# Deep Learning for NLP, part II
## Stanford ICME Summer workshop 2021

Instructor: **Afshine Amidi**

18-20 August 2021

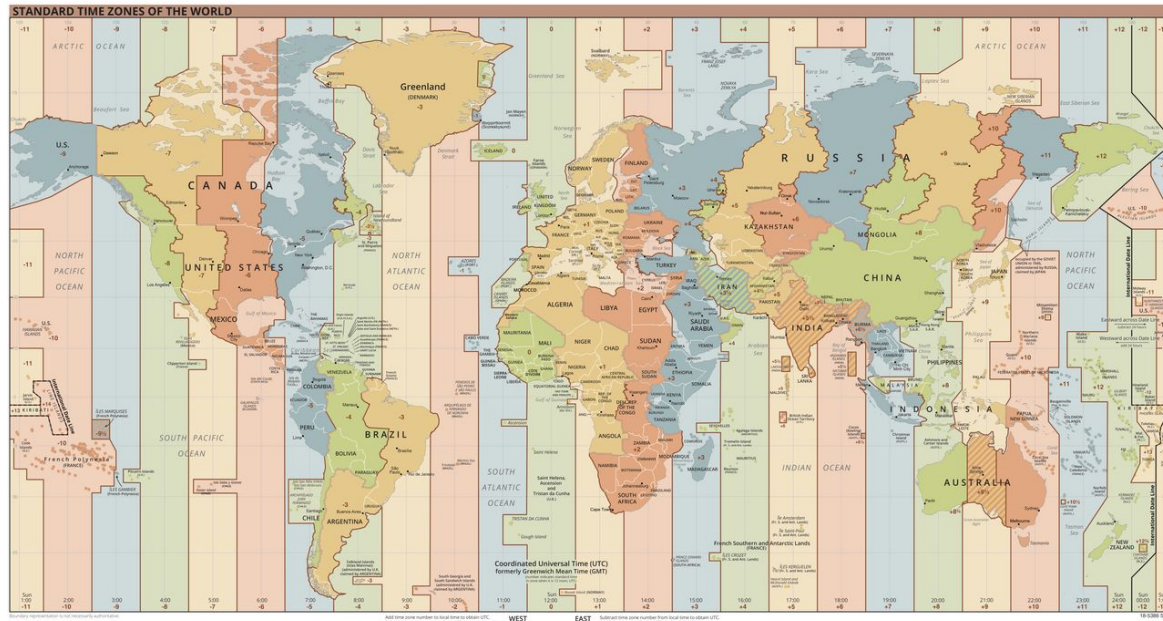# Teaching staff



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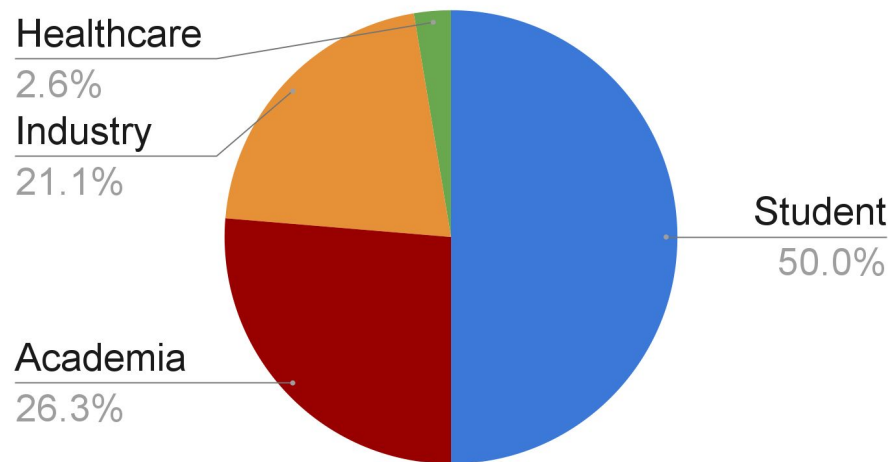


**Sam**, TA
Stanford ('21)
Peerlift, Iris Labs
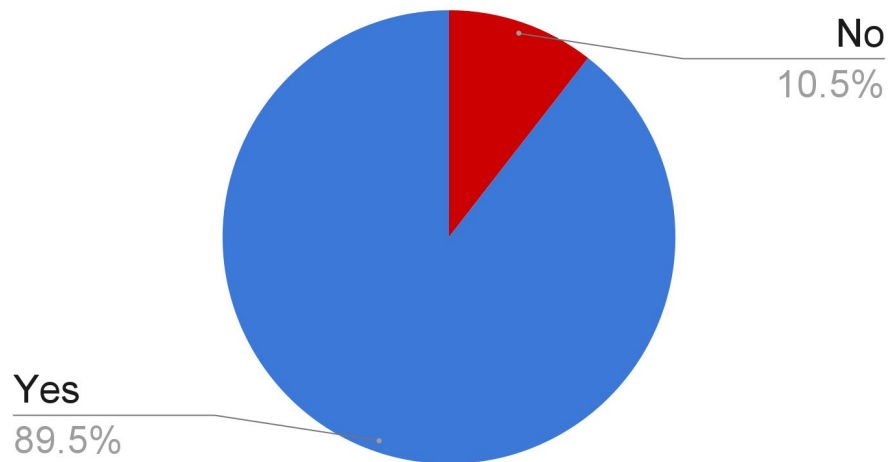
# Poll results

**What is (approximately) your timezone?**

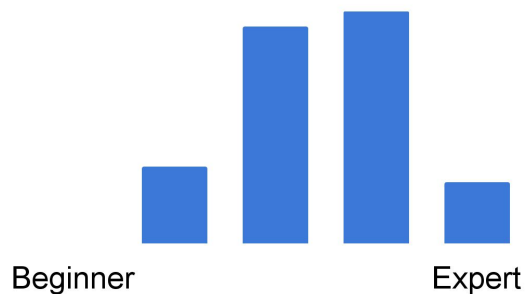# Poll results

**Participant category**

# Poll results

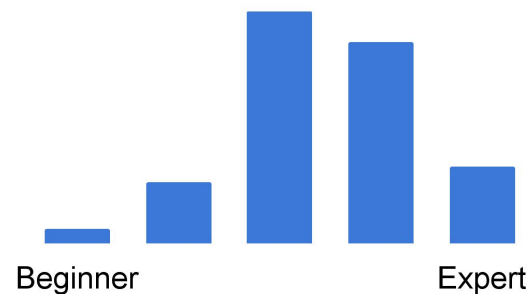**Have you attended the first part of the ICME Workshop NLP series?**

No
10.5%

Yes
89.5%

# Poll results

**Familiarity level**

# Anonymized quotes from feedback section

"*Some **context on evolution of NLP** will be super helpful*"

# Anonymized quotes from feedback section

*"Some **context on evolution of NLP** will be super helpful"*

*"Start from **medium level** and then go upwards to **higher difficulty**"*

# Anonymized quotes from feedback section

*"Some **context on evolution of NLP** will be super helpful"*

*"Start from **medium level** and then go upwards to **higher difficulty**"*

*"Hope to get the **summary of the materials** (including additional articles/books) and **links to them** to have better understanding. Hope to **try BERT in practice** (in Python notebooks). [...]"*

*I am very interested in **hands on** experience. I hope this session will help to start running my first deep learning model.*
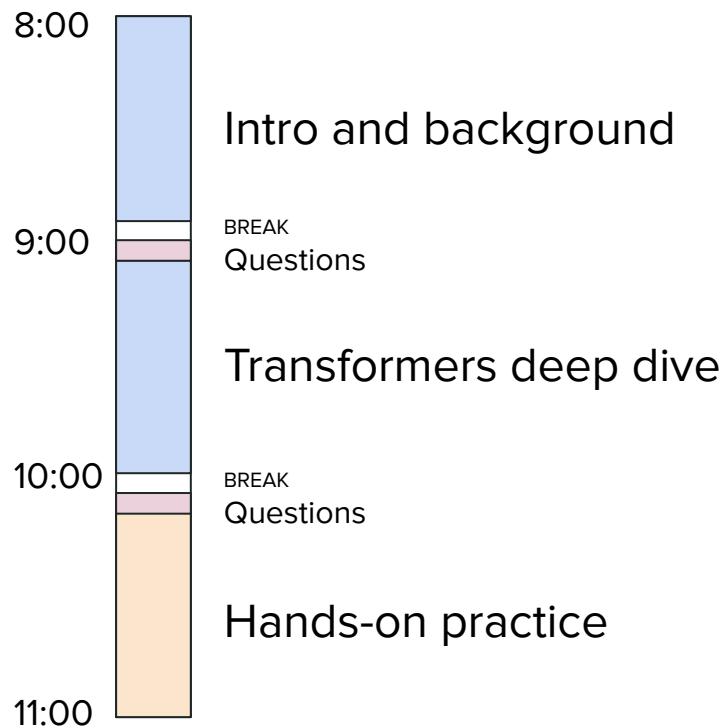
*"[...] My hope is we **dig into the code** and the **details** of running a program"*

*"Among the given topics, I am more interested in the application areas of **sentiment extraction**"*
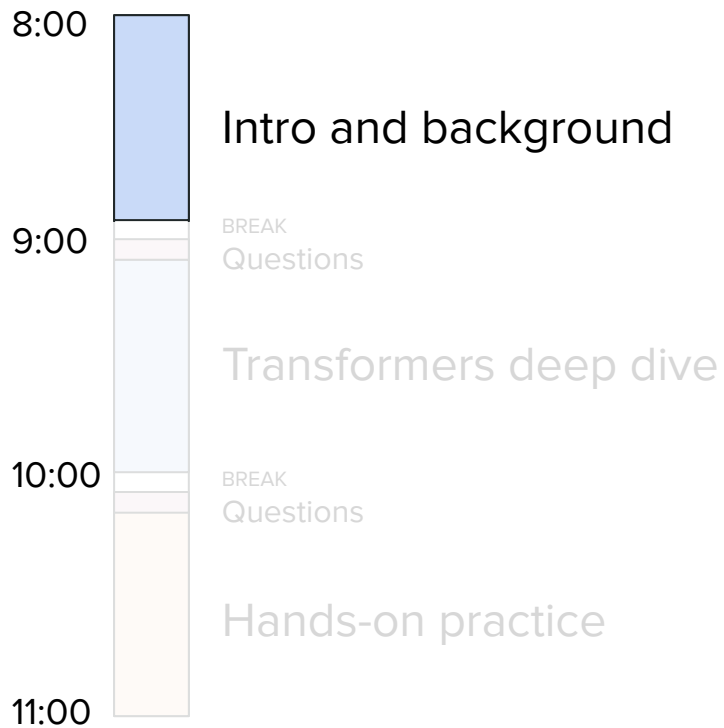
# Logistics

- **Two half-days**
  - Wednesday 8/18, 8am - 11am PT
  - Friday 8/20, 8am - 11am PT
- **Hands-on** format
  - ~2/3 slides
  - ~1/3 code via Colab
- **Questions**
  - Preferably ask questions via Ed
  - Pause from time to time to answer questions
  - After each break, dedicated time for Q&A
- **Homework** between the 2 days
  - Apply concepts in a practical use case
  - Completely optional, but recommended
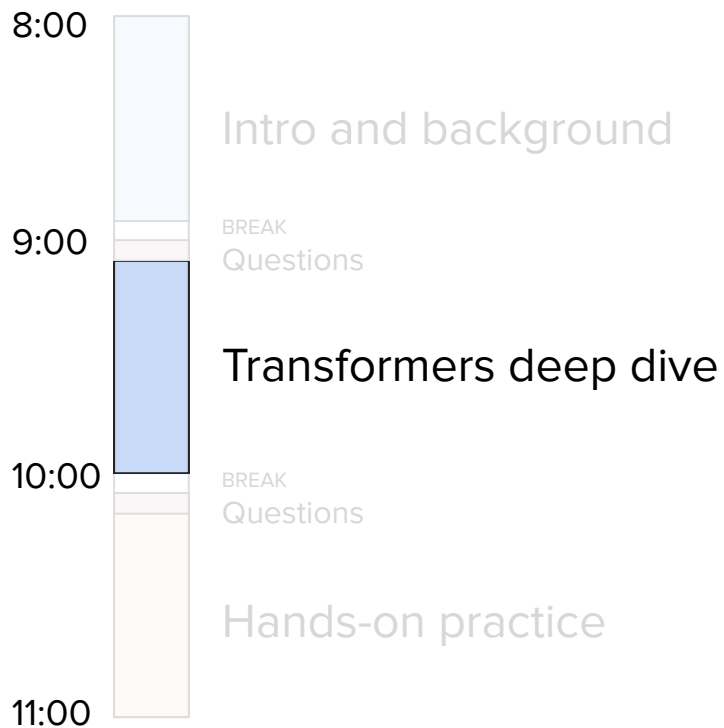
# Tentative schedule for today

8:00

Intro and background

9:00 — BREAK
Questions

Transformers deep dive

10:00 — BREAK
Questions

Hands-on practice

11:00

*All times are in Pacific Time (UTC-7)*

# Tentative schedule for today

8:00

**Intro and background**

9:00
BREAK
Questions

Transformers deep dive

10:00
BREAK
Questions

Hands-on practice

11:00

*All times are in Pacific Time (UTC-7)*

# Tentative schedule for today

8:00

Intro and background

9:00

BREAK
Questions

Transformers deep dive

10:00

BREAK
Questions

Hands-on practice

11:00

*All times are in Pacific Time (UTC-7)*

# Tentative schedule for today

8:00

Intro and background

9:00
BREAK
Questions

Transformers deep dive

10:00
BREAK
Questions

Hands-on practice

11:00

*All times are in Pacific Time (UTC-7)*

# Tentative schedule for today

8:00

Intro and background

9:00      **BREAK**
       Questions

Transformers deep dive

10:00      **BREAK**
       Questions

Hands-on practice

11:00

*All times are in Pacific Time (UTC-7)*

# Deep Learning for NLP, part II

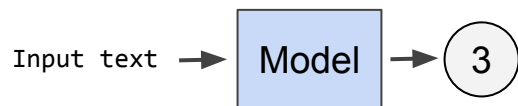Stanford ICME Summer workshop 2021

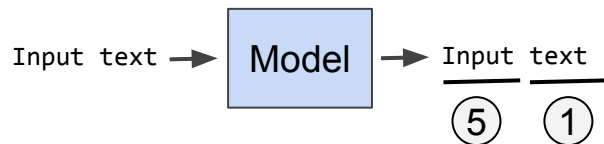**Motivation and setup**

Background

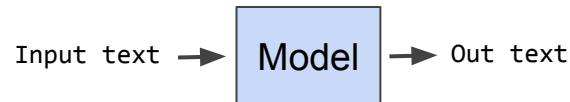Transformers

BERT

Conclusion

# 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

# NLP task: Sentiment Extraction

This teddy bear is SO CUTE! → Model → +

## Datasets
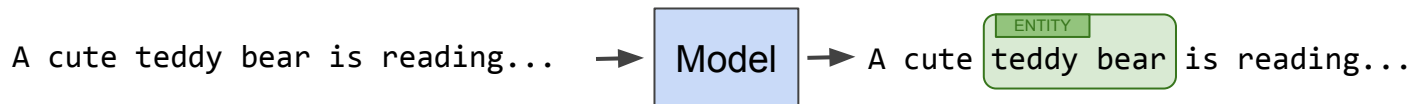
Amazon reviews                IMDB critiques                Twitter

## 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 `ENTITY` teddy bear` `is reading...`
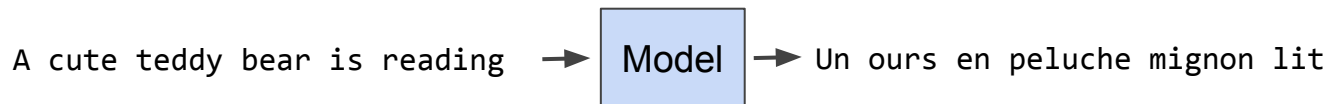
## Datasets

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

## Evaluation metrics

- Accuracy
- Precision
- Recall
- F1 score

at a token level, per entity type

# NLP task: Machine Translation

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

# Standardized benchmark for NLP

**GLUE**: **G**eneral **L**anguage **U**nderstanding **E**valuation

| | | | |
|---|---|---|---|
| Grammatical correctness<br><br>*CoLA* | Paraphrase<br><br>*MRPC* | Similarity<br><br>*QQP, STS-B* | Common sense<br><br>*WNLI* |
| Entailment<br><br>*RTE, MNLI* | Sentiment Extraction<br><br>*SST-2* | Question Answering<br><br>*QNLI* | **Glue score** |

# Disclaimer before starting: many abbreviations….

BLEU

ELMo

BPE

WNLI

MLM

BERT

CoNLL

MRPC

GPT

EM

LSTM

PoS

T5

QA

ROUGE

GLUE

WMT

NER

F1

GloVe

C4

MT

ACL

mT5

PPL

GRU

SQuAD

WP

EMNLP

SP

NLG

METEOR

# ...but don't worry!

BERT, DistilBERT, ALBERT, T5, mT5, GPT

**Transformer-based models**

LSTM, GRU, GloVe, ELMo, BPE, WP, SP

**Some techniques**

ACL, EMNLP, WMT, CoNLL

**Conferences**

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

**Tasks**

MNLI, WNLI, C4, SQuAD, GLUE, MRPC

**Datasets**

F1, PPL, ROUGE, BLEU, METEOR, EM

**Metrics**

# Deep Learning for NLP, part II

Stanford ICME Summer workshop 2021

Motivation and setup

**Background**

Transformers

BERT

Conclusion

# High-level timeline

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

Theoretical foundations

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

**2013**      Word2vec
**2017**      Transformers

Lots of data, growing
computing power
Fast iterations on ideas

# Word representations

**Motivation**

Naive (one-hot) encoding



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

# Word representations

**Motivation**

Naive (one-hot) encoding



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

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

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

Learned embedding



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

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

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

# Word2vec

## 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...

# Word2vec

## Architecture

output                                   size V

hidden                                  size d

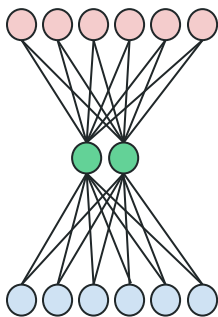input                                     size V

# Word2vec

**Example with left context window = 1**

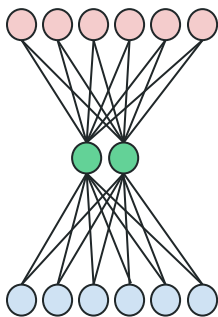A `cute` teddy bear is reading

`A` cute teddy bear is reading

# Word2vec

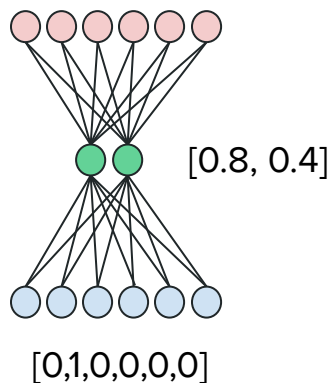**Example with left context window = 1**

A cute teddy bear is reading

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

A cute teddy bear is reading

# Word2vec

**Example with left context window = 1**
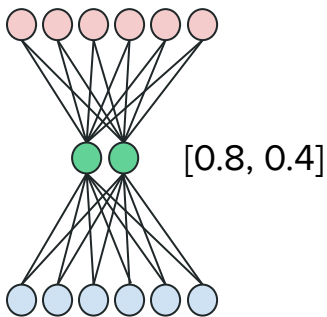
A `cute` teddy bear is reading



[0.2, 0.9]

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

`A` cute teddy bear is reading

# Word2vec

**Example with left context window = 1**

A cute teddy bear is reading
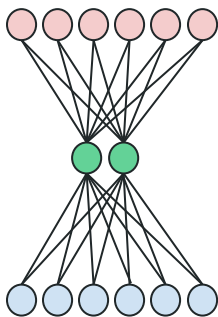
[0.2, 0.4, 0.1, 0.1, 0.1, 0.1]
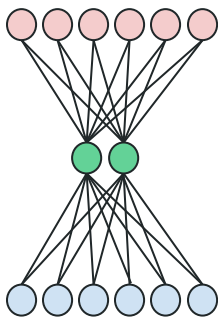
[0.2, 0.9]

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

A cute teddy bear is reading

# Word2vec

**Example with left context window = 1**

A cute `teddy bear` is reading



A `cute` teddy bear is reading
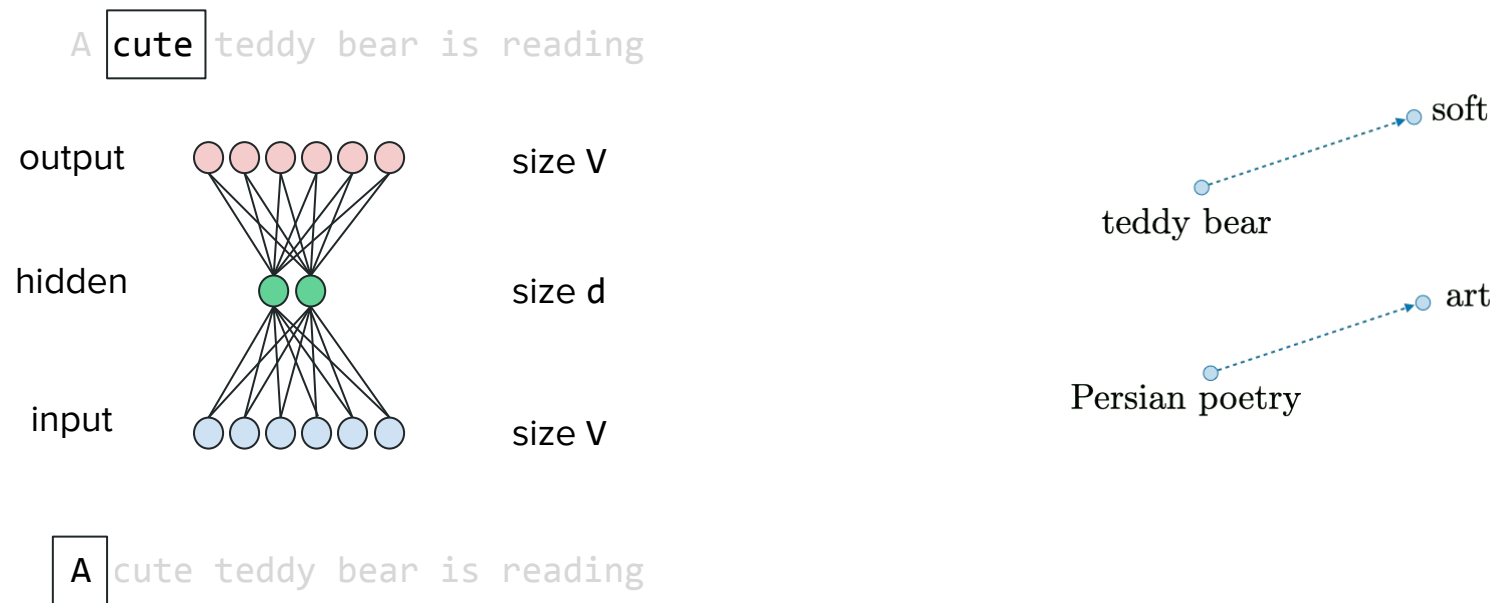
# Word2vec

**Example with left context window = 1**

A cute teddy bear is reading

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

A cute teddy bear is reading

# Word2vec

**Example with left context window = 1**

A cute teddy bear is reading

[0.8, 0.4]

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

A cute teddy bear is reading

# Word2vec

**Example with left context window = 1**



A cute teddy bear is reading

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

[0.8, 0.4]

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

A cute teddy bear is reading

# Word2vec

**Example with left context window = 1**

A cute teddy bear is reading

A cute teddy bear is reading
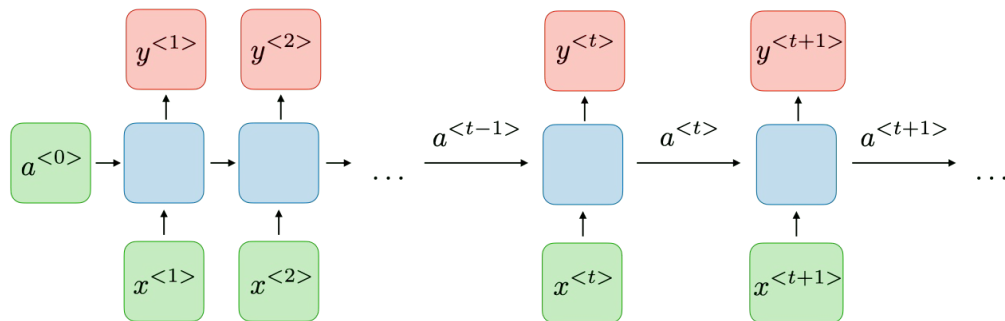
# Word2vec

**Example with left context window = 1**

A cute teddy bear is reading

A cute teddy bear is reading

# Word2vec

## Example with left context window = 1

A |cute| teddy bear is reading

output          size V

hidden          size d

input           size V

|A| cute teddy bear is reading

soft

teddy bear

art

Persian poetry

# Recurrent Neural Networks (RNNs)

## Overview

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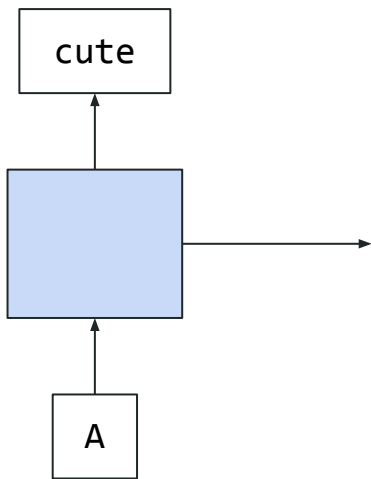- Class of neural networks where connections form a temporal sequence

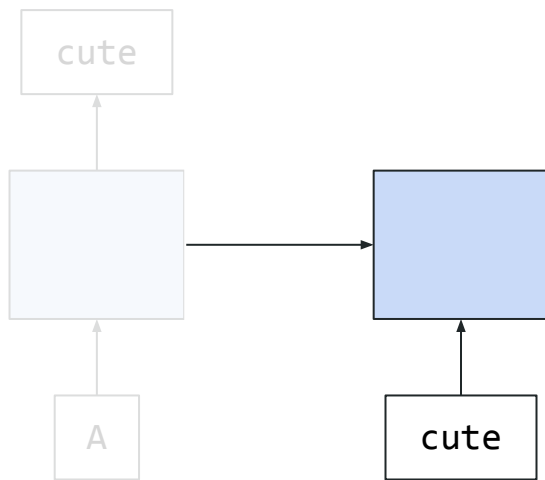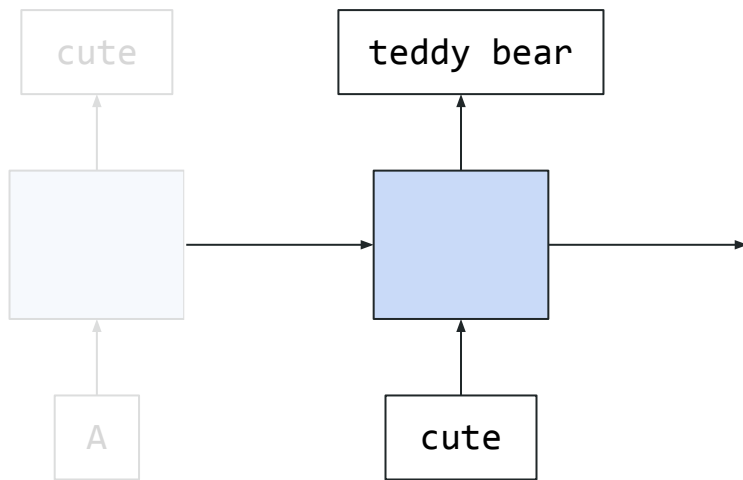## General form

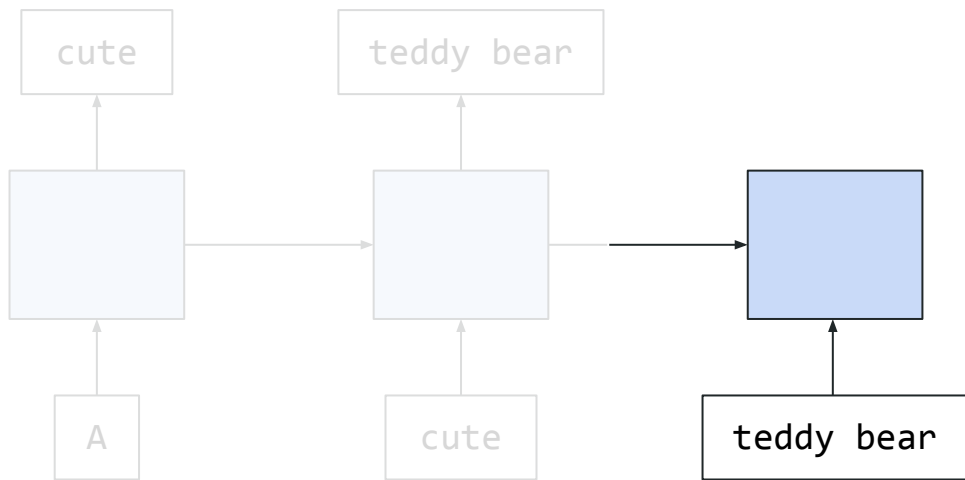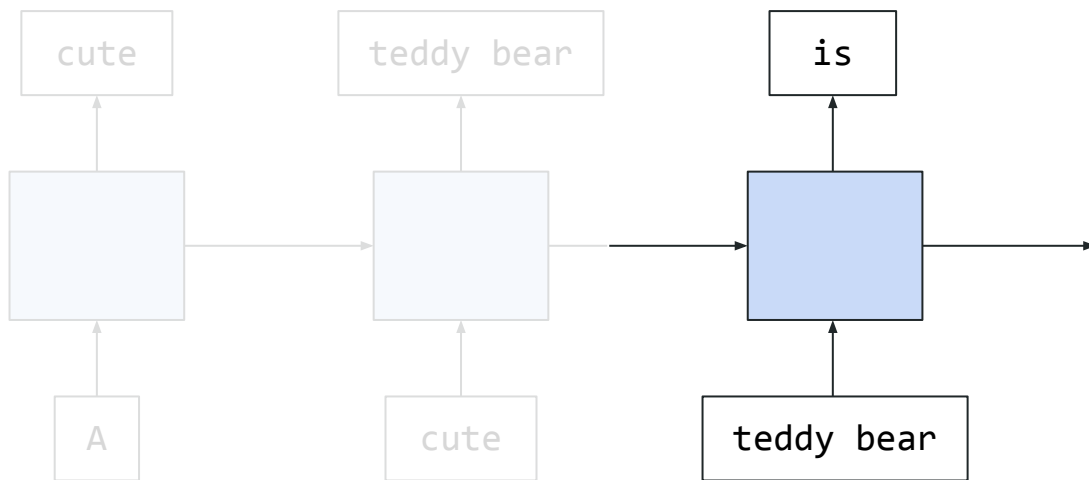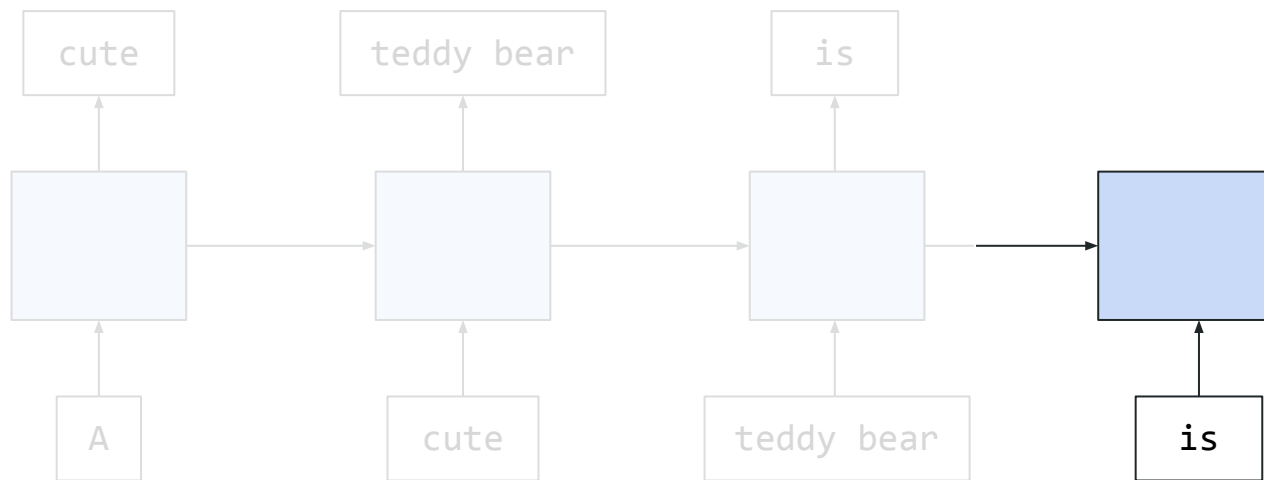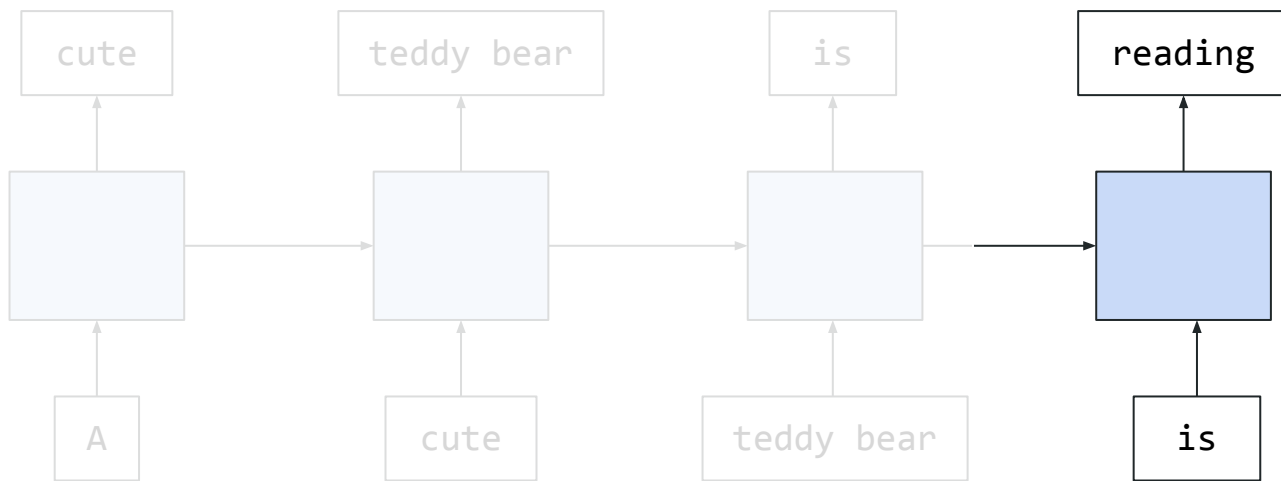# Recurrent Neural Networks (RNNs)

A

# Recurrent Neural Networks (RNNs)

# Recurrent Neural Networks (RNNs)
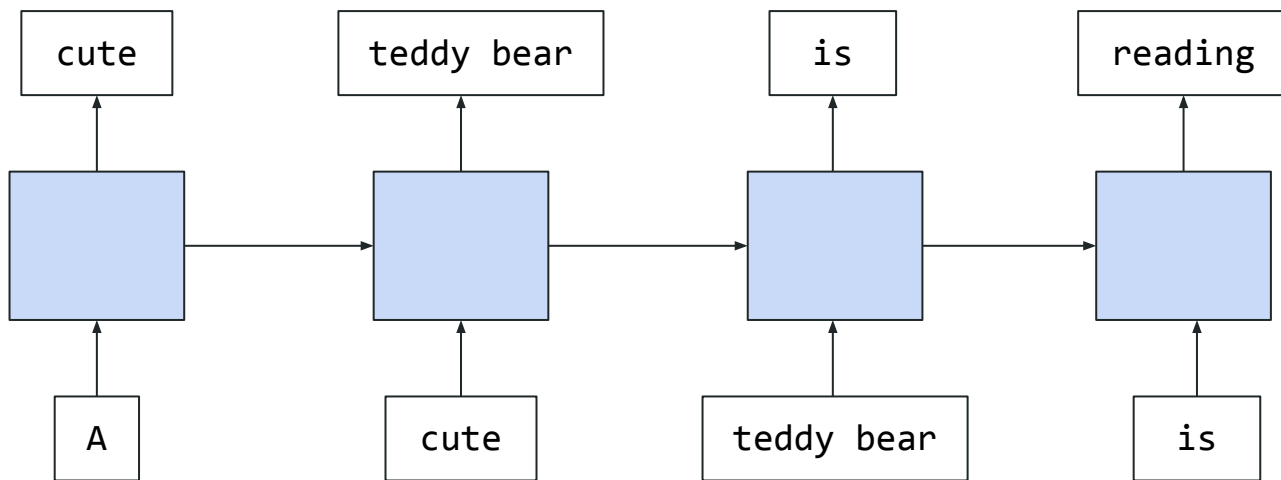
# Recurrent Neural Networks (RNNs)

# Recurrent Neural Networks (RNNs)

# Recurrent Neural Networks (RNNs)

# Recurrent Neural Networks (RNNs)

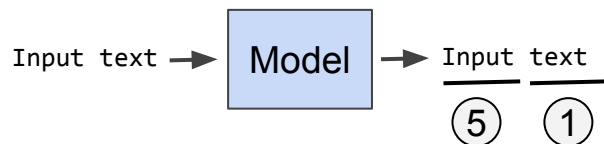# Recurrent Neural Networks (RNNs)

## Classification

Input text → Model → (3)

Sentiment



Opinion

## "Multi"-classification

Input text → Model → Input text (5) (1)

Tags



Text

## Generation

Input text → Model → Out text

Translation



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

## General form

# Summary of main methods (non-exhaustive list)

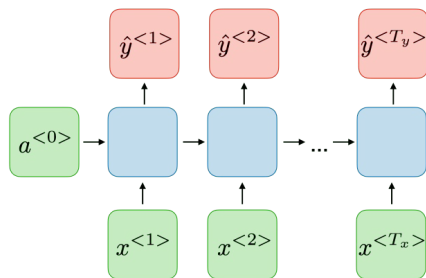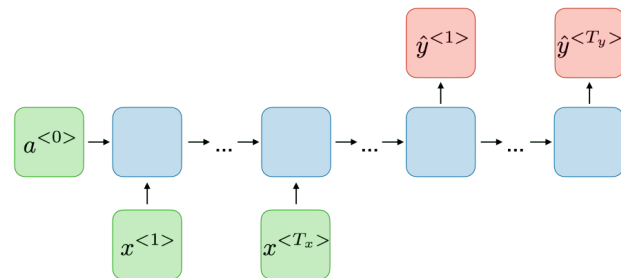| Method | Pros | Cons |
|---|---|---|
| **Word2vec**<br><br>e.g. CBOW, Skip-gram | ● Very simple, yet powerful<br>● Intuitive embeddings | ● Word order does not count<br>● Embeddings not context aware |
| **Recurrent Neural Networks**<br><br>e.g. traditional RNN, LSTM | ● Word order matters<br>● State-of-the-art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

# Break + questions

# Deep Learning for NLP, part II

Stanford ICME Summer workshop 2021

# 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.*

# 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.
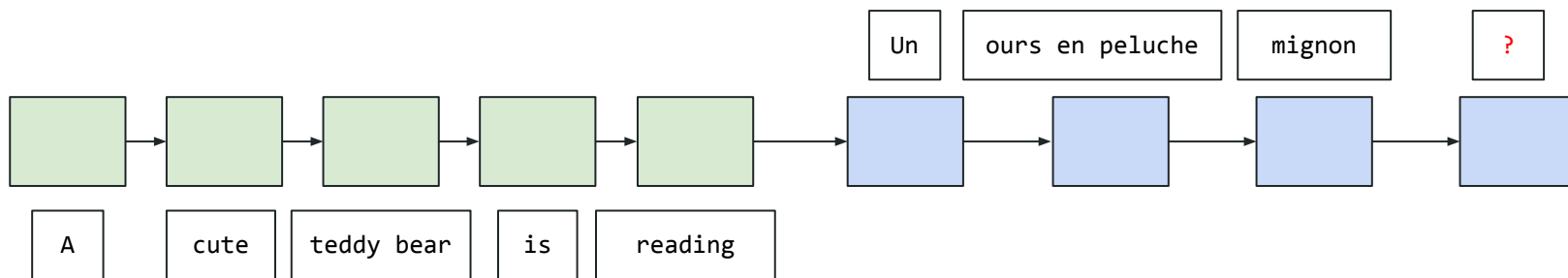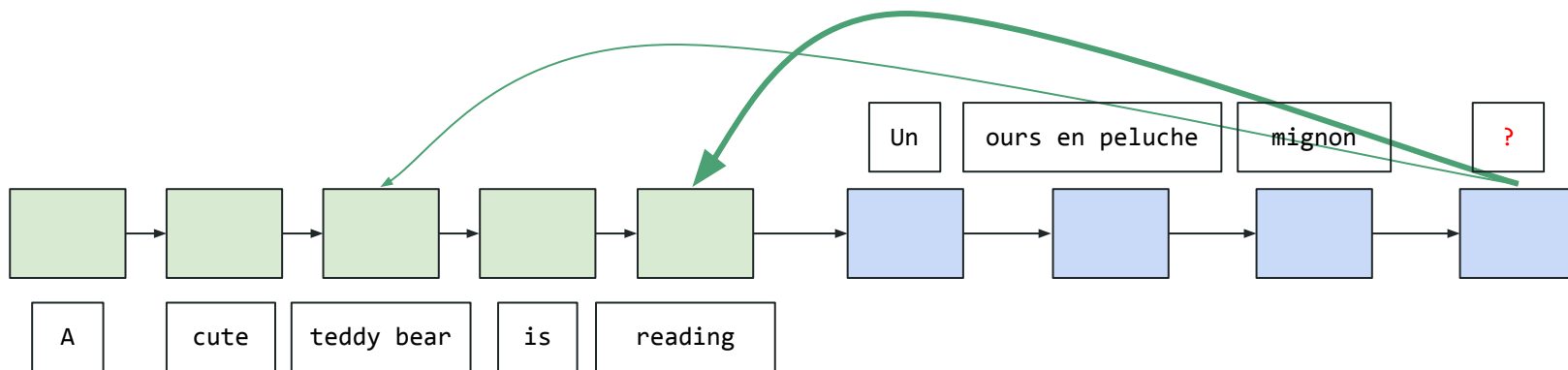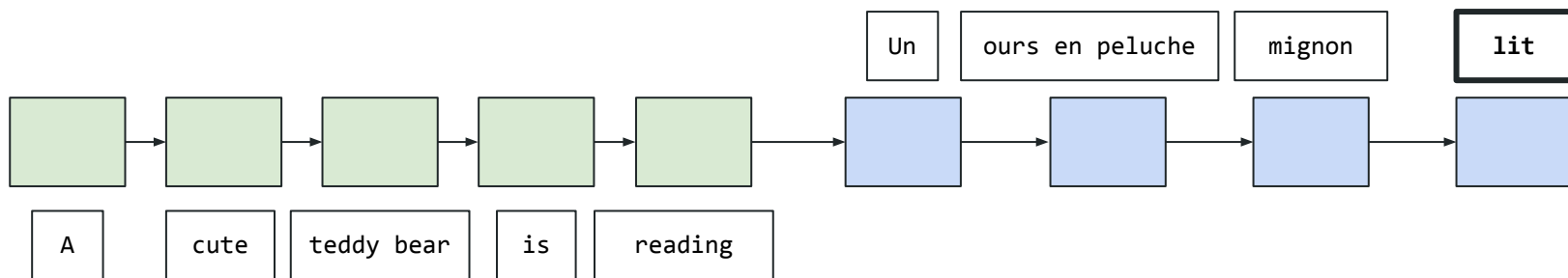
# History of attention

- 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.

# Overview of the Transformer

- Introduced in the 2017 paper "Attention is All You Need"
- Relies on the self-attention mechanism
- Encoder/decoder parts that are used in a lot of models
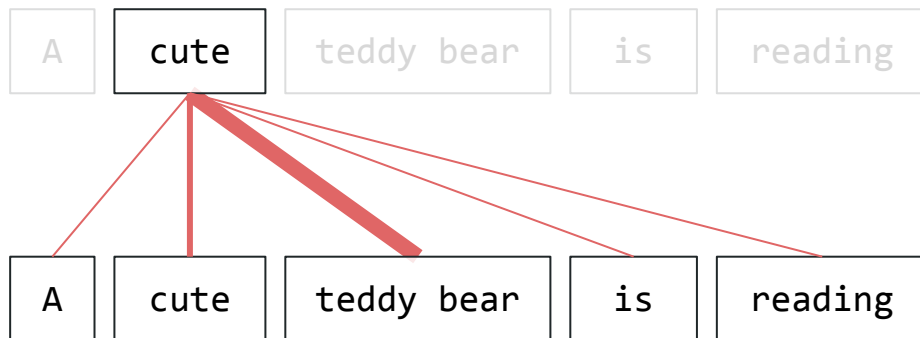- State of the art results on machine translation tasks

# Overview of the Transformer

- Introduced in the 2017 paper "Attention is All You Need"
- Relies on the self-attention mechanism
- Encoder/decoder parts that are used in a lot of models
- State of the art results on machine translation tasks

# Attention mechanism



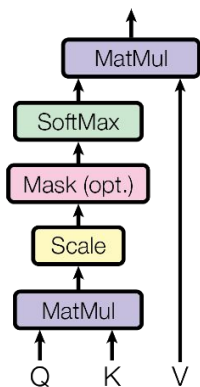Scaled Dot-Product Attention

Multi-Head Attention

- **Q**uery, **K**ey, **V**alue

- Computationally efficient with matrices

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

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

# Transformer architecture



- **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

# … with a trick!

## Positional encoding

Idea:
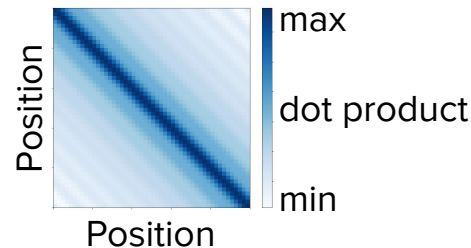
- From a convolutional seq2seq 2017 paper
- Add **position information** to inputs
- Can be either learned or hardcoded



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.*
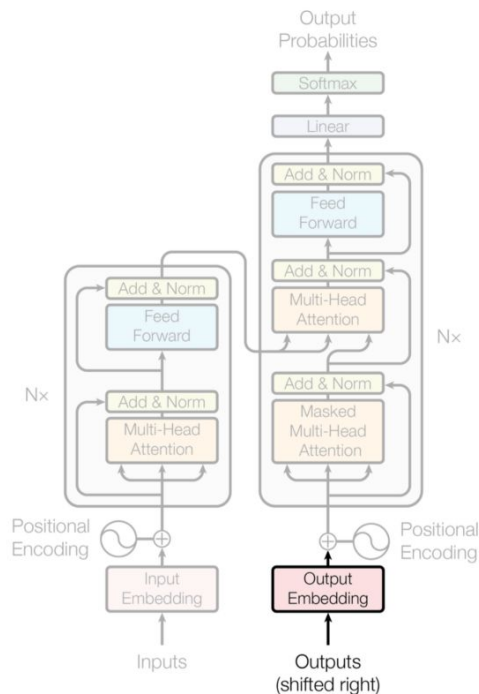
# Decoder



## 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
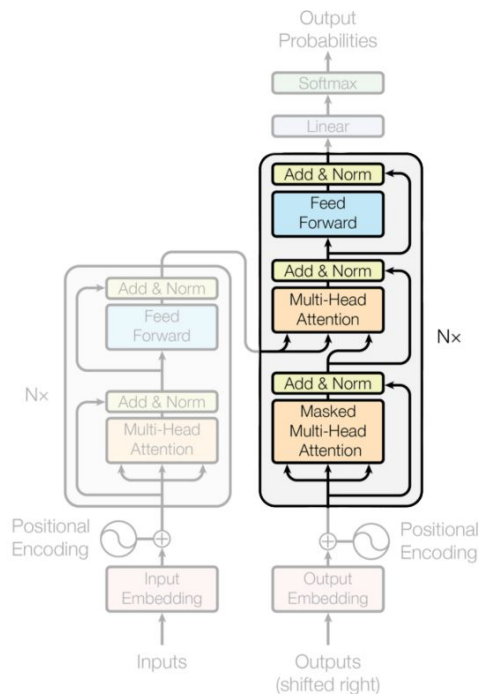- d_model: embedding dimensions

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

# 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

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

# 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

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

# Stitching all the pieces together with an example

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

**embedding**

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

⊕ **position embedding**

embedding

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

**position-aware embedding**

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

# Stitching all the pieces together with an example



**position-aware embeddings**

| [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] |

# Stitching all the pieces together with an example

**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

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

# Stitching all the pieces together with an example

**encoder**

| Wq | Wk | Wv |

position-aware
embeddings
matrix

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

# Stitching all the pieces together with an example

Q

Wq

Wk

Wv

encoder

PAUSE

position-aware
embeddings
matrix

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

# Stitching all the pieces together with an example

**Q**

[BOS]

A

cute

teddy bear

is

reading

.

[EOS]

# Stitching all the pieces together with an example

$Q$

[BOS]

A

cute

teddy bear

is

reading

.

[EOS]

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

$K^T$

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

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

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

$$\vdots$$

$$QK^T$$

$$QK^T$$

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

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

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

$$\vdots$$

**V**

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

# Stitching all the pieces together with an example

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

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

$$\vdots$$

$$QK^T V$$

# 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 $\langle q, k \rangle$

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

**Q**

**K**

**V**

**Wq**

**Wk**

**Wv**

**encoder**

position-aware
embeddings
matrix

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

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

**Wp**

**Q**

**K**

**V**

**Wq**

**Wk**

**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]

# 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
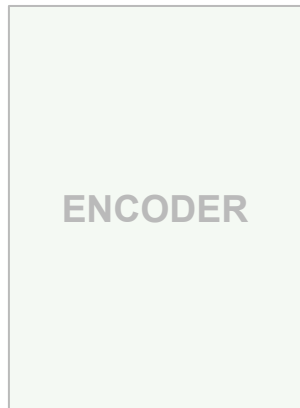
A cute teddy bear
   is reading.

# Stitching all the pieces together with an example

encoded
embedding

ENCODER

A cute teddy bear
   is reading.

encoded
embedding

ENCODER

A cute teddy bear
   is reading.

[BOS]

# Stitching all the pieces together with an example

**encoded embedding**
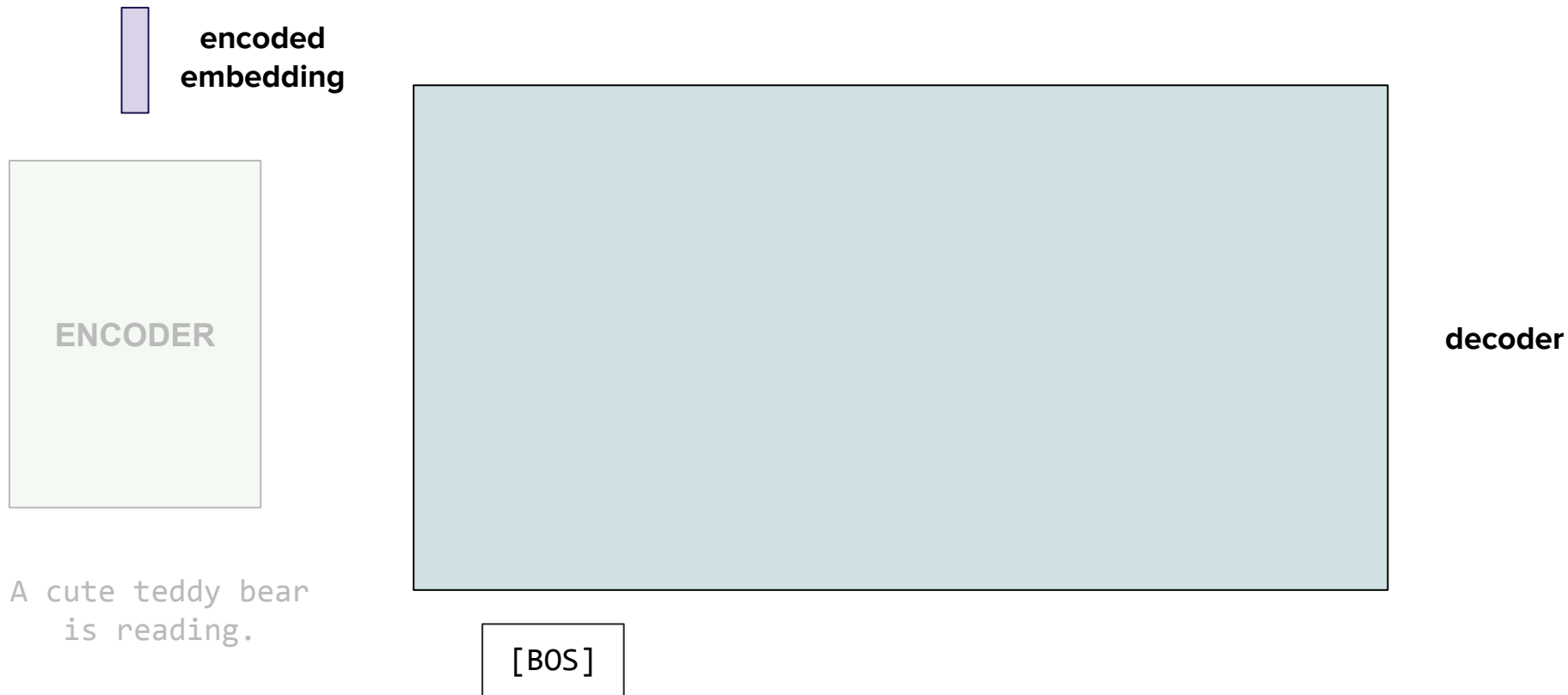
**ENCODER**

A cute teddy bear
is reading.

**decoder**

[BOS]

# Stitching all the pieces together with an example

encoded embedding

ENCODER

A cute teddy bear is reading.

decoder

Self-attention layer

[BOS]

# Stitching all the pieces together with an example

**encoded embedding**

**ENCODER**

A cute teddy bear
is reading.

**decoder**

Self-attention layer

[BOS]

# Stitching all the pieces together with an example

encoded
embedding

ENCODER

A cute teddy bear
is reading.

decoder

Encoder - Decoder
attention layer

Self-attention layer

[BOS]

# Stitching all the pieces together with an example

**encoded embedding**

**ENCODER**

A cute teddy bear
is reading.

**decoder**

Encoder - Decoder
attention layer

Self-attention layer

[BOS]

**encoded embedding**

**ENCODER**

A cute teddy bear
is reading.

**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**

A cute teddy bear
   is reading.

**decoder**

| Softmax layer |
| --- |

| Feed forward network |
| --- |

| |
| --- |

| Encoder - Decoder attention layer |
| --- |

| |
| --- |

| Self-attention layer |
| --- |

[BOS]

encoded
embedding

Un

Softmax layer

Feed forward network

Encoder - Decoder
attention layer

Self-attention layer

ENCODER

decoder

A cute teddy bear
is reading.

[BOS]

# Stitching all the pieces together with an example

Un

encoded embedding

ENCODER

A cute teddy bear is reading.

DECODER

[BOS]  Un

# Stitching all the pieces together with an example

encoded
embedding

Un

ours en peluche

ENCODER

DECODER

A cute teddy bear
is reading.

[BOS]

Un

# Stitching all the pieces together with an example

encoded
embedding

Un     ours en peluche

ENCODER

DECODER

A cute teddy bear
is reading.

[BOS]     Un     ours en peluche

# Stitching all the pieces together with an example

encoded embedding

| Un | ours en peluche | mignon |

**ENCODER**

**DECODER**

A cute teddy bear is reading.

| [BOS] | Un | ours en peluche |

encoded embedding

| Un | ours en peluche | mignon |

ENCODER

DECODER

A cute teddy bear
is reading.

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

# Stitching all the pieces together with an example

# Stitching all the pieces together with an example

encoded embedding

| Un | ours en peluche | mignon | lit |

**ENCODER**

**DECODER**

A cute teddy bear is reading.

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

# Stitching all the pieces together with an example

encoded
embedding

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

**ENCODER**

**DECODER**

A cute teddy bear
is reading.

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

# Stitching all the pieces together with an example

encoded
embedding

🇫🇷 Un ours en peluche mignon lit.

**ENCODER**

**DECODER**

🇬🇧 A cute teddy bear is reading.

# See you on Friday!