# Attention Models

- o Neural Machine Translation

- o Text Summation

- o Question Answering

- Neural Machine Translation

  - o Introduction to Neural Machine Translation

  - o Seq2Seq model and its shortcomings

  - o Solution for the information bottleneck

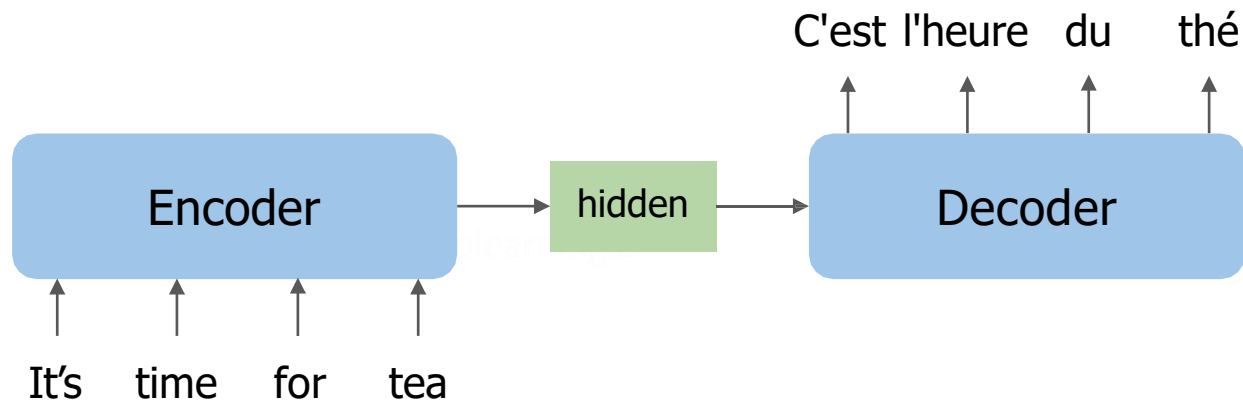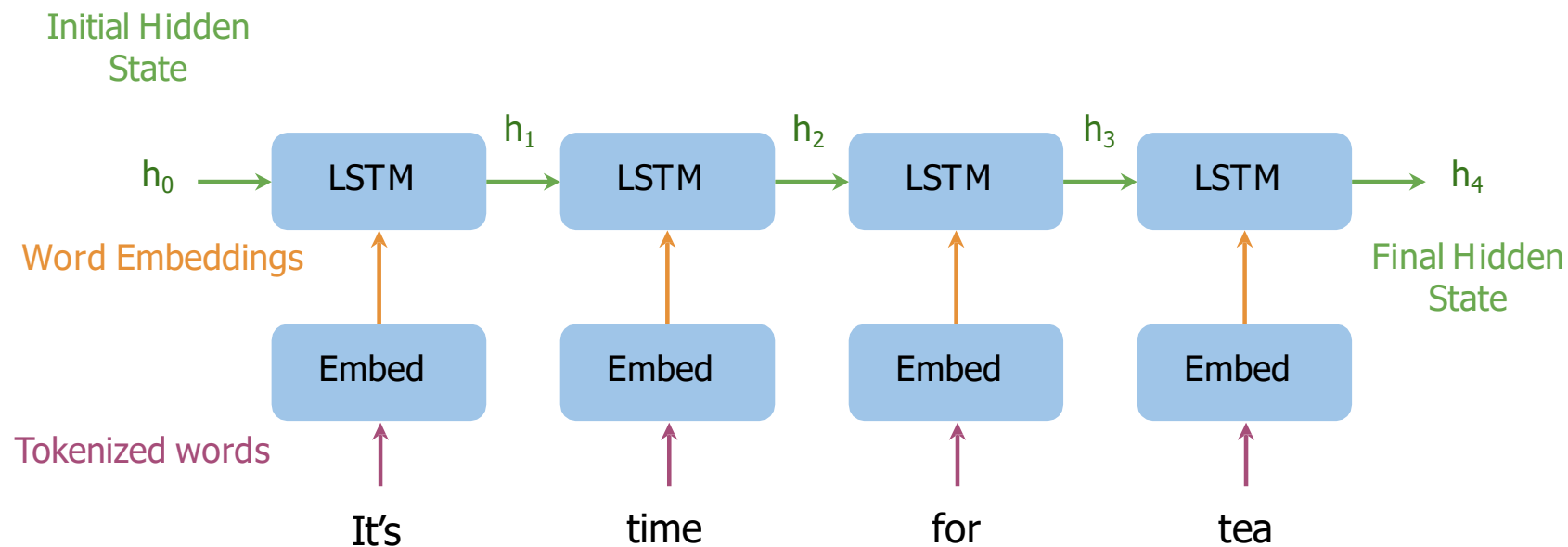It's time for tea   ⟶   C'est l'heure du thé

U              F

# Seq2Seq model

- Introduced by Google in 2014

- Maps variable-length sequences to fixed-length memory

- Inputs and outputs can have different lengths

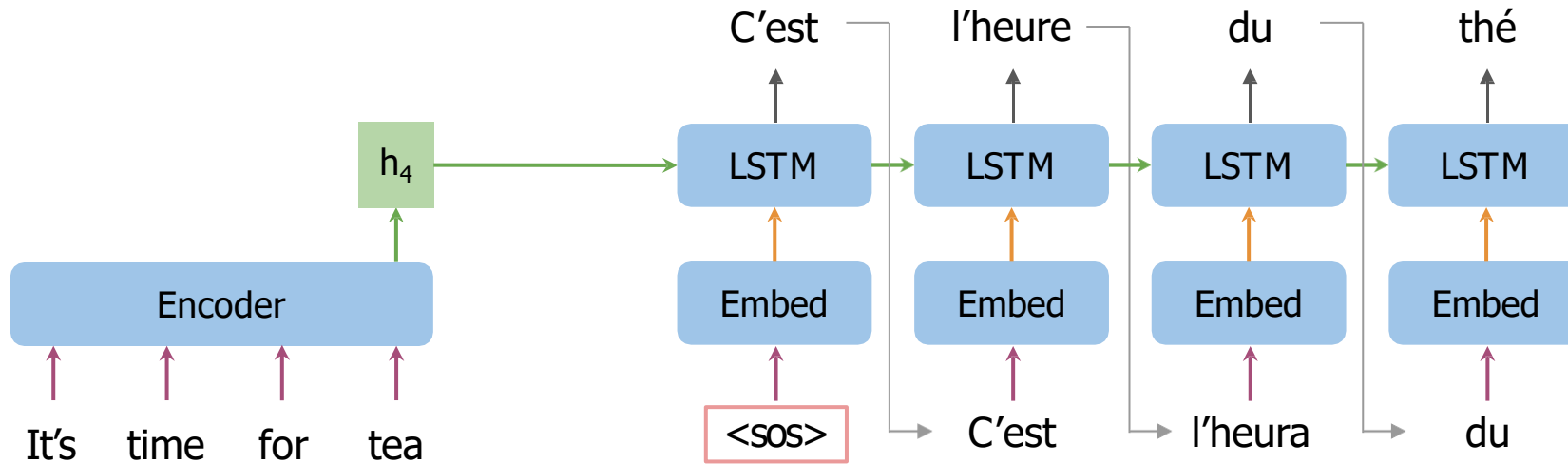- LSTMs and GRUs to avoid vanishing    and exploding gradient problems

# Seq2Seq encoder



- typically consists of an embedding layer and an LSTM module with one or more layers
- transforms words tokenized first into a vector for input to the LSTM module
- the LSTM module receives inputs from the embedding layer, as well as the hidden states from the previous step
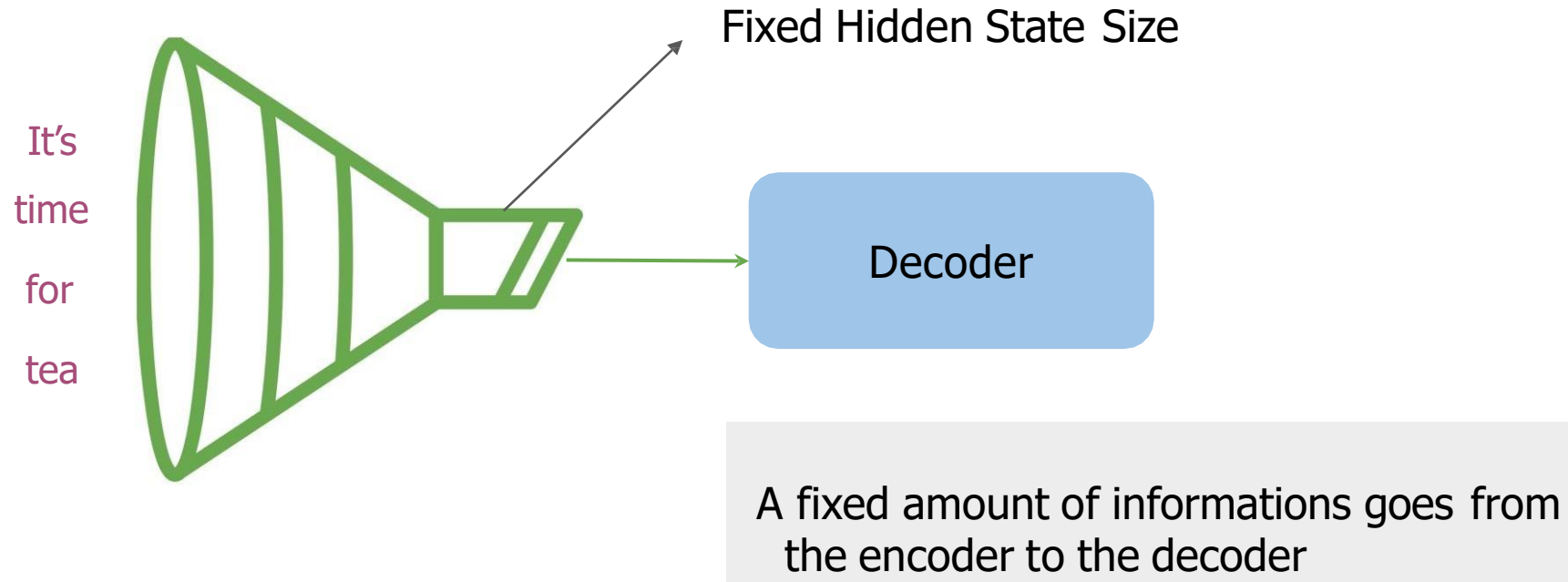
# Seq2Seq decoder



- The decoder is constructed similarly with an embedding layer and an LSTM layer.

- the output word of a step as the input word for the next step.

- pass the LSTM hidden state to the next step.

# The information bottleneck

Fixed Hidden State Size

It's
time
for
tea

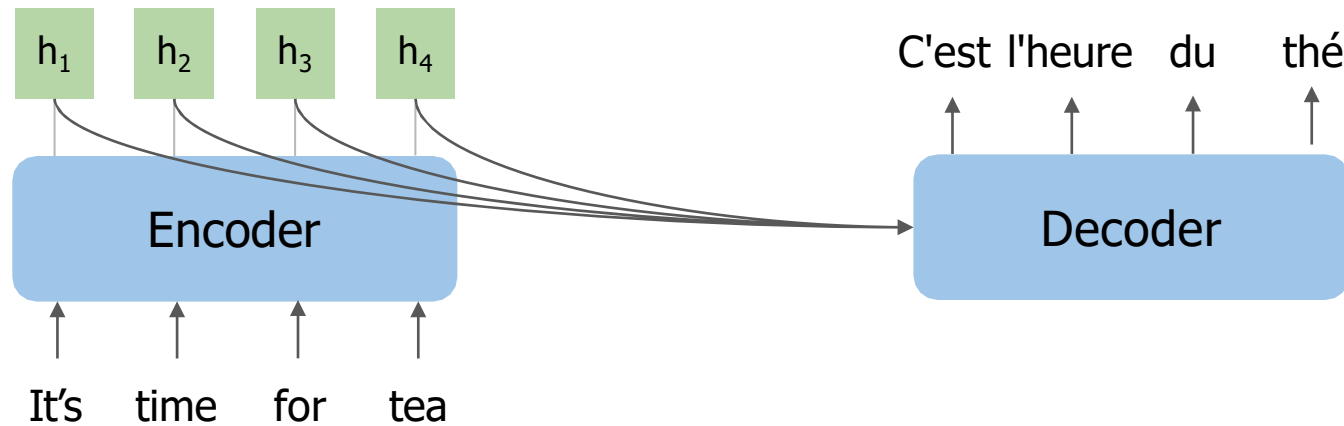Decoder

A fixed amount of informations goes from the encoder to the decoder

- Variable-length sentences + fixed-length memory

=>As sequence size increases, model performance decreases

# Use all the encoder hidden states?



- Solution: focus attention in the right place



The model can focus on specific hidden states at every step

# Seq2Seq model with attention

**Dzmitry Bahdanau**
Jacobs University Bremen, Germany

**KyungHyun Cho**     **Yoshua Bengio***
Université de Montréal

- Performance

Greater BLEU is better



Seq2Seq with Attention

Traditional Seq2Seq Models

# Traditional seq2seq models

# How to use all the hidden states?

# How to use all the hidden states?

# The attention layer in more depth

$$e_{ij} = a(s_{i-1}, h_j)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

$h_j$

$s_{i-1}$

| Feedforward Network | $\rightarrow$ $e_{ij}$ $\rightarrow$ | Softmax | $\rightarrow$ | $\alpha_{ij}$ |

Weights used for the sum of hidden states

Learnable parameters

$$c_i = \sum_{j=1}^{T_x} \boxed{\alpha_{ij}} \boxed{h_j}$$

Context Vector is an expected value

$\alpha_{i1}h_1 \;+\; \alpha_{i2}h_2 \;+\; \alpha_{i3}h_3 \;+\cdots+\; \alpha_{iM}h_M \longrightarrow c_i$

# Queries, Keys, Values and Attention

- Queries, Keys, Values

| Query | | Key | Value |
|-------|---|-----|-------|
| l'heure | → | It's | [0.5, 0.2, -1.2, …, ] |
| | → | time | [0.2, -0.7, 0.9, …, ] |
| | → | for | [1.3, 0.3, 0.8, …, ] |
| Similarity is used in for weighted sum | → | tea | [-0.4, 0.6, -1.1, …, ] |

# Scaled dot-product attention

**Similarity Between Q and K**

$$\mathrm{softmax}\left(\frac{QK^{\top}}{\sqrt{d_k}}\right)V$$

**Weights for the weighted sum**

**Scale using the root of the key vector size**

**Weighted sum of values $V$**

Just two matrix multiplications and a Softmax!

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q          K          V

Queries  Keys      Values

(Vaswani et al., 2017)

# Alignment Weights

Keys

it's    time    for    tea

c'est

l'heure

Queries

du

thé

Similar words have large weights



- Flexible attention
- Works for languages with different grammar structures!

# Summary

- Attention is a layer that lets a model focus on what's important

- Queries, Values, and Keys are used for information retrieval inside the  Attention layer

- Works for languages with very different grammatical structures

# Setup for machine translation

- Data in machine translation

| English | French |
|---------|--------|
| I am hungry! | J'ai faim! |
| ... | ... |
| I watched the soccer game. | J'ai regardé le match de football. |

**Attention!** (pun intended) Assignment dataset is not as squeaky-clean as this example and contains some Spanish translations.

# Machine translation setup

- Use pre-trained vector embeddings

- Otherwise, initially represent words with a one-hot vectors

- Keep track of index mappings with word2ind and ind2word dictionaries

- Add end of sequence tokens:        <EOS>

- Pad the token vectors with zeros

# Preparing to Translate to English

ENGLISH SENTENCE:

Both the ballpoint and the mechanical pencil in the series are equipped with a special mechanism: when the twist mechanism is activated, the lead is pushed forward.

TOKENIZED VERSION OF THE ENGLISH SENTENCE:

[4546    4    11358  362    8    4    23326    20104    1745    8210    9641    5    6
4  3103    31  2767    30    13  914  4797    64    196    4    22474    5  4797    16
24864    86    2    4    1060  16          6413  1138    3    1    0    0    0    0    0    0    0    0
0    0    0    0    0    0    0    0    0    0    0    0    0    0    Padding    0    0    0]

Padding

<EOS>

FRENCH TRANSLATION:

Le stylo à bille et le porte-mine de la série sont équipés d'un mécanisme spécial: lorsque le mécanisme de torsion est activé, le plomb est poussé vers l'avant.

TOKENIZED VERSION OF THE FRENCH TRANSLATION:

[7    29587    9    18240    8    7    420    5    3440    2    6    156    39    7941    14    19
5548    2648    562    7    5548    2    23194    18    20114    1    7    5695    18    8865    149
12  137    1    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0    0
0    0    0]

Padding

# Teacher Forcing

- o Training for NMT

- o Teacher forcing

- Traditional seq2seq models

Outputs

| C'est | l'heure | du | thé |
|-------|---------|-----|-----|

| Encoder | Decoder | Decoder | Decoder | Decoder |

It's  time  for  tea

<sos>  C'est  l'heure  du

Inputs

# Training seq2seq models

Encoder

It's time for tea



≠ **Fluffy**

**Targets**  C'est  l'heura  du  thé
Yes!  No...  Even worse  lol, no

**Output** ✓C'est  ✗un  ✗chat  ✗duveteux

Decoder  Decoder  Decoder  Decoder

<sos>  C'est  ✗un  ✗chat

Errors from early steps propagate

# Teacher Forcing

Improves training performance

**Targets**

C'est    l'heura    du    café

**Output**

Yes!    No...    Yes!    No, but not bad

✓C'est    ✗un    ✓ du    ✗duveteux

It's    time    for    tea

Encoder → Decoder → Decoder → Decoder → Decoder

<sos>    C'est    l'heura    du

**Correct sequence of words as input (shifted right)**

# Neural Machine Translation Model with Attention

o How everything fits together

o NMT model in detail

- NMT Model



The decoder has to pass the hidden state to the Attention Mechanism

Difficult to implement, so a **pre-attention decoder** is introduced.

# Neural Machine Translation Model

# Neural Machine Translation Model

# BLEU Score

- **Bi**Lingual **E**valuation **U**nderstudy

- Compares candidate translations to reference (human) translations The closer to **1**, the better

0                    1

| Candidate | I | I | am | I | |
|---|---|---|---|---|---|
| Reference 1 | Younes | said | I | am | hungry |
| Reference 2 | He | said | I | am | hungry |

How many words from the **candidate** appear in the **reference** translations?

25

# BLEU Score

| Candidate | I | I | am | I |
|---|---|---|---|---|
| **Reference 1** | Younes | said | I | am | hungry |
| **Reference 2** | He | said | I | am | hungry |

Count: $\dfrac{1+1+1+1}{4} = 1$

A model that always outputs common words will do great!

- BLEU Score (Modified)

| Candidate | I | I | am | I |
|---|---|---|---|---|
| **Reference 1** | Younes | said | | | hungry |
| **Reference 2** | He | said | | | |

Count: $\dfrac{1+1}{4} = 0.5$

Better than the previous implementation version!

# ROUGE

o **R**ecall-**O**riented **U**nderstudy for **G**isting **E**valuation

o Compares candidates with reference (human) translations

o Multiple versions for this metric

- ROUGE-N

| Candidate | I | I | am | I | |
|---|---|---|---|---|---|
| **Reference 1** | Younes | said | I | am | hungry |
| **Reference 2** | He | said | I | am | hungry |

How many words from the **reference** appear in the **candidate** translations?

# ROUGE-N

| | | | | | |
|---|---|---|---|---|---|
| **Candidate** | I | I | am | I | |
| **Reference 1** | Younes | said | I | am | hungry |
| **Reference 2** | He | said | I | am | hungry |

$$\text{Count 1: } \frac{1+1}{5} = 0.4 \qquad \text{Count 2: } \frac{1+1}{5} = 0.4$$

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \longrightarrow F1 = 2 \times \frac{\text{BLEU} \times \text{ROUGE-N}}{\text{BLEU} + \text{ROUGE-N}}$$

$$F1 = 2 \times \frac{0.5 \times 0.4}{0.5 + 0.4} = \frac{4}{9} \approx 0.44$$

# Sampling and Decoding

- Random sampling

- Temperature in sampling

- Greedy decoding

- **Seq2Seq model**

| Words | de | la | le | et | à | ... |
|-------|------|------|-----|-------|------|-----|
| $P(w_i)$ | 0.02 | 0.04 | 0.1 | 0.005 | 0.08 | ... |

```
Softmax
  ↑
Dense
  ↑
Decoder
  ↑
<SOS>
```

Probability distribution over words in target language

# Greedy decoding

o Selects the most probable word at each step

o But the best word at each step may not be the best for longer sequences...

o Can be fine for shorter sequences, but limited by inability to look further down the sequence

J'ai faim.

I am __hungry__.

I am, am, am, am...

• Random sampling

| am | full | hungry | I | the |
|------|------|--------|------|------|
| 0.05 | 0.3 | 0.15 | 0.25 | 0.25 |

Often a little too random for accurate translation!

Solution: Assign more weight to more probable words, and less weight to less probable words.

# Temperature

- Can control for more or less randomness in predictions

- Lower temperature setting : More confident, conservative network

- Higher temperature setting : More excited, random network

# Beam search

Most probable translation **is not** the one with the most probable word at each step

| Solution

Calculate probability of multiple possible sequences

Beam search

- Beam search decoding
  - Probability of multiple possible sequences at each step
  - Beam width B determines number of sequences you keep
  - Until all B most probable sequences end with <EOS>

Beam search with **B=1** is **greedy decoding**.

# Beam search example

B = 2

**P(w₁ | "<sos>")**

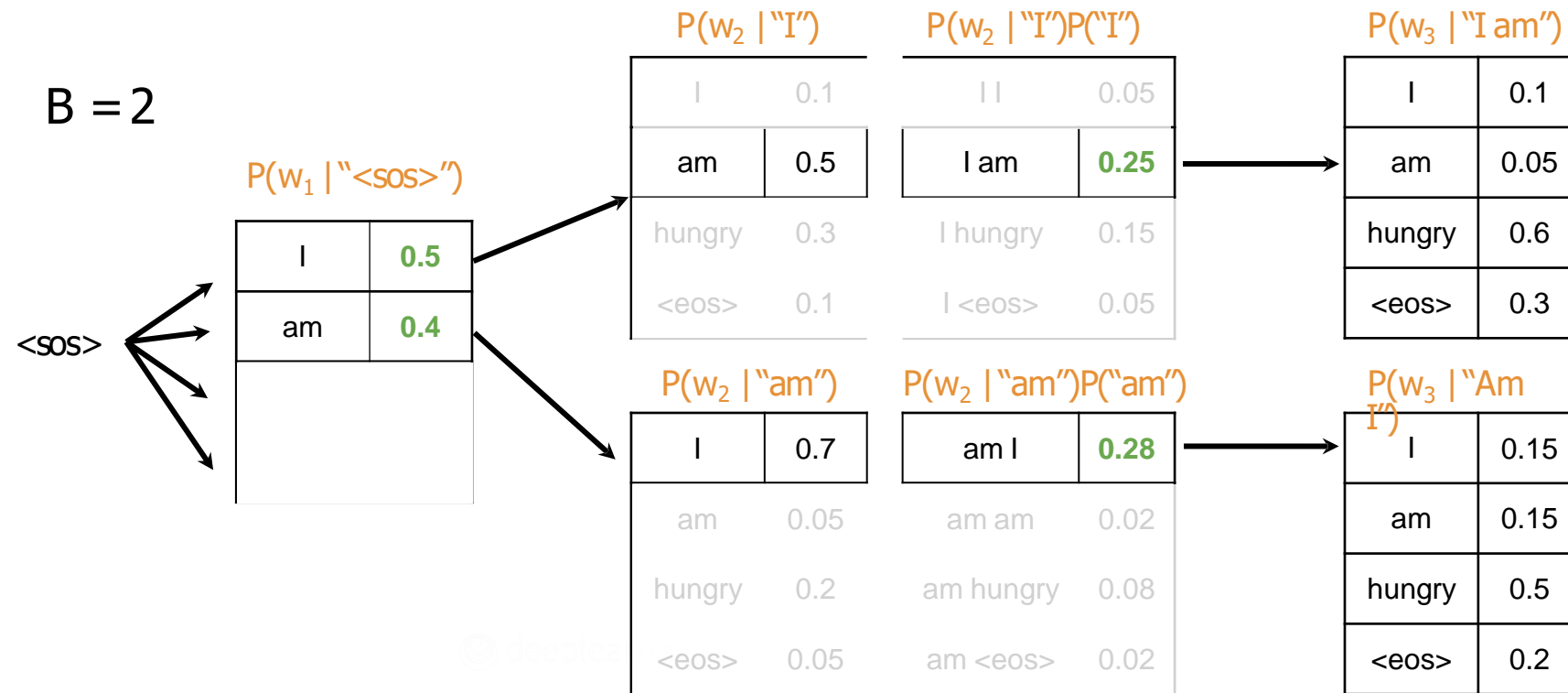| I | **0.5** |
|---|---|
| am | **0.4** |
| | |

<sos>

**P(w₂ | "I")**

| I | 0.1 |
|---|---|
| am | 0.5 |
| hungry | 0.3 |
| <eos> | 0.1 |

**P(w₂ | "I")P("I")**

| I I | 0.05 |
|---|---|
| I am | **0.25** |
| I hungry | 0.15 |
| I <eos> | 0.05 |

**P(w₃ | "I am")**

| I | 0.1 |
|---|---|
| am | 0.05 |
| hungry | 0.6 |
| <eos> | 0.3 |

**P(w₂ | "am")**

| I | 0.7 |
|---|---|
| am | 0.05 |
| hungry | 0.2 |
| <eos> | 0.05 |

**P(w₂ | "am")P("am")**

| am I | **0.28** |
|---|---|
| am am | 0.02 |
| am hungry | 0.08 |
| am <eos> | 0.02 |

**P(w₃ | "Am I")**

| I | 0.15 |
|---|---|
| am | 0.15 |
| hungry | 0.5 |
| <eos> | 0.2 |

# Beam search decoding

I

$P(w_2|\text{``I''})$

```
          → c → Decoder → Decoder
                  ↑          ↑
               <SOS>         I
```

$P(w_1|\text{``<sos>''})$

Select B most probable words →

am

$P(w_2|\text{``am''})$

<SOS>

```
          → c → Decoder → Decoder
                  ↑          ↑
               <SOS>        am
```

B model runs

# Problems with beam search

Penalizes long sequences, so you should normalize by the sentence length

Computationally expensive and consumes a lot of memory

- Minimum Bayes Risk (MBR)
  - Generate several candidate translations
  - Assign a similarity to every pair using a similarity score (such as ROUGE!)
  - Select the sample with the highest average similarity

# Minimum Bayes Risk (MBR)

$$\arg\max_{E} \frac{1}{n} \sum_{E'} \text{ROUGE}(E, E')$$

**Find the candidate translation that maximizes**

**Compare with every other candidate**

**ROUGE score between pair of candidates**

- Example: MBR Sampling

$\text{ROUGE}(C_1, C_2)$

$\text{ROUGE}(C_1, C_3)$

$\text{ROUGE}(C_1, C_4)$

Compute average ROUGE

$$R_1 = \frac{1}{3} \sum_{i \neq 1} \text{ROUGE}(C_1, C_i)$$

Repeat for every candidate

Select the candidate with the highest average