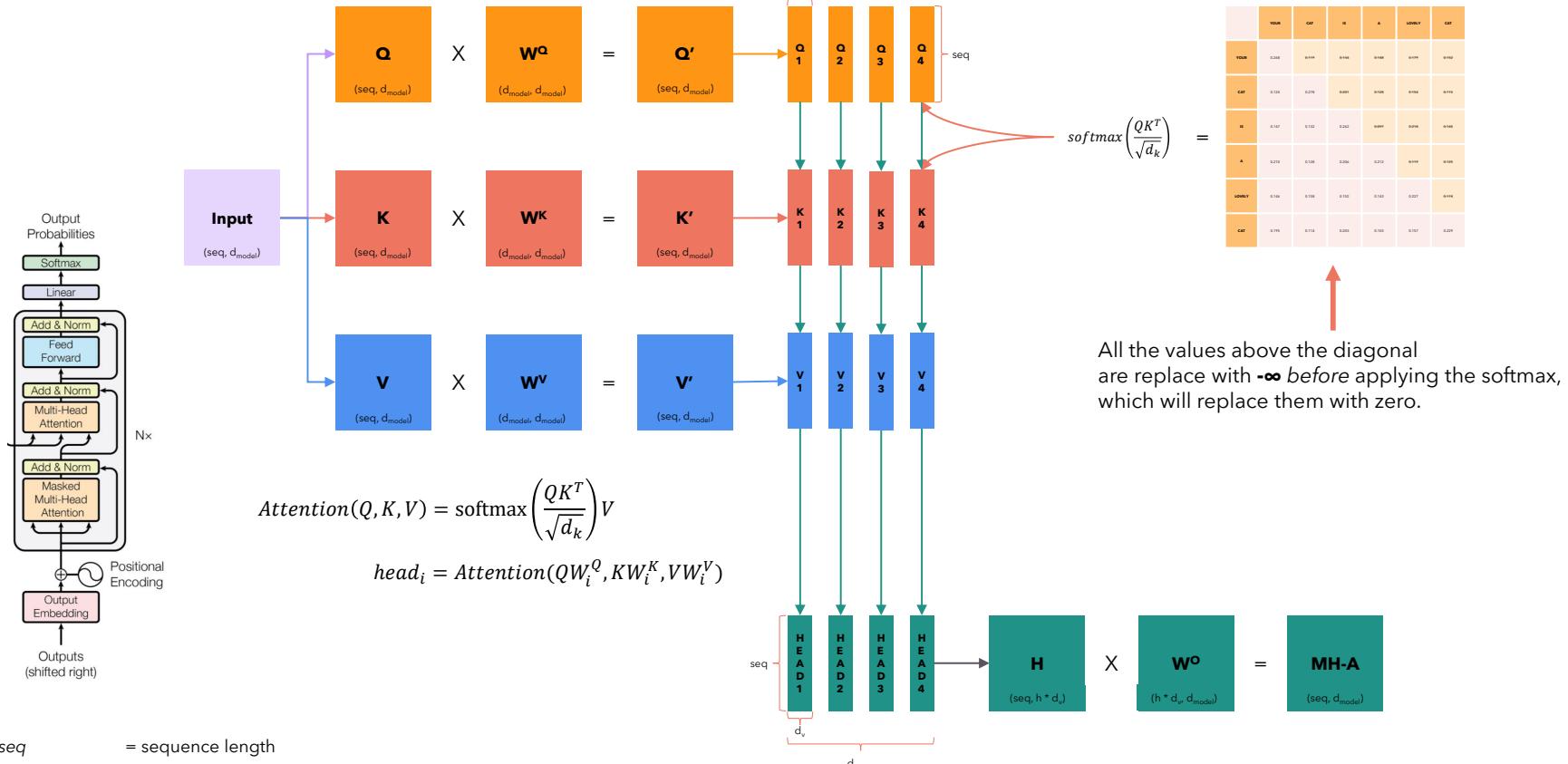


# Decoder

What is Masked Multi-Head Attention?



# Decoder



# Training



I love you very much



Ti amo molto



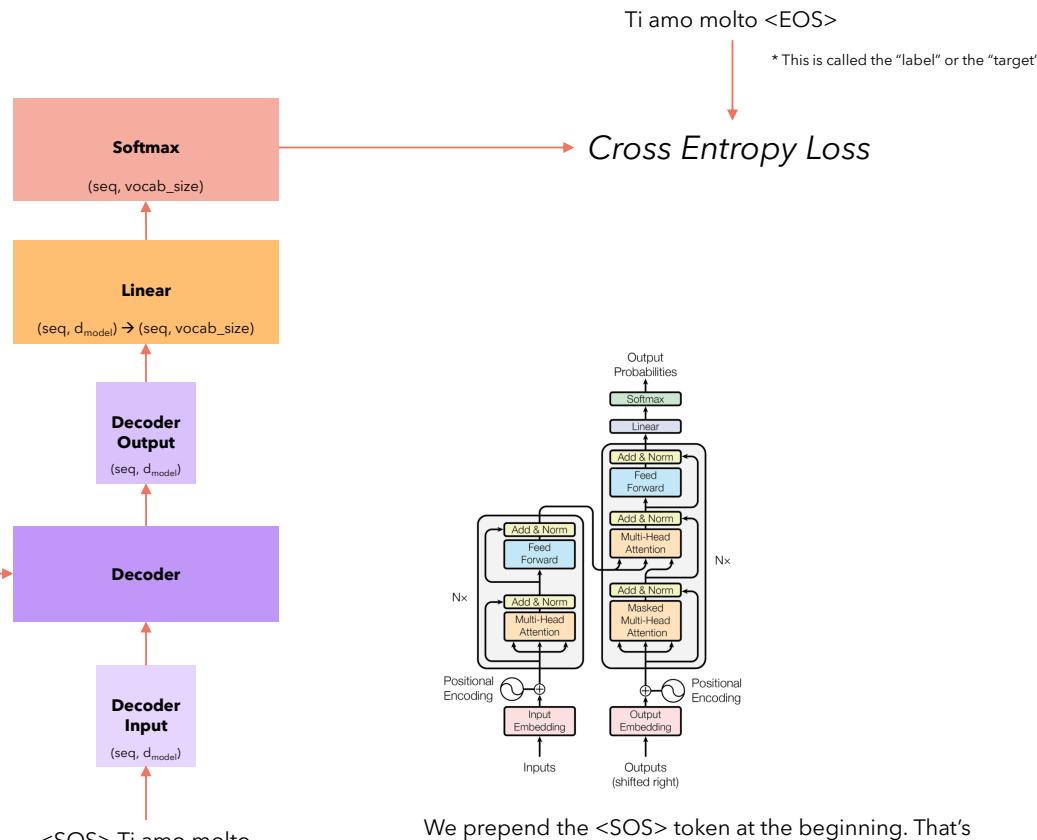
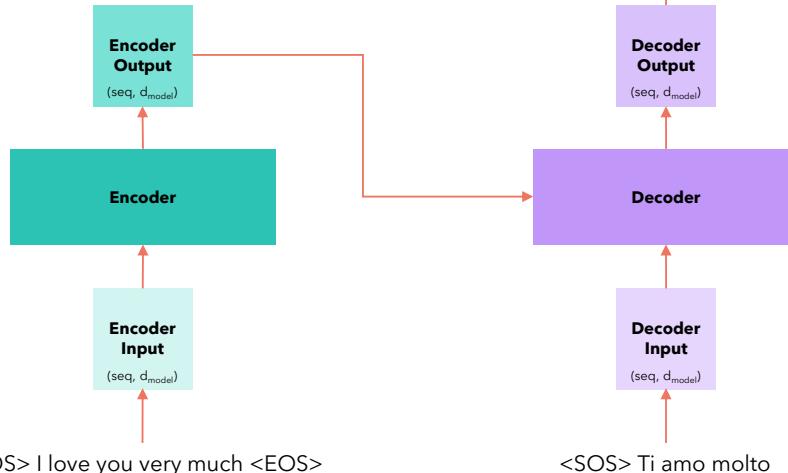
# Training

## Training

Time Step = 1

**It all happens in one time step!**

The encoder outputs, for each word a vector that not only captures its meaning (the embedding) or the position, but also its interaction with other words by means of the multi-head attention.

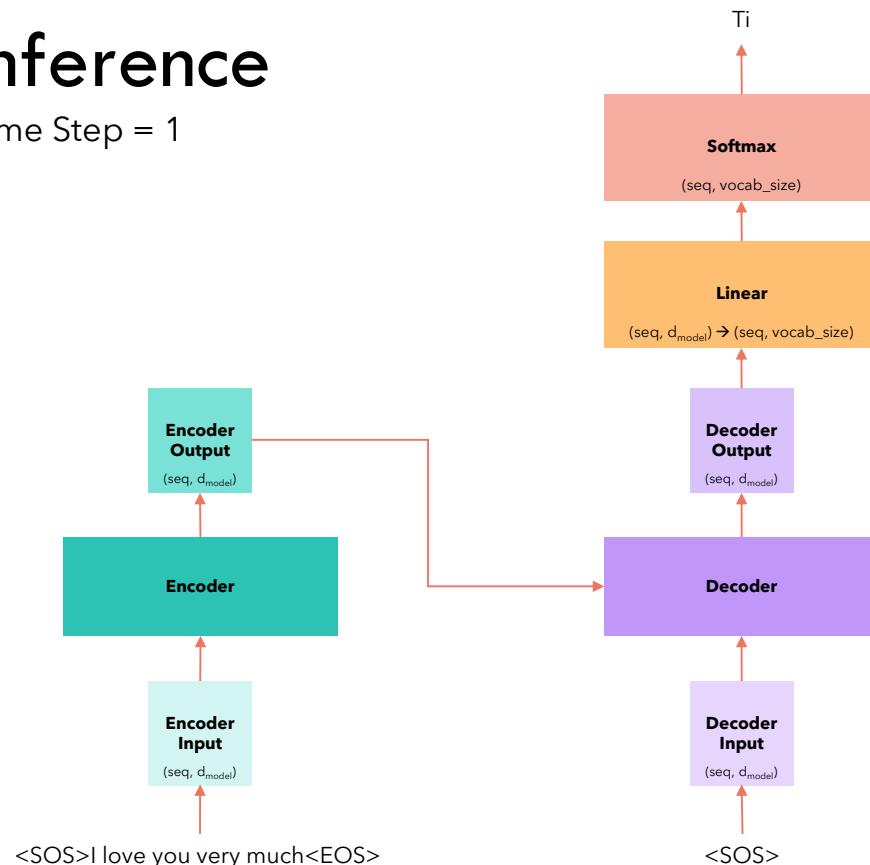


We prepend the **<SOS>** token at the beginning. That's why the paper says that the decoder input is shifted right.

# Inference

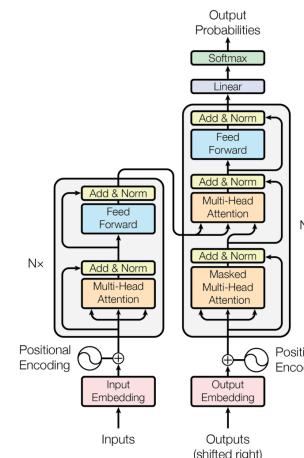
## Inference

Time Step = 1



We select a token from the vocabulary corresponding to the position of the token with the maximum value.

The output of the last layer is commonly known as **logits**



\* Both sequences will have same length thanks to padding

# Inference

## Inference

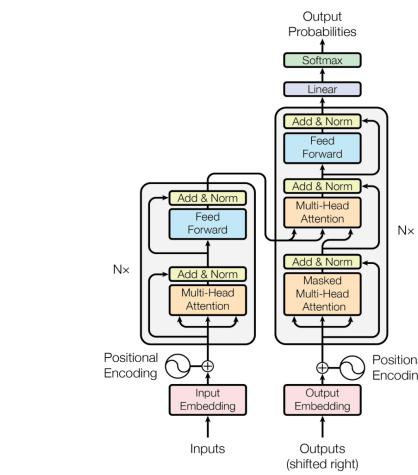
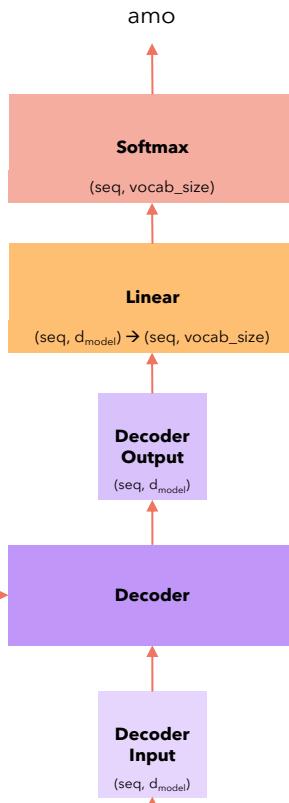
Time Step = 2

Use the encoder output from the first time step

<SOS>I love you very much<EOS>

<SOS> ti

Since decoder input now contains **two** tokens, we select the softmax corresponding to the second token.

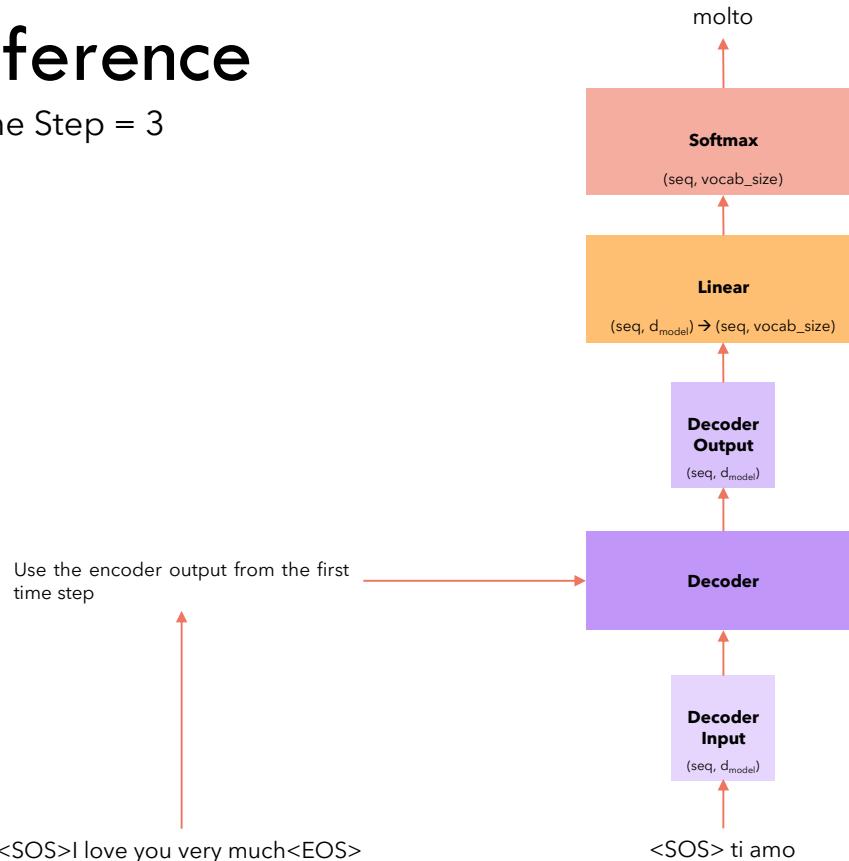


Append the previously output word to the decoder input

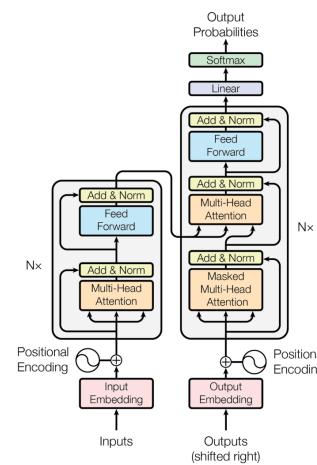
# Inference

## Inference

Time Step = 3



Since decoder input now contains **three** tokens, we select the softmax corresponding to the third token.

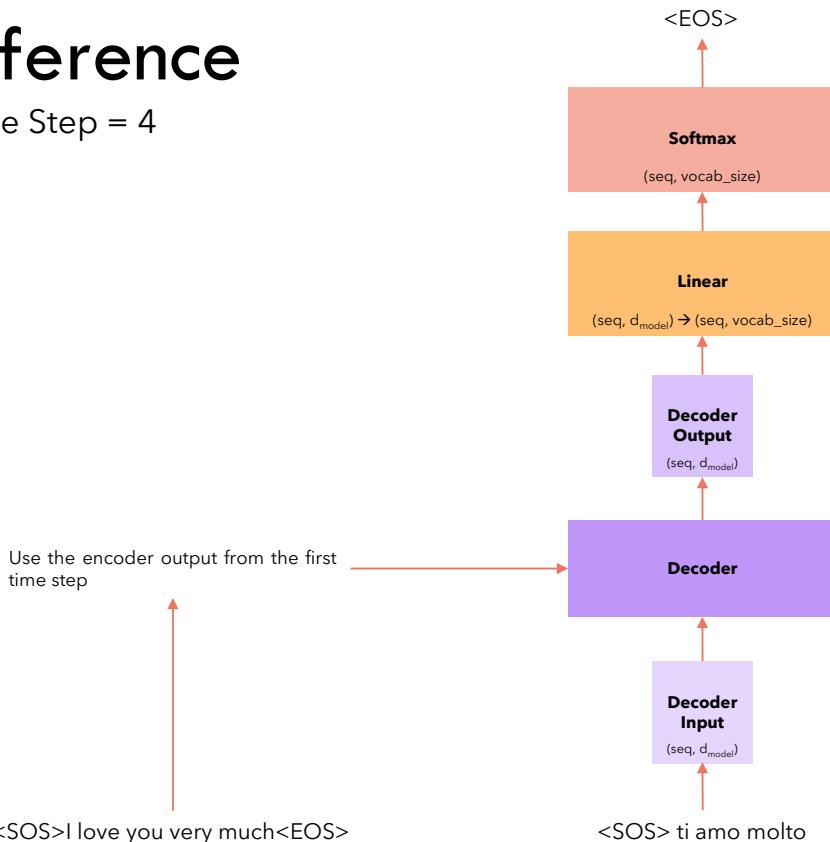


Append the previously output word to the decoder input

# Inference

## Inference

Time Step = 4



Since decoder input now contains **four** tokens, we select the softmax corresponding to the fourth token.

