

# CME 295: Transformers & Large Language Models



Afshin Amidi & Shervine Amidi



# Teaching staff



**Afshine and Shervine**

# Teaching staff



**Afshine**  
Centrale Paris ('16), MIT ('17)  
Uber, Google, Netflix



**Shervine**  
Centrale Paris ('16), Stanford ('19)  
Uber, Google, Netflix

# Welcome to CME 295!

## Goals.

1. Understand how **Transformers** work and how they **relate** to **LLMs**
2. Learn how **LLMs** are **trained** and used in various **applications**

# Welcome to CME 295!

## Goals.

1. Understand how **Transformers** work and how they **relate** to **LLMs**
2. Learn how **LLMs** are **trained** and used in various **applications**

## Audience.

- Interested in LLMs
  - Career goal
  - Personal project
  - "AI literacy" or curiosity
- Prerequisite: machine learning basics, linear algebra

# Logistics

## Date & time.

- Fridays from 3:30pm to 5:20pm
- Thornton 110

# Logistics

## Date & time.

- Fridays from 3:30pm to 5:20pm
- Thornton 110

## Details about the class.

- 2 units
- Letter or Credit/No credit
- Lectures are recorded
- Midterm (50% grade), scheduled on October 24th
- Final exam (50% grade), week of December 8th (exact date TBD)

# Material

**Class website.** [cme295.stanford.edu](http://cme295.stanford.edu)

- Contains syllabus & logistics
- Slides and recordings will be posted there

# Material

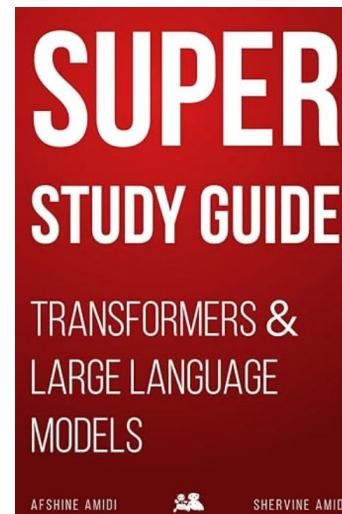
**Class website.** [cme295.stanford.edu](https://cme295.stanford.edu)

- Contains syllabus & logistics
- Slides and recordings will be posted there

**Class textbook.**

**Super Study Guide:**  
Transformers &  
Large Language Models

Book: <https://superstudy.guide>



# Material

CME 295 – TRANSFORMERS & LARGE LANGUAGE MODELS

<https://cme295.stanford.edu>

## VIP Cheatsheet:

### Transformers & Large Language Models

Afshin AMIDI and Shervine AMIDI

March 23, 2025

This VIP cheatsheet gives an overview of what is in the 'Super Study Guide: Transformers & Large Language Models' book, which contains ~600 illustrations over 250 pages and goes into the following concepts in depth. You can find more details at <https://superstudy.guide>.

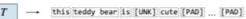
#### 1 Foundations

##### 1.1 Tokens

□ Definition – A token is an indivisible unit of text, such as a word, subword or character, and is part of a predefined vocabulary.

Remark: The unknown token `[UNK]` represents unknown pieces of text while the padding token `[PAD]` is used to fill empty positions to ensure consistent input sequence lengths.

□ Tokenizer – A tokenizer  $T$  divides text into tokens of an arbitrary level of granularity.

this teddy bear is reaaaally cute → 

Here are the main types of tokenizers:

Type	Pros	Cons	Illustration
Word	<ul style="list-style-type: none"><li>• Easy to interpret</li><li>• Short sequence</li></ul>	<ul style="list-style-type: none"><li>• Large vocabulary size</li><li>• Word variations not handled</li></ul>	
Subword	<ul style="list-style-type: none"><li>• Word roots leveraged</li><li>• Intuitive embeddings</li></ul>	<ul style="list-style-type: none"><li>• Increased sequence length</li><li>• Tokenization more complex</li></ul>	
Character	<ul style="list-style-type: none"><li>• No out-of-vocabulary concerns</li><li>• Small vocabulary size</li></ul>	<ul style="list-style-type: none"><li>• Much longer sequence length</li><li>• Patterns hard to interpret because too low-level</li></ul>	
Byte			

Remark: Byte-Pair Encoding (BPE) and Unigram are commonly-used subword-level tokenizers.

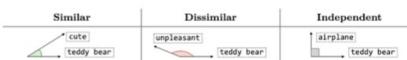
##### 1.2 Embeddings

□ Definition – An embedding is a numerical representation of an element (e.g. token, sentence) and is characterized by a vector  $\mathbf{x} \in \mathbb{R}^n$ .

□ Similarity – The cosine similarity between two tokens  $t_1, t_2$  is quantified by:

$$\text{similarity}(t_1, t_2) = \frac{t_1 \cdot t_2}{\|t_1\| \|t_2\|} = \cos(\theta) \in [-1, 1]$$

The angle  $\theta$  characterizes the similarity between the two tokens:



Remark: Approximate Nearest Neighbors (ANN) and Locality Sensitive Hashing (LSH) are methods that approximate the similarity operation efficiently over large databases.

#### 2 Transformers

##### 2.1 Attention

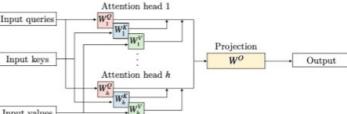
□ Formula – Given a query  $q$ , we want to know which key  $k$  the query should pay 'attention' to with respect to the associated value  $v$ .



Attention can be efficiently computed using matrices  $Q, K, V$  that contain queries  $q$ , keys  $k$  and values  $v$  respectively, along with the dimension  $d_k$  of keys:

$$\text{attention} = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

□ MHA – A Multi-Head Attention (MHA) layer performs attention computations across multiple heads, then projects the result in the output space.



It is composed of  $h$  attention heads as well as matrices  $W^Q, W^K, W^V$  that project the input to obtain query  $Q$ , keys  $K$  and values  $V$ . The projection is done using matrix  $W^O$ .

Remark: Grouped-Query Attention (GQA) and Multi-Query Attention (MQA) are variations of MHA that reduce computational overhead by sharing keys and values across attention heads.

##### 2.2 Architecture

□ Overview – Transformer is a landmark model relying on the self-attention mechanism and is composed of encoders and decoders. Encoders compute meaningful embeddings of the input that are then used by decoders to predict the next token in the sequence.

Link to PDF: <https://github.com/afshinea/stanford-cme-295-transformers-large-language-models>

# Class communications

## **Canvas.**

- Announcements
- Class discussions *via Ed*

# Class communications

Canvas.

- Announcements
- Class discussions via Ed

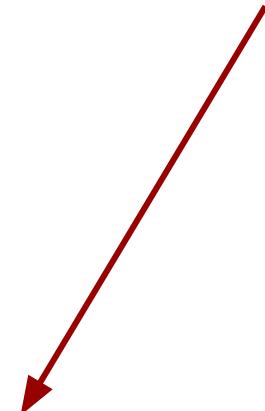
Any other general **inquiries / questions**, email us:

- [cme295-aut2526-staff@lists.stanford.edu](mailto:cme295-aut2526-staff@lists.stanford.edu)
- [afshine@stanford.edu](mailto:afshine@stanford.edu), [shervine@stanford.edu](mailto:shervine@stanford.edu)

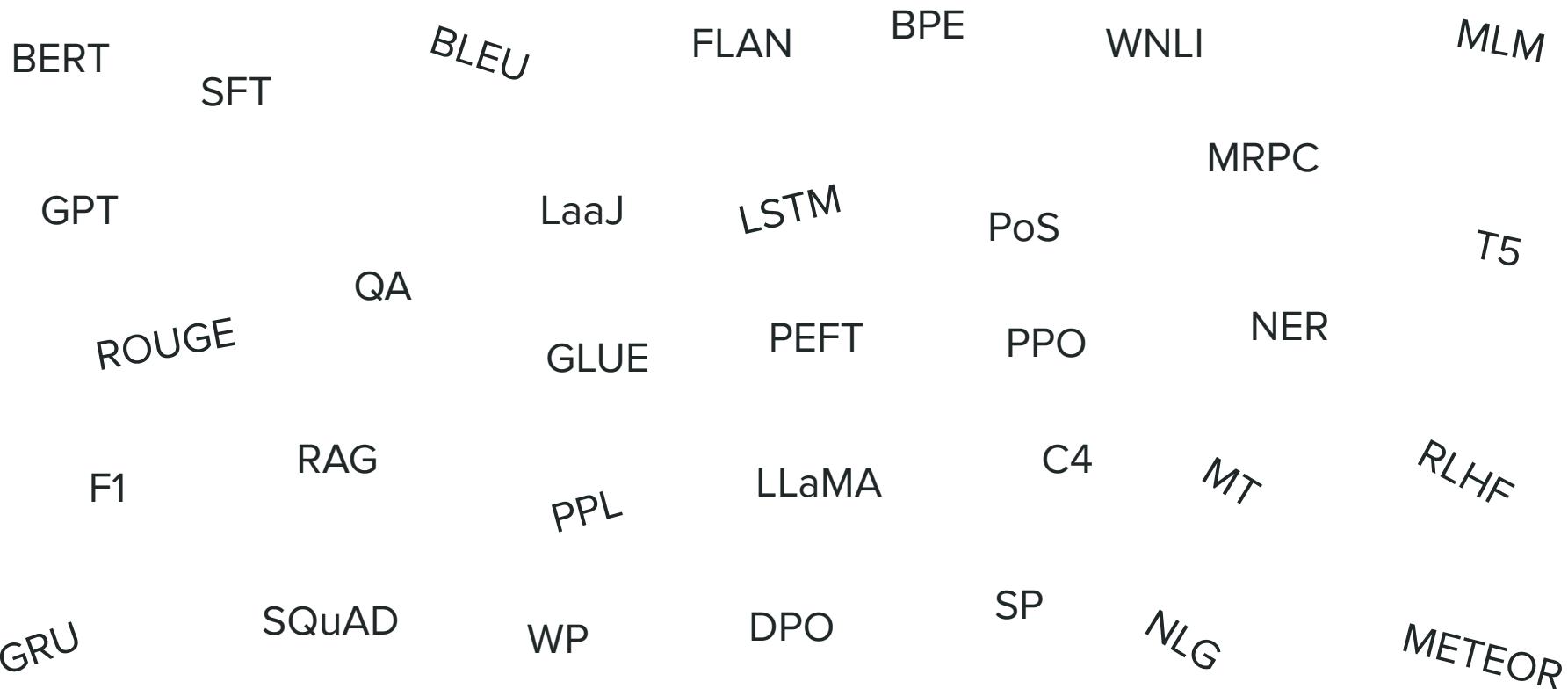
# Example of a slide in this class

Explanation of a CME 295 concept

Source & suggested reading, if interested



# Disclaimer before starting: many abbreviations....



# ...but don't worry!

BERT, T5, GPT, LLaMA

**Transformer-based models**

SFT, PEFT, FLAN, RL, RM, RLHF, PPO, DPO

**Training strategies**

MNLI, WNLI, C4, SQuAD, GLUE, MRPC

**Datasets**

LSTM, GRU, GloVe, BPE, CoT, ToT, SC, RAG

**Misc architectures & techniques**

NER, PoS, MLM, NSP, MT, QA, NLG

**Tasks**

F1, PPL, ROUGE, BLEU, METEOR, LaaJ, WER

**Metrics**



# CME 295

## Transformers & Large Language Models

NLP overview

Tokenization

Word representation

RNNs

Self-attention mechanism

Transformer architecture

End-to-end example

---

# NLP tasks overview

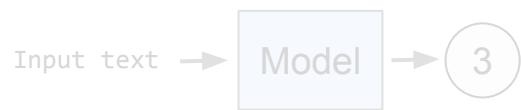
## Classification



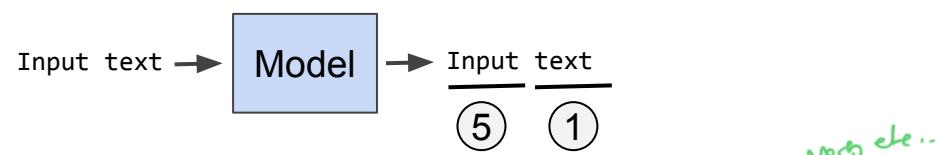
- Sentiment extraction
- Intent detection
- Language detection
- Topic modeling

# NLP tasks overview

## Classification



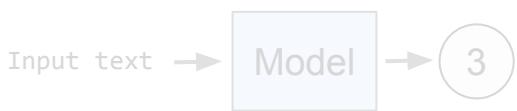
## “Multi”-classification



- Sentiment extraction
- Intent detection
- Language detection
- Topic modeling
- Part of speech tagging
- ✓ Named entity recognition
  - understanding noun, verb etc..
  - labelling some specific word (text)
- Dependency parsing
- Constituency parsing

# NLP tasks overview

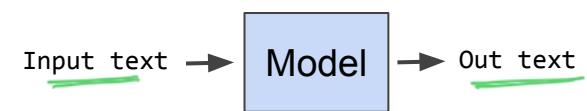
## Classification



## “Multi”-classification



## Generation

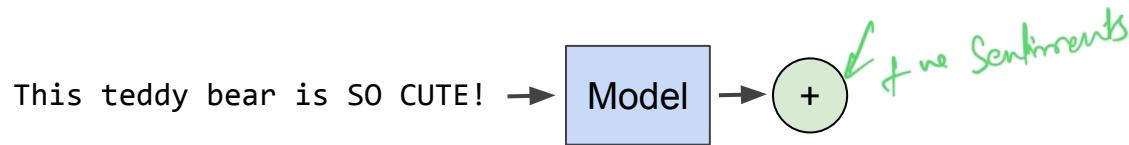


- Sentiment extraction
- Intent detection
- Language detection
- Topic modeling

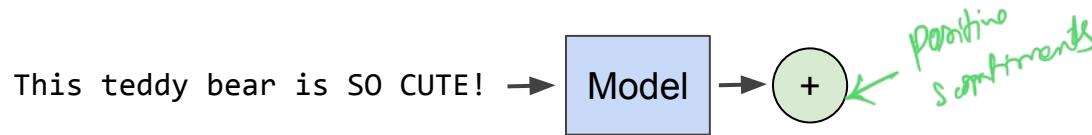
- Part of speech tagging
- Named entity recognition
- Dependency parsing
- Constituency parsing

- Machine translation
- Question answering
- Summarization
- Text generation

# NLP task: Sentiment Extraction



# NLP task: Sentiment Extraction



## Datasets

 Amazon reviews

 IMDB critiques

 Twitter

## Evaluation metrics

- Accuracy → % of observations that were correctly predicted?
- Precision → % of predicted positive that were correct?
- Recall → % of actually positive that were correct?
- F1 score → score that is a function of precision and recall

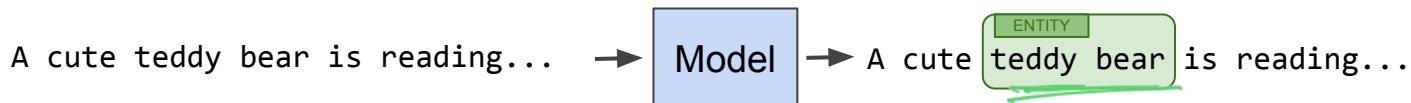
[used in case of  
imbalanced dataset]

# NLP task: Named Entity Recognition

multi classified form



# NLP task: Named Entity Recognition



## Datasets



Annotated Reuters newspaper (CoNLL-2003, CoNLL++)

## Evaluation metrics

- Accuracy
- Precision at a token level, per entity type
- Recall
- F1 score

# NLP task: Machine Translation



# NLP task: Machine Translation

Teddy bear → *Un ours*

A cute teddy bear is reading → Model → Un ours en peluche mignon lit

## Datasets

🇺🇸🇫🇷 WMT'14 English-French

🇺🇸🇩🇪 WMT'14 English-German

## Evaluation metrics

*bilingual Evaluation under study // How well translation stands w.r.t reference text*

- 🔴 BLEU → quality of text translated, similar to “precision”
- 🟡 ROUGE → quality of text generated, similar to “recall”  
*↳ suite of metrics*
- 🟢 Perplexity → quantifies how ‘surprised’ the model is to see some words together  
*(lower the better)*

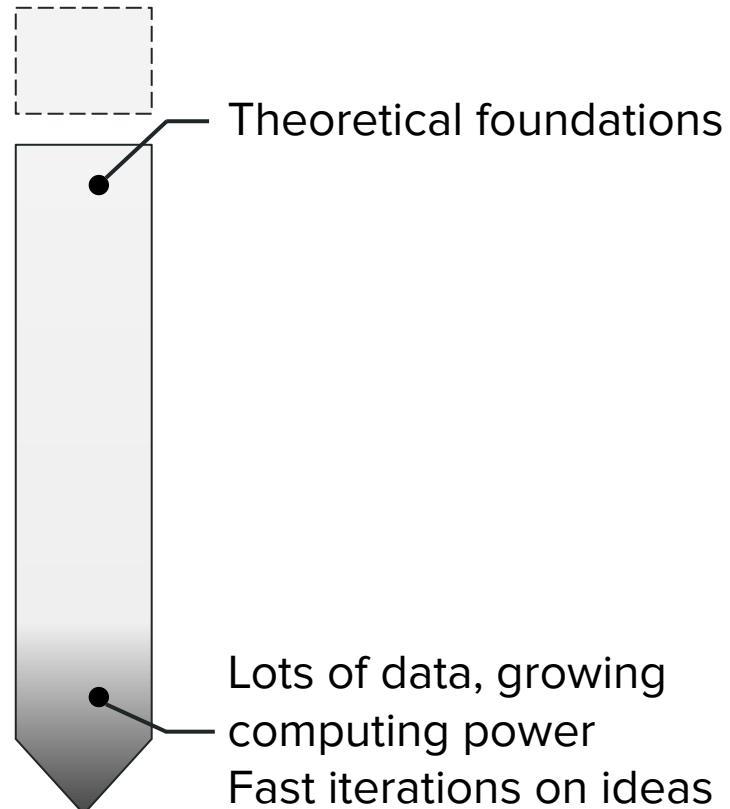
# High-level timeline

**1980s** Recurrent neural networks (RNNs)

**1997** Long short-term memory (LSTM)

**2013**  
**2017**  
**2020s**

Word2vec  
Transformers  
Large Language Models





# CME 295

## Transformers & Large Language Models

NLP overview

**Tokenization**

Word representation

RNNs

Self-attention mechanism

Transformer architecture

End-to-end example

---

# Tokenization

↓  
Putting the sentence to  
the form of model can  
understand.

i.e. pass the sentence  
in such a way a  
model can understand the  
text being provided.

A cute teddy bear is reading.

# Tokenization

A cute teddy bear is reading.

arbitrary

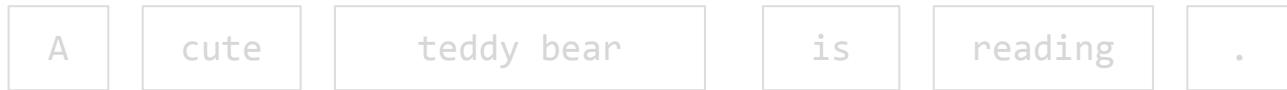


unit of texts  
called as "tokens"  
and method of  
separating out the tokens  
from the sentences which is  
going to pass to the model  
is called "Tokeniser"



# Tokenization

A cute teddy bear is reading.



Each word will be represented as a token.  
but need to separate them in meaningful way.  
(concerning word level tokenization)  
can end up with words which looks similar  
but considered as different tokens  
for that need to compute embeddings for  
similar but different words tokens  
and somehow the embedding becomes  
similar. ex: bear / bears  
run / runs

} having two different entity  
but having similar embedding

# Tokenization

↓  
to get rid of the limitation on word level

tokenid is ee sub-word

↓ leveraging roots of words

A cute teddy bear is reading.

bear / bears

play / playing

read / reading

cute

teddy bear

is

reading

A

cute

teddy

bear

is

reading

A

cute

teddy

bear

is

reading

A

cute

teddy

bear

is

reading

A

cute

ted

##dy

bear

is

read

###ng

.

pros: leveraging the root of words  
cons: longer sequence length

↑  
sub-word

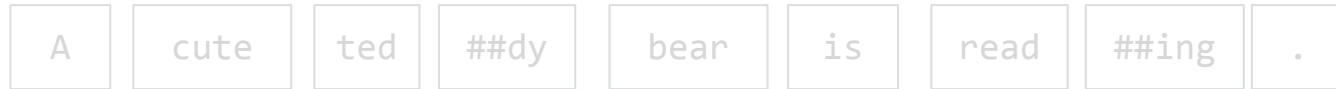
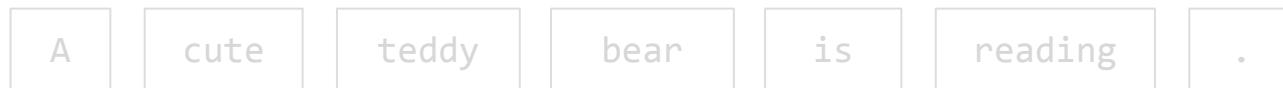
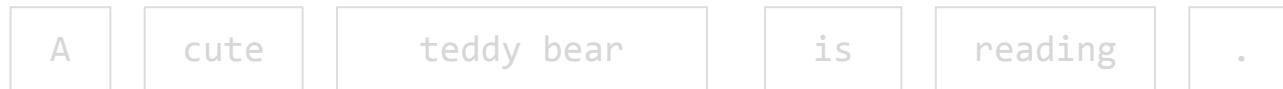
but cons of doing sub words  
Sequence will be longer

\* Complexity of large language model is — Sequence length - the more token will be passed to model more time it will take to Run/process.

# Tokenization

↓  
character level tokenization  
but again having much much  
larger length of sequence.  
and hard to represent the letters in  
character level tokenization

A cute teddy bear is reading.



# Tokenization summary

Method	Pros	Cons
<b>Word-level</b> <p><i>↑ super simple way of tokenizing.</i></p> <p><i>(dividing texts into arbitrary units)</i></p>	<ul style="list-style-type: none"> <li>Simple</li> <li>Interpretable</li> </ul> <p><i>cons times this</i></p>	<ul style="list-style-type: none"> <li>Risk of OOV <i>(OOV: out of vocabulary) i.e. when tokenization is being done and kept into training sets but some token which is not seen in training set but arising for inference set, but it is not available in training set go, just word token will be called as "out of vocabulary"</i></li> <li>Does not leverage knowledge of root <i>// leveraging root of words</i></li> </ul>
<b>Subword-level</b> <p><i>this is basically based on the case of Compositional understanding Problem: it is leveraging</i></p> <p>e.g. WordPiece, BPE</p>	<ul style="list-style-type: none"> <li>Leverages <u>common prefixes and suffixes</u></li> <li>Learned from the data</li> </ul>	<ul style="list-style-type: none"> <li>Risk of OOV, though less than word-level</li> </ul>
<b>Character-level</b> <p><i>↓ it might be considered based for language translation or multilingual problem statements; bcz it can improve the grammatical understanding of language.</i></p>	<ul style="list-style-type: none"> <li>Small chance of OOV</li> <li>Robust to casing and misspellings</li> </ul>	<ul style="list-style-type: none"> <li>Makes computations slower</li> <li>Embeddings not interpretable</li> </ul>



# CME 295

## Transformers & Large Language Models

NLP overview

Tokenization

**Word representation**

RNNs

Self-attention mechanism

Transformer architecture

End-to-end example

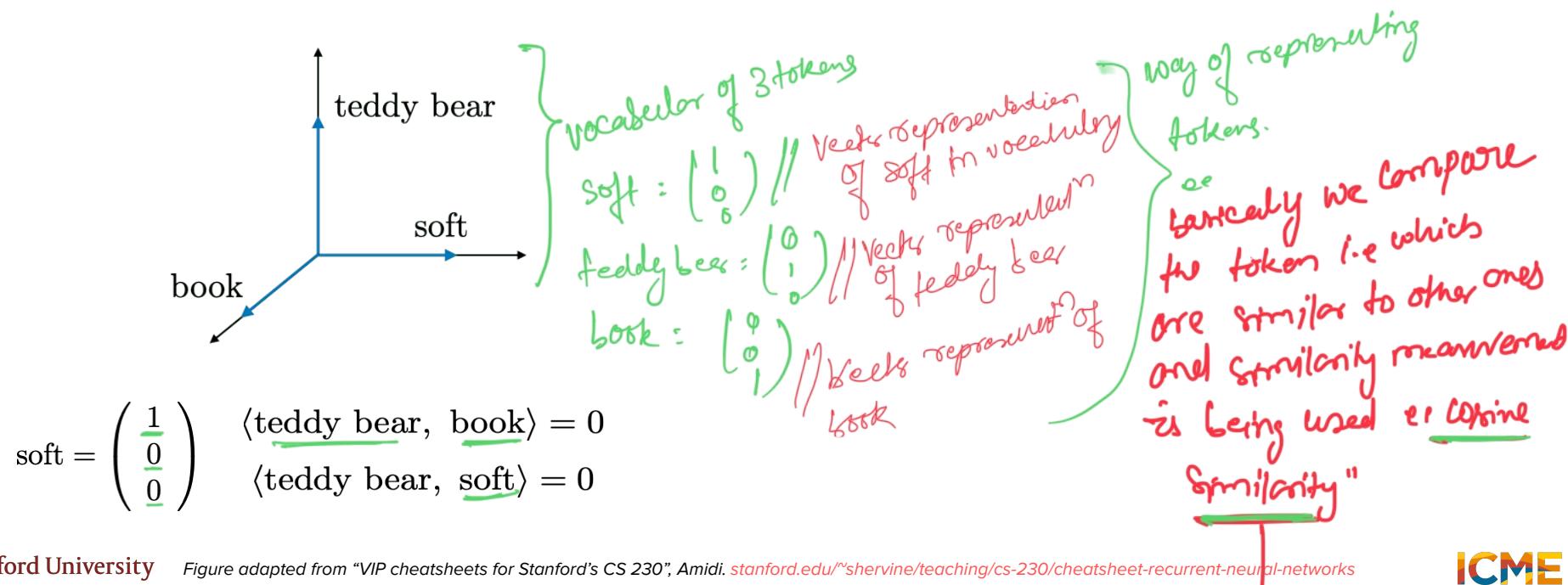
---

# Token representations

## Motivation

in the form of model to understand the sentences, we need to have representations of token — called “Token representation”

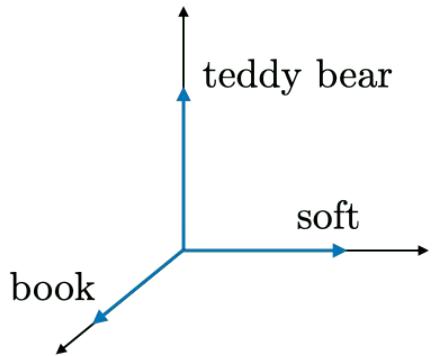
### Naive (one-hot) encoding



# Token representations

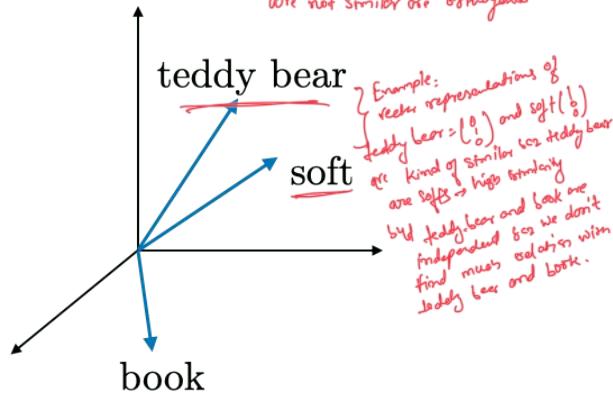
## Motivation

### Naive (one-hot) encoding



$$\text{soft} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad \langle \text{teddy bear}, \text{book} \rangle = 0$$
$$\langle \text{teddy bear}, \text{soft} \rangle = 0$$

### Learned embedding



$$\text{soft} = \begin{pmatrix} 0.95 \\ 0.32 \\ 0.01 \end{pmatrix} \quad \langle \text{teddy bear}, \text{book} \rangle \sim 0$$
$$\langle \text{teddy bear}, \text{soft} \rangle \sim 1$$

so lexically similarity measurement is used for an intuitively what angle the vector is being made in m-dimensional space and vectors from another dimension will be checked using cosine similarity - which say how vector are similar/opposite/independent

but the problem is if we represent tokens in one-hot encoding fashion, end up all the vectors are orthogonal to one-another so, basically what we want is tokens which are similar will be having high similarity, creating are not similar are orthogonal"

# Word2vec

## Overview

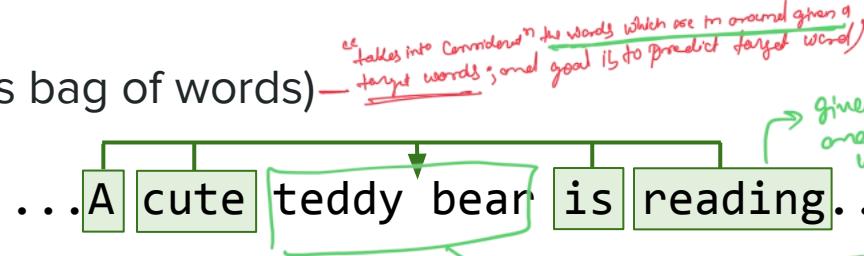
till here, we get to know that one-hot encoding is not great bce needs  
 can be orthogonal in m-dimensions using one-hot encoding &  
 basically sub-word tokenization learnt from data / Embeddings &  
 Word2Vec being used bcz it make sense in term of embedding.

- Neural network with a **proxy task** over billions of words worth of text
- Learns an embedding layer

because it already haven't kept in such a way → what words will be predicted after  
what word is already being kept in Vector Space which are Skip-gram, and Continuous  
word of bags  
(i.e. taking the leverage of texts we have  
and try to predict something part of text  
based on "context")

## Proxy tasks

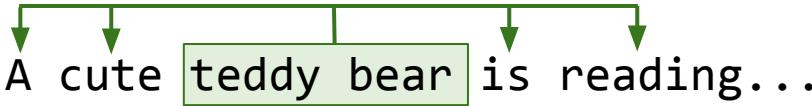
- CBOW (continuous bag of words)



given the around words,  
and predict the target  
word. like here - teddy  
bear has to be predicted  
based on input given  
as "A cute — —  
is reading"

- Skip-gram

go to the target words and  
try to predict the words  
around it.



opposite of CBOW  
target to predicted around  
word given "teddy" to

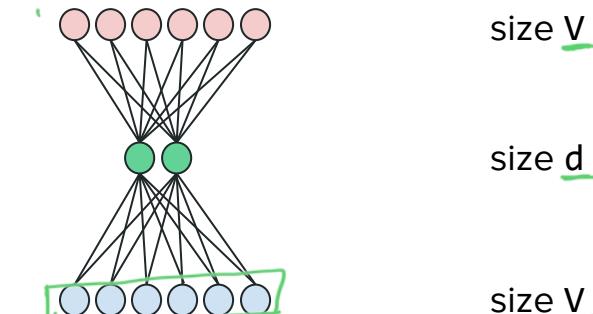
# Word2vec

Architecture → Predicting the next word

output

hidden

input



Input is given as vector of word  
and with the help of hidden state; output  
will be predicted. (Vocabulary is being  
given)

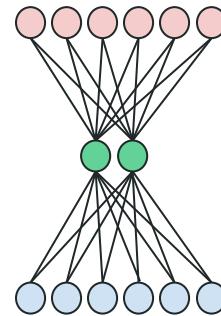
↓ basically learning the “word represented”  
through this task – try to consider  
as input and predict the  
next words.

So, it would be better to have  
a model which can understand  
the context being input to predict  
target words/around so, word2vec  
has introduced.

# Word2vec

## Example with predicting next word

A **cute** teddy bear is reading



↓ this is text ↓ it has go through some processing  
for model to understand

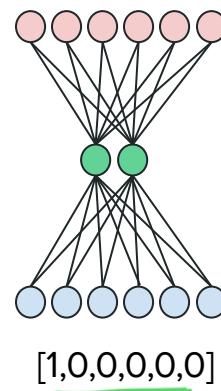
A **cute** teddy bear is reading

i.e. Embedding like one-hot encoding  
to predict the relationship b/w words for prediction  
using cosine similarity in n-dimension of vector space

# Word2vec

## Example with predicting next word

A cute teddy bear is reading



[1,0,0,0,0,0]

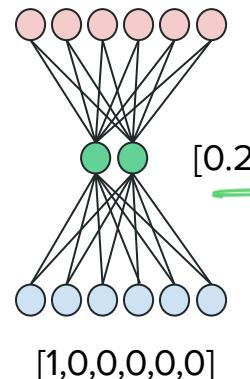
Embedding layer one-hot Encoding to vocabulary.

A cute teddy bear is reading

# Word2vec

## Example with predicting next word

A **cute** teddy bear is reading



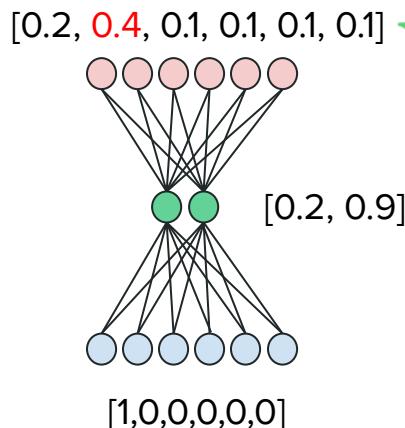
// One-hot encoding is passed through network;  
and weight are stored (after so much of  
mathematical calculation)  
known as "facing".  
this is the optimal weight  
stored

A **cute** teddy bear is reading

# Word2vec

## Example with predicting next word

A **cute** teddy bear is reading



finding the probability of words  
which will be highest probability  
will be optimal word for word  
"A" which is [1,0,0,...] had the  
embedding but cute might have [0,1,0,0]  
so that in n-dimensions of vector space, the  
cosine similarity has given the highest  
similarity value for prediction of next  
word.

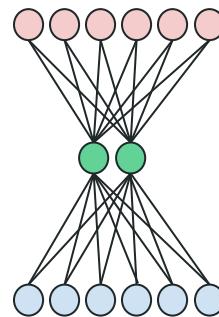
A **cute** teddy bear is reading

# Word2vec

## Example with predicting next word

A cute **teddy bear** is reading

↓  
some repeats with next word  
and it Continously goes on - remember  
it so much of computation

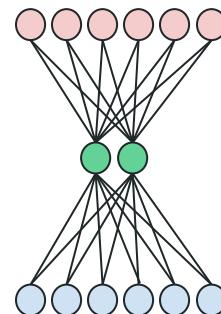


A **cute** teddy bear is reading

# Word2vec

## Example with predicting next word

A cute **teddy bear** is reading



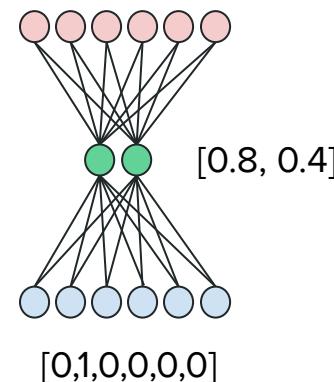
[0,1,0,0,0,0]

A **cute** teddy bear is reading

# Word2vec

## Example with predicting next word

A cute **teddy bear** is reading



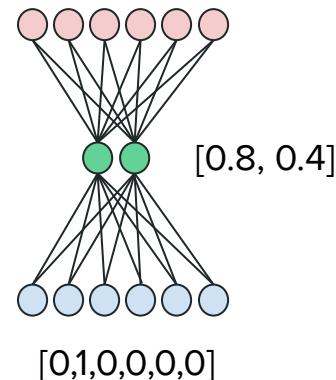
A **cute** teddy bear is reading

# Word2vec

## Example with predicting next word

A cute **teddy bear** is reading

[0.2, 0.2, 0.2, 0.1, **0.2**, 0.1]

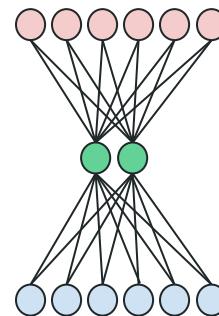


A **cute** teddy bear is reading

# Word2vec

## Example with predicting next word

A cute teddy bear **is** reading

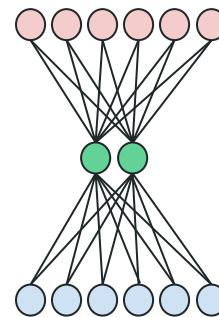


A cute **teddy bear** is reading

# Word2vec

## Example with predicting next word

A cute teddy bear is **reading**

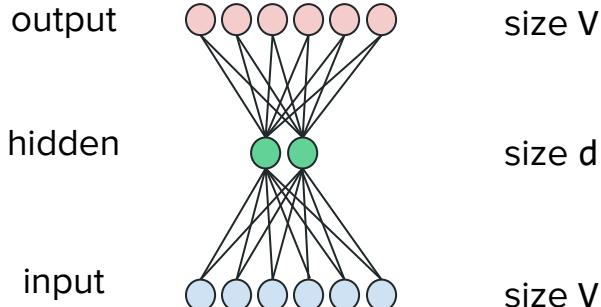


A cute teddy bear **is** reading

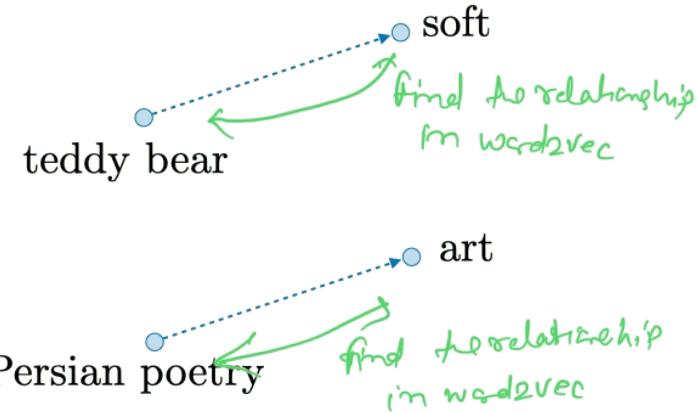
# Word2vec

after a lot of computation for predicting next word is bit hectic  
but this is how a main/ state of art architecture is being prepared  
for complex task.

A **cute** teddy bear is reading



A **cute** teddy bear is reading





# CME 295

## Transformers & Large Language Models

NLP overview

Tokenization

Word representation

RNNs

Self-attention mechanism

Transformer architecture

End-to-end example

---

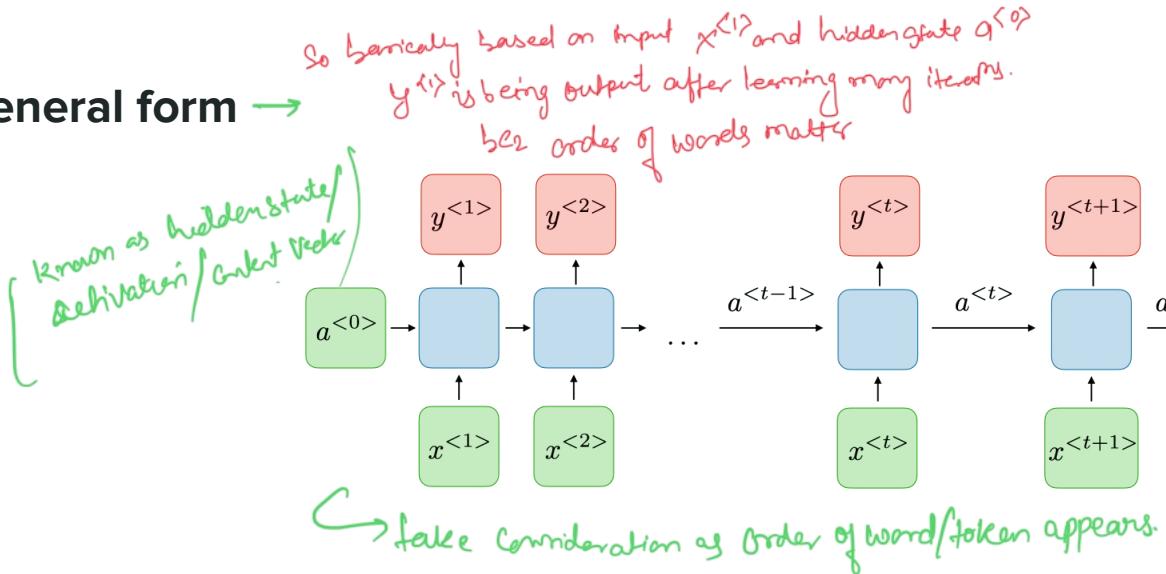
# Recurrent Neural Networks (RNNs)

## Overview

- First introduced in the 80s
- Class of neural networks where connections form a temporal sequence

So, we have a naive and state of art method to predict the next words/arounds but we wanted to make it complicated by not only predicts the words but sentences/pieces of sent/absolute relation which is given in puts i.e. we need to find the best relation word inputs → which does by Sequential model known as "RNN"

## General form →

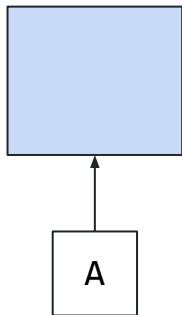


It basically captures the sequential nature of how text appears.  
(Instead of predicting words at a time, it also keeps hidden representation of the sentence and considers tokens one at a time.)

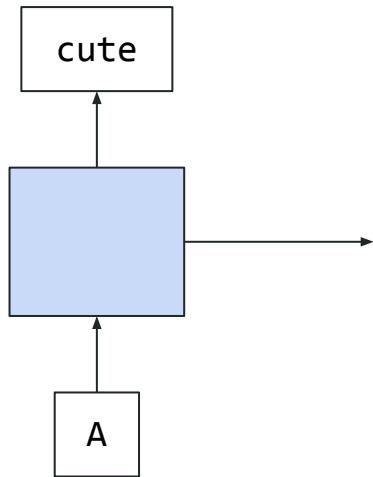
# Recurrent Neural Networks (RNNs)

A

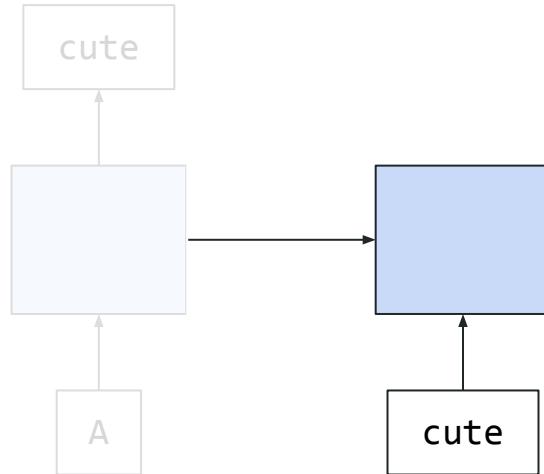
# Recurrent Neural Networks (RNNs)



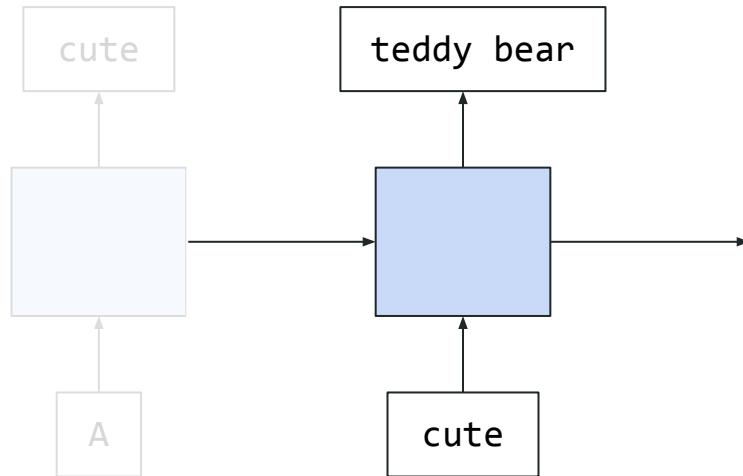
# Recurrent Neural Networks (RNNs)



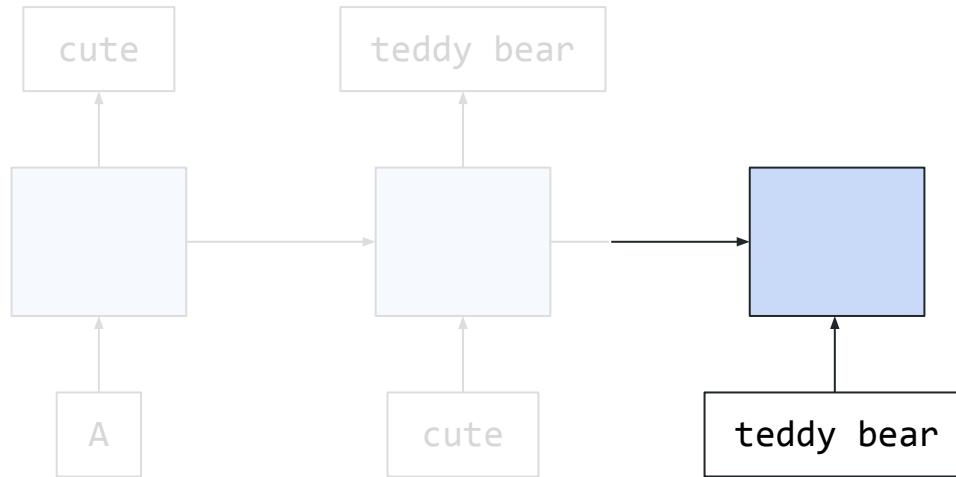
# Recurrent Neural Networks (RNNs)



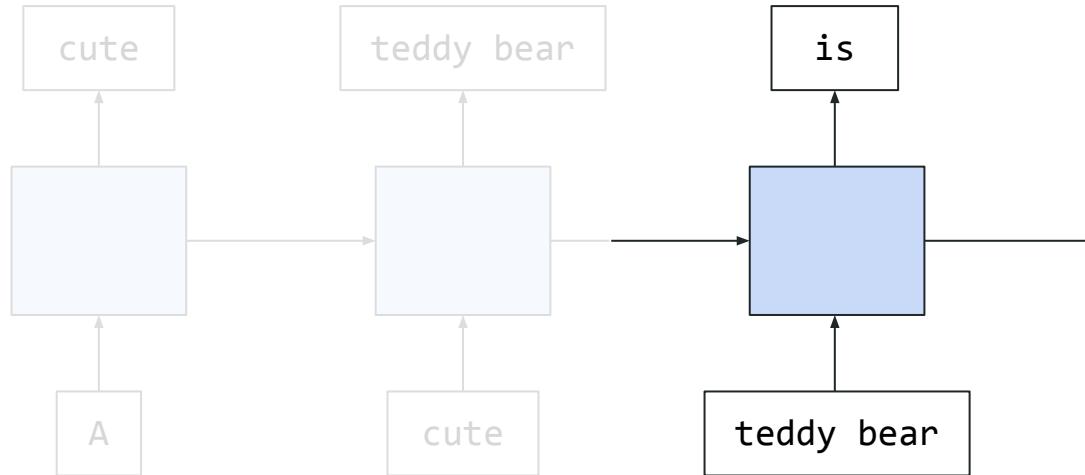
# Recurrent Neural Networks (RNNs)



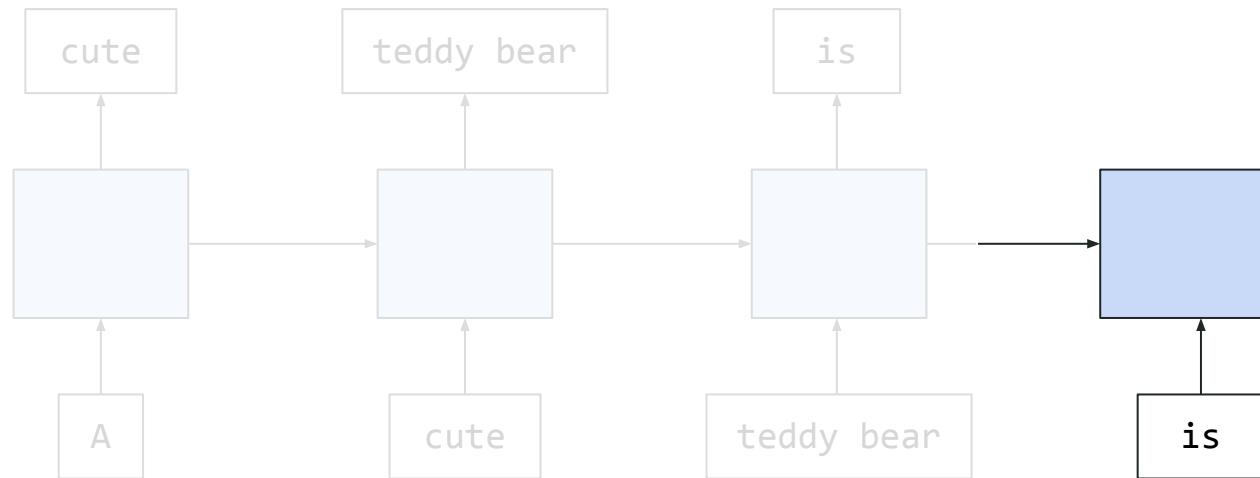
# Recurrent Neural Networks (RNNs)



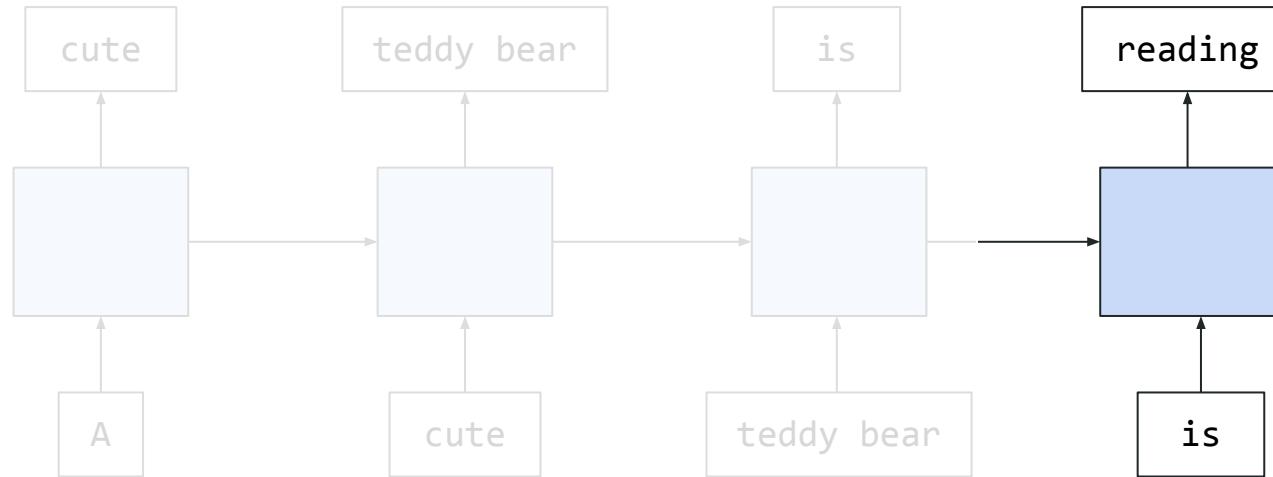
# Recurrent Neural Networks (RNNs)



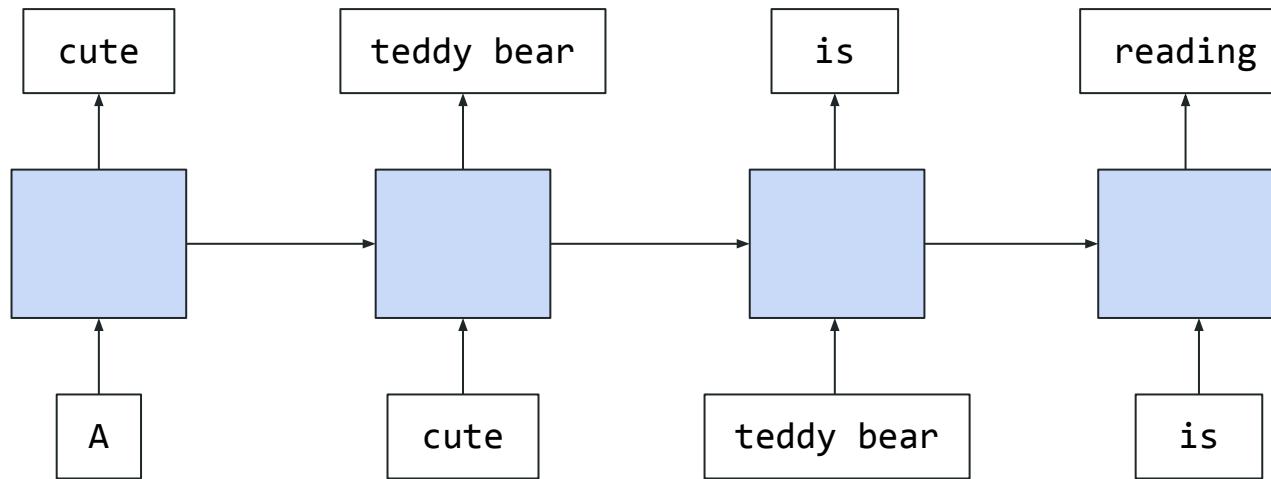
# Recurrent Neural Networks (RNNs)



# Recurrent Neural Networks (RNNs)



# Recurrent Neural Networks (RNNs)



# Recurrent Neural Networks (RNNs)

Cons of RNN is it keeps word representation based on input for expected output is on hidden state but how many hidden states will be it could be unmanageable so, it does have issue of "Long range dependencies" that's why LSTM comes

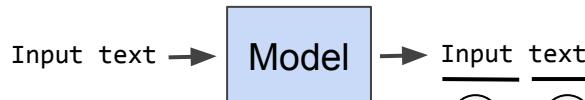
Generation  
based on problem statement or interest

Input text → Model → Out text

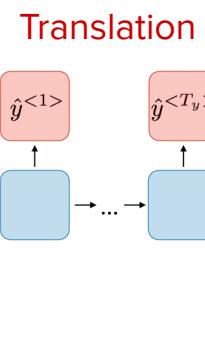
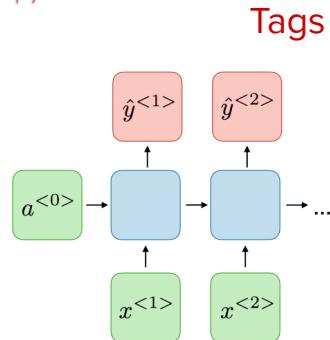
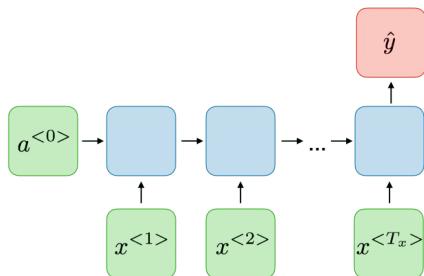
## Classification

## "Multi"-classification

## Generation



Sentiment  
{ sentiment over sentence/paragraph/review }



Cons of RNN is

it keeps word representation based on input for expected output is on hidden state but how many hidden states will be it could be unmanageable so, it does have issue of "Long range dependencies" that's why LSTM comes

Generation

based on problem statement or interest

Input text → Model → Out text

5      1

may be this for token of interest; like given input as sentences and predicted sentence or paragraph.

Opinion

Text

Source

# Long Short-Term Memory (LSTM)

## Overview

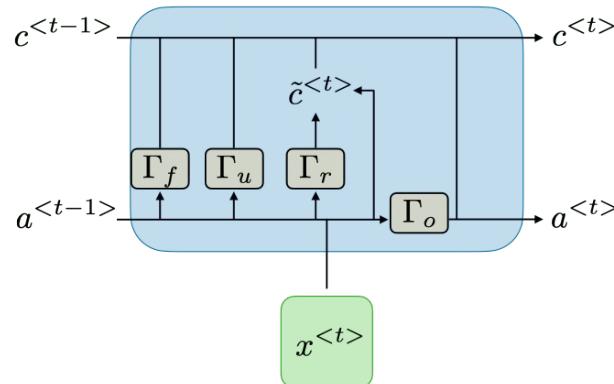
- Introduced in “Long short-term memory” (1997)
- Uses a more structured approach in the cell’s hidden state

*ways to keep track of “important” to remember on the top of hidden state (like RNN)  
basically LSTM keeps the information of important (memory) + hidden state of next predicted wsg*

*↓  
bcz it also need more computation now  
bcz impaired memory also have to “keep” and comes with vanishing gradient problem.*

*multiple gradients has to be kept and also it will have higher risk in memory.*

## General form



# Summary of main methods (non-exhaustive list)

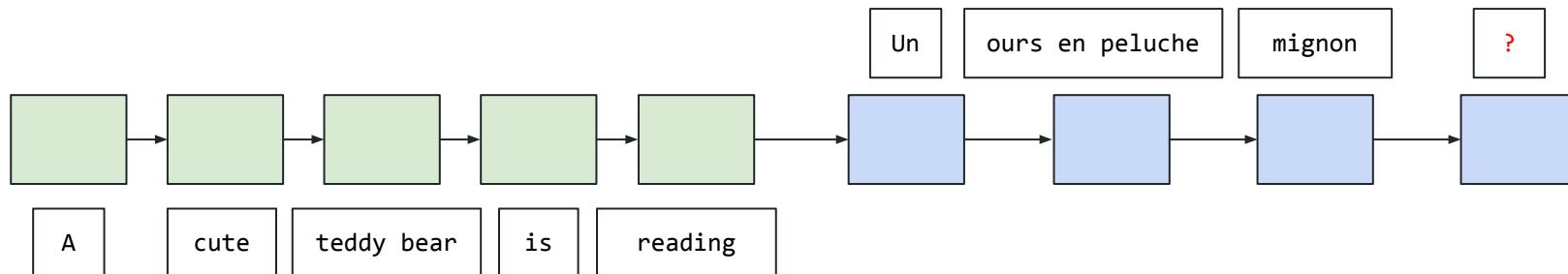
Method	Pros	Cons
<b>Word2vec</b> e.g. CBOW, Skip-gram	<ul style="list-style-type: none"><li>• <u>Very simple, yet powerful</u></li><li>• <u>Intuitive embeddings</u></li></ul>	<ul style="list-style-type: none"><li>• <u>Word order does not count</u></li><li>• <u>Embeddings not context aware</u></li></ul>
<b>Recurrent Neural Networks</b> e.g. traditional RNN, LSTM	<ul style="list-style-type: none"><li>• <u>Word order matters</u></li><li>• <u>State-of-the-art results</u></li></ul>	<ul style="list-style-type: none"><li>• <u>Vanishing gradient problem</u></li><li>• <u>Slow computations</u></li></ul>

# History of attention

- ← bez having problem of memorizing the words  
(memory)
- ↓ tries to find direct link b/w what we are trying to predict  
and something from past. It solve the problem of long-range dependencies
- Introduced in 2014
  - Translation tasks had a real issue with long-term dependencies
  - Seq2seq unable to "remember" what input sentence was saying

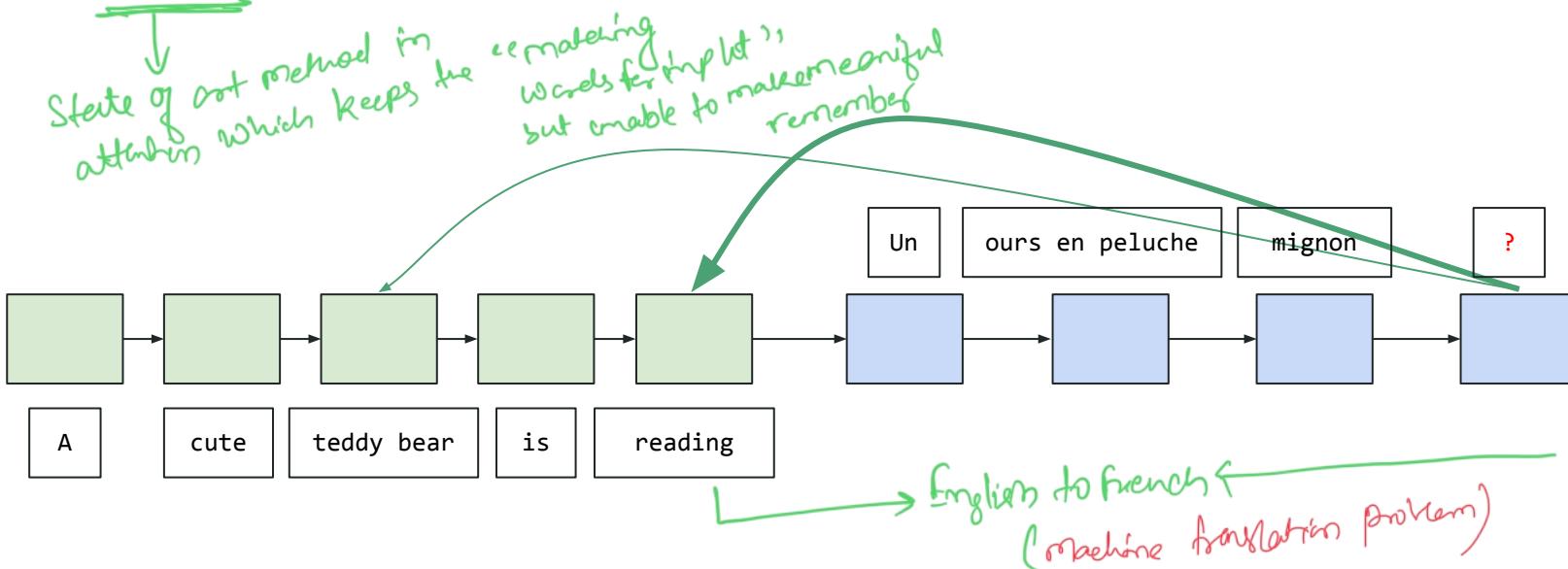
# History of attention

- Introduced in 2014
- Translation tasks had a real issue with long-term dependencies
- Seq2seq unable to "remember" what input sentence was saying



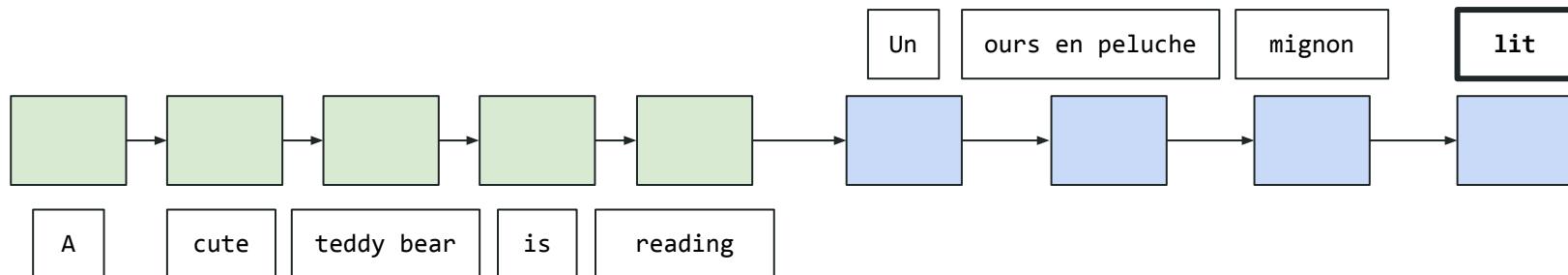
# History of attention

- Introduced in 2014
- Translation tasks had a real issue with long-term dependencies
- Seq2seq unable to "remember" what input sentence was saying



# History of attention

- Introduced in 2014
- Translation tasks had a real issue with long-term dependencies
- Seq2seq unable to "remember" what input sentence was saying





# Transformers & Large Language Models

NLP overview

Tokenization

Word representation

RNNs

**Self-attention mechanism**

Transformer architecture

End-to-end example

---

# Overview of the Transformer

- Introduced in the 2017 paper "**Attention is All You Need**"
- Relies on the **self-attention** mechanism
- State-of-the-art results on machine translation tasks

transformer also solve the problem of long-term dependency, memory, matching  
Content pair".

it says "move away from sequential processing  
of text and let the model have direct  
connection between all parts of texts at once"  
i.e. each token will be connected with all

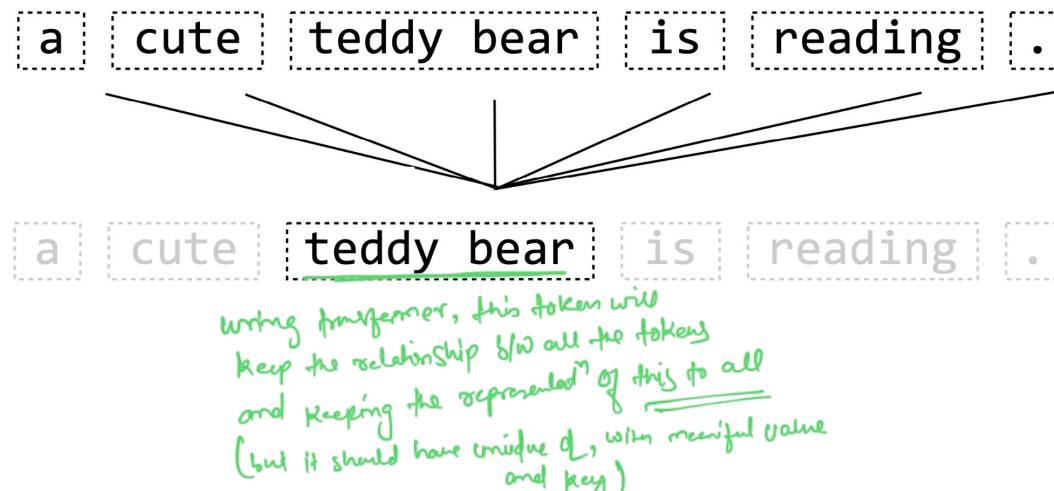
# Overview of the Transformer

- Introduced in the 2017 paper "**Attention is All You Need**"
- Relies on the **self-attention** mechanism
- State-of-the-art results on machine translation tasks

a cute teddy bear is reading .

# Overview of the Transformer

- Introduced in the 2017 paper "Attention is All You Need"
- Relies on the **self-attention** mechanism
- State-of-the-art results on machine translation tasks



# Attention mechanism

## Concept of Query, Key, Value

In the example,  
what are tokens are  
similar to teddy  
bear?

Is it a query → writing query tries to  
identify the similar key  
which should have been  
closer to already provided  
key-value pair dataset.

a cute **teddy bear** is reading .

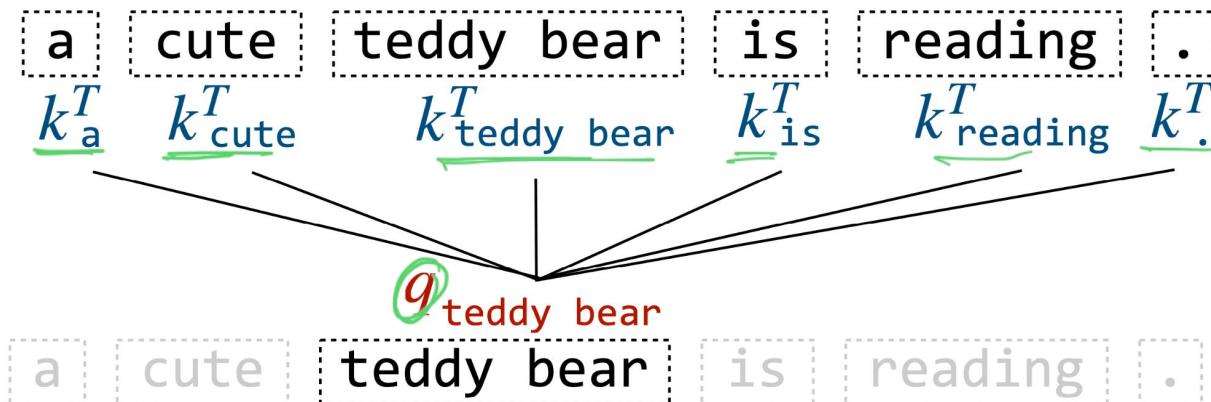
# Attention mechanism

Concept of **Query**, **Key**, **Value**

$q_{\text{teddy bear}}$   
a cute teddy bear is reading .

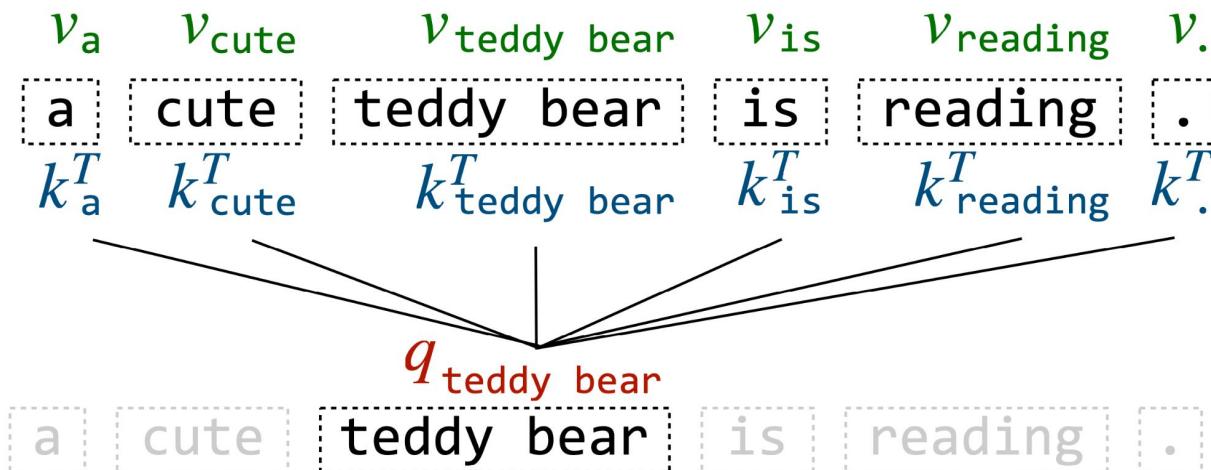
# Attention mechanism

Concept of **Query**, **Key**, **Value**



# Attention mechanism

Concept of **Query**, **Key**, **Value**



# Attention mechanism

benefit of using self-attention mechanism  
bcz input can be represented using matrix  
format and GPU loves matrix so, computation  
is not the problem.

Efficient computations with matrices:

$$\text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

this gives the value in form of probability  
which says find  $e^p$  based on given  
query, which value will be more important  
smaller important will have smaller weight  
and higher important will have higher weight

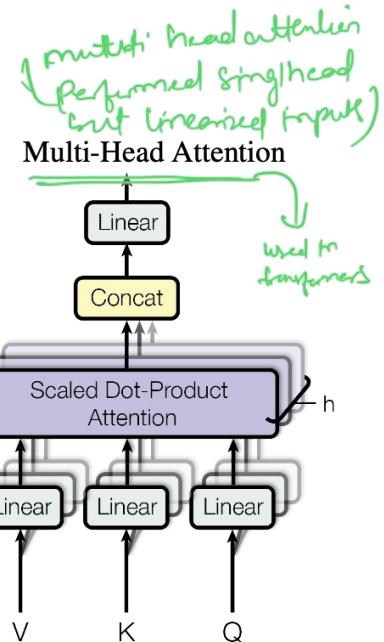
# Attention mechanism

Efficient computations with matrices:

$$\text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$



i.e. any condition  
you want to  
put as per  
your interest





# Transformers & Large Language Models

NLP overview

Tokenization

Word representation

RNNs

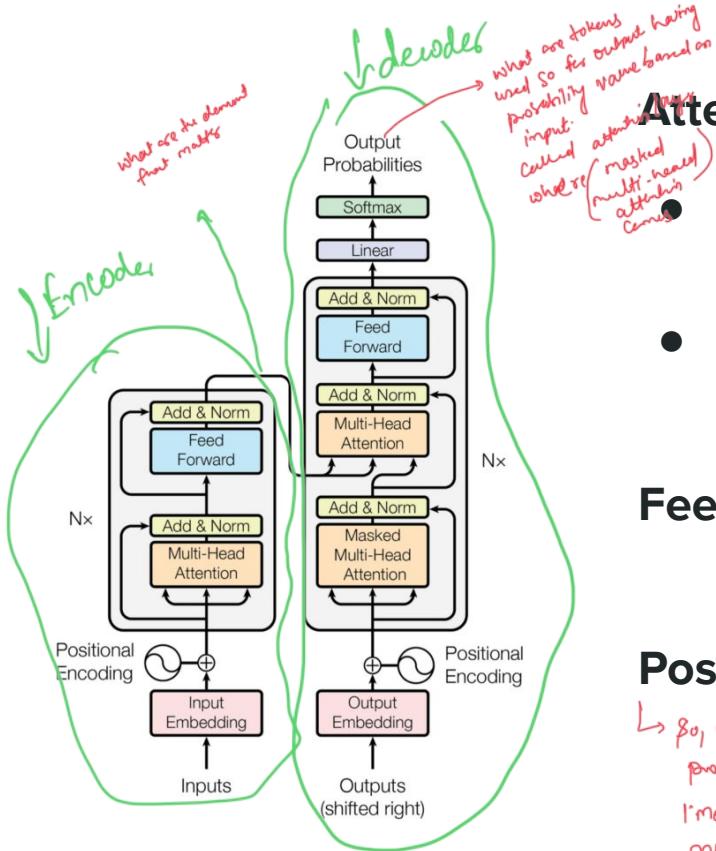
Self-attention mechanism

**Transformer architecture**

End-to-end example

---

# Transformer architecture



## Attention layer (MHA)

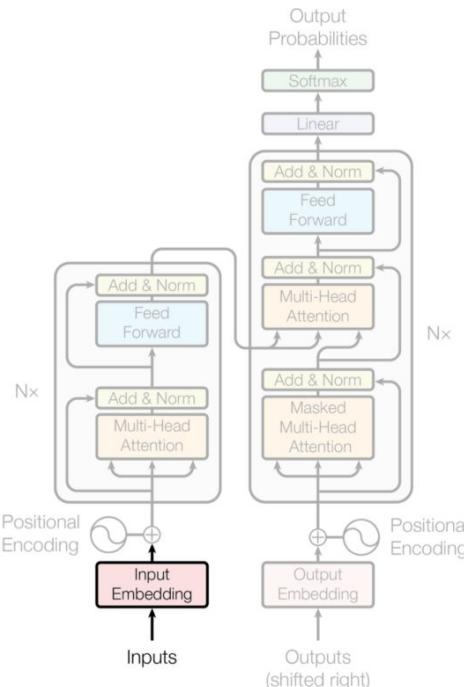
- Self-attention (Encoder-Encoder, Decoder-Decoder)
- Encoder-Decoder attention layer

## Feed Forward Neural Network (FFNN)

## Positional Encoding (PE)

↳ so, we have the direct link b/w one token to all which solves the problem of long-term dependencies but order did not care. I mean, in LSTM at least order were maintained but here not, so to get rid of these problems, we have positional embedding which takes care of ordering w/o tokens

# Input



## Overview

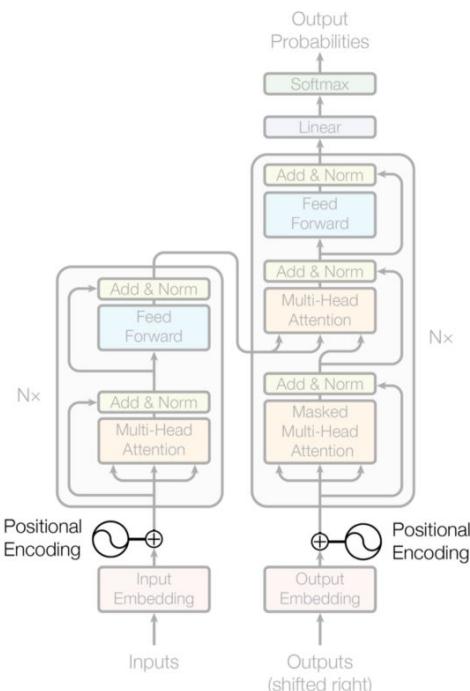
- Text is "tokenized"
- Learned embeddings for tokens

dividing into arbitrary units to learn the  
links all tokens to all  
and then learn the embeddings

## Parameters

- $V$ : vocabulary size
- $d_{model}$ : embedding dimensions

# ... with a trick!

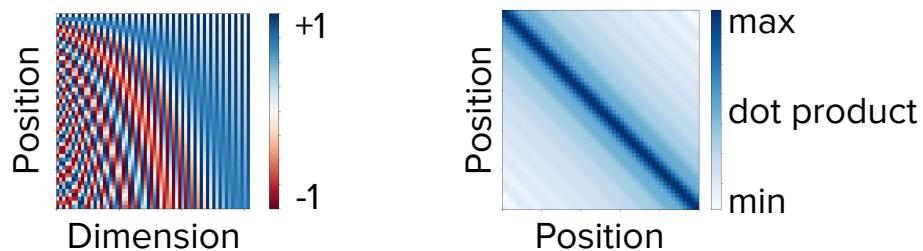


## Positional encoding

Idea:

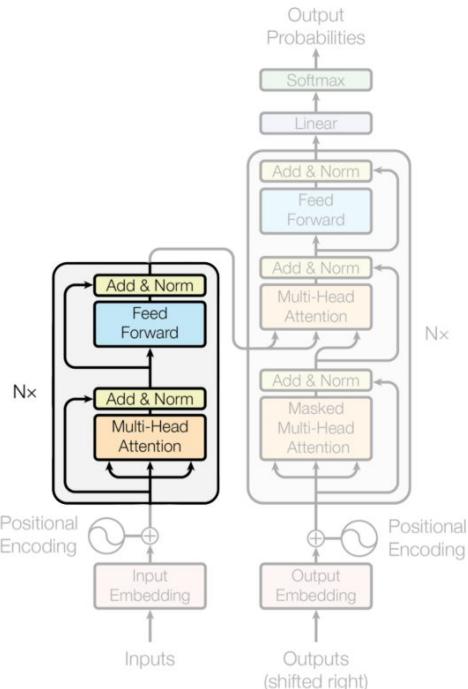
- Add **position information** to inputs
- Can be either learned or hardcoded

gives a position information for the embedding  
of token is being provided as input for  
relationship b/w all tokens to all.



Goal: let model understand relative input position

# Encoder



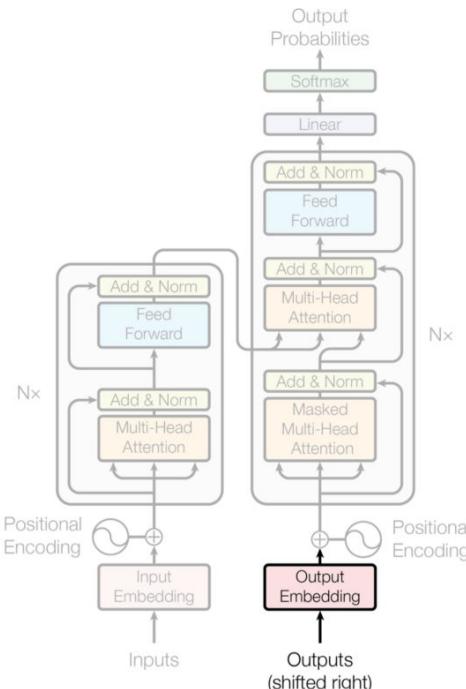
## Overview

- Encoder-Encoder attention / self-attention
- Feed Forward Neural Network
- Normalization layer

## Parameters

- ✓ N: layers stacked
- ✓ h: number of attention heads
- ✓ d\_FF, d\_key, d\_value: sub-layer dimension
- ✓ d\_model: embedding dimensions

# Output "shifted right"



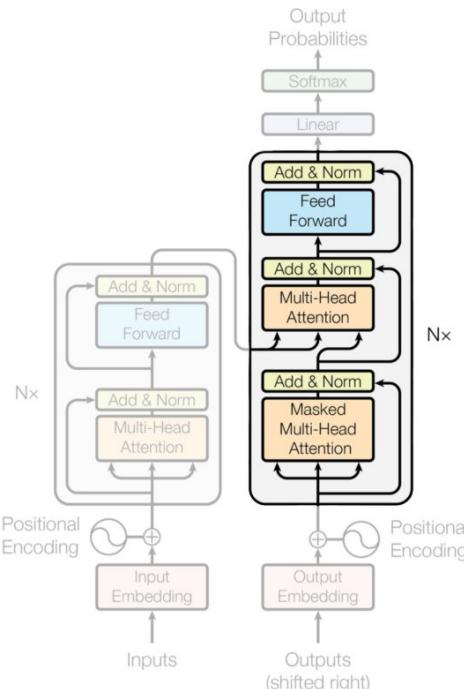
## Overview

- Learned embeddings for output tokens
- In practice, will start with [BOS] during translation

## Parameters

- $V$ : vocabulary size
- $d_{\text{model}}$ : embedding dimensions

# Decoder



## Overview

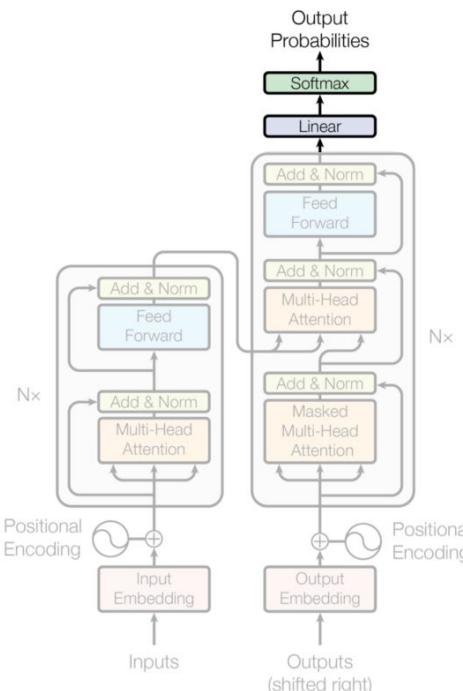
- Decoder-Decoder attention / self-attention
- Encoder-Decoder attention
- Feed Forward Neural Network
- Normalization layer

## Parameters

- $N$ : layers stacked
- $h$ : number of attention heads
- $d_{FF}$ ,  $d_{key}$ ,  $d_{value}$ : sub-layer dimension
- $d_{model}$ : embedding dimensions

# Output

## Overview

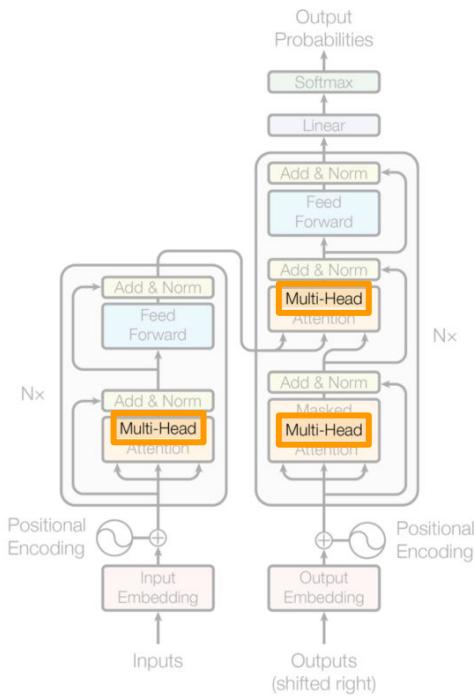


- Linear projection
- Classification problem that outputs probability of belonging to a class, where class = word

## Parameters

- $V$ : vocabulary size
- $d_{model}$ : embedding dimensions

# Computational tricks



## Multi-head attention

Idea:

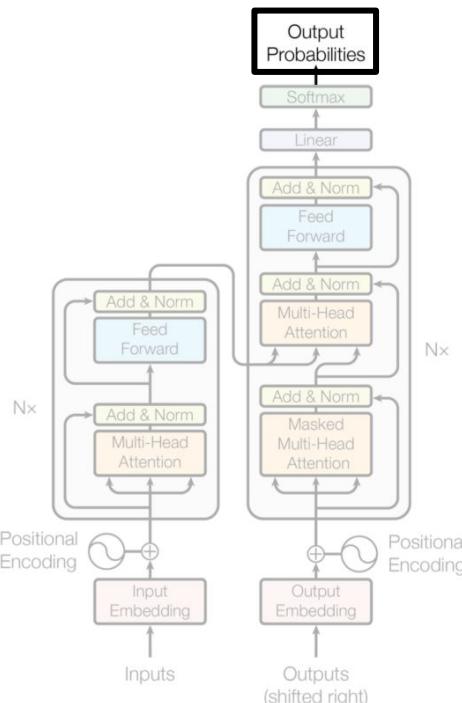
- Run **multiple self-attention layers in parallel**

Benefits:

- Enables the model to capture different attention features in parallel
- Comparison: multiple filters of a convolutional layer in computer vision



# Computational tricks



## Label smoothing //

In NLP, when word to predict some texts  
then it has kind of multiple choices so  
label smoothing tries to fit the problem of  
multiple choices and keeping one choice  
more particular or certain

Idea:

- **2015 vision paper:** overconfidence is bad
- Introduce **noise** in true labels

$$q(k|x) = \delta_{k,y} \rightarrow q'(k|x) = (1 - \epsilon)\delta_{k,y} + \epsilon u(k)$$

Benefits:

- General technique that prevents overfitting
- Improves accuracy and BLEU score



# Transformers & Large Language Models

NLP overview

Tokenization

Word representation

RNNs

Self-attention mechanism

Transformer architecture

**End-to-end example**

---

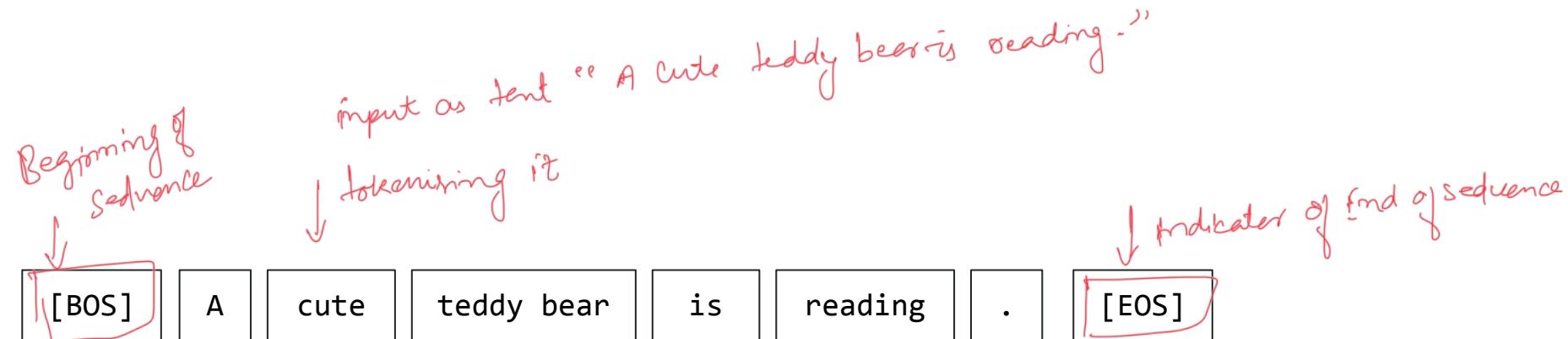
# Stitching all the pieces together with an example

A cute teddy bear is reading.

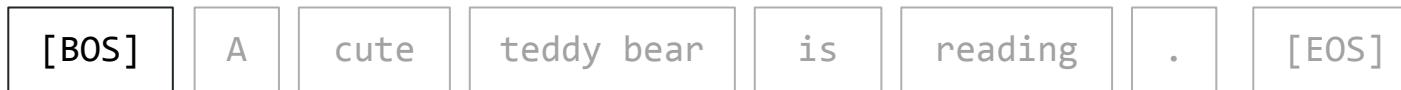
# Stitching all the pieces together with an example

A      cute      teddy bear      is      reading      .

# Stitching all the pieces together with an example



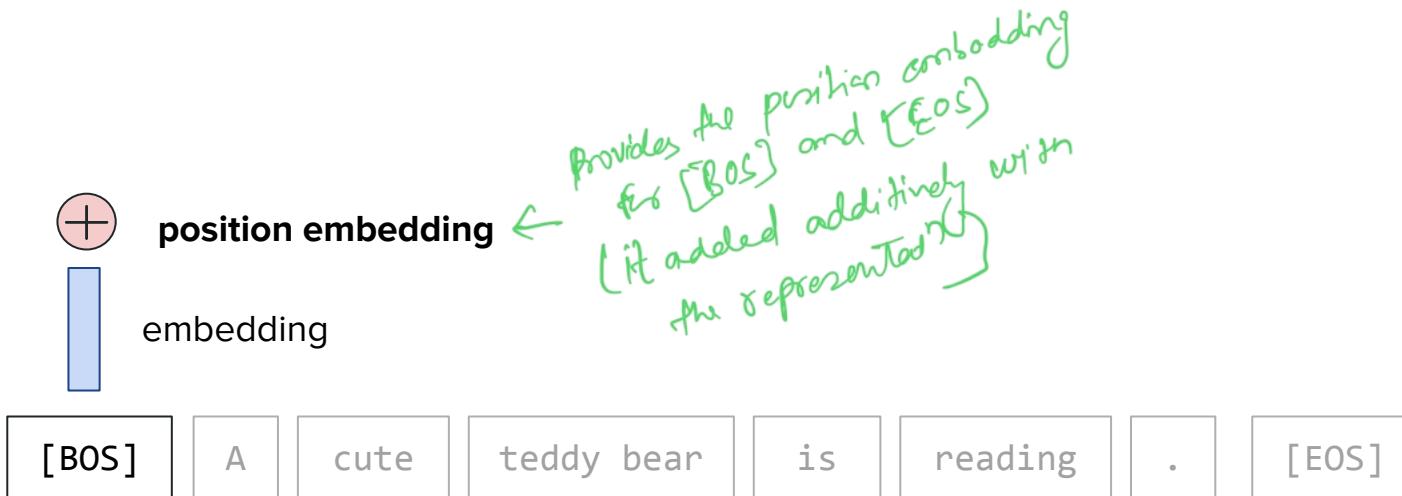
# Stitching all the pieces together with an example



# Stitching all the pieces together with an example



# Stitching all the pieces together with an example



# Stitching all the pieces together with an example



position-aware embedding

← bez using position Embedding (based on cosine, sine)  
we know where [BOS] will add with representation

[BOS]

A

cute

teddy bear

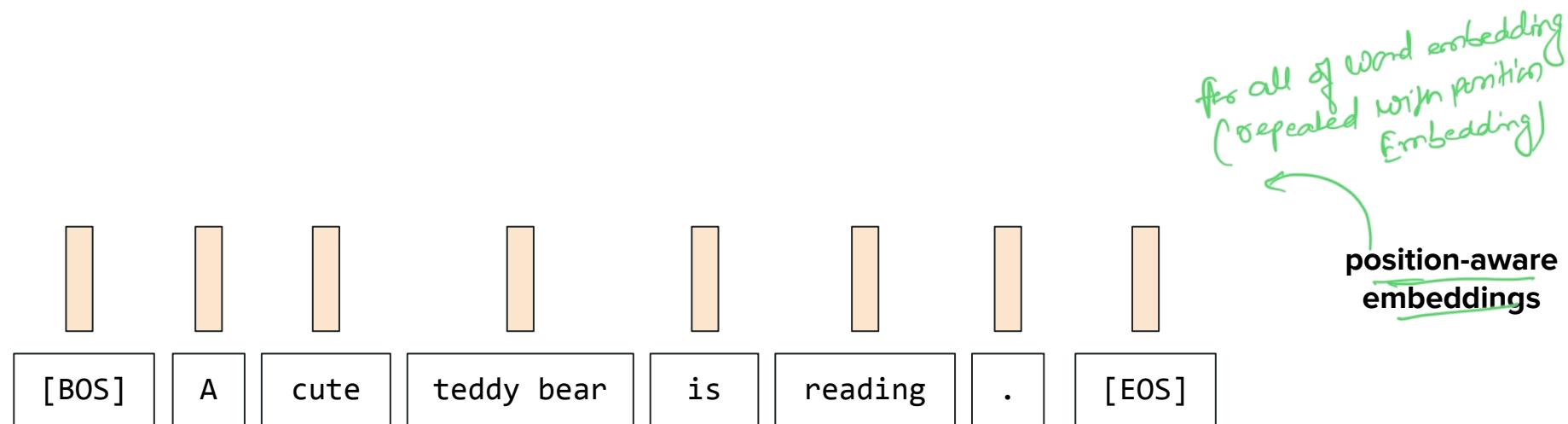
is

reading

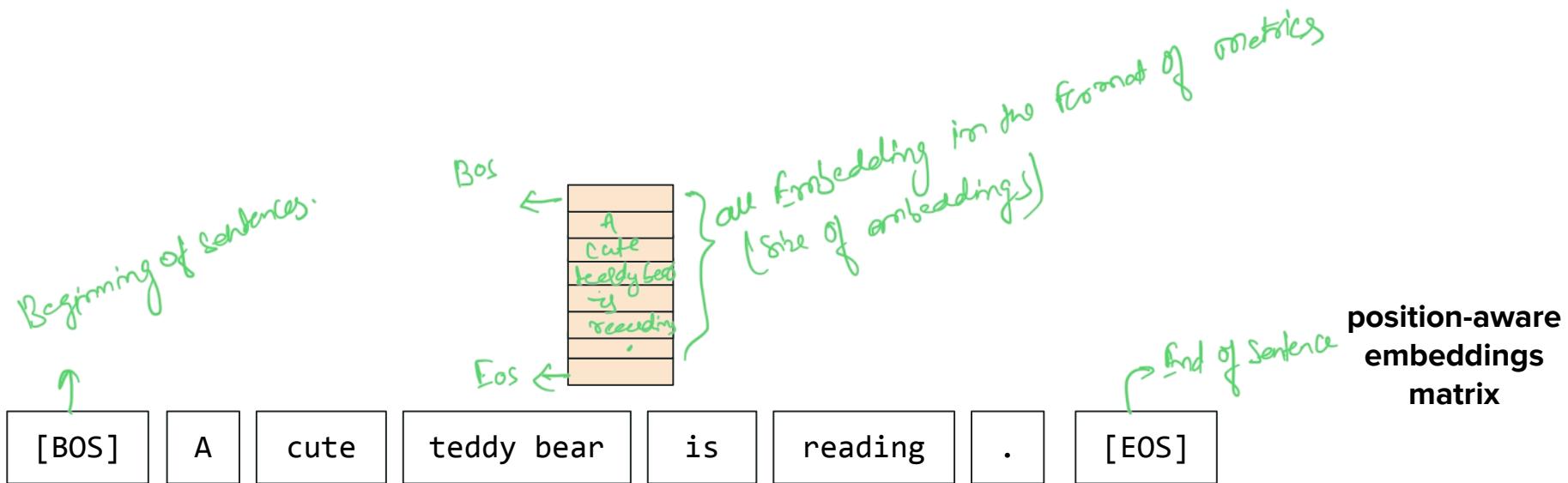
.

[EOS]

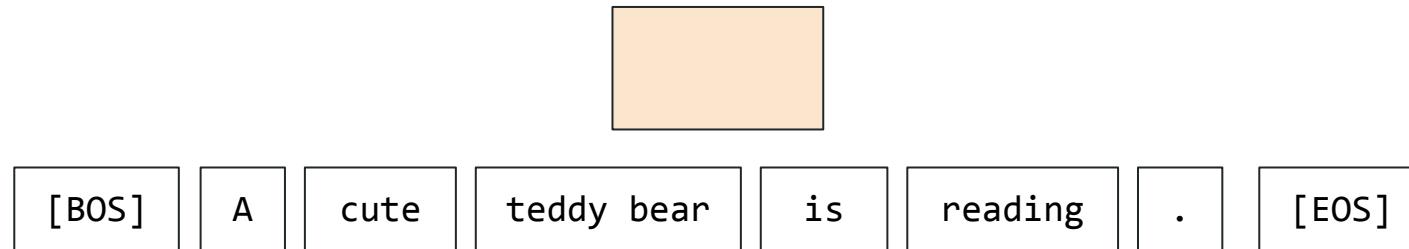
# Stitching all the pieces together with an example



# Stitching all the pieces together with an example

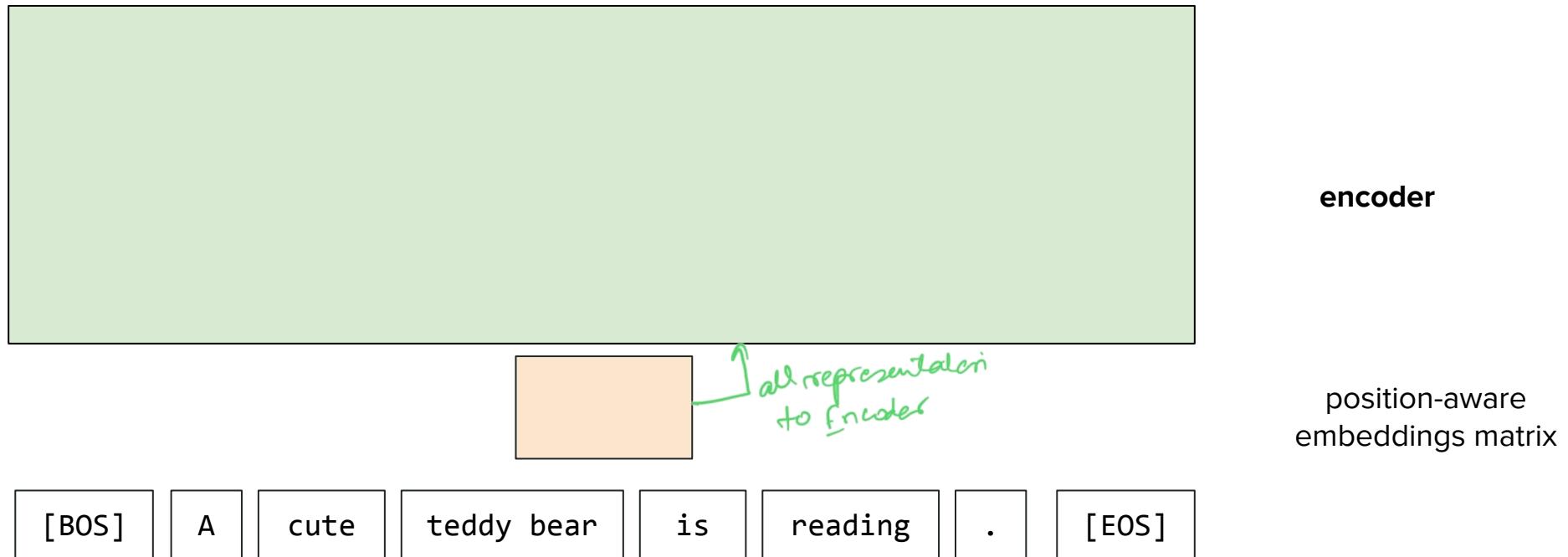


# Stitching all the pieces together with an example

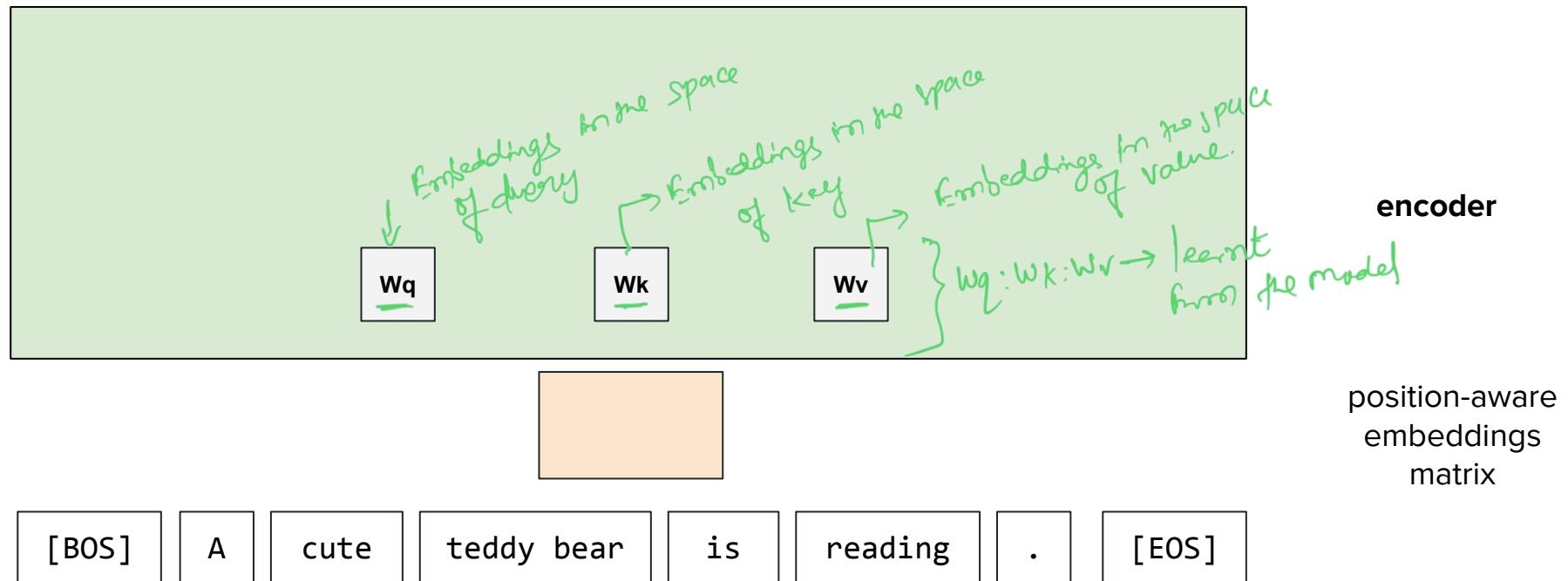


**position-aware  
embeddings  
matrix**

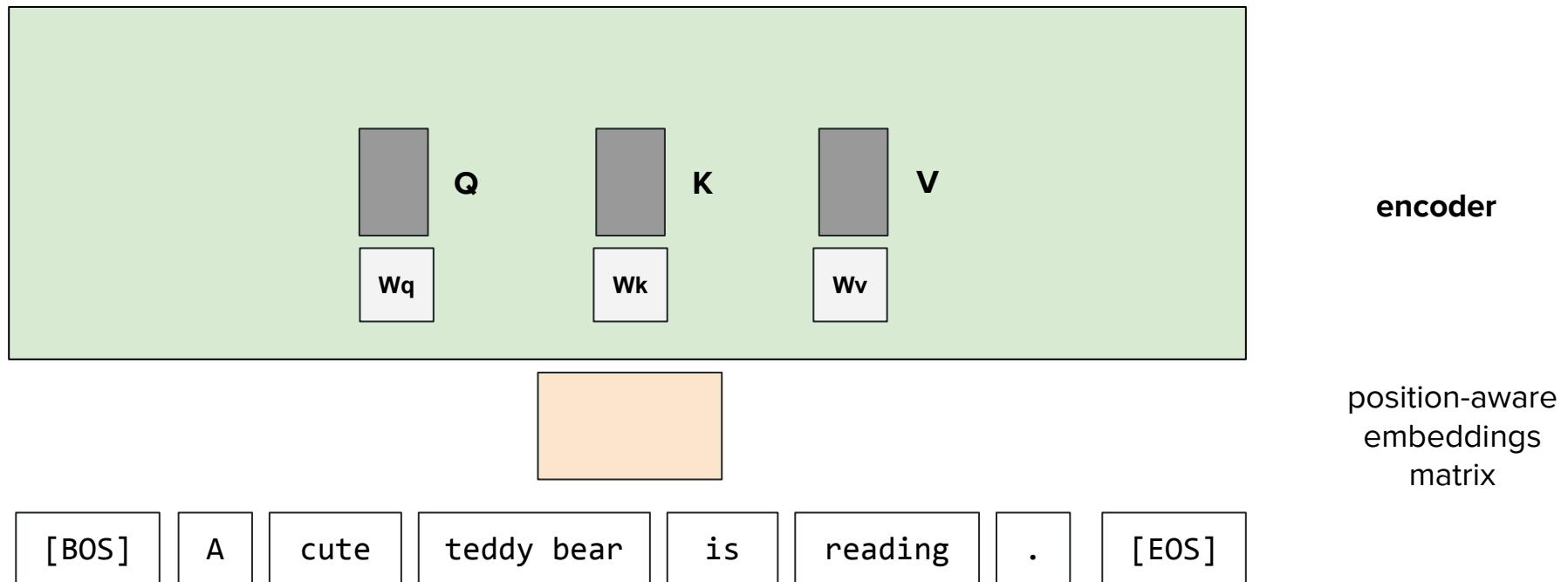
# Stitching all the pieces together with an example



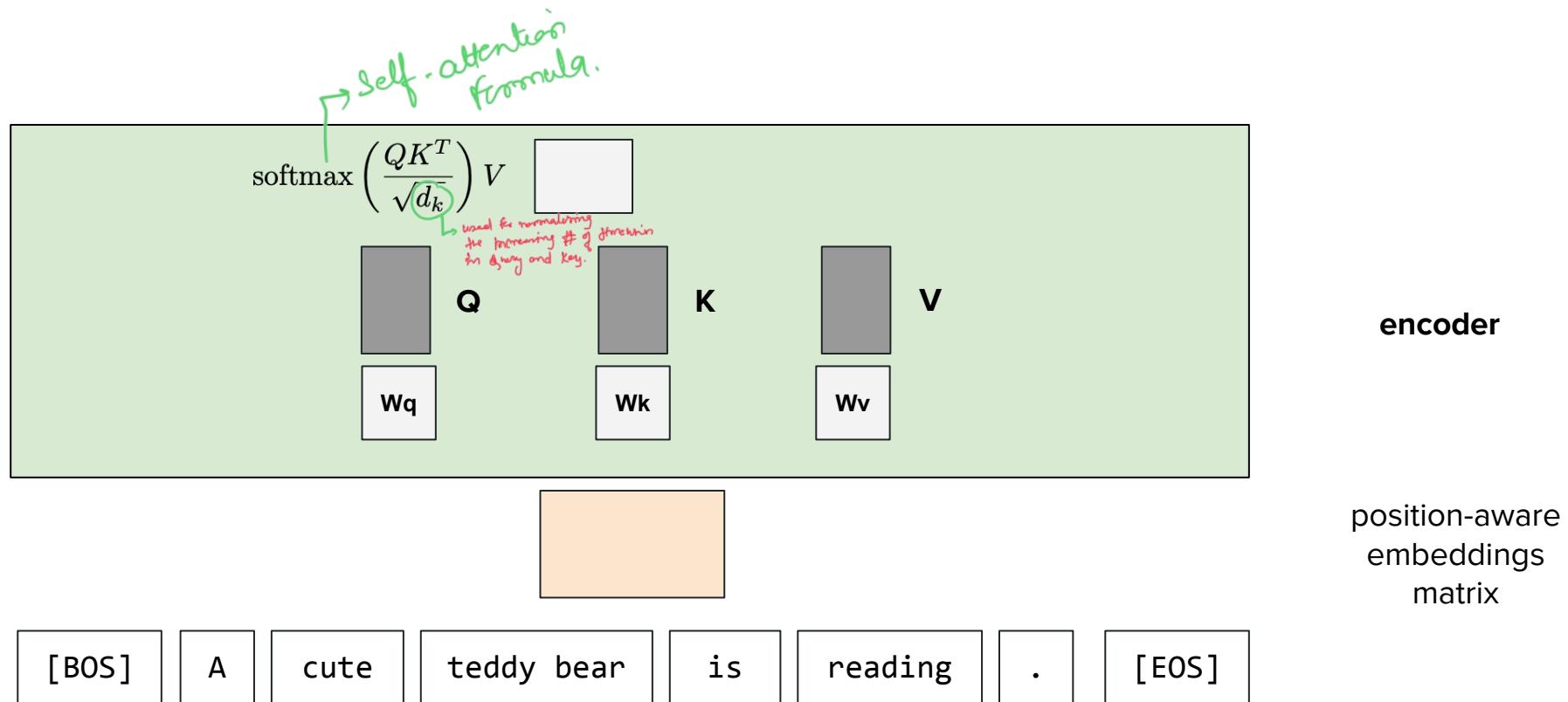
# Stitching all the pieces together with an example



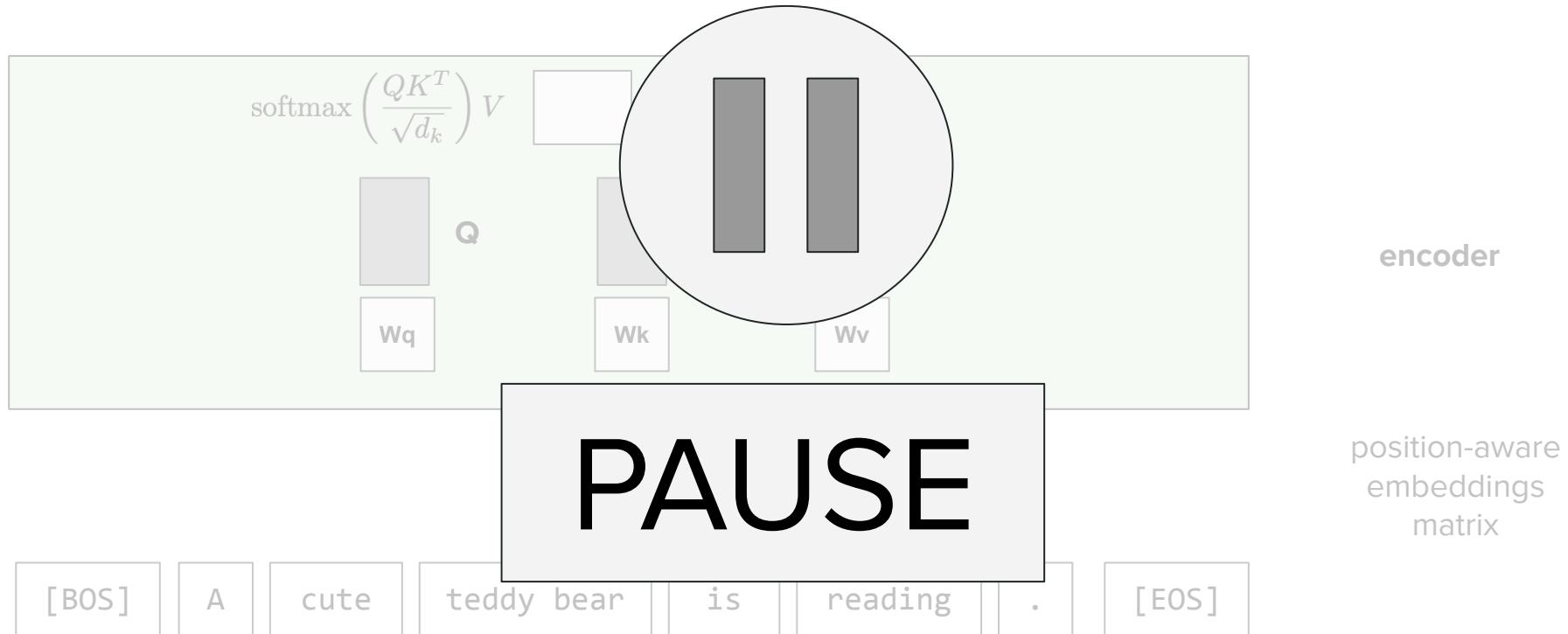
# Stitching all the pieces together with an example



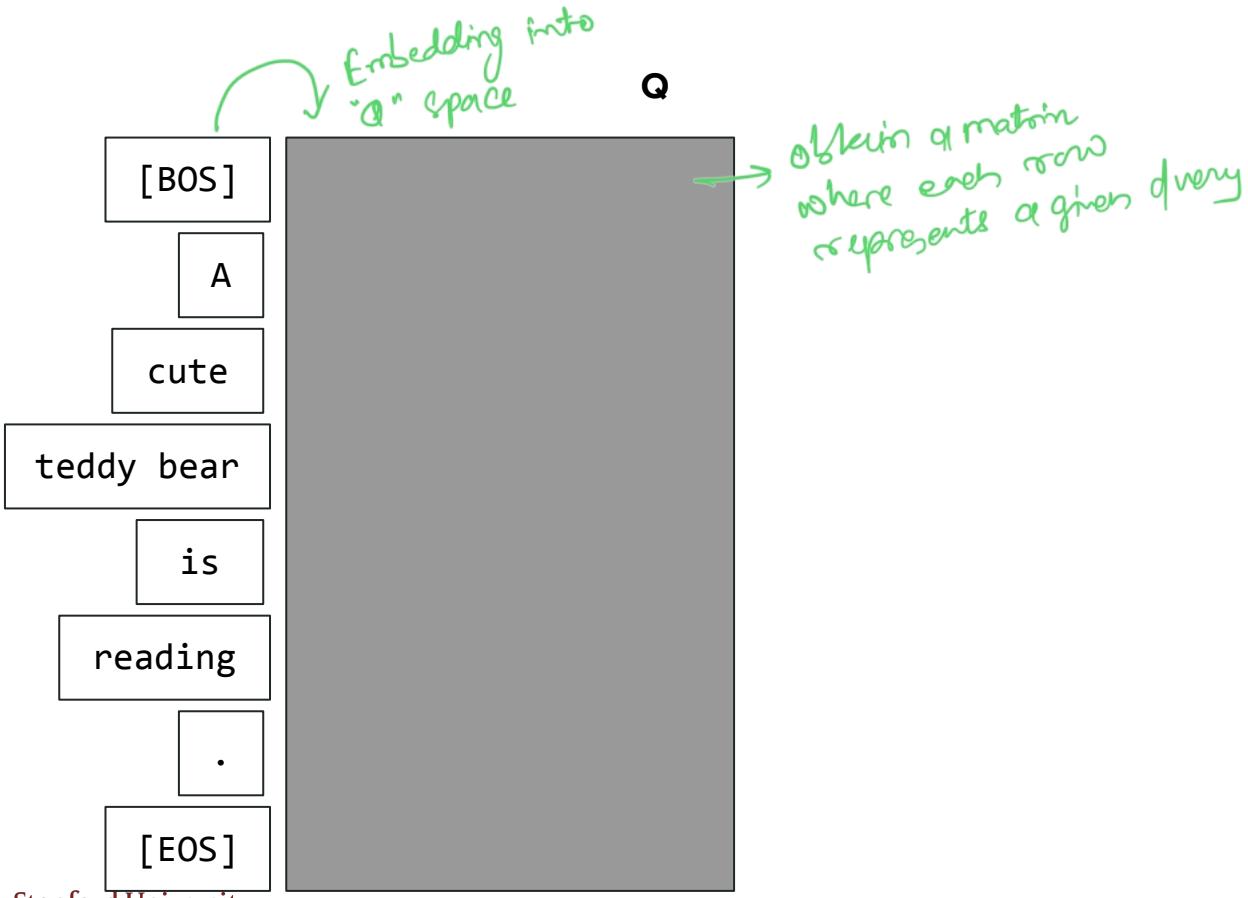
# Stitching all the pieces together with an example



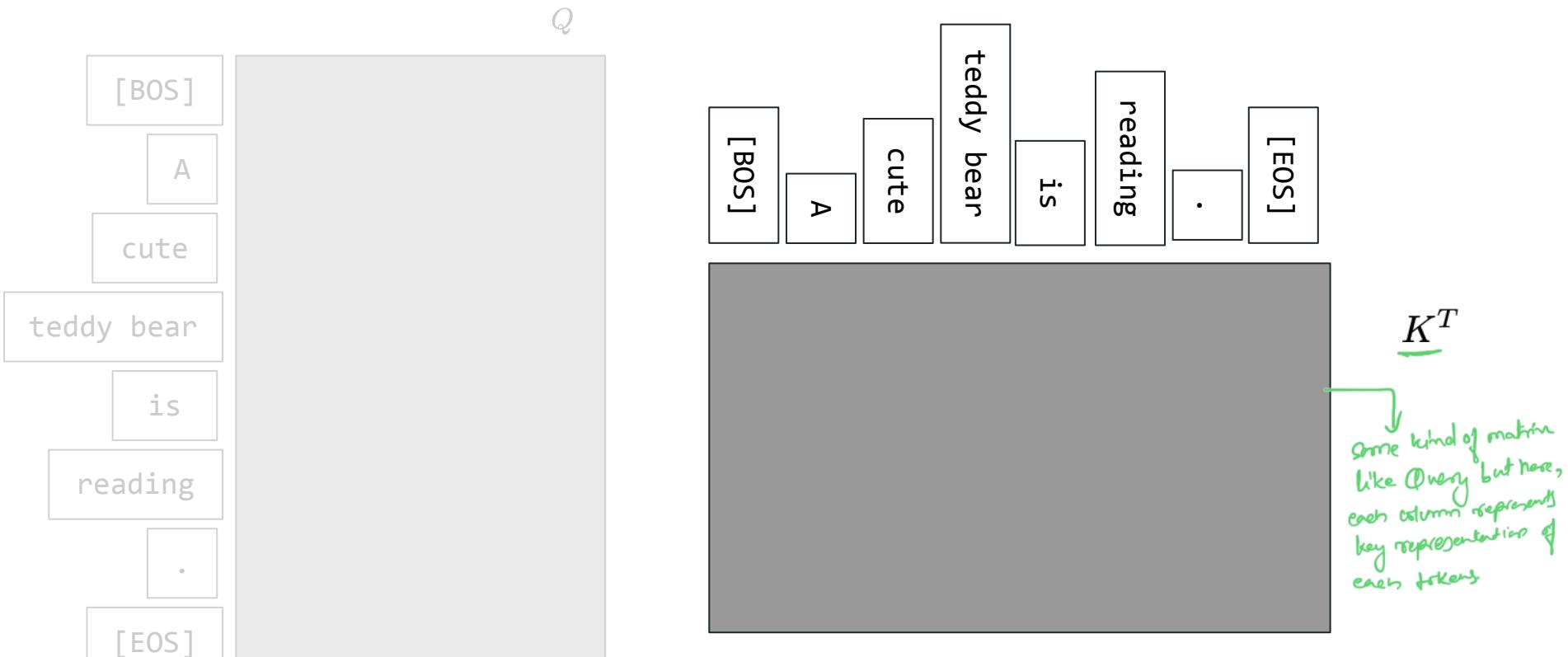
# Stitching all the pieces together with an example



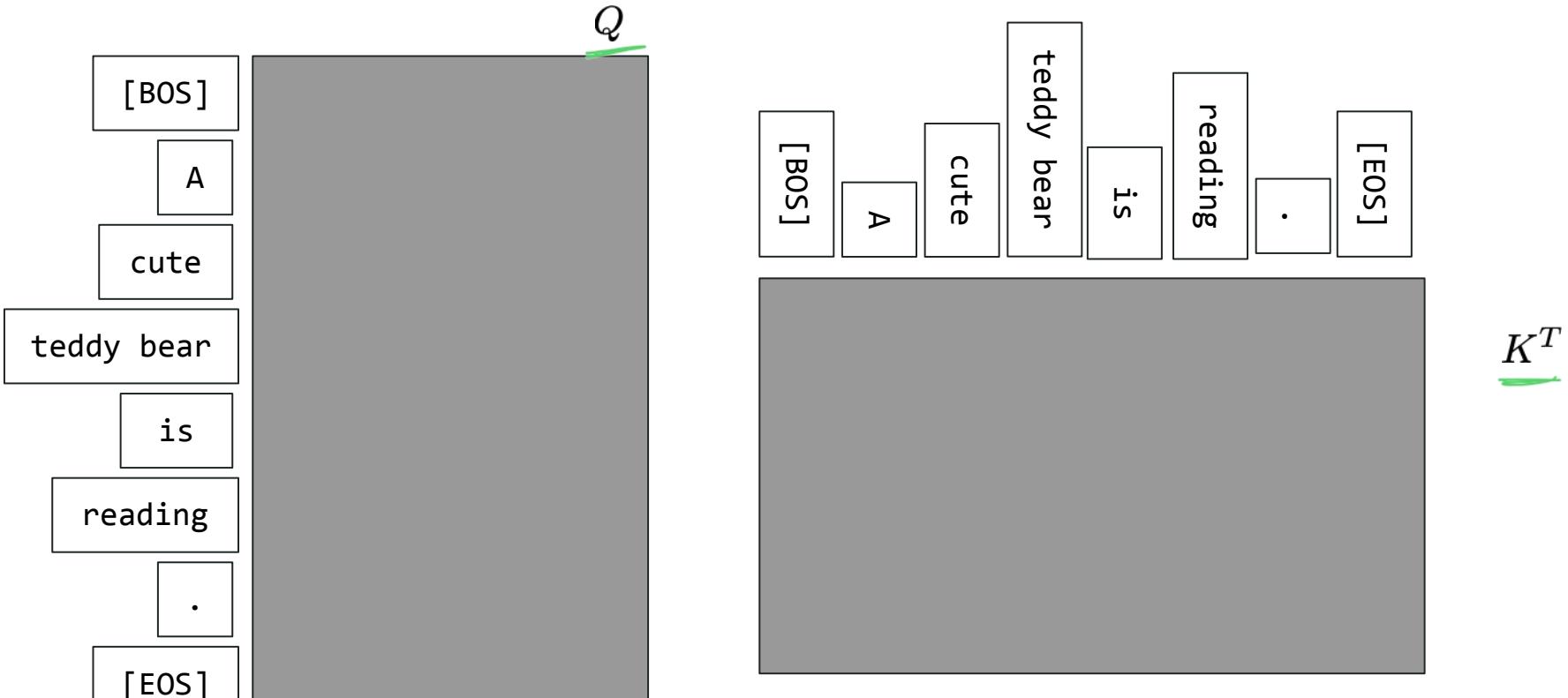
# Stitching all the pieces together with an example



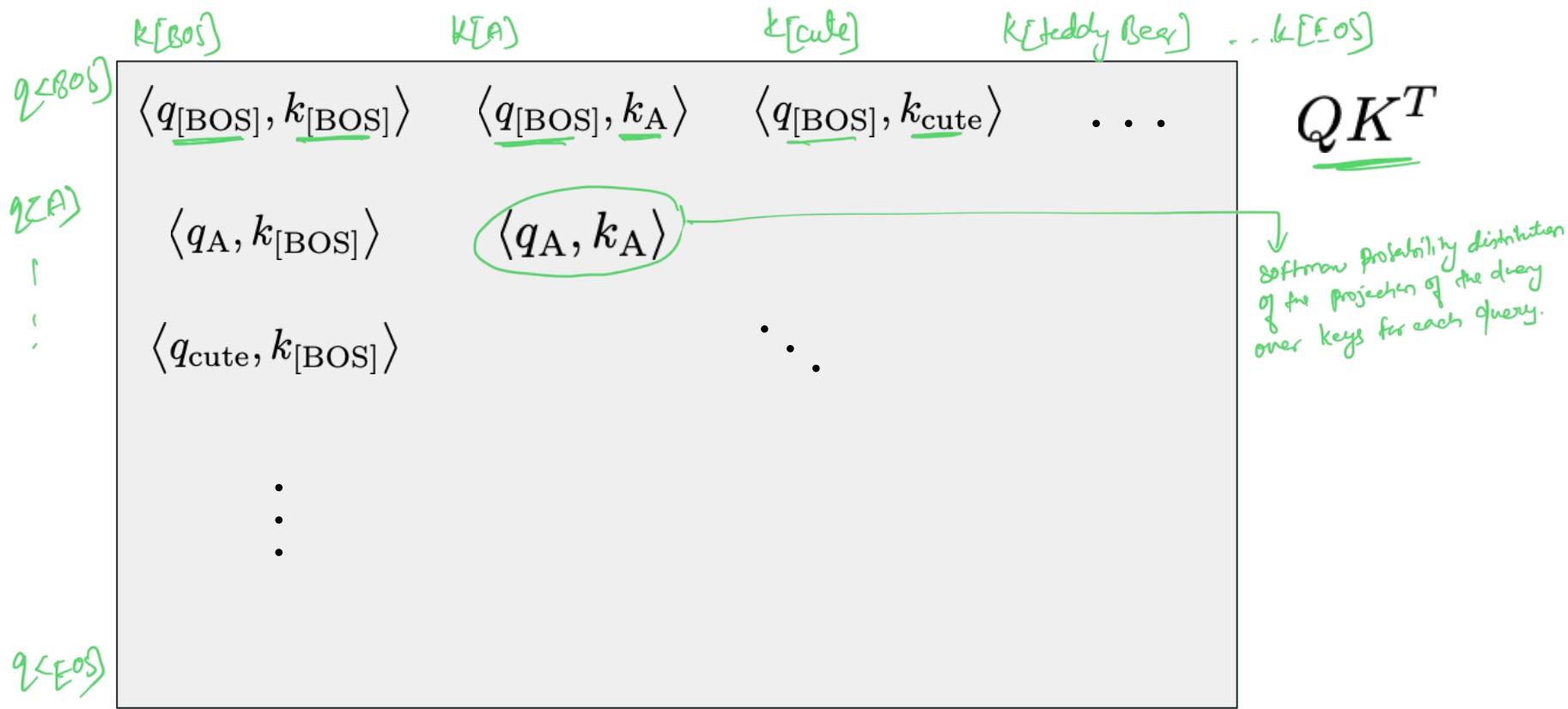
# Stitching all the pieces together with an example



# Stitching all the pieces together with an example



# Stitching all the pieces together with an example



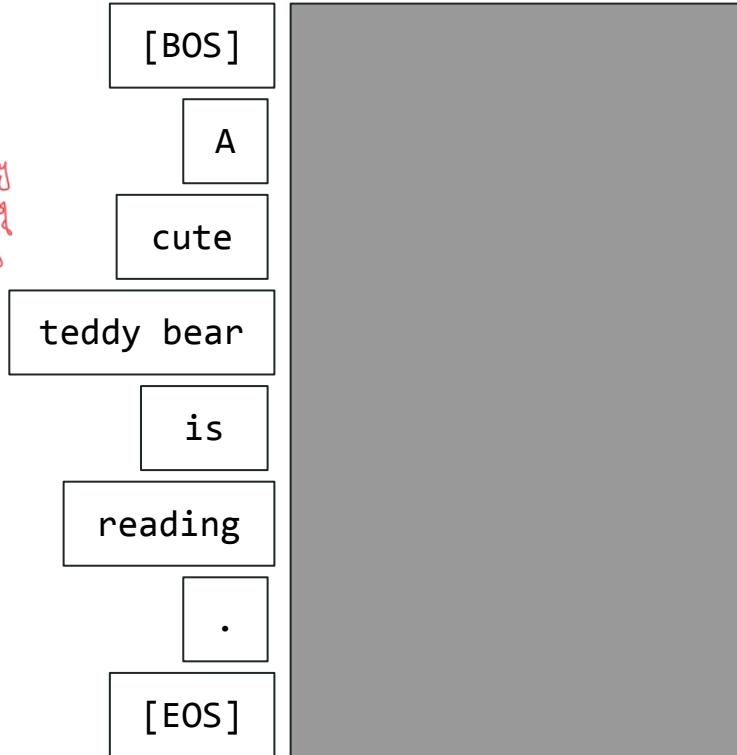
# Stitching all the pieces together with an example

V

$QK^T$

$\langle q_{[\text{BOS}]}, k_{[\text{BOS}]} \rangle$	$\langle q_{[\text{BOS}]}, k_A \rangle$	$\langle q_{[\text{BOS}]}, k_{\text{cute}} \rangle$	$\dots$
$\langle q_A, k_{[\text{BOS}]} \rangle$	$\langle q_A, k_A \rangle$		
$\langle q_{\text{cute}}, k_{[\text{BOS}]} \rangle$		$\ddots$	
$\vdots$			

\* softmax  
of all of the  
query and key  
combinations  
embeddings



# Stitching all the pieces together with an example

$$\begin{array}{c} \text{QK}^T \\ \downarrow \quad \downarrow \quad \downarrow \\ \langle q_{[\text{BOS}]}, k_{[\text{BOS}]} \rangle v_{[\text{BOS}]} + \langle q_{[\text{BOS}]}, k_A \rangle v_A + \langle q_{[\text{BOS}]}, k_{\text{cute}} \rangle v_{\text{cute}} + \dots \\ \\ \langle q_A, k_{[\text{BOS}]} \rangle v_{[\text{BOS}]} + \langle q_A, k_A \rangle v_A + \langle q_A, k_{\text{cute}} \rangle v_{\text{cute}} + \dots \\ \\ \vdots \end{array}$$

$QK^TV$

Each row represents the weight sum of each value with each query.

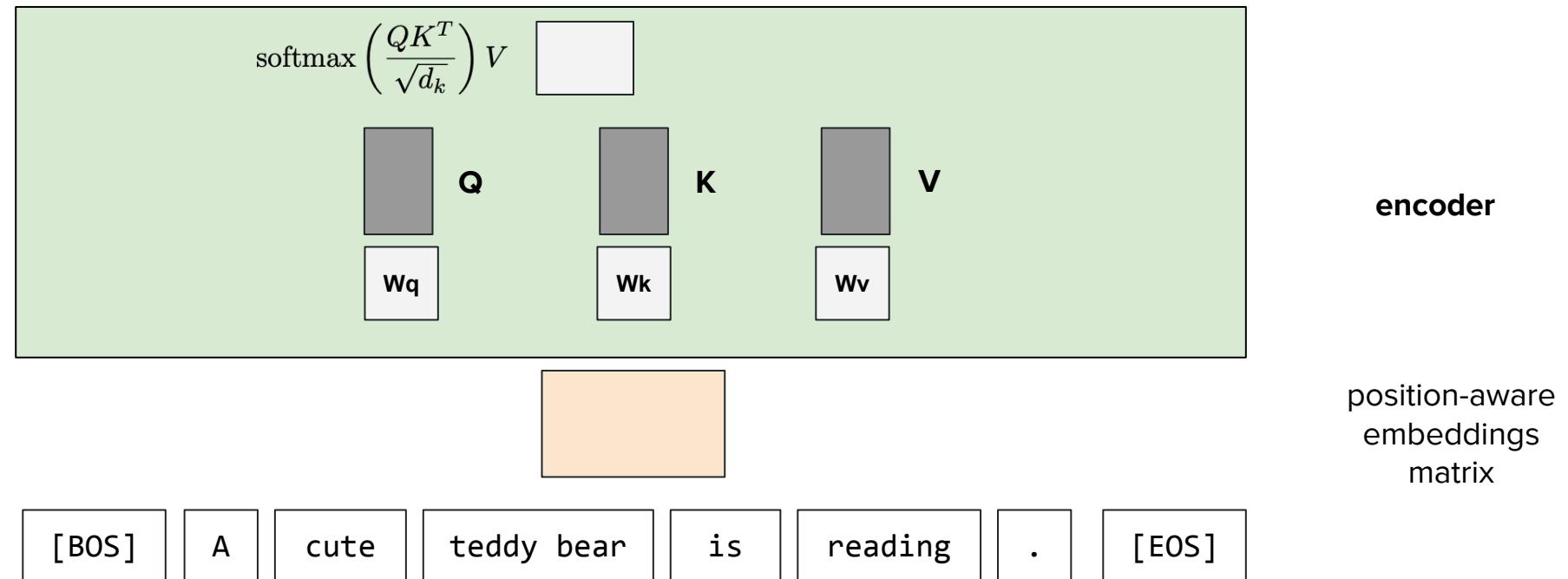
# Stitching all the pieces together with an example

$$\text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

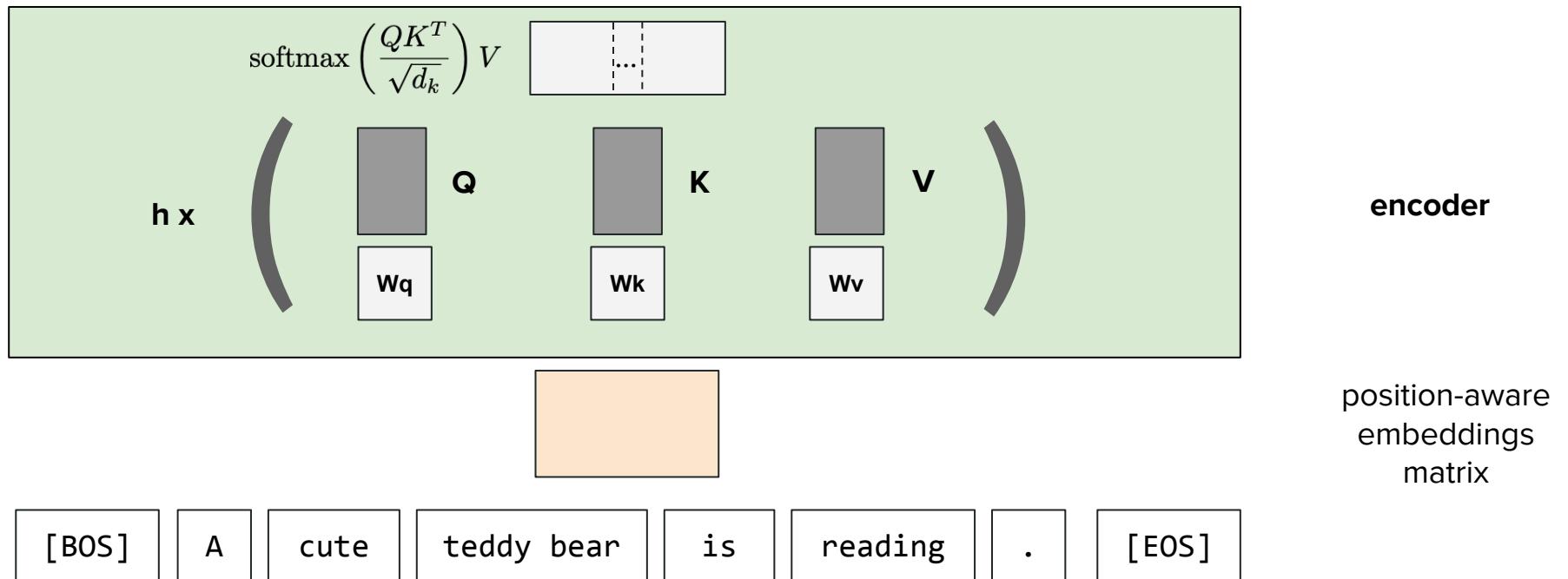
weighted average of values

with weights being a function of  $\langle q, k \rangle$

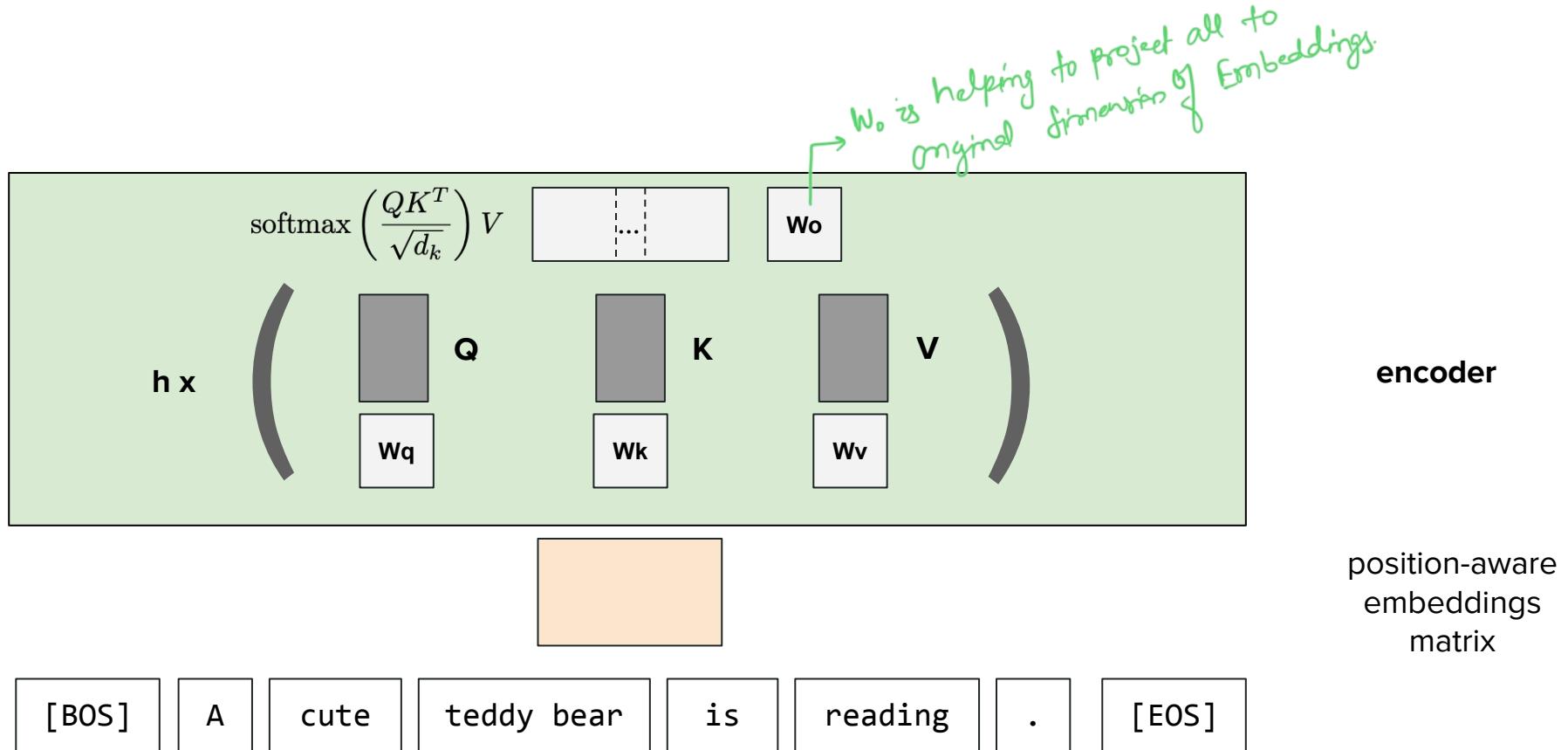
# Stitching all the pieces together with an example



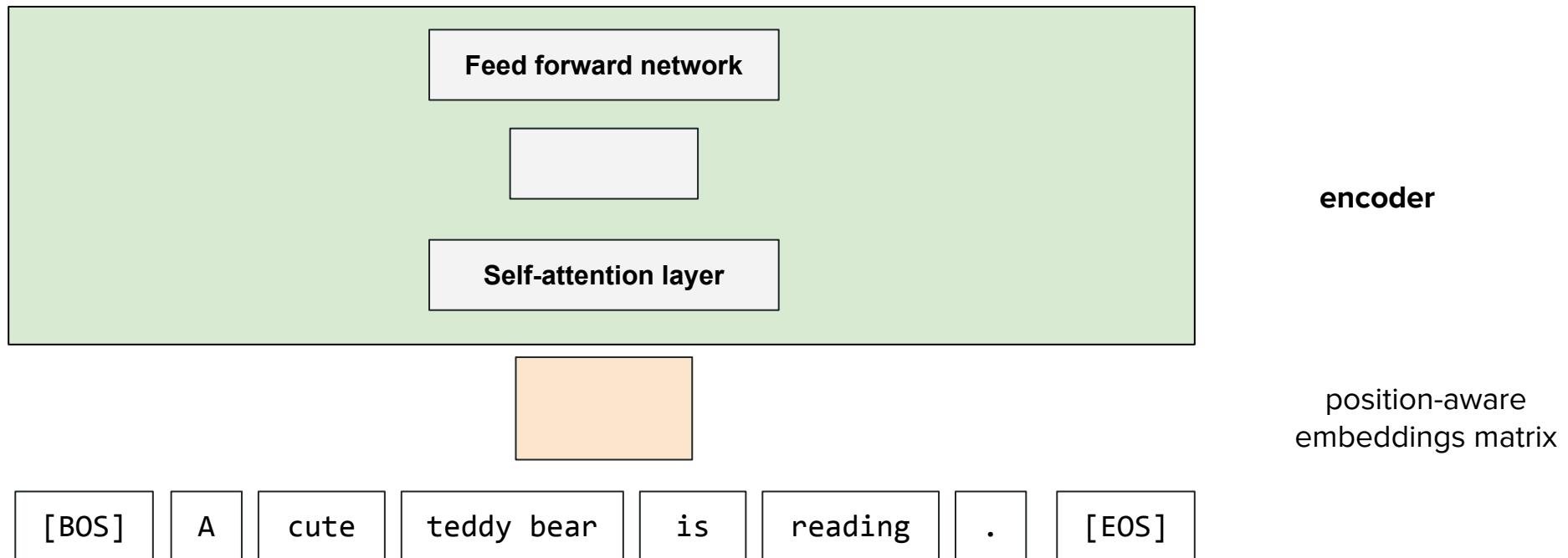
# Stitching all the pieces together with an example



# Stitching all the pieces together with an example



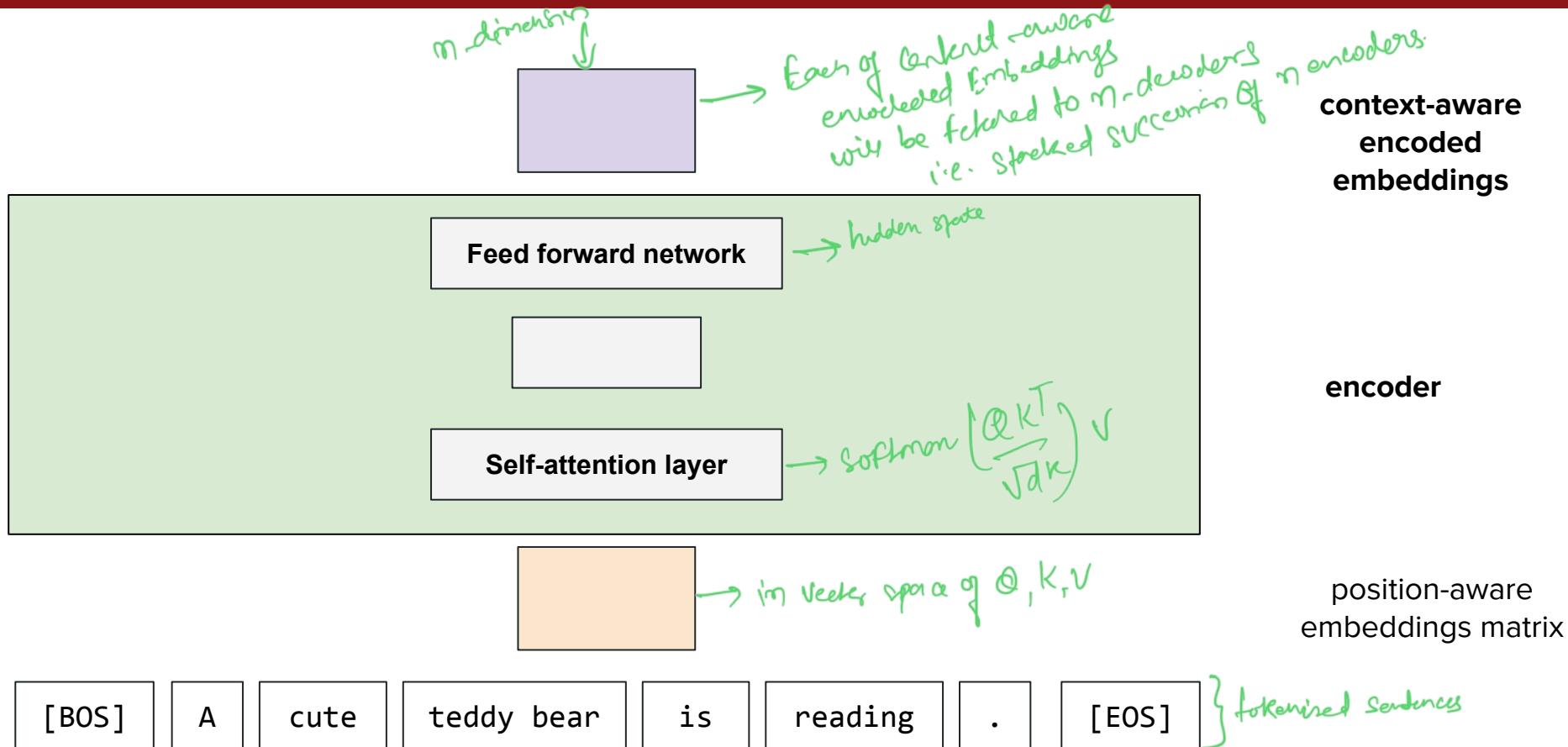
# Stitching all the pieces together with an example



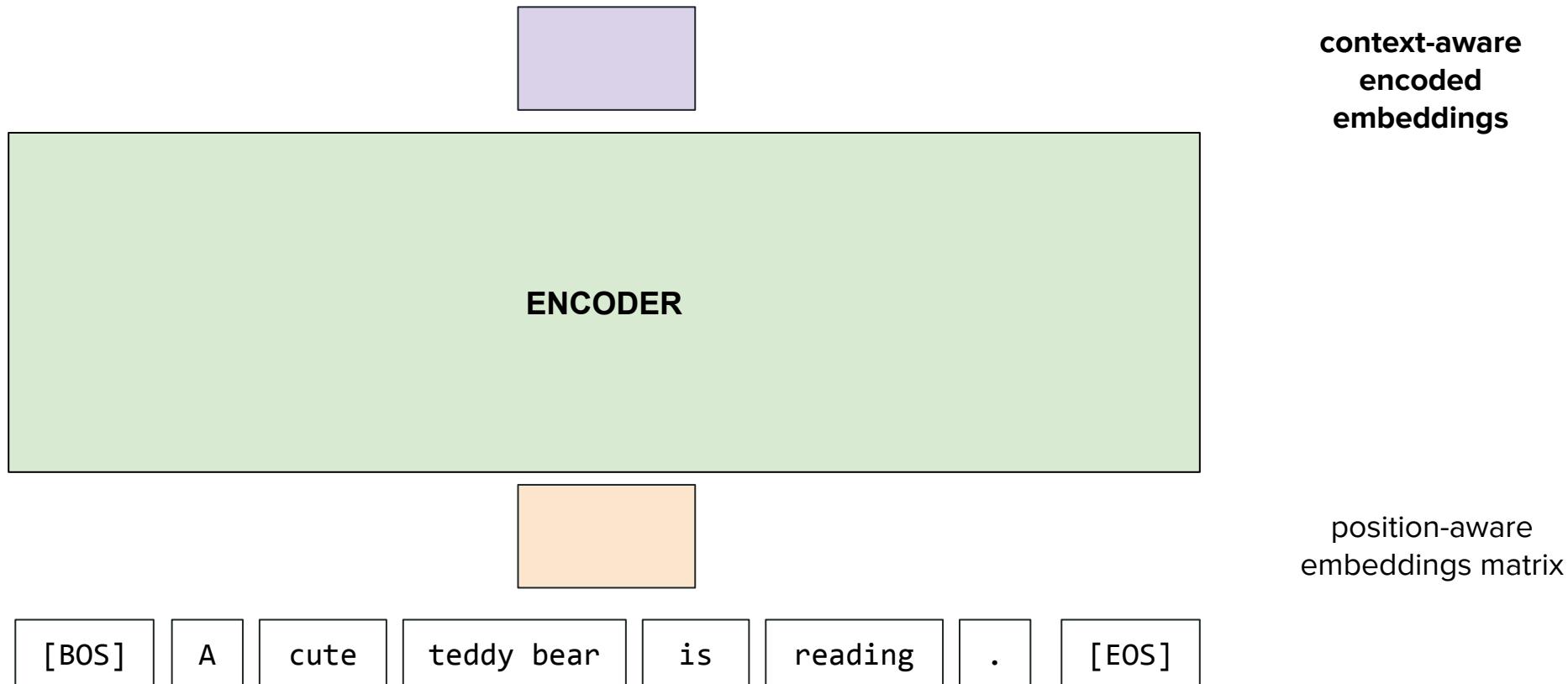
encoder

position-aware  
embeddings matrix

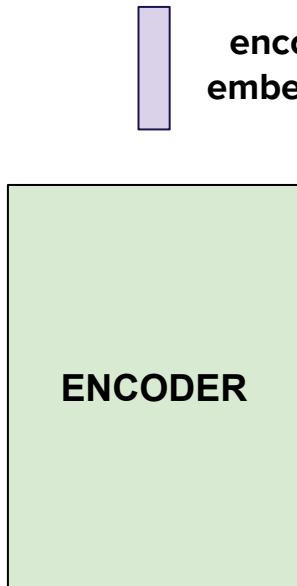
# Stitching all the pieces together with an example



# Stitching all the pieces together with an example

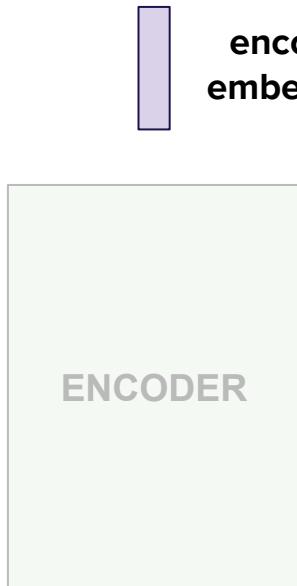


# Stitching all the pieces together with an example

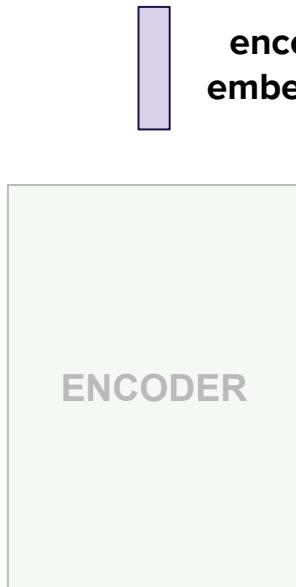


A cute teddy bear  
is reading.

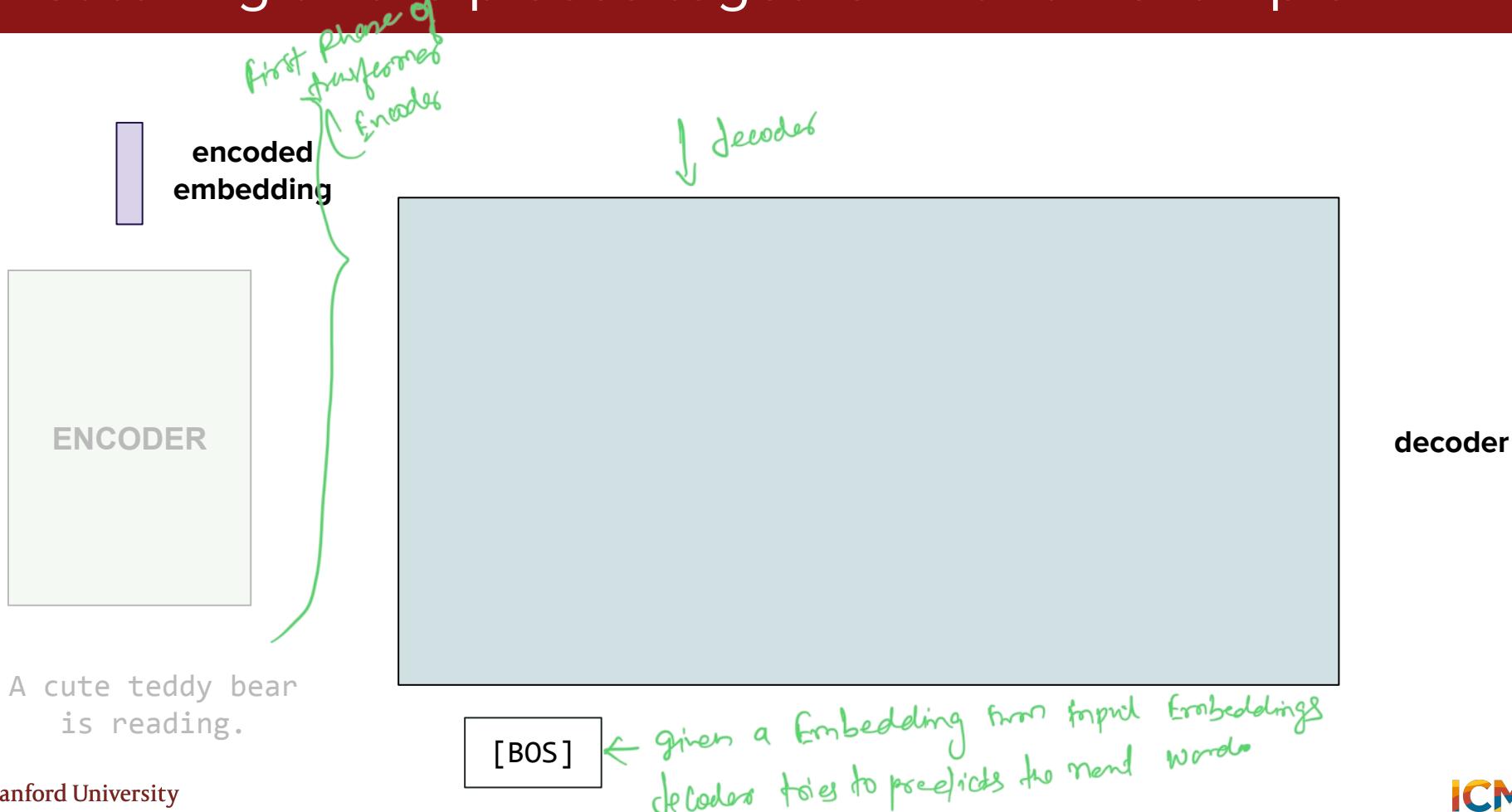
# Stitching all the pieces together with an example



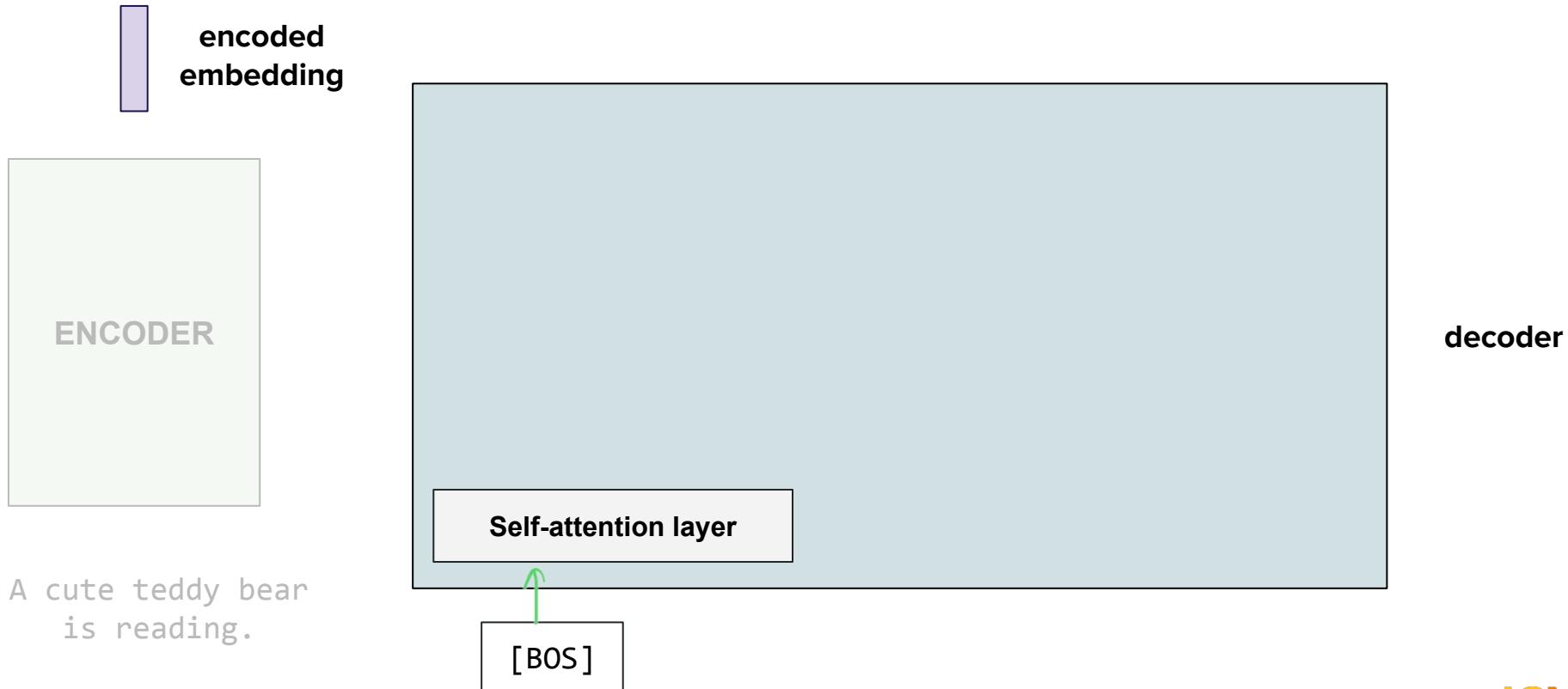
# Stitching all the pieces together with an example



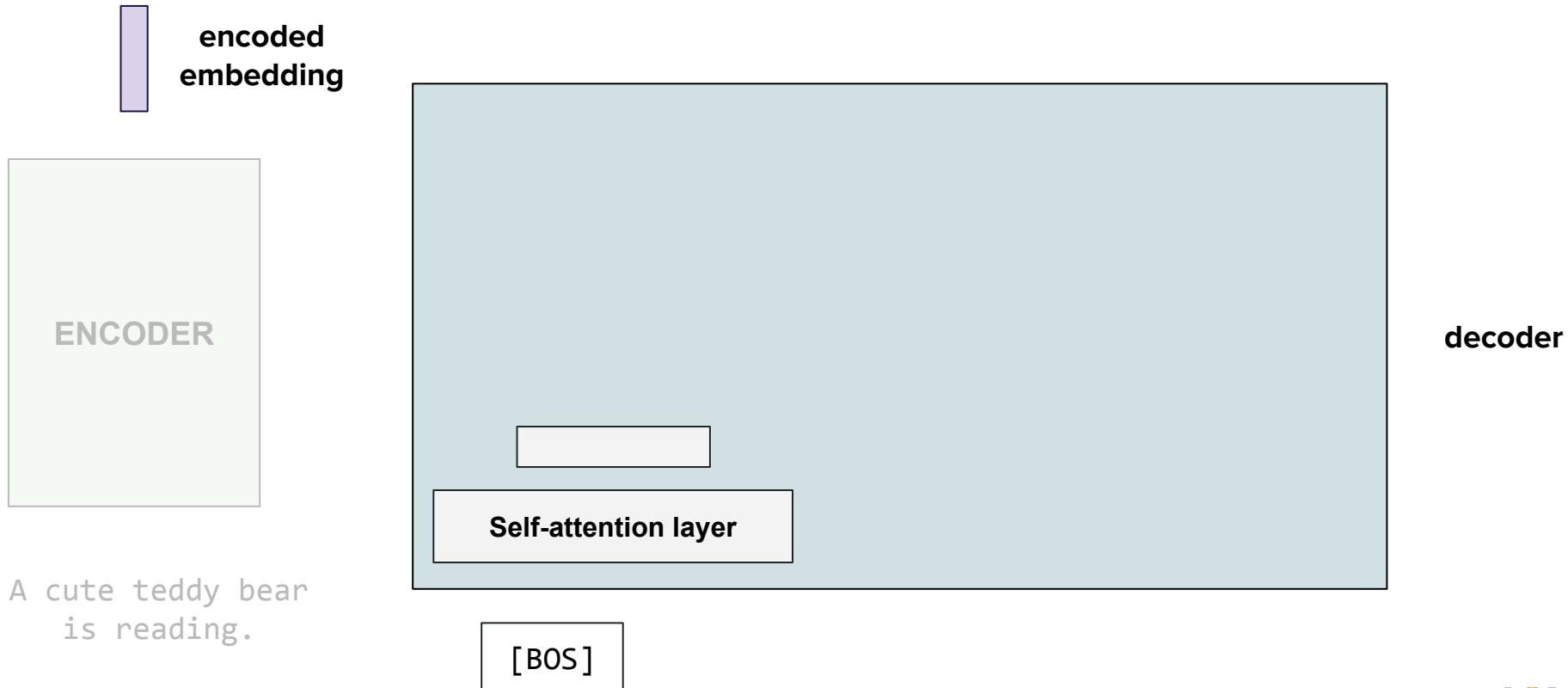
# Stitching all the pieces together with an example



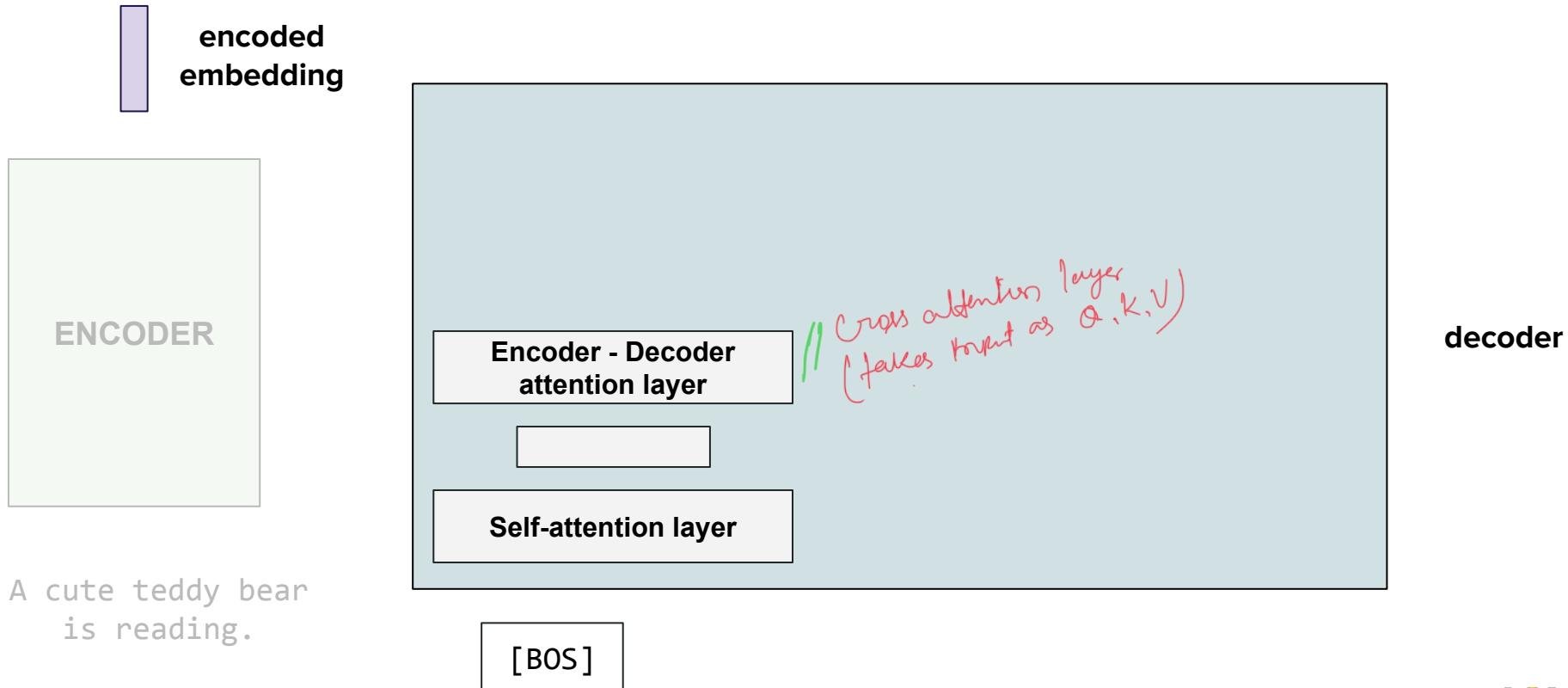
# Stitching all the pieces together with an example



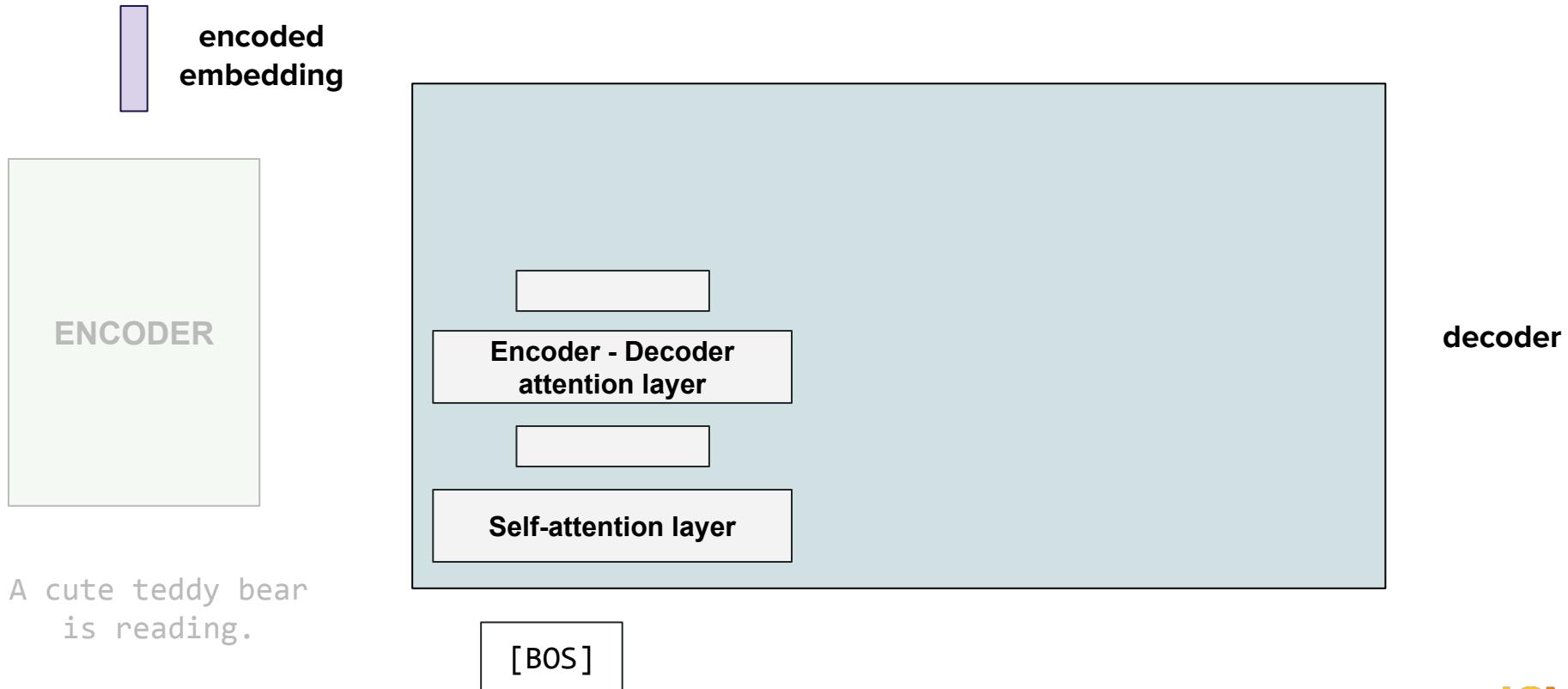
# Stitching all the pieces together with an example



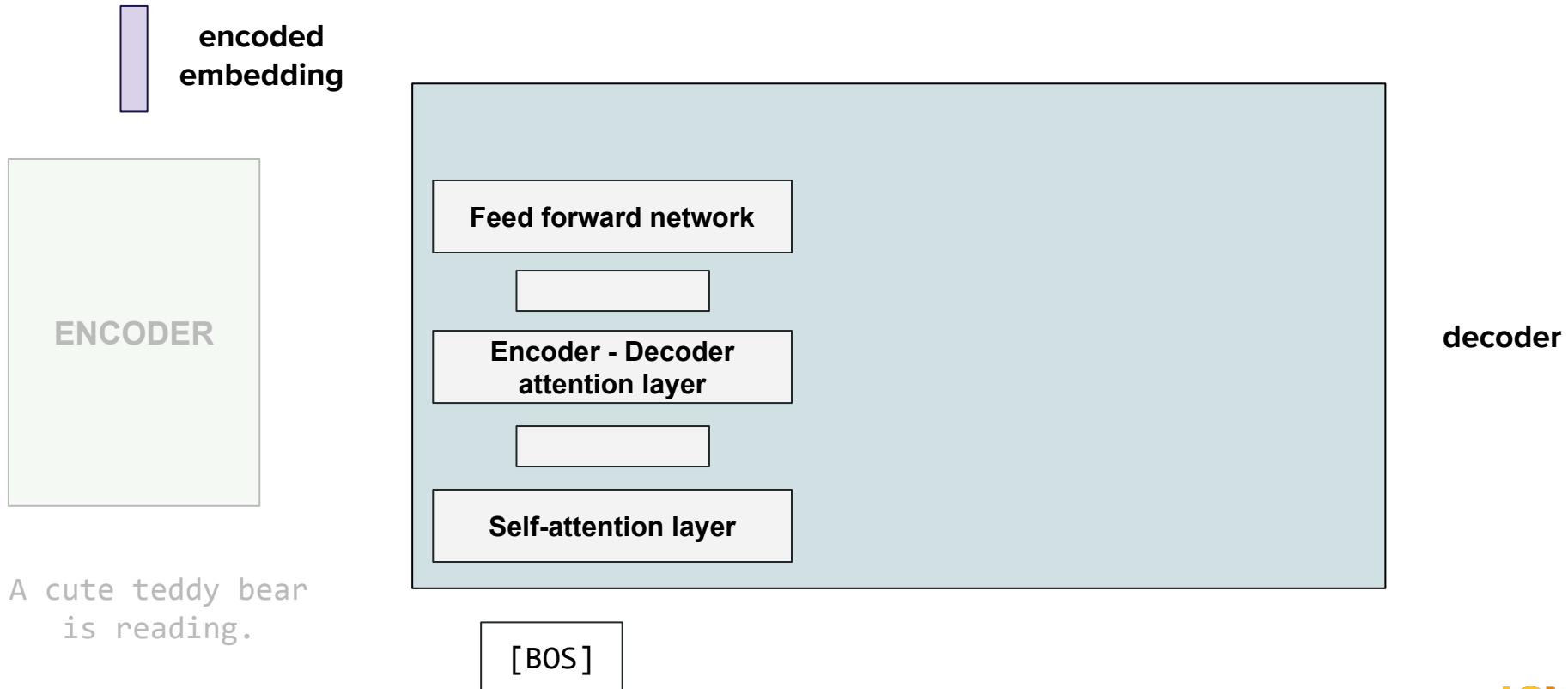
# Stitching all the pieces together with an example



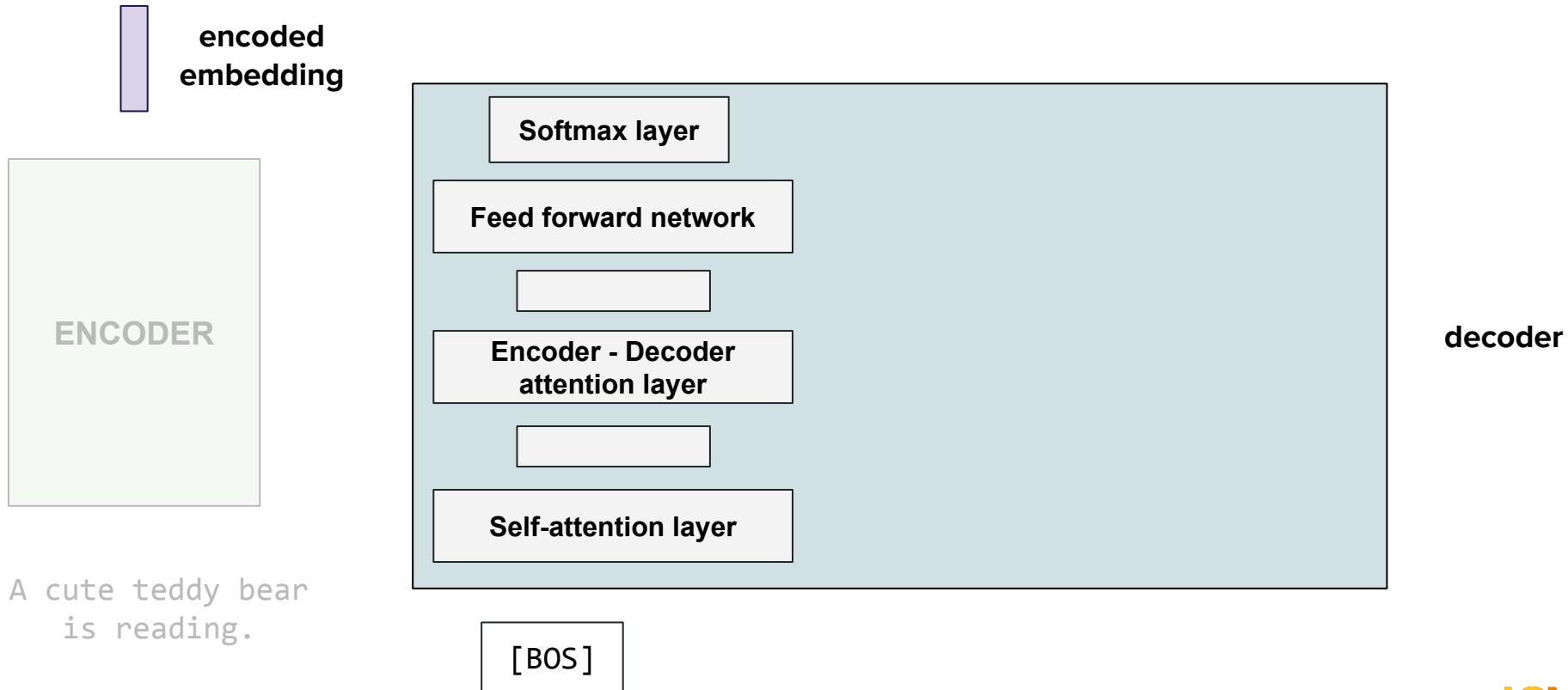
# Stitching all the pieces together with an example



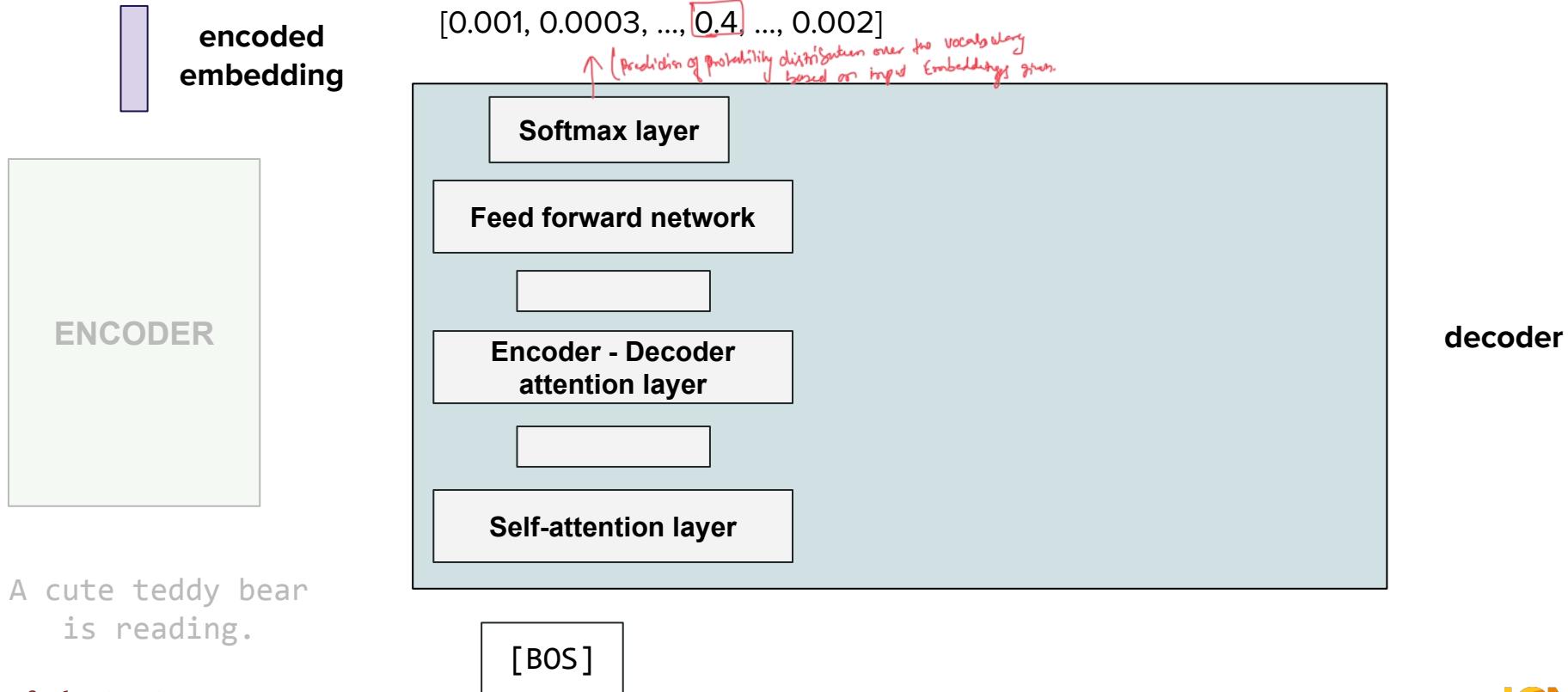
# Stitching all the pieces together with an example



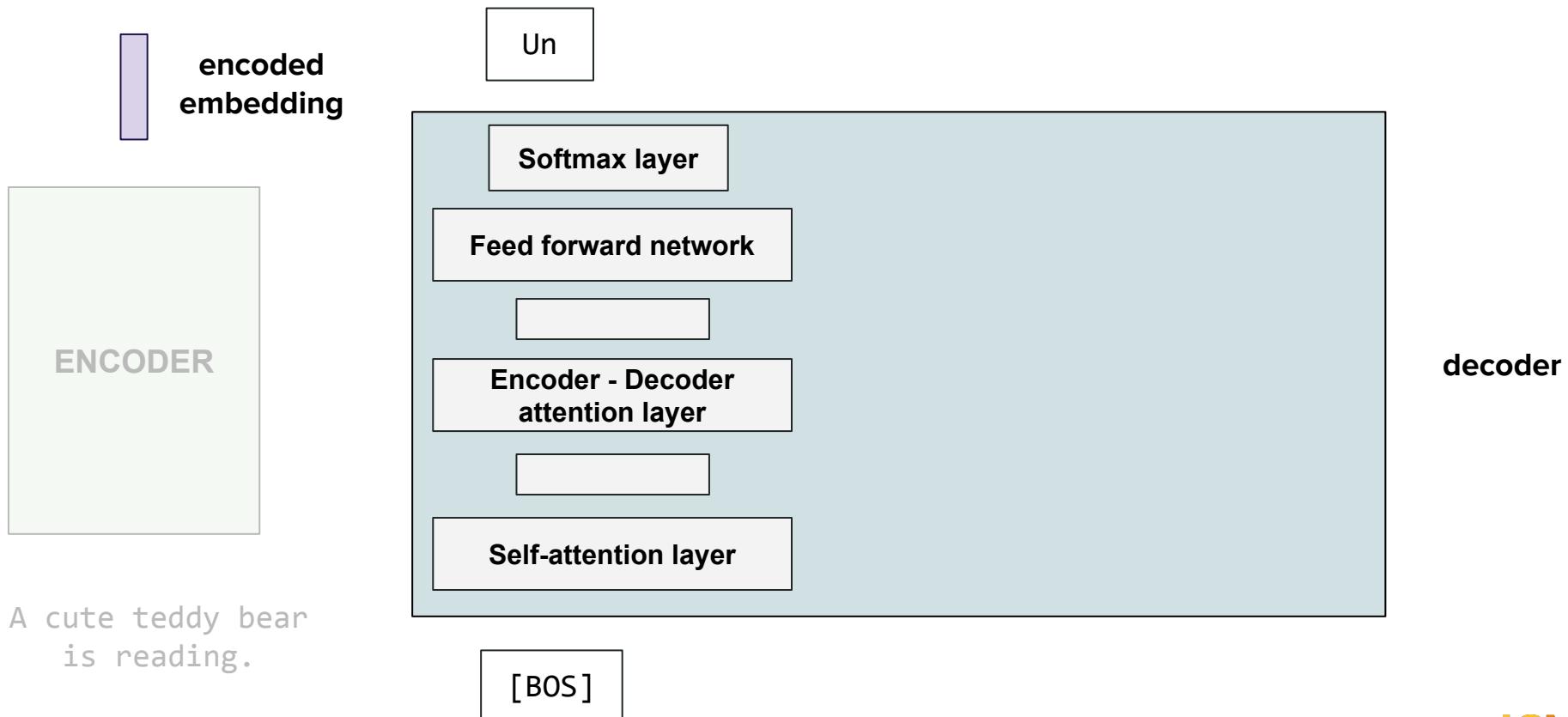
# Stitching all the pieces together with an example



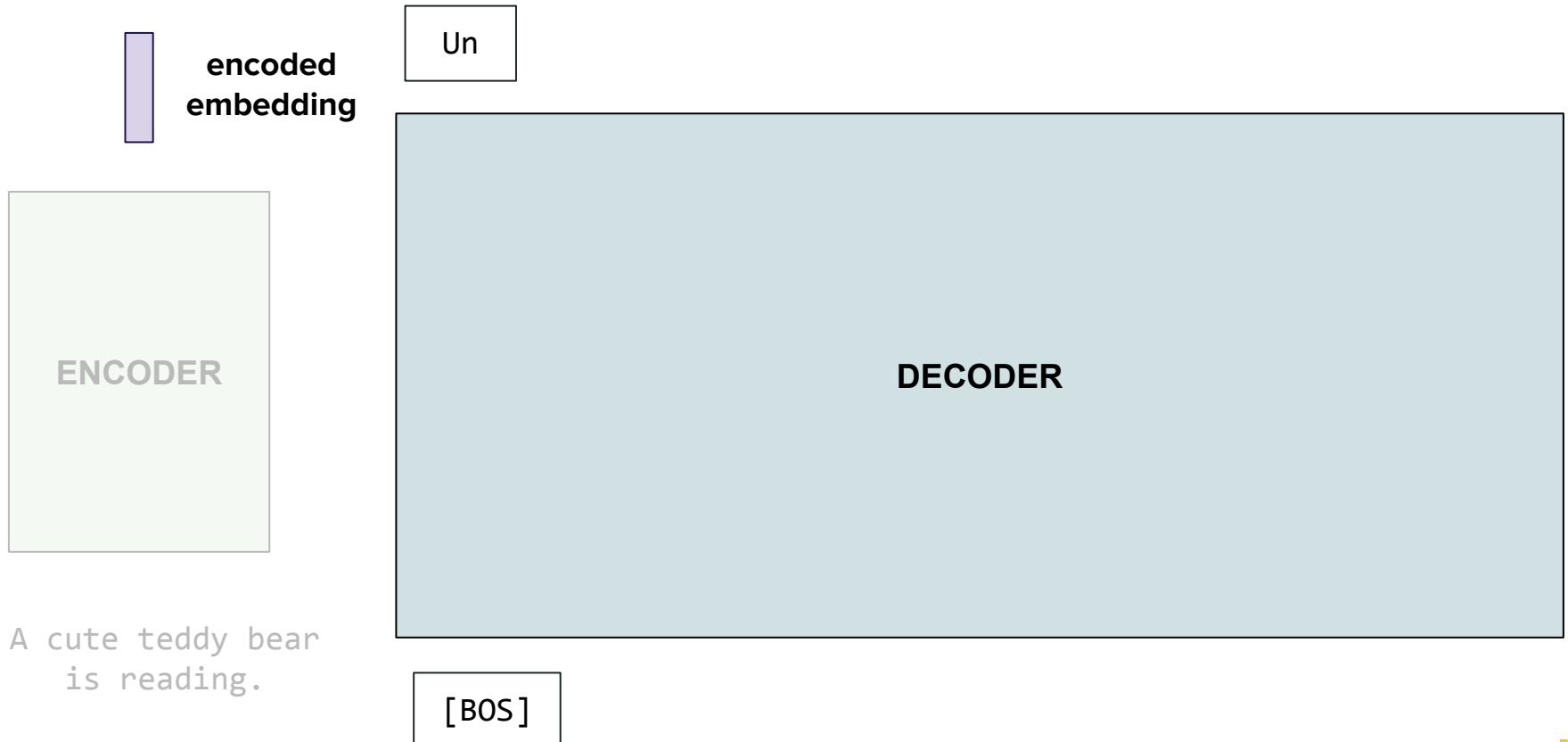
# Stitching all the pieces together with an example



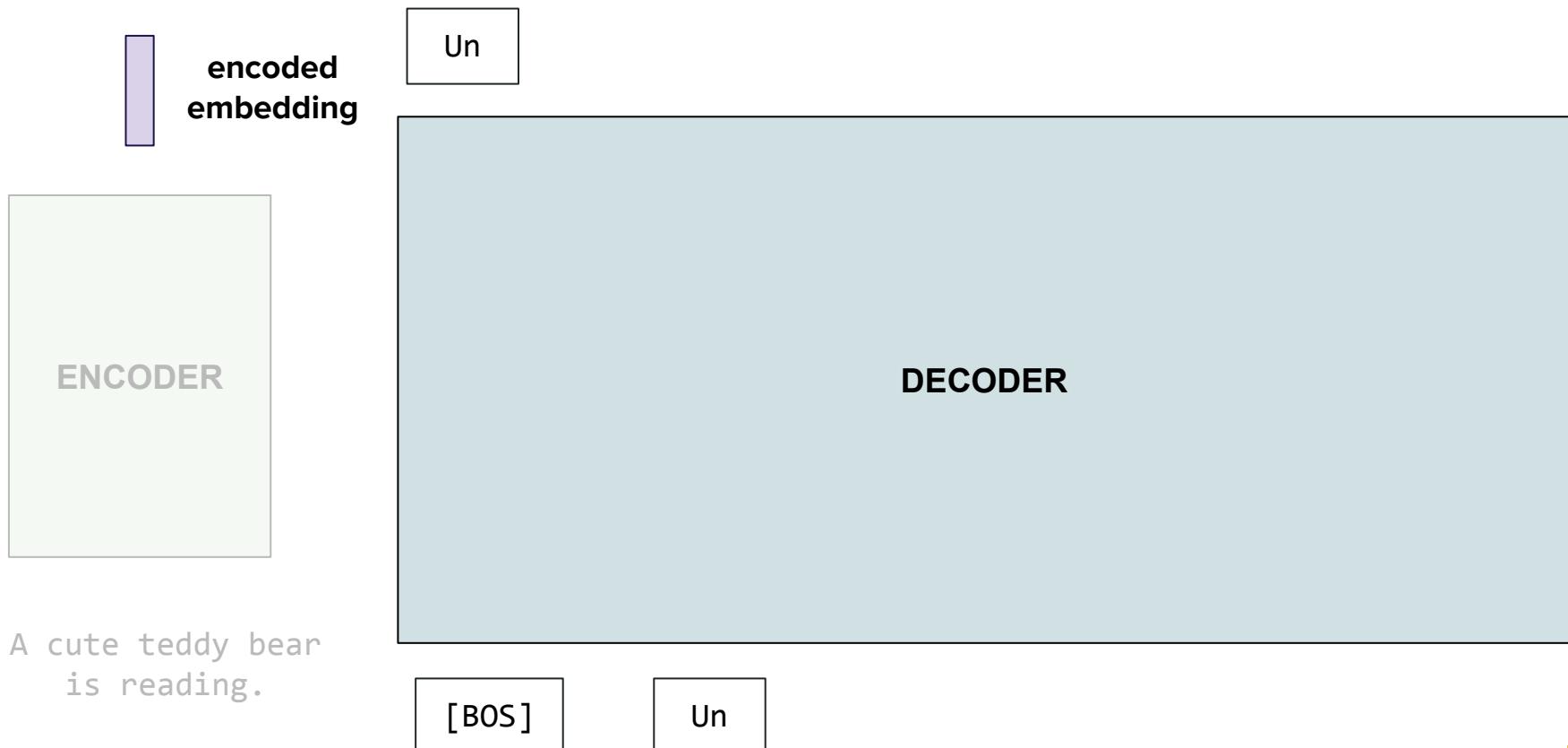
# Stitching all the pieces together with an example



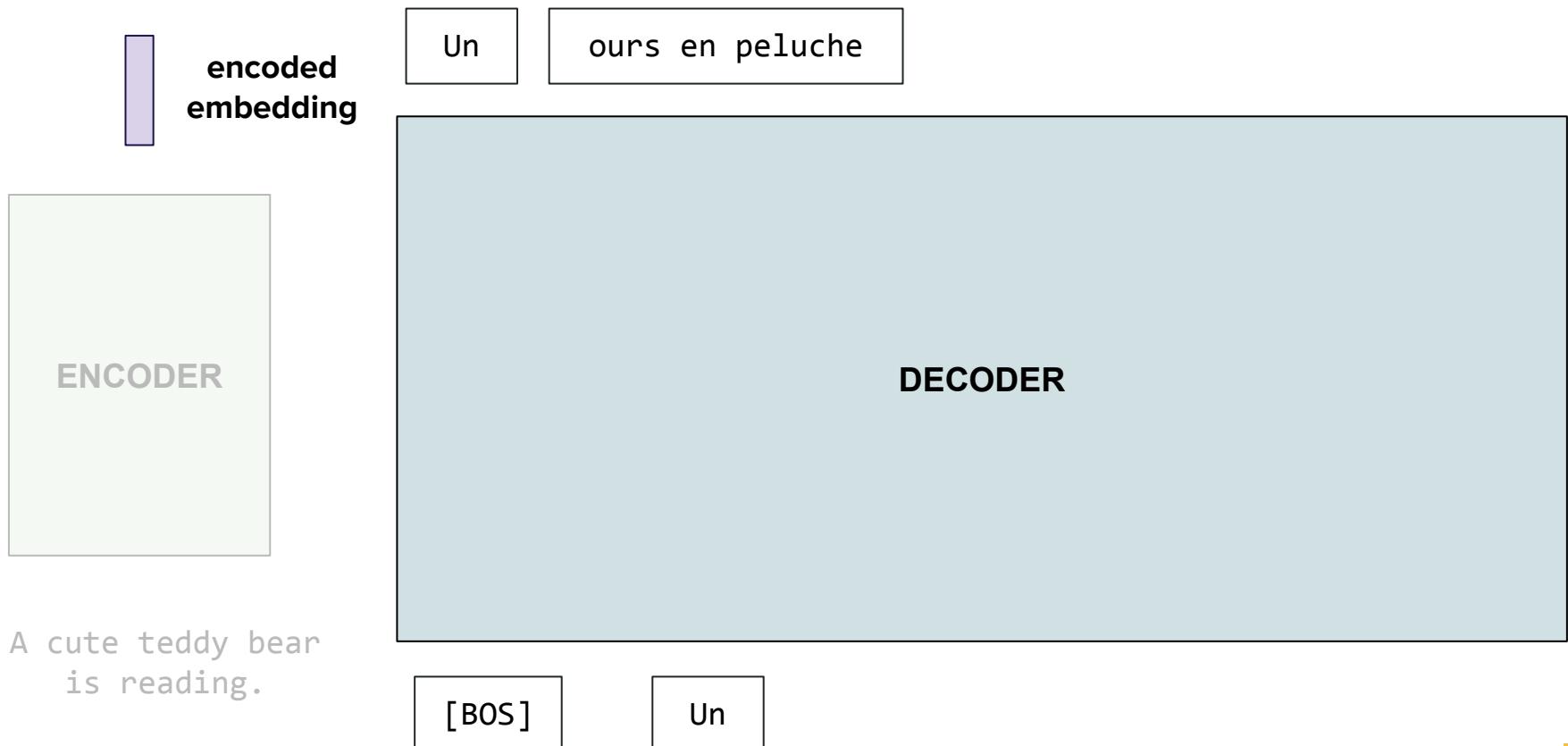
# Stitching all the pieces together with an example



# Stitching all the pieces together with an example

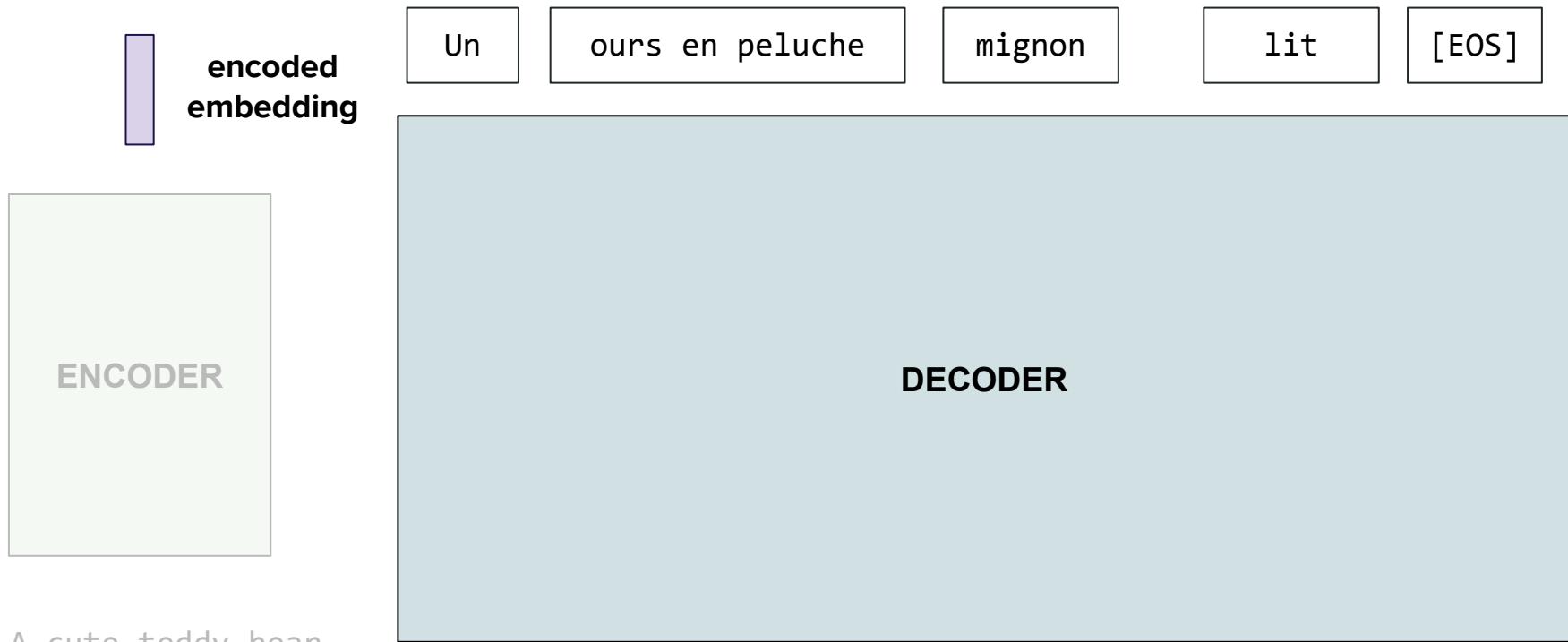


# Stitching all the pieces together with an example



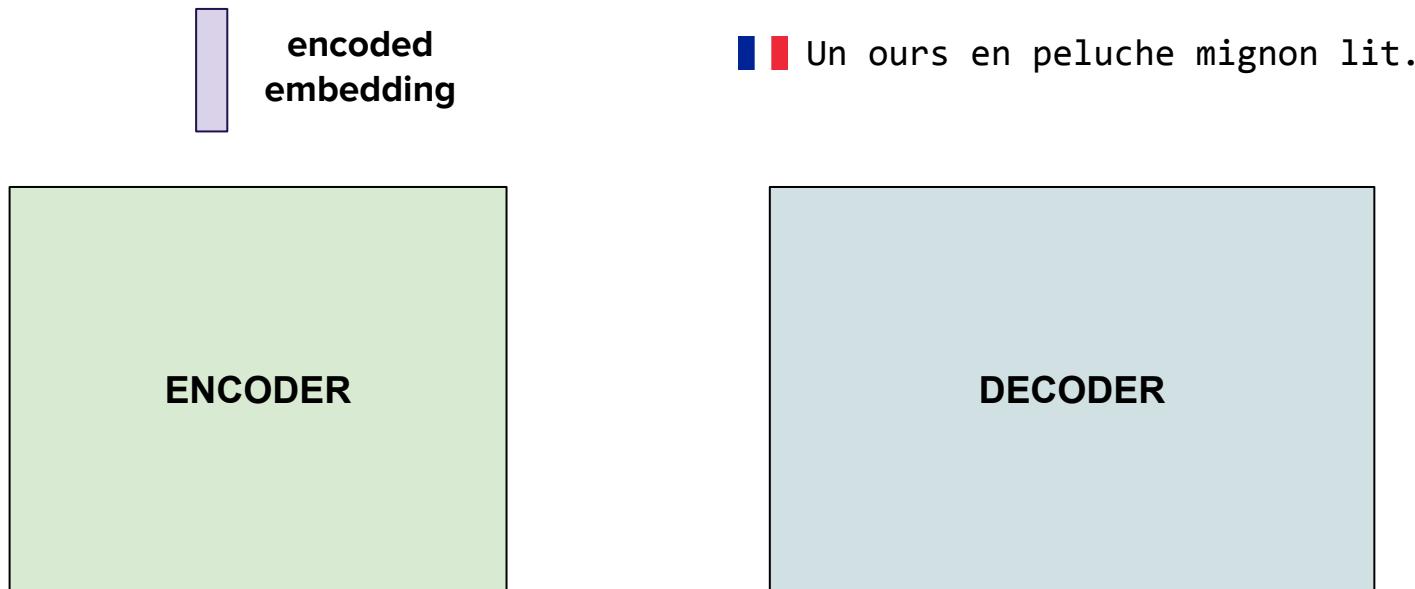
# Stitching all the pieces together with an example

Machine translation summary.



A cute teddy bear  
is reading.

# Stitching all the pieces together with an example



🇺🇸 A cute teddy bear is reading.

Thank you for your attention!

---