# CME 295: Transformers &
# Large Language Models

**Afshine Amidi** & **Shervine Amidi**

**Afshine** and **Shervine**

# Teaching staff



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



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

ICME

# 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**

ICME

# 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

ICME

# 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)

ICME

# Material

**Class website**. cme295.stanford.edu

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

# Material

**Class website**. 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

ICME

# Material



**VIP Cheatsheet**

(also translated in 11 languages so far)

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

Stanford University

**Canvas**.

- Announcements
- Class discussions via Ed

ICME

# Class communications

**Canvas**.

- Announcements
- Class discussions via Ed

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

- cme295-aut2526-staff@lists.stanford.edu
- afshine@stanford.edu, shervine@stanford.edu

ICME

# Example of a slide in this class

Explanation of a CME 295 concept

Source & suggested reading, if interested

Stanford University

*Some paper*

ICME

# Disclaimer before starting: many abbreviations….

BERT

SFT

BLEU

FLAN

BPE

WNLI

MLM

MRPC

GPT

LaaJ

LSTM

PoS

T5

QA

ROUGE

GLUE

PEFT

PPO

NER

F1

RAG

PPL

LLaMA

C4

MT

RLHF

GRU

SQuAD

WP

DPO

SP

NLG

METEOR

ICME

# ...but don't worry!

BERT, T5, GPT, LLaMA

**Transformer-based models**

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

**Misc architectures &  techniques**

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

**Training strategies**

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

**Tasks**

MNLI, WNLI, C4, SQuAD, GLUE, MRPC

**Datasets**

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

**Metrics**

# CME 295 Transformers & Large Language Models

**Stanford University**

## NLP overview

Tokenization

Word representation

RNNs

Self-attention mechanism

Transformer architecture

End-to-end example

**Classification**

Input text → Model → 3

- Sentiment extraction

- Intent detection

- Language detection

- Topic modeling

ICME

**Classification**

**"Multi"-classification**

Input text → Model → 3

Input text → Model → Input text
                        5   1

- Sentiment extraction

- Intent detection

- Language detection

- Topic modeling

- Part of speech tagging

- Named entity recognition

- Dependency parsing

- Constituency parsing

ICME

# NLP tasks overview

## Classification

Input text → Model → ③

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

## "Multi"-classification

Input text → Model → Input text
                      ⑤      ①

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

## Generation

Input text → Model → Out text

- Machine translation
- Question answering
- Summarization
- Text generation

ICME

This teddy bear is SO CUTE! → Model → ( + )

ICME

# NLP task: Sentiment Extraction

This teddy bear is SO CUTE! → Model → +

## 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

# NLP task: Named Entity Recognition

A cute teddy bear is reading... → Model → A cute [teddy bear]ENTITY is reading...

# NLP task: Named Entity Recognition

```
A cute teddy bear is reading...  →  Model  →  A cute [ENTITY teddy bear] is reading...
```

## Datasets

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

## Evaluation metrics

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

ICME

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

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

## Datasets

🇺🇸🇫🇷 WMT'14 English-French          🇺🇸🇩🇪 WMT'14 English-German

## Evaluation metrics

- BLEU ➡ quality of text translated, similar to "precision"

- ROUGE ➡ quality of text generated, similar to "recall"

- Perplexity ➡ quantifies how 'surprised' the model is to see some words together

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

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

Theoretical foundations

**2013**     Word2vec
**2017**     Transformers
**2020s**     Large Language Models

Lots of data, growing
computing power
Fast iterations on ideas

A cute teddy bear is reading.

Stanford University

ICME

A cute teddy bear is reading.

**arbitrary**   | A | cute | teddy bear | is | reading | . |

A cute teddy bear is reading.

| A | cute | teddy bear | is | reading | . |

**word**

| A | cute | teddy | bear | is | reading | . |

A cute teddy bear is reading.

| A | cute | teddy bear | | is | reading | . |

| A | cute | teddy | bear | is | reading | . |

**sub-word** | A | cute | ted | ##dy | bear | is | read | ##ing | . |

ICME

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 | ##ing | . |

| A | _ | c | u | t | e | _ | t | e | d | d | y | _ | b | e | a | r | _ | i | s | _ | r | e | a | d | i | n | g | . |

# Tokenization summary

| Method | Pros | Cons |
| --- | --- | --- |
| **Word-level** | • Simple<br>• Interpretable | • Risk of OOV<br>• Does not leverage knowledge of root |
| **Subword-level**<br><br>e.g. WordPiece, BPE | • Leverages common prefixes and suffixes<br>• Learned from the data | • Risk of OOV, though less than word-level |
| **Character-level** | • Small chance of OOV<br>• RoBUsT tO CASinG anD MIspeliNGs | • Makes computations slower<br>• Embeddings not interpretable |

ICME

**Motivation**

Naive (one-hot) encoding



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

## Motivation

Naive (one-hot) encoding

Learned embedding



$$\text{soft} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

$\langle \text{teddy bear, book} \rangle = 0$

$\langle \text{teddy bear, soft} \rangle = 0$

$$\text{soft} = \begin{pmatrix} 0.95 \\ 0.32 \\ 0.01 \end{pmatrix}$$

$\langle \text{teddy bear, book} \rangle \sim 0$

$\langle \text{teddy bear, soft} \rangle \sim 1$

## Overview

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

## Proxy tasks

- CBOW (continuous bag of words)

...A cute teddy bear is reading...

- Skip-gram

...A cute teddy bear is reading...

**Architecture**

output                 size V

hidden                size d

input                  size V

# Word2vec

**Example with predicting next word**

A [cute] teddy bear is reading



[A] cute teddy bear is reading

**Example with predicting next word**

A cute teddy bear is reading



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

A cute teddy bear is reading

**Example with predicting next word**

A [cute] teddy bear is reading

[0.2, 0.9]

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

[A] cute teddy bear is reading

# Word2vec

**Example with predicting next word**

A `cute` teddy bear is reading

[0.2, 0.4, 0.1, 0.1, 0.1, 0.1]

[0.2, 0.9]

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

`A` cute teddy bear is reading

# Word2vec

**Example with predicting next word**

A cute teddy bear is reading



A cute teddy bear is reading

ICME

**Example with predicting next word**

A cute `teddy bear` is reading

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

A `cute` teddy bear is reading

**Example with predicting next word**

A cute teddy bear is reading

[0.8, 0.4]

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

A cute teddy bear is reading

ICME

# Word2vec

**Example with predicting next word**

A cute [teddy bear] is reading

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

[0.8, 0.4]

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

A [cute] teddy bear is reading

ICME

# Word2vec

**Example with predicting next word**

A cute teddy bear $\boxed{\texttt{is}}$ reading



A cute $\boxed{\texttt{teddy bear}}$ is reading
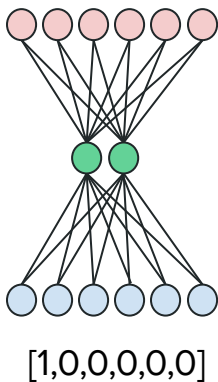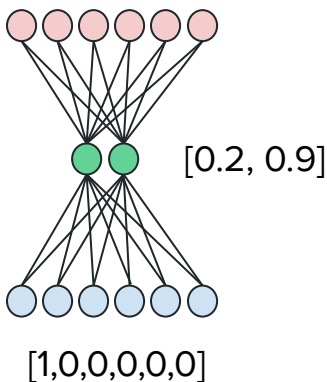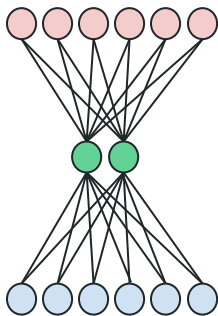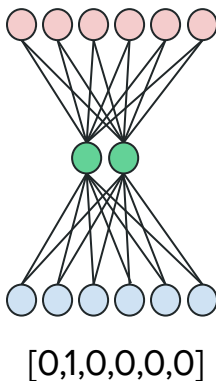
ICME

# Word2vec

**Example with predicting next word**

A cute teddy bear is `reading`



A cute teddy bear `is` reading

A cute teddy bear is reading

output        size V

hidden        size d

input         size V

A cute teddy bear is reading

teddy bear → soft

Persian poetry → art

# CME 295
# Transformers & Large Language Models

Stanford University

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

## General form



*Figure adapted from "VIP cheatsheets for Stanford's CS 230", Amidi.* stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks

A

# Recurrent Neural Networks (RNNs)

# Recurrent Neural Networks (RNNs)

# Recurrent Neural Networks (RNNs)

## Classification

Input text $\rightarrow$ Model $\rightarrow$ (3)

Sentiment



Opinion

## "Multi"-classification

Input text $\rightarrow$ Model $\rightarrow$ $\underline{\text{Input}}$ $\underline{\text{text}}$
                                          (5)    (1)

Tags



Text

## Generation

Input text $\rightarrow$ Model $\rightarrow$ Out text

Translation



Source

*Figure adapted from "VIP cheatsheets for Stanford's CS 230", Amidi. stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks*

ICME

# Long Short-Term Memory (LSTM)

## Overview

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

## General form



*Figure adapted from "VIP cheatsheets for Stanford's CS 230", Amidi. stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks*

# Summary of main methods (non-exhaustive list)

| Method | Pros | Cons |
|---|---|---|
| **Word2vec**<br><br>e.g. CBOW, Skip-gram | • Very simple, yet powerful<br>• Intuitive embeddings | • Word order does not count<br>• Embeddings not context aware |
| **Recurrent Neural Networks**<br><br>e.g. traditional RNN, LSTM | • Word order matters<br>• State-of-the-art results | • Vanishing gradient problem<br>• Slow computations |

ICME

# 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

*"Neural Machine Translation by Jointly Learning to Align and Translate", Bahdanau et al., 2014.*

ICME

# 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



*"Neural Machine Translation by Jointly Learning to Align and Translate", Bahdanau et al., 2014.*

# 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



*"Neural Machine Translation by Jointly Learning to Align and Translate", Bahdanau et al., 2014.*

- Introduced in 2014

- Translation tasks had a real issue with long-term dependencies

- Seq2seq unable to "remember" what input sentence was saying

| | | | | | Un | ours en peluche | mignon | lit |
|---|---|---|---|---|---|---|---|---|

| A | cute | teddy bear | is | reading | | | | |
|---|---|---|---|---|---|---|---|---|

# Transformers & Large Language Models

Stanford University

ICME

# 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

ICME

- 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    .
```

- 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

Concept of **Q**uery, **K**ey, **V**alue

a   cute   **teddy bear**   is   reading   .

Concept of **Q**uery, **K**ey, **V**alue

$q_{\text{teddy bear}}$

a cute **teddy bear** is reading .

Concept of **Q**uery, **K**ey, **V**alue



*"Super Study Guide: Transformers & Large Language Models", by Amidi, 2024*

Concept of **Q**uery, **K**ey, **V**alue



*"Super Study Guide: Transformers & Large Language Models", by Amidi, 2024*

Efficient computations with matrices:

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

# Attention mechanism

Efficient computations with matrices:

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



Scaled Dot-Product Attention

Multi-Head Attention

*Figure adapted from "Attention Is All You Need", Vaswani et al., 2017.*

**Attention layer** (MHA)

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

**Feed Forward Neural Network** (FFNN)

**Positional Encoding** (PE)

# Input



## Overview

- Text is "tokenized"
- Learned embeddings for tokens

## Parameters

- V: vocabulary size
- `d_model`: embedding dimensions

**Positional encoding**

Idea:

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



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

*Figure adapted from "Attention Is All You Need", Vaswani et al., 2017.*

# Output "shifted right"



## Overview

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

## Parameters

- V: vocabulary size
- `d_model`: embedding dimensions

*Figure adapted from "Attention Is All You Need", Vaswani et al., 2017.*

# 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

## Overview

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

## Parameters

- `V`: vocabulary size
- `d_model`: embedding dimensions

Figure adapted from *"Attention Is All You Need"*, Vaswani et al., 2017.

**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



Output Probabilities
Softmax
Linear
Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention
Nx
Add & Norm
Feed Forward
Nx
Add & Norm
Multi-Head Attention
Add & Norm
Masked Multi-Head Attention
Positional Encoding
Positional Encoding
Input Embedding
Output Embedding
Inputs
Outputs (shifted right)

**Label smoothing**

Idea:

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

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

Benefits:

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

*Figure adapted from "Attention Is All You Need", Vaswani et al., 2017.*

ICME

# Transformers & Large Language Models

Stanford University

ICME

A cute teddy bear is reading.

Stanford University

# Stitching all the pieces together with an example

| A | cute | teddy bear | is | reading | . |

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

ICME

# Stitching all the pieces together with an example

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

ICME

embedding

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

Stanford University

ICME

position embedding

embedding

[BOS]  A  cute  teddy bear  is  reading  .  [EOS]

**position-aware embedding**

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

**position-aware embeddings**

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

**position-aware embeddings matrix**

[BOS]   A   cute   teddy bear   is   reading   .   [EOS]

# Stitching all the pieces together with an example

position-aware
embeddings
matrix

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

ICME

encoder

position-aware
embeddings matrix

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

# Stitching all the pieces together with an example



encoder

position-aware embeddings matrix

| Wq | Wk | Wv |

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

ICME

encoder

position-aware
embeddings
matrix

Q    K    V

Wq    Wk    Wv

[BOS]   A   cute   teddy bear   is   reading   .   [EOS]

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

**Q**

**K**

**V**

**encoder**

Wq

Wk

Wv

position-aware
embeddings
matrix

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

ICME

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

**Q**

**Wq**      **Wk**      **Wv**

**encoder**

position-aware
embeddings
matrix

[BOS]    A    cute    teddy bear    is    reading    .    [EOS]

PAUSE

# Stitching all the pieces together with an example
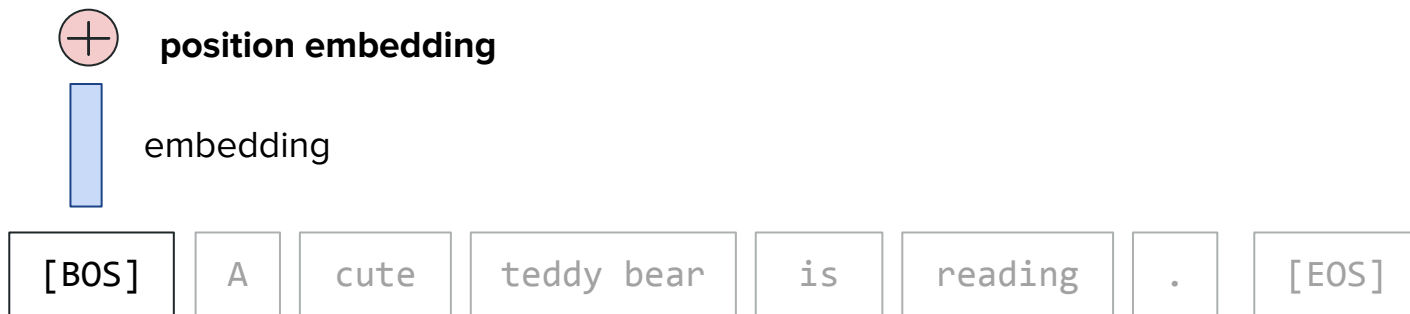


Q

[BOS]

A

cute

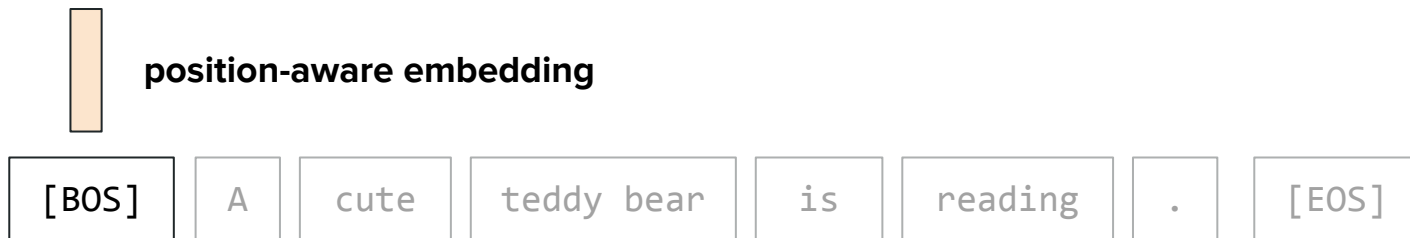teddy bear

is

reading

.

[EOS]

ICME

$Q$

[BOS]

A

cute

teddy bear

is

reading

.

[EOS]

[BOS]  A  cute  teddy bear  is  reading  .  [EOS]

$K^T$

Stanford University

ICME

$Q$

[BOS]

A

cute

teddy bear

is

reading

.

[EOS]

[BOS]

A

cute

teddy bear

is

reading

.

[EOS]

$K^T$

ICME

$$\begin{array}{cccc} \langle q_{[\text{BOS}]}, k_{[\text{BOS}]} \rangle & \langle q_{[\text{BOS}]}, k_A \rangle & \langle q_{[\text{BOS}]}, k_{\text{cute}} \rangle & \cdots \\[1em] \langle q_A, k_{[\text{BOS}]} \rangle & \langle q_A, k_A \rangle & & \\[1em] \langle q_{\text{cute}}, k_{[\text{BOS}]} \rangle & & \ddots & \\[1em] \vdots & & & \end{array} \qquad QK^T$$

$QK^T$

$$\langle q_{[\text{BOS}]}, k_{[\text{BOS}]} \rangle \quad \langle q_{[\text{BOS}]}, k_{\text{A}} \rangle \quad \langle q_{[\text{BOS}]}, k_{\text{cute}} \rangle \quad \dots$$

$$\langle q_{\text{A}}, k_{[\text{BOS}]} \rangle \qquad \langle q_{\text{A}}, k_{\text{A}} \rangle$$

$$\langle q_{\text{cute}}, k_{[\text{BOS}]} \rangle \qquad\qquad \ddots$$

$$\vdots$$

**V**

[BOS]

A

cute

teddy bear

is

reading

.

[EOS]

ICME

$$\langle q_{\text{[BOS]}}, k_{\text{[BOS]}} \rangle \, v_{\text{[BOS]}} + \langle q_{\text{[BOS]}}, k_{\text{A}} \rangle \, v_{\text{A}} + \langle q_{\text{[BOS]}}, k_{\text{cute}} \rangle \, v_{\text{cute}} + \ldots$$

$$\langle q_{\text{A}}, k_{\text{[BOS]}} \rangle \, v_{\text{[BOS]}} + \langle q_{\text{A}}, k_{\text{A}} \rangle \, v_{\text{A}} + \langle q_{\text{A}}, k_{\text{cute}} \rangle \, v_{\text{cute}} + \ldots$$

$$\vdots$$

$$QK^{T}V$$

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

weighted average of values

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

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

**Q**  **K**  **V**

Wq  Wk  Wv

**encoder**

position-aware
embeddings
matrix

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \qquad \boxed{\;|...|\;}$$

**h x** $\left(\begin{array}{ccc} \boxed{\phantom{Q}}\; \mathbf{Q} & \boxed{\phantom{K}}\; \mathbf{K} & \boxed{\phantom{V}}\; \mathbf{V} \\ \boxed{\text{Wq}} & \boxed{\text{Wk}} & \boxed{\text{Wv}} \end{array}\right)$     **encoder**

**position-aware embeddings matrix**

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |
|---|---|---|---|---|---|---|---|

Stanford University

# Stitching all the pieces together with an example



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

Wo

h x ( Q Wq K Wk V Wv )

encoder

position-aware embeddings matrix

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |

# Stitching all the pieces together with an example
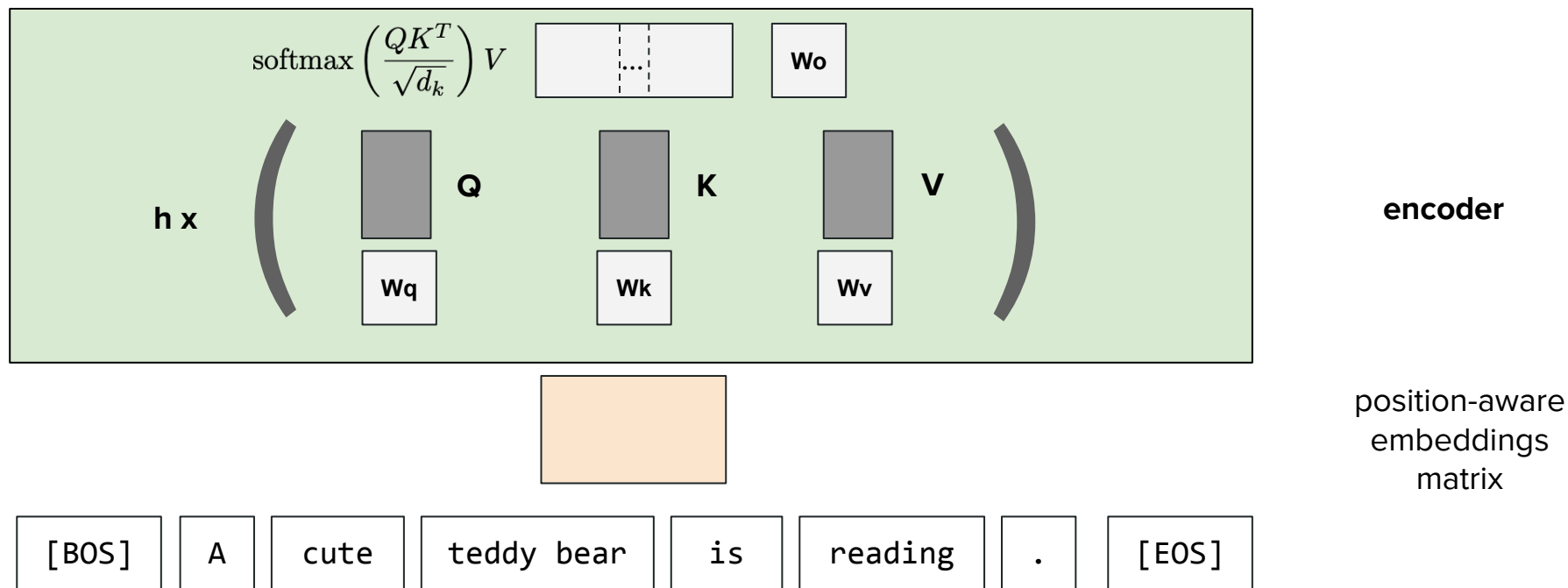
# Stitching all the pieces together with an example

**context-aware encoded embeddings**

| Feed forward network |
| --- |

**encoder**

| Self-attention layer |
| --- |

position-aware embeddings matrix

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |
| --- | --- | --- | --- | --- | --- | --- | --- |

ICME

# Stitching all the pieces together with an example



**context-aware encoded embeddings**

**ENCODER**

position-aware embeddings matrix

| [BOS] | A | cute | teddy bear | is | reading | . | [EOS] |
|---|---|---|---|---|---|---|---|

**encoded embedding**

**ENCODER**

```
A cute teddy bear
   is reading.
```

ICME

encoded embedding

ENCODER

A cute teddy bear is reading.

Stanford University

encoded embedding

ENCODER

A cute teddy bear
   is reading.

[BOS]

ICME

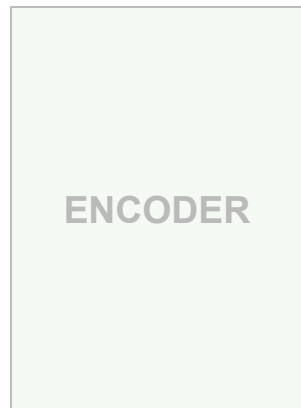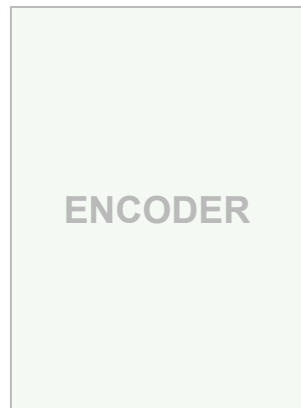# Stitching all the pieces together with an example

encoded
embedding

ENCODER

A cute teddy bear
is reading.

decoder

[BOS]

ICME

**encoded embedding**

**ENCODER**

A cute teddy bear is reading.

**decoder**

**Self-attention layer**

[BOS]

Stanford University

ICME

**encoded embedding**

**ENCODER**

**decoder**

**Self-attention layer**

A cute teddy bear is reading.

[BOS]

Stanford University

ICME

# Stitching all the pieces together with an example

encoded embedding

ENCODER

decoder

Encoder - Decoder attention layer

Self-attention layer

A cute teddy bear is reading.

[BOS]

ICME

**encoded embedding**

**ENCODER**

A cute teddy bear is reading.

**decoder**

Encoder - Decoder attention layer

Self-attention layer

[BOS]

ICME

# Stitching all the pieces together with an example

encoded embedding

**Feed forward network**

**Encoder - Decoder attention layer**
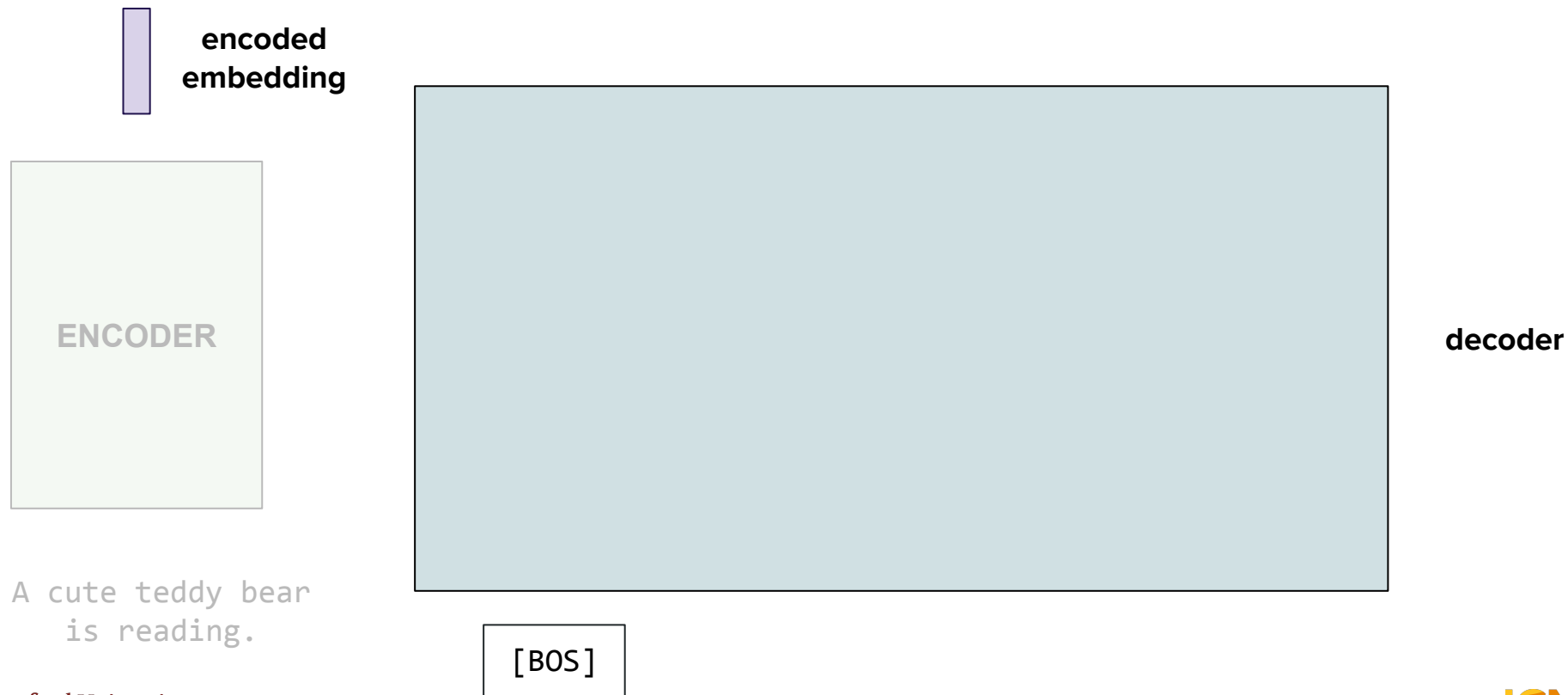
**Self-attention layer**

ENCODER

decoder

A cute teddy bear is reading.

[BOS]

ICME

# Stitching all the pieces together with an example

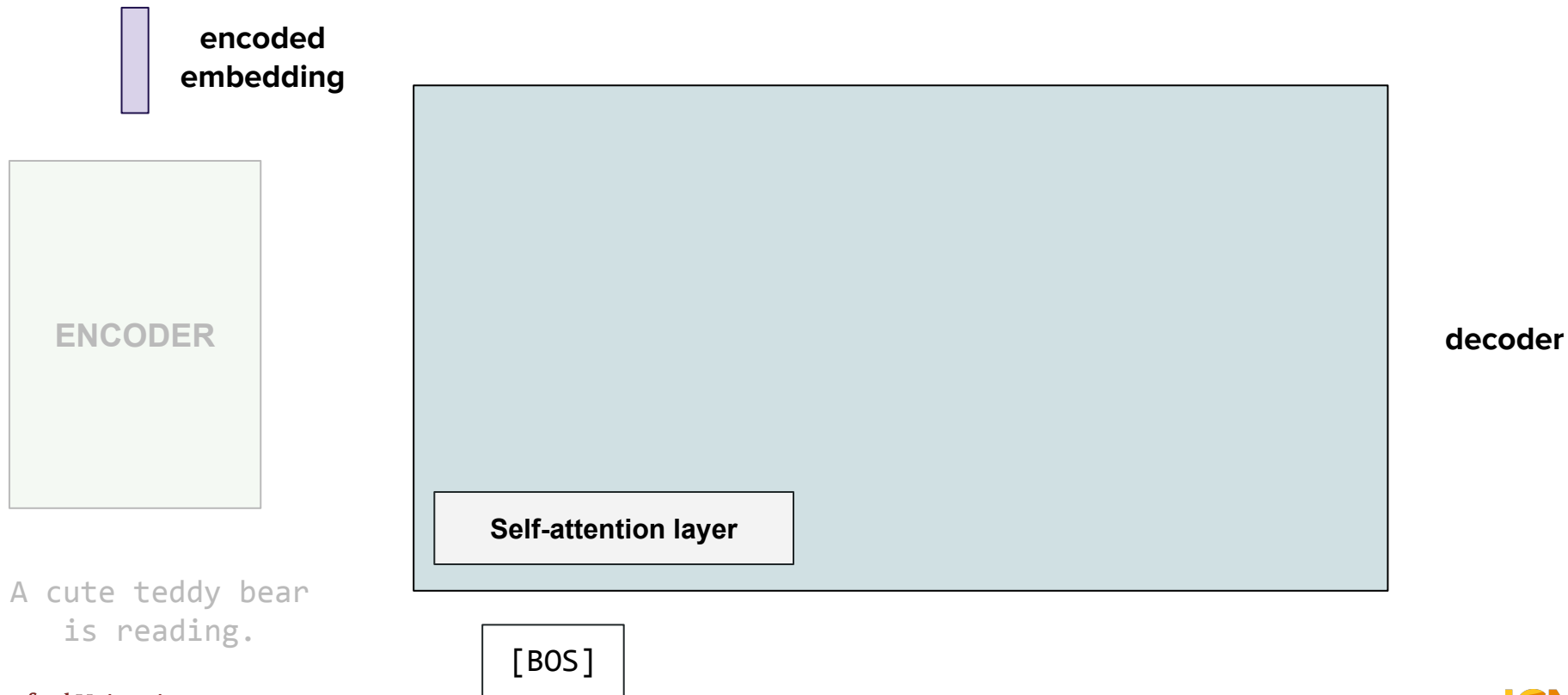

encoded embedding
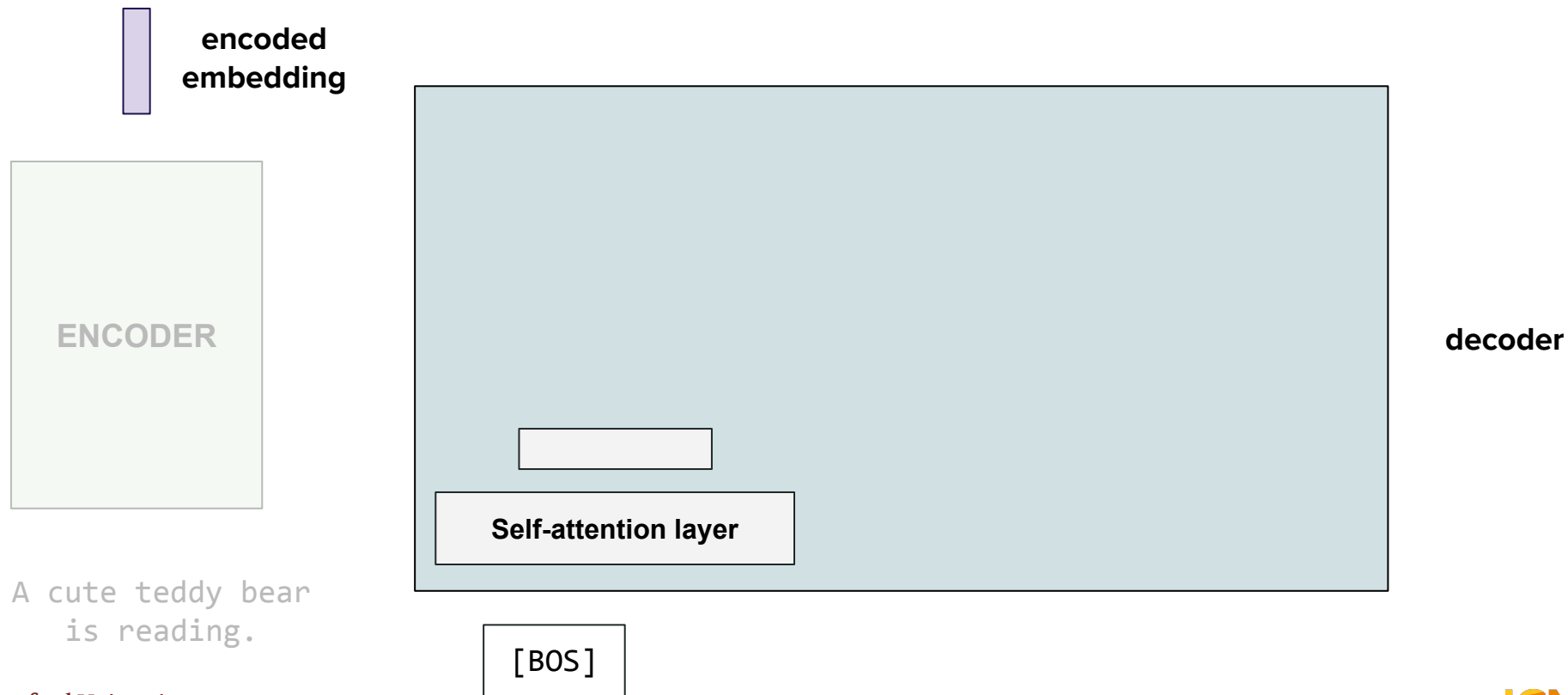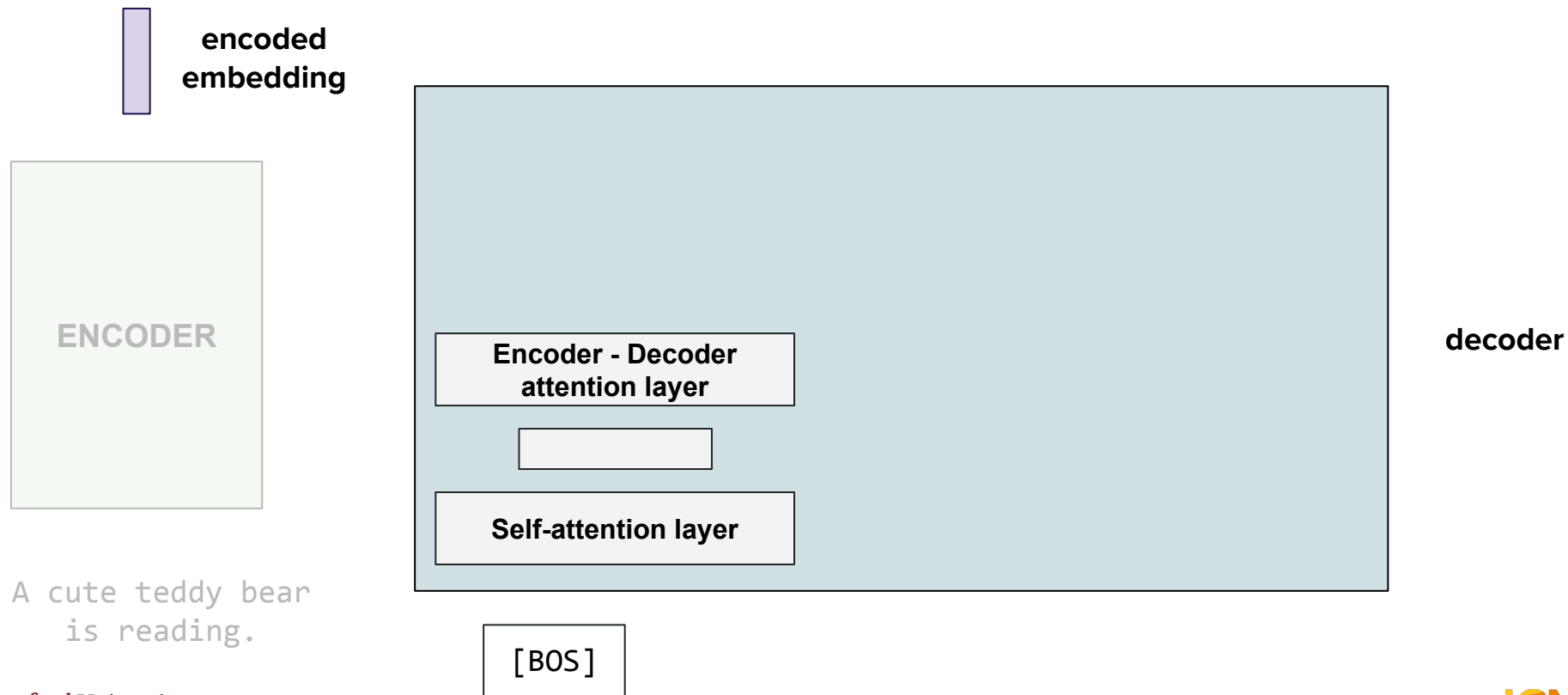
ENCODER

A cute teddy bear is reading.

Softmax layer

Feed forward network

Encoder - Decoder attention layer

Self-attention layer

decoder

[BOS]

# Stitching all the pieces together with an example

encoded embedding

[0.001, 0.0003, ..., 0.4, ..., 0.002]

ENCODER

**Softmax layer**

**Feed forward network**

**Encoder - Decoder attention layer**

**Self-attention layer**

decoder

A cute teddy bear is reading.

[BOS]

ICME

# Stitching all the pieces together with an example

encoded embedding

ENCODER

A cute teddy bear
is reading.

Un

decoder

Softmax layer

Feed forward network

Encoder - Decoder
attention layer

Self-attention layer

[BOS]

ICME

encoded embedding

Un

ENCODER

DECODER

A cute teddy bear is reading.

[BOS]

Stanford University

ICME

encoded embedding

Un

ENCODER

DECODER

A cute teddy bear is reading.

[BOS]

Un

ICME

# Stitching all the pieces together with an example

encoded embedding

Un   |   ours en peluche

ENCODER

DECODER

A cute teddy bear is reading.

[BOS]   |   Un

ICME

**encoded embedding**

| Un | ours en peluche | mignon | lit | [EOS] |

**ENCODER**

**DECODER**

A cute teddy bear is reading.

| [BOS] | Un | ours en peluche | mignon | lit |

ICME

# Stitching all the pieces together with an example



encoded embedding

🇫🇷 Un ours en peluche mignon lit.

**ENCODER**

**DECODER**

🇺🇸 A cute teddy bear is reading.

# Thank you for your attention!