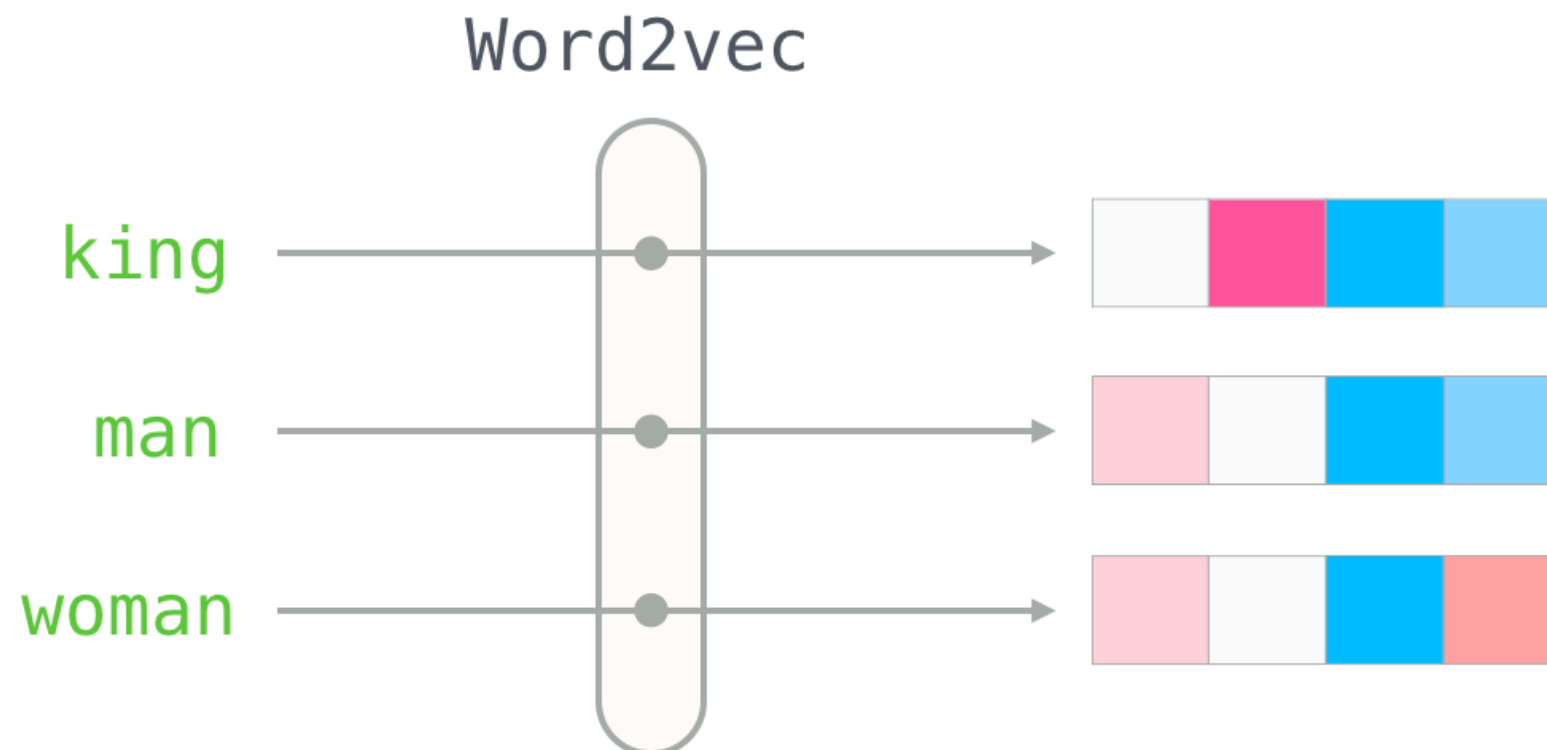# Language Models

Boris Zubarev

@bobazooba

# Word Embeddings
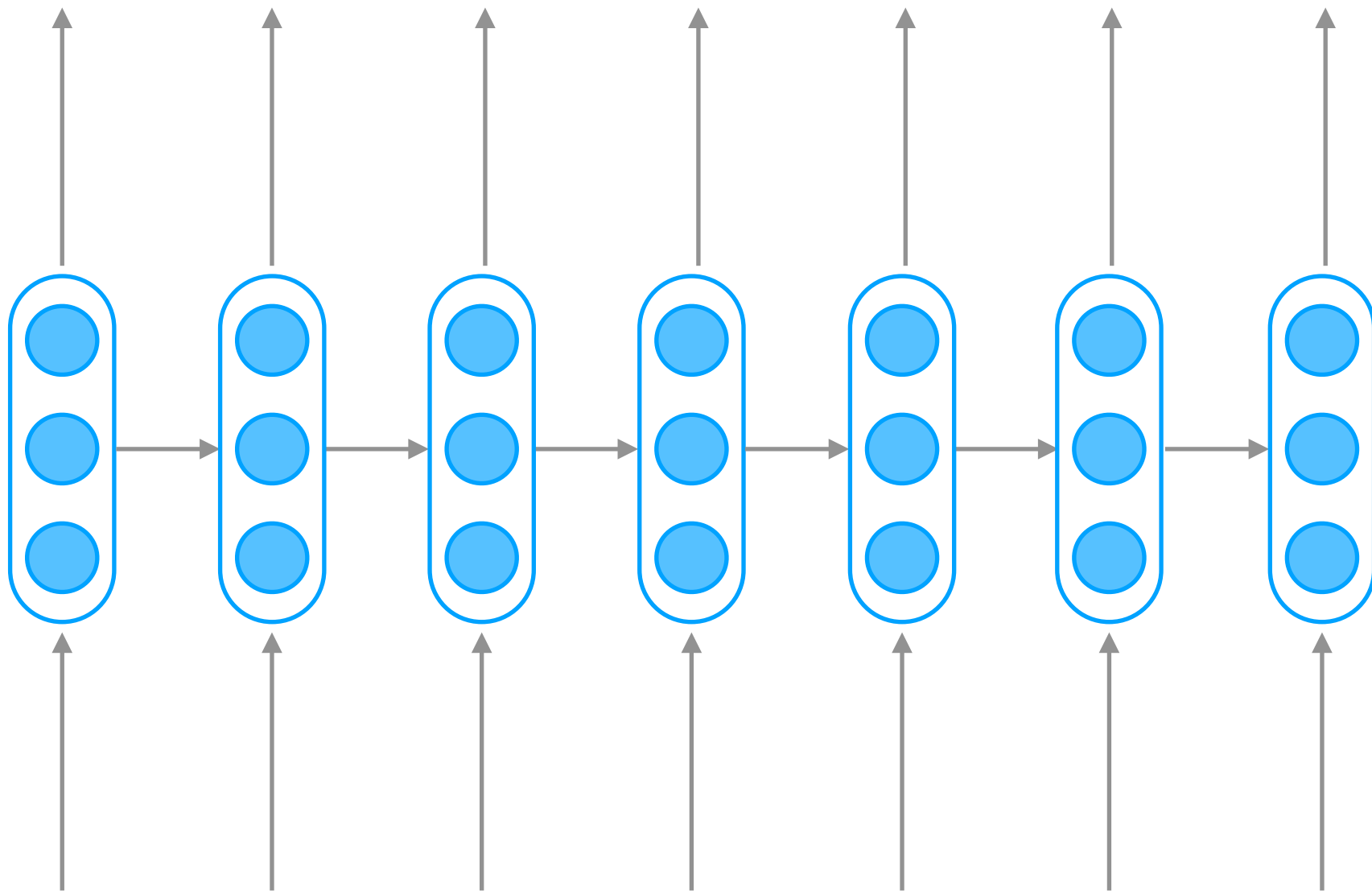
**w2v, glove, etc**

- Just key—value storage at inference

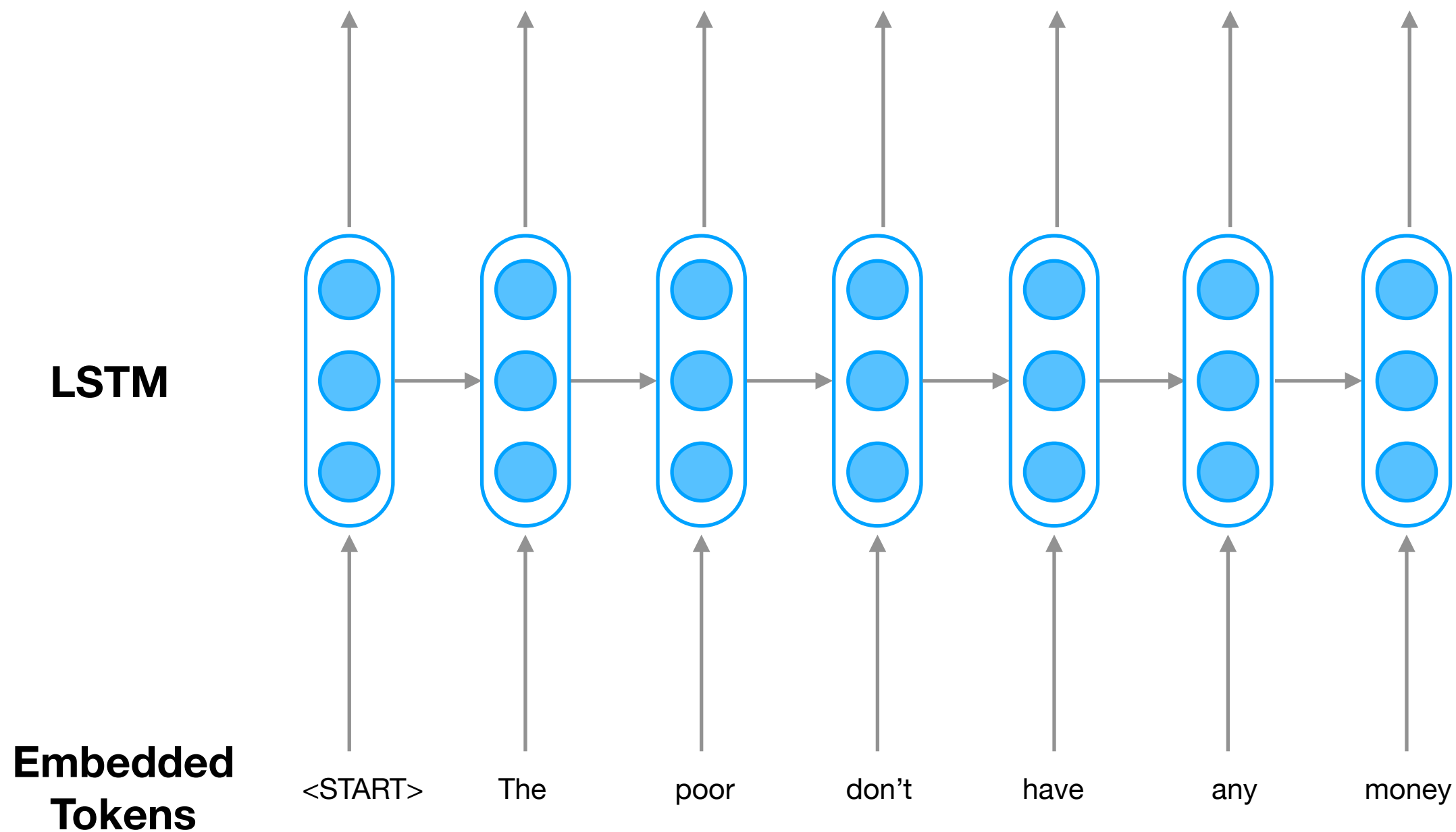- Don't change from relationships with other words in the current text

Word2vec

king

man

woman

# Language Model

**LSTM**

# Language Model

**LSTM**

**Embedded Tokens**

<START> The poor don't have any money

# Language Model



**Softmax**

The     poor     don't     have     any     money     <END>

**LSTM**

**Embedded Tokens**

<START>     The     poor     don't     have     any     money

# Language Model

## Inference

**LSTM**

**Embedded Tokens**    <START>

# Language Model

**Inference**

**Softmax**

**LSTM**

**Embedded Tokens**

<START>

# Language Model

## Inference

**Softmax**     The

**LSTM**

**Embedded Tokens**    <START>

# Language Model

**Inference**

**Softmax**                    The

**LSTM**

**Embedded
Tokens**       <START>        The

# Language Model

**Inference**

**Softmax**

The

**LSTM**

**Embedded Tokens**

<START>       The

# Language Model

**Inference**

**Softmax**

The      poor

**LSTM**

**Embedded Tokens**

<START>      The

# Language Model

**Inference**

**Softmax**

The       poor

**LSTM**

**Embedded Tokens**

&lt;START&gt;      The      poor

# Language Model

**Inference**

**Softmax**

The         poor

**LSTM**

**Embedded Tokens**

&lt;START&gt;        The        poor

# Language Model

**Inference**

**Softmax**    The        poor        don't

**LSTM**

**Embedded
Tokens**    <START>    The        poor

# Language Model

**Inference**

# Language Model

**Inference**

**Softmax**

The        poor        don't

**LSTM**

**Embedded Tokens**

&lt;START&gt;      The      poor      don't

# Language Model

## Inference

**Softmax**      The        poor        don't        have

**LSTM**

**Embedded Tokens**    <START>        The        poor        don't

# Language Model

**Inference**

**Softmax**

The      poor      don't      have

**LSTM**

**Embedded Tokens**

&lt;START&gt;      The      poor      don't      have

# Language Model

**Inference**

# Language Model

**Inference**

**Softmax**

The          poor          don't          have          any

**LSTM**



**Embedded Tokens**

<START>          The          poor          don't          have          any

# Language Model

**Inference**

# Language Model

**Inference**

**Softmax**

The      poor      don't      have      any      money

**LSTM**

**Embedded Tokens**

&lt;START&gt;      The      poor      don't      have      any      money

# Language Model

## Inference



**Softmax**    The    poor    don't    have    any    money

**LSTM**

**Embedded Tokens**    <START>    The    poor    don't    have    any    money

# Language Model

## Inference



Softmax: The   poor   don't   have   any   money   <END>

LSTM

Embedded Tokens: <START>   The   poor   don't   have   any   money

# Language Model

## Training

# Language Model

# Language Model

## Training

**Tying**

**Softmax**

| The | poor | don't | have | any | money | <END> |

**LSTM**

**Embedded Tokens**

| <START> | The | poor | don't | have | any | money |

# Teacher Forcing

**Softmax**   The    poor    don't    have    any    money    <END>



**LSTM**

**Embedded Tokens**   <START>    The    poor    don't    have    any    money

# Teacher Forcing

**Softmax**  The     poor     don't     have     any     money     <END>

**LSTM**

**Embedded Tokens**  <START>     The     poor     don't     have     any     money

# Language Model

<START> The poor don't have any money <END>

# Language Model

**X**  &lt;START&gt;

**Y**    The

# Language Model

**X**   &lt;START&gt;   The

**Y**   poor

# Language Model

**X**    &lt;START&gt;    The    poor

**Y**    don't

# Language Model

**X**     &lt;START&gt;     The          poor          don't

**Y**                                                    have

# Language Model

**X**    &lt;START&gt;    The    poor    don't    have

**Y**                                                            any

# Language Model

**X**        &lt;START&gt;      The        poor        don't        have        any

**Y**                                                                money

# Language Model

**X**      <START>      The      poor      don't      have      any      money

**Y**                                                                <END>

# Language Model

**X**             &lt;START&gt;      The        poor        don't        have        any        money

**Y**                                                   &lt;END&gt;

**Sum of Losses**

# Transfer Learning

## CNN Intuition

# NLP LM Transfer Learning

**LSTM**

# NLP LM Transfer Learning



**Frozen LSTM Language Model**

# NLP LM Transfer Learning



Frozen LSTM Language Model

Tokenized embedded text

# NLP LM Transfer Learning

**Contextualized word embeddings**

**Frozen LSTM Language Model**

**Tokenized embedded text**

# NLP LM Transfer Learning

Your typical classifier

**Contextualized word embeddings**

**Frozen LSTM Language Model**

**Tokenized embedded text**

# NLP LM Transfer Learning

**Featuring Model**

Your typical classifier

Contextualized word embeddings

**Frozen LSTM Language Model**

Tokenized embedded text

# ELMo

# ELMo

# ELMo

# ELMo

# ELMo

**Source**  &lt;START&gt;  The  poor  don't  have  any  money  &lt;END&gt;

# ELMo

**Source**  &lt;START&gt;   The   poor   don't   have   any   money   &lt;END&gt;

**Forward**

**Backward**

# ELMo

| **Source** | <START> | The | poor | don't | have | any | money | <END> |
|------------|---------|-----|------|-------|------|-----|-------|-------|
| **Forward** | <START> | | | | | | | |
| **Backward** | <END> | | | | | | | |

# ELMo

**Source**   &lt;START&gt;   The   poor   don't   have   any   money   &lt;END&gt;

**Forward**   &lt;START&gt;   The

**Backward**   &lt;END&gt;   money

# ELMo

**Source**   &lt;START&gt;     The       poor      don't      have       any       money      &lt;END&gt;

**Forward**   &lt;START&gt;     The       poor

**Backward**   &lt;END&gt;     money      any

# ELMo

| **Source** | <START> | The | poor | don't | have | any | money | <END> |
|---|---|---|---|---|---|---|---|---|
| **Forward** | <START> | The | poor | don't | | | | |
| **Backward** | <END> | money | any | have | | | | |

# ELMo

**Source**   <START>   The   poor   don't   have   any   money   <END>

**Forward**   <START>   The   poor   don't   have

**Backward**   <END>   money   any   have   don't

# ELMo

| **Source** | <START> | The | poor | don't | have | any | money | <END> |
|---|---|---|---|---|---|---|---|---|
| **Forward** | <START> | The | poor | don't | have | any | | |
| **Backward** | <END> | money | any | have | don't | poor | | |

# ELMo

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Source** | <START> | The | poor | don't | have | any | money | <END> |
| **Forward** | <START> | The | poor | don't | have | any | money | |
| **Backward** | <END> | money | any | have | don't | poor | The | |

# ELMo

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Source** | <START> | The | poor | don't | have | any | money | <END> |
| **Forward** | <START> | The | poor | don't | have | any | money | <END> |
| **Backward** | <END> | money | any | have | don't | poor | The | <START> |

# ELMo

# ELMo

# ELMo

# ELMo

## Results



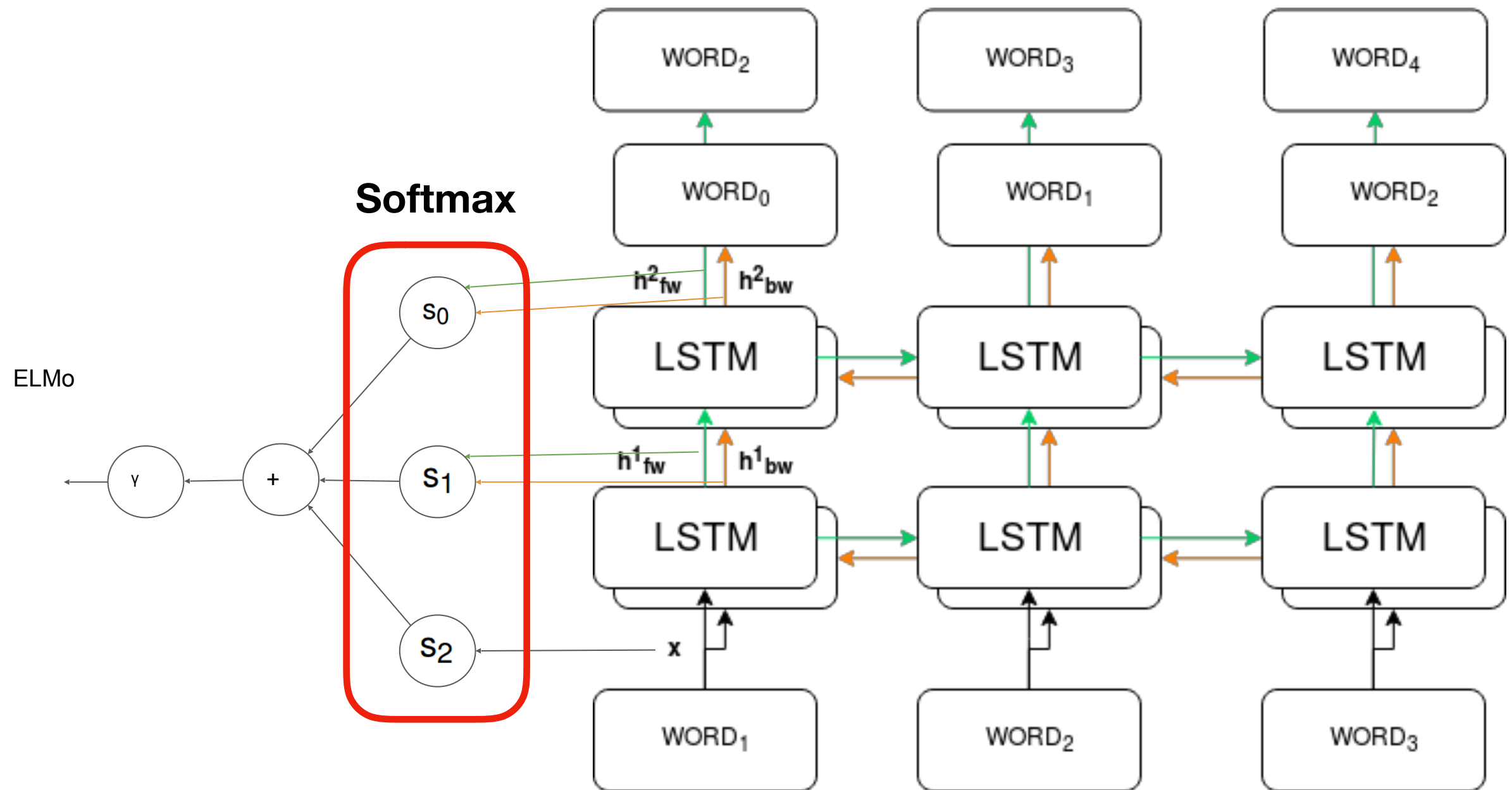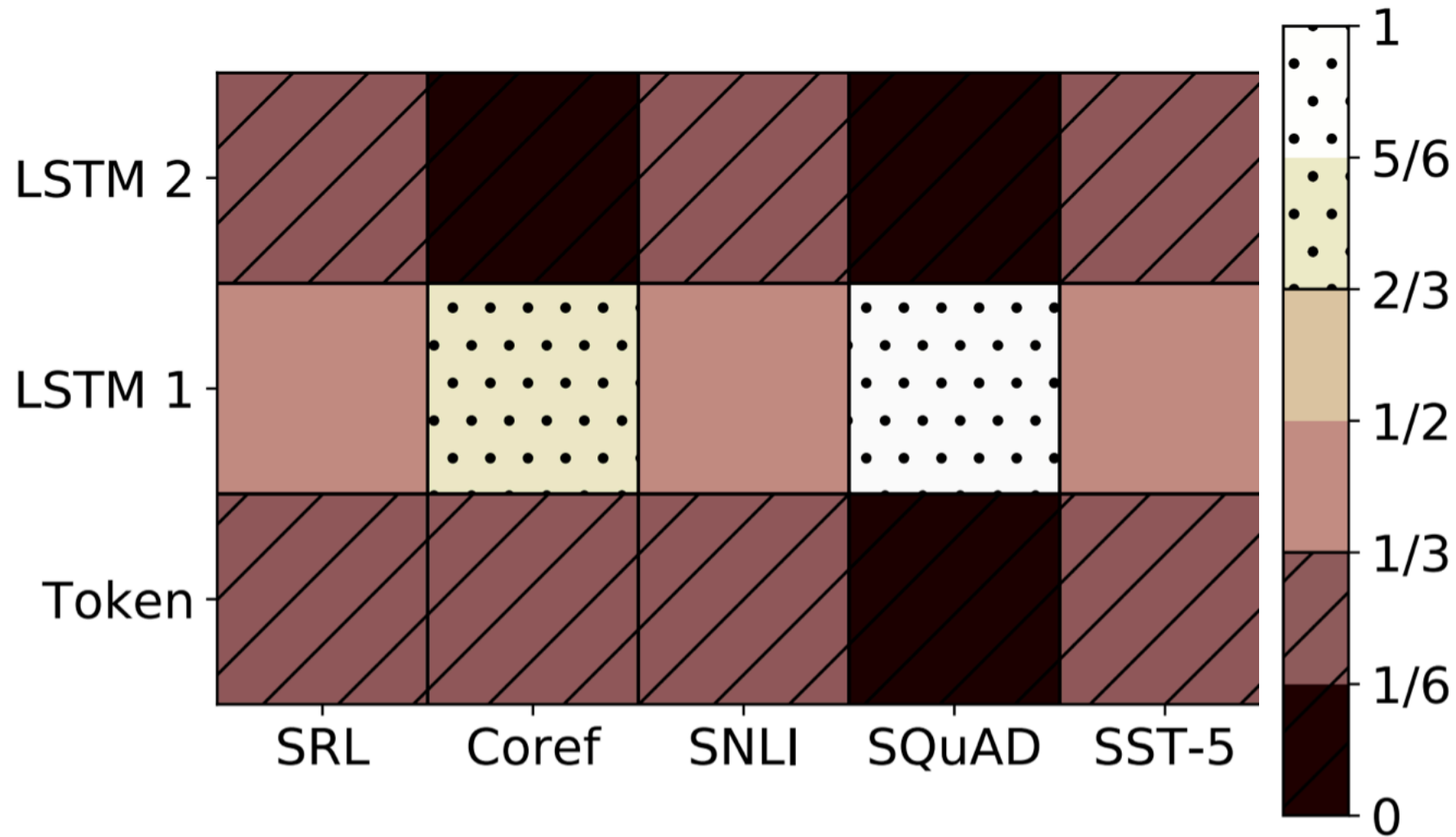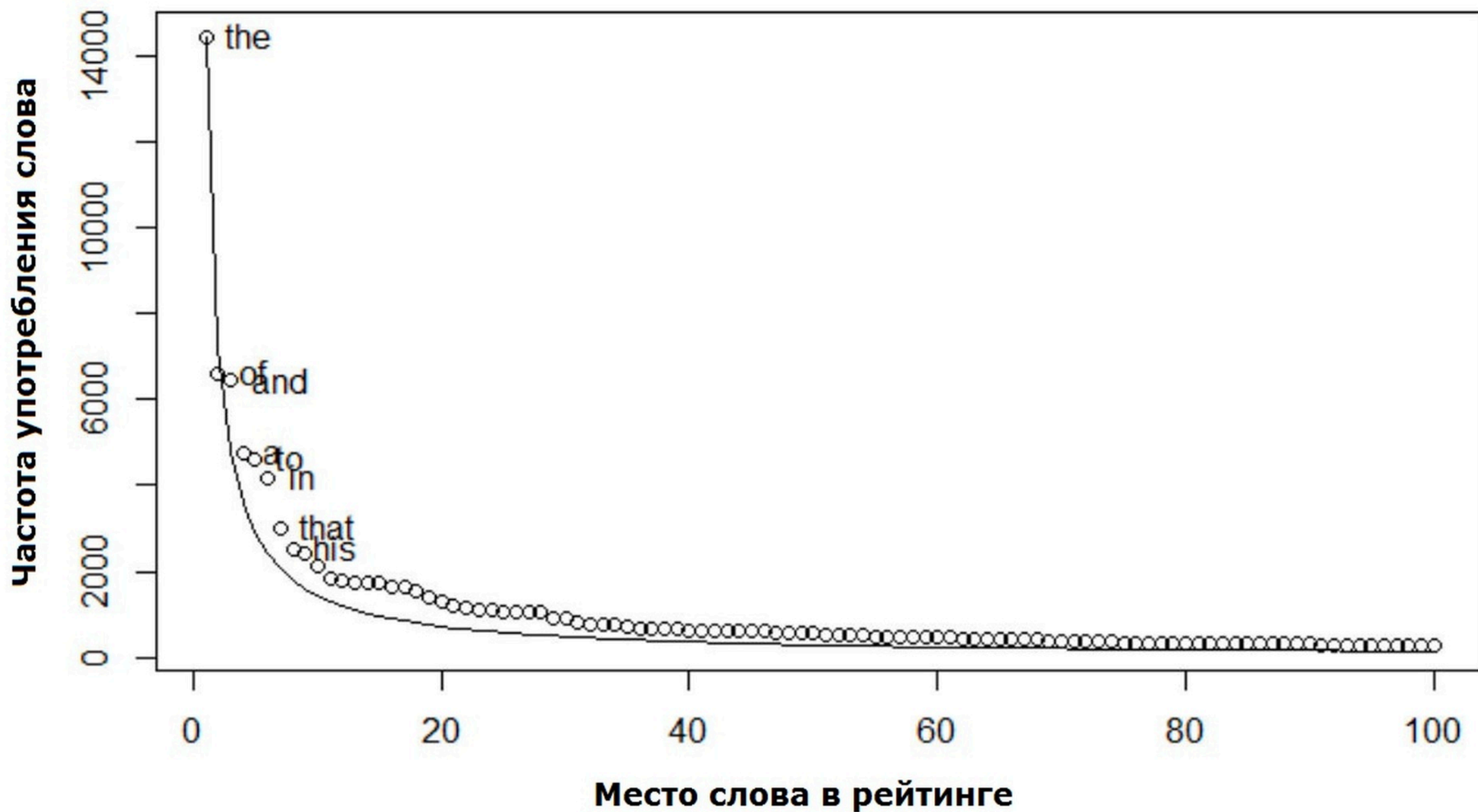Figure 2: Visualization of softmax normalized biLM layer weights across tasks and ELMo locations. Normalized weights less then 1/3 are hatched with horizontal lines and those greater then 2/3 are speckled.

# Tokenization

## Закон Ципфа

# Tokenization

**Word level**

+   Small text length
- Big vocabulary size
-       OOV

**Character level**

-      Long text
+ Small vocabulary size
+   Almost no OOV

# Tokenization

**Word level**

+ Small text length
- Big vocabulary size
- OOV

**Character level**

- Long text
+ Small vocabulary size
+ Almost no OOV

1. Word level

i'm a second year student in an ivy league school ->

["i'm", 'a', 'second', 'year', 'student', 'in', 'an', 'ivy', 'league', 'school']

2. Character level

['i', "'", 'm', ' ', 'a', ' ', 's', 'e', 'c', 'o', 'n', 'd', ' ', 'y', 'e', 'a', 'r', ' ', 's', 't', 'u', 'd', 'e', 'n', 't', ' ', 'i', 'n', ' ', 'a', 'n', ' ', 'i', 'v', 'y', ' ', 'l', 'e', 'a', 'g', 'u', 'e', ' ', 's', 'c', 'h', 'o', 'o', 'l']

# BPE

I saw a girl with a telescope. ->

['__I', '__saw', '__a', '__girl', '__with', '__a', '__', 'te', 'le', 's', 'c', 'o', 'pe', '.']

опубликовано видео убитого саудовского журналиста джамаля хашкуджи ->

['__опубликовано', '__видео', '__убитого', '__саудов', 'ского', '__журналиста',
'__джама', 'ля', '__ха', 'шку', 'джи']

# BPE

---

**Algorithm 1** Learn BPE operations

```python
import re, collections

def get_stats(vocab):
  pairs = collections.defaultdict(int)
  for word, freq in vocab.items():
    symbols = word.split()
    for i in range(len(symbols)-1):
      pairs[symbols[i],symbols[i+1]] += freq
  return pairs

def merge_vocab(pair, v_in):
  v_out = {}
  bigram = re.escape(' '.join(pair))
  p = re.compile(r'(?<!\S)' + bigram + r'(?!\S)')
  for word in v_in:
    w_out = p.sub(''.join(pair), word)
    v_out[w_out] = v_in[word]
  return v_out

vocab = {'l o w </w>' : 5, 'l o w e r </w>' : 2,
         'n e w e s t </w>':6, 'w i d e s t </w>':3}
num_merges = 10
for i in range(num_merges):
  pairs = get_stats(vocab)
  best = max(pairs, key=pairs.get)
  vocab = merge_vocab(best, vocab)
  print(best)
```

---

| | | |
|---|---|---|
| r · | → | r· |
| l o | → | lo |
| lo w | → | low |
| e r· | → | er· |

Figure 1: BPE merge operations learned from dictionary {'low', 'lowest', 'newer', 'wider'}.

# BPE

```
import re, collections

def get_stats(vocab):
  pairs = collections.defaultdict(int)
  for word, freq in vocab.items():
    symbols = word.split()
    for i in range(len(symbols)-1):
      pairs[symbols[i],symbols[i+1]] += freq
  return pairs

def merge_vocab(pair, v_in):
  v_out = {}
  bigram = re.escape(' '.join(pair))
  p = re.compile(r'(?<!\S)' + bigram + r'(?!\S)')
  for word in v_in:
    w_out = p.sub(''.join(pair), word)
    v_out[w_out] = v_in[word]
  return v_out

vocab = {'l o w </w>' : 5, 'l o w e r </w>' : 2,
         'n e w e s t </w>':6, 'w i d e s t </w>':3}
num_merges = 10
for i in range(num_merges):
  pairs = get_stats(vocab)
  best = max(pairs, key=pairs.get)
  vocab = merge_vocab(best, vocab)
  print(best)
```

|       |               |       |
|-------|---------------|-------|
| r ·   | $\rightarrow$ | r·    |
| l o   | $\rightarrow$ | lo    |
| lo w  | $\rightarrow$ | low   |
| e r·  | $\rightarrow$ | er·   |

Figure 1: BPE merge operations learned from dictionary {'low', 'lowest', 'newer', 'wider'}.

- learning
  - word:freq : {low:5, lowest:2, newer:6, wider:3}
  - marge & count
    1. 'r' '</w>' : 9  → marge'r</w>'
    2. 'e' 'r</w>' : 9 →marge'er</w>'
    3. 'l' 'o' : 7        →marge'lo'
    4. 'lo' 'w' : 7      →marge'low'

  → OOV : 'lower' segmented 'low er</w>'

# BPE

**Algorithm 1** Learn BPE operations

```python
import re, collections

def get_stats(vocab):
  pairs = collections.defaultdict(int)
  for word, freq in vocab.items():
    symbols = word.split()
    for i in range(len(symbols)-1):
      pairs[symbols[i],symbols[i+1]] += freq
  return pairs

def merge_vocab(pair, v_in):
  v_out = {}
  bigram = re.escape(' '.join(pair))
  p = re.compile(r'(?<!\S)' + bigram + r'(?!\S)')
  for word in v_in:
    w_out = p.sub(''.join(pair), word)
    v_out[w_out] = v_in[word]
  return v_out

vocab = {'l o w </w>' : 5, 'l o w e r </w>' : 2,
         'n e w e s t </w>':6, 'w i d e s t </w>':3}
num_merges = 10
for i in range(num_merges):
  pairs = get_stats(vocab)
  best = max(pairs, key=pairs.get)
  vocab = merge_vocab(best, vocab)
  print(best)
```

|       |               |      |
|-------|---------------|------|
| r ·   | $\rightarrow$ | r·   |
| l o   | $\rightarrow$ | lo   |
| lo w  | $\rightarrow$ | low  |
| e r·  | $\rightarrow$ | er·  |

Figure 1: BPE merge operations learned from dictionary {'low', 'lowest', 'newer', 'wider'}.

- learning
  - word:freq : {low:5, lowest:2, newer:6, wider:3}
  - marge & count
    1. 'r' '</w>' : 9   → marge'r</w>'
    2. 'e' 'r</w>' : 9  → marge'er</w>'
    3. 'l' 'o' : 7       → marge'lo'
    4. 'lo' 'w' : 7      → marge'low'

  → OOV : 'lower' segmented 'low er</w>'

**Vocabulary sizes:**
5000, 10000, 15000, …, 50000

# BPE

## SentencePiece

SentencePiece is an unsupervised text tokenizer and detokenizer mainly for Neural Network-based text generation systems where the vocabulary size is predetermined prior to the neural model training. SentencePiece implements **subword units** (e.g., **byte-pair-encoding (BPE)** [Sennrich et al.]) and **unigram language model** [Kudo.]) with the extension of direct training from raw sentences. SentencePiece allows us to make a purely end-to-end system that does not depend on language-specific pre/postprocessing.

**This is not an official Google product.**

## YouTokenToMe

YouTokenToMe is an unsupervised text tokenizer focused on computational efficiency. It currently implements fast Byte Pair Encoding (BPE) [Sennrich et al.]. Our implementation is much faster in training and tokenization than both fastBPE and SentencePiece. In some test cases, it is 90 times faster. Check out our benchmark results.

Key advantages:

- Multithreading for training and tokenization
- The algorithm has `O(N)` complexity, where `N` is the length of training data
- Highly efficient implementation in C++
- Python wrapper and command-line interface

# Sequence to sequence

Source sentence

# Sequence to sequence
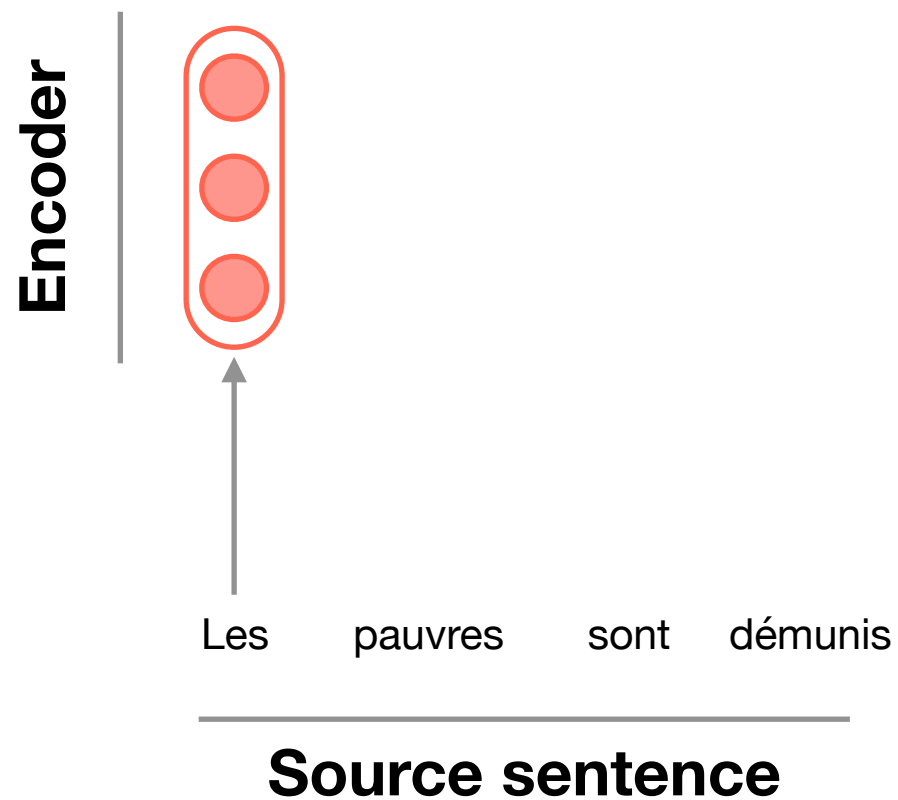
Les     pauvres    sont   démunis
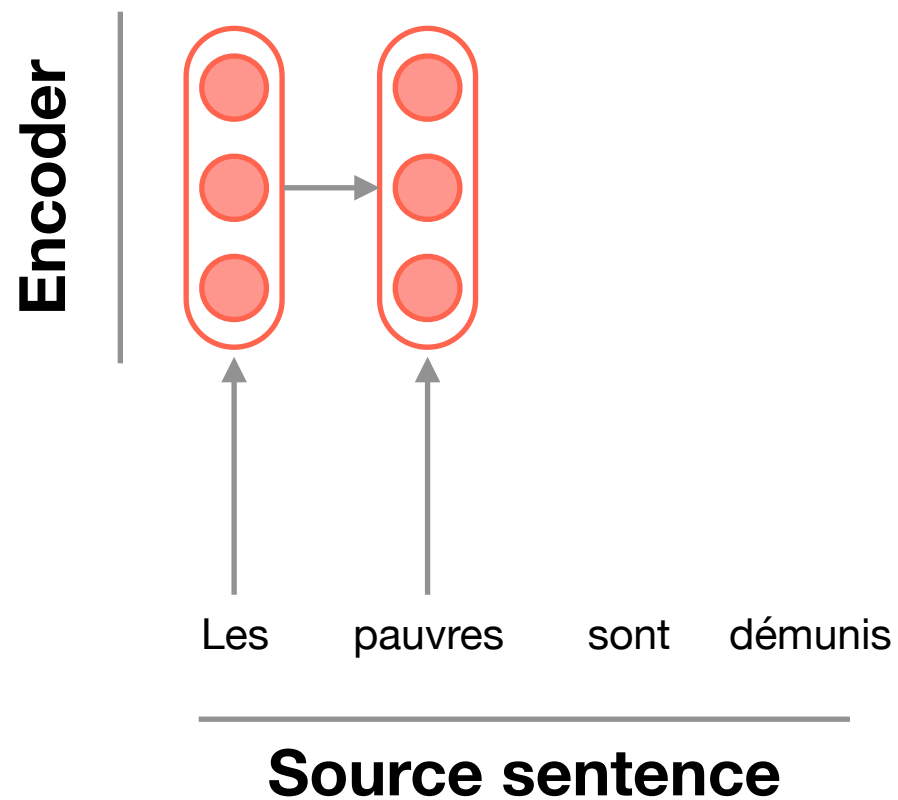
---

**Source sentence**

# Sequence to sequence
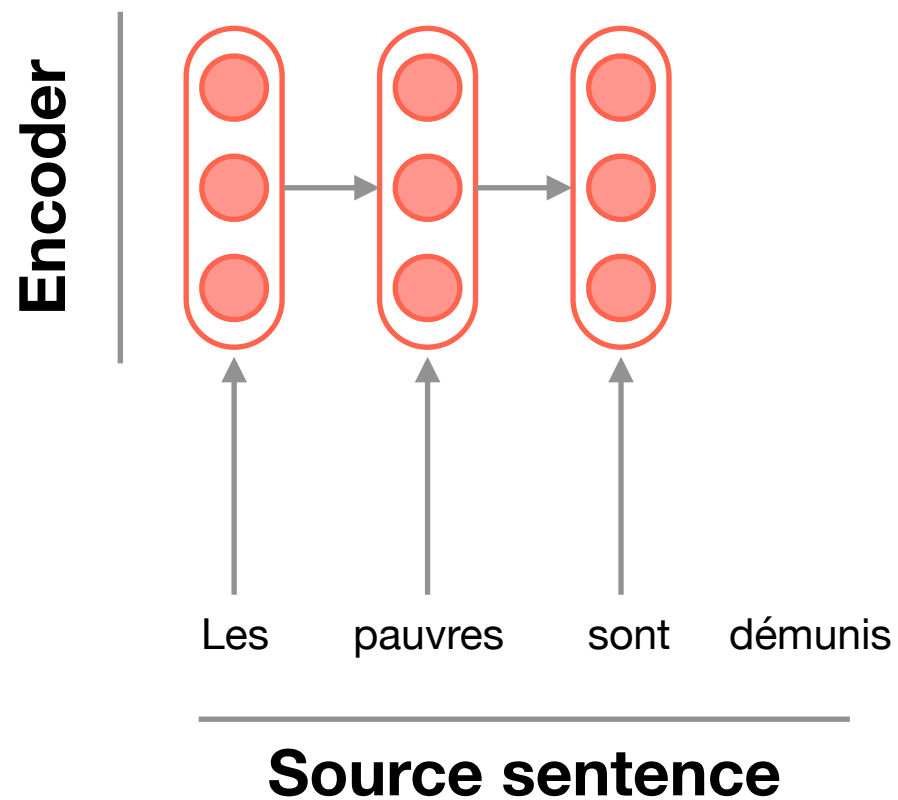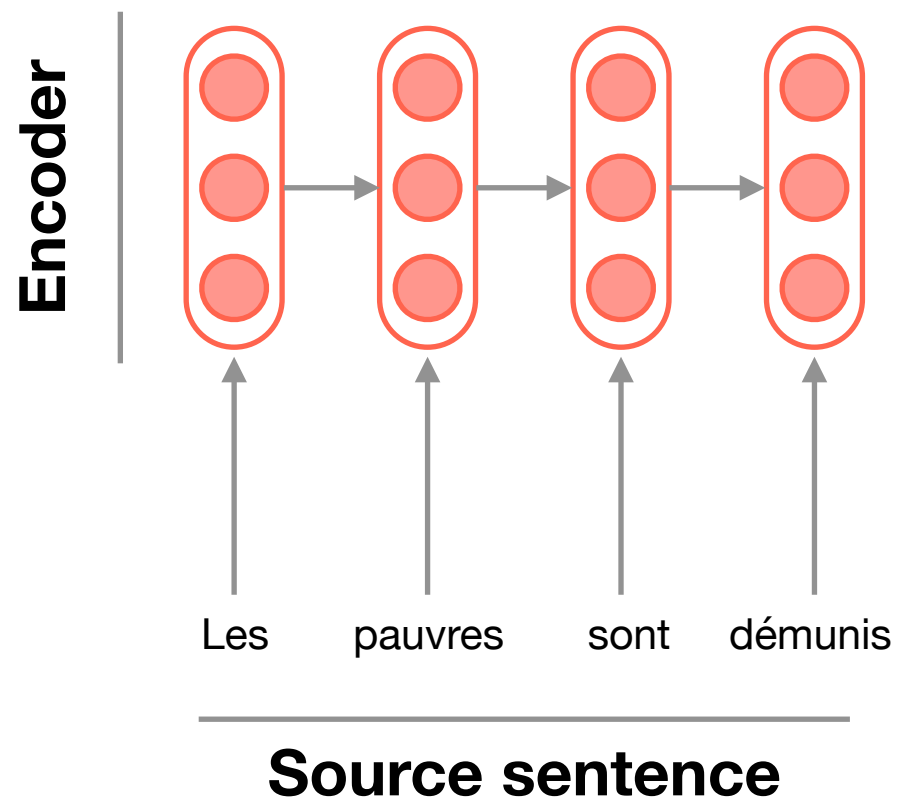
**Encoder**

Les    pauvres    sont    démunis

**Source sentence**

# Sequence to sequence

Encoder

Source sentence

Les    pauvres    sont    démunis

# Sequence to sequence

**Encoder**

Les     pauvres     sont     démunis

**Source sentence**

# Sequence to sequence

**Encoder**

**Source sentence**

Les   pauvres   sont   démunis

# Sequence to sequence

# Sequence to sequence

**Encoded source sentence**



**Encoder**

Les  pauvres  sont  démunis

**Source sentence**

# Sequence to sequence

**Inference**

**Encoded source sentence**

**Encoder**

**Decode**

Les   pauvres   sont   démunis

**Source sentence**

# Sequence to sequence

## Inference

**Encoded source sentence**



Encoder

Decode

Les    pauvres    sont    démunis            <START>

**Source sentence**

# Sequence to sequence

## Inference

**Encoded source sentence**

**Target sentence**

**Encoder**

The

**Decode**

Les    pauvres    sont    démunis

&lt;START&gt;
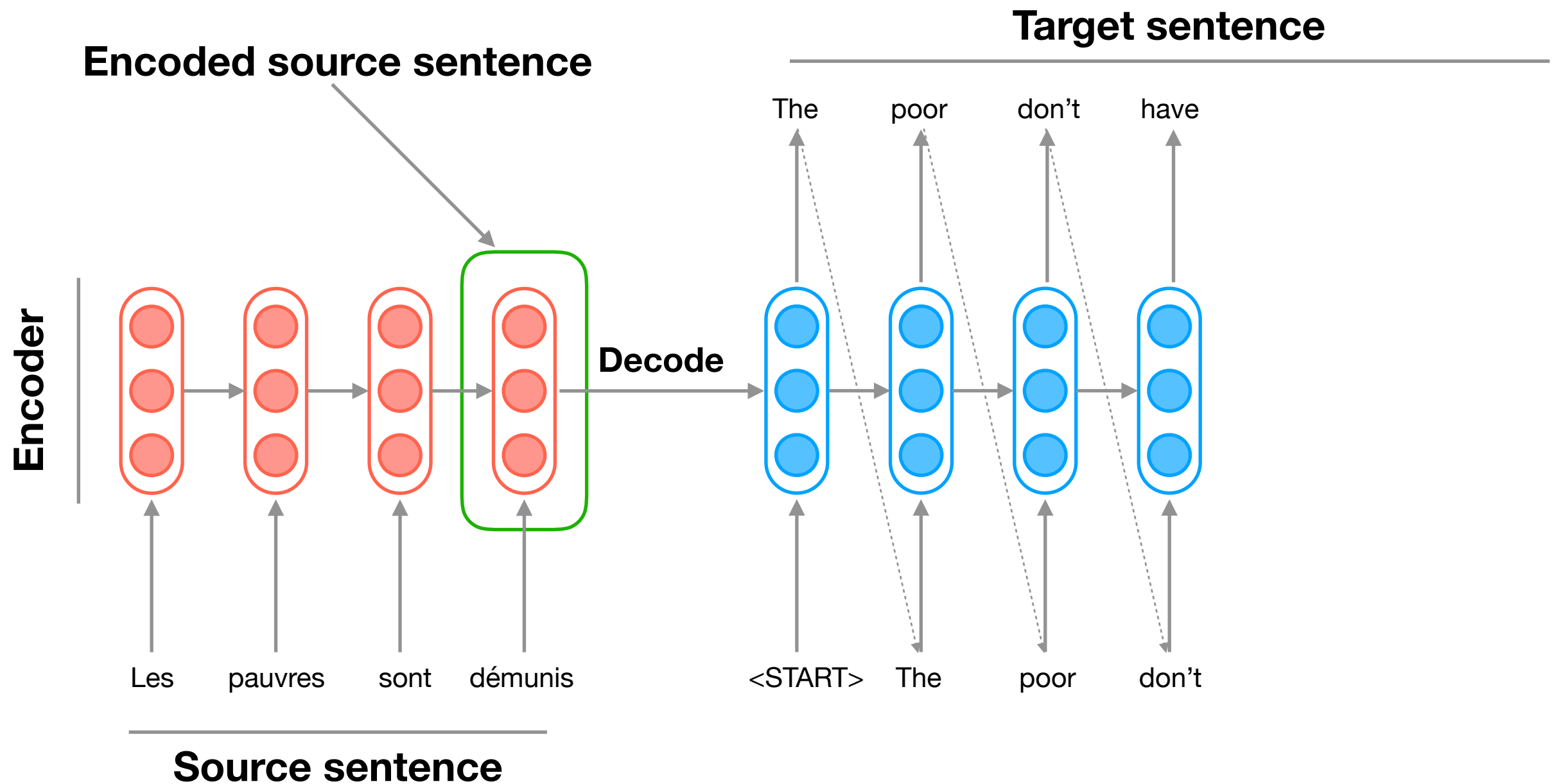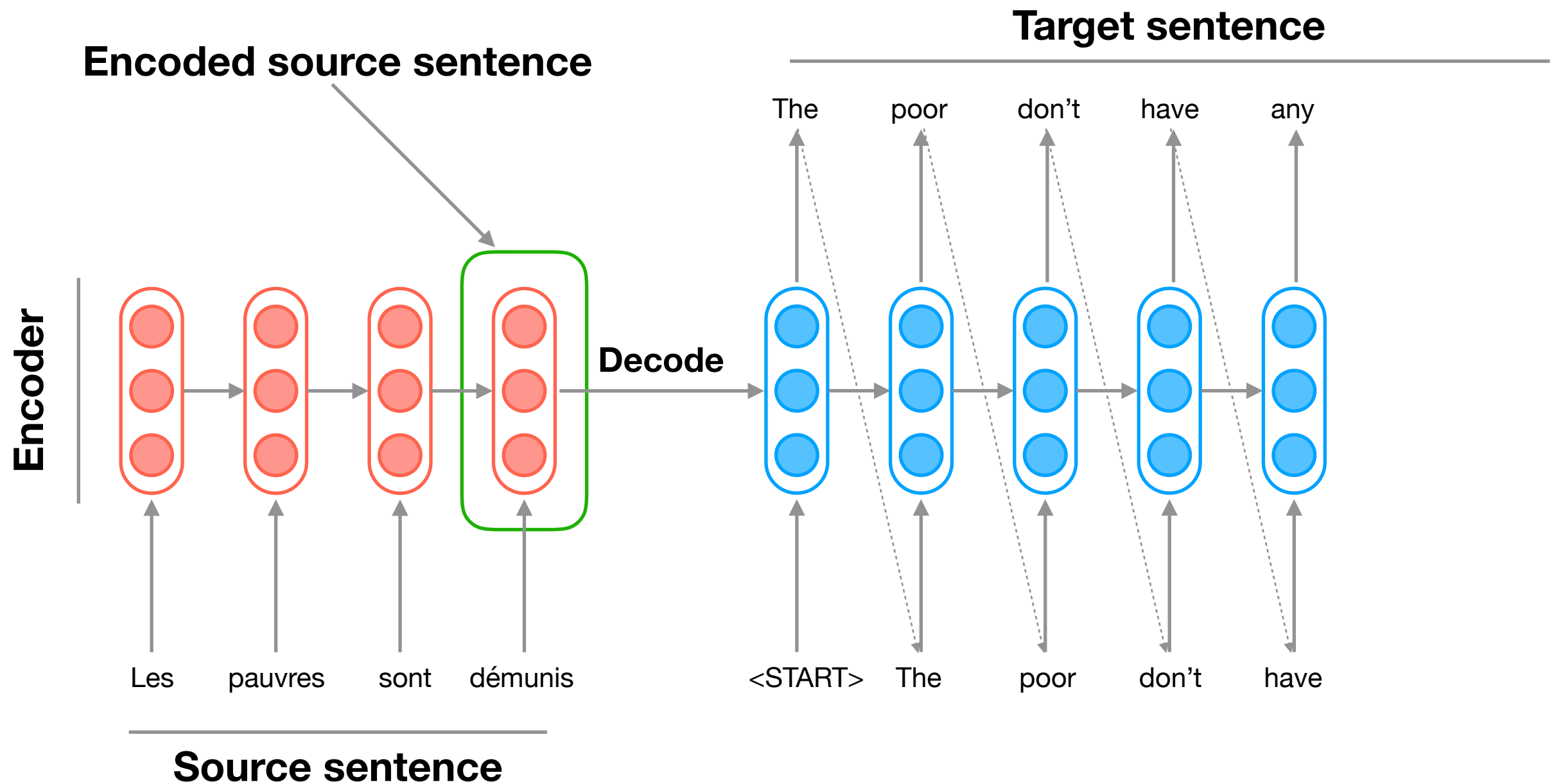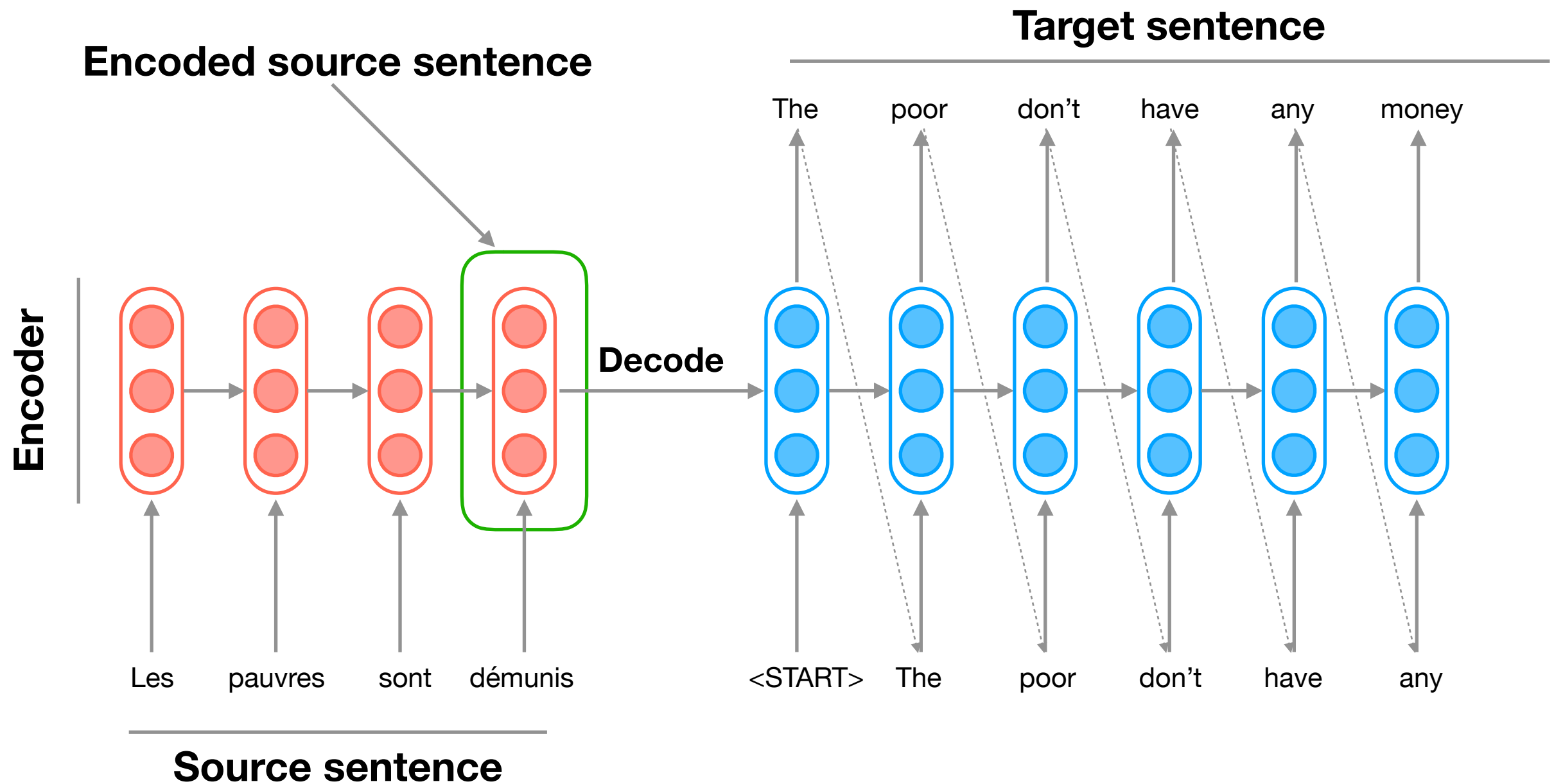
**Source sentence**

# Sequence to sequence

## Inference

# Sequence to sequence

## Inference

# Sequence to sequence

## Inference



**Encoded source sentence**

**Target sentence**

**Encoder**

The    poor    don't    have

Decode

Les    pauvres    sont    démunis

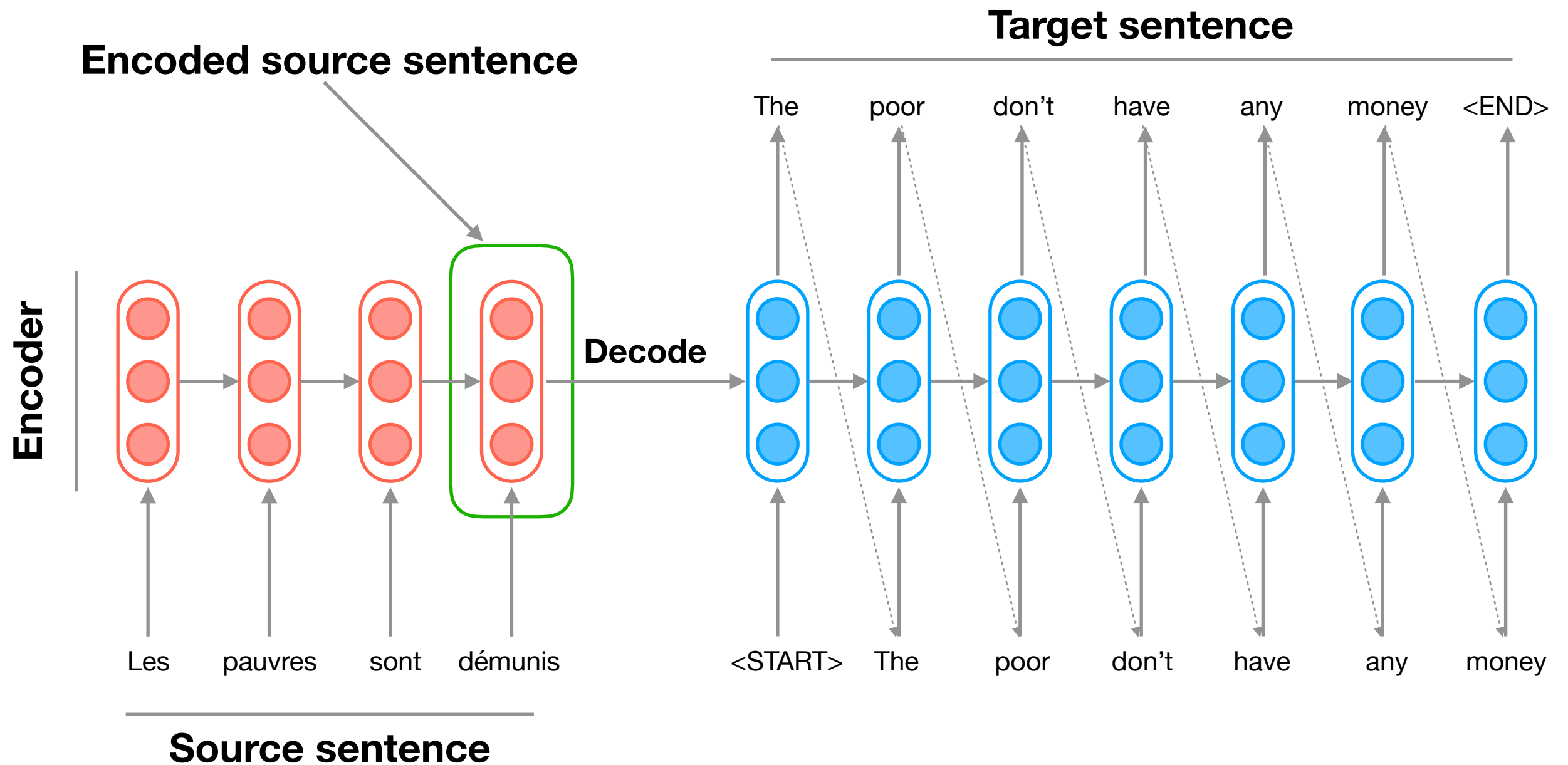<START>    The    poor    don't

**Source sentence**

# Sequence to sequence

## Inference

# Sequence to sequence

## Inference

# Sequence to sequence

## Inference

# Sequence to sequence

## Inference



Encoder

Encoded source sentence

Bottleneck!

Decode

Source sentence: Les pauvres sont démunis

Target sentence: The poor don't have any money <END>

<START> The poor don't have any money
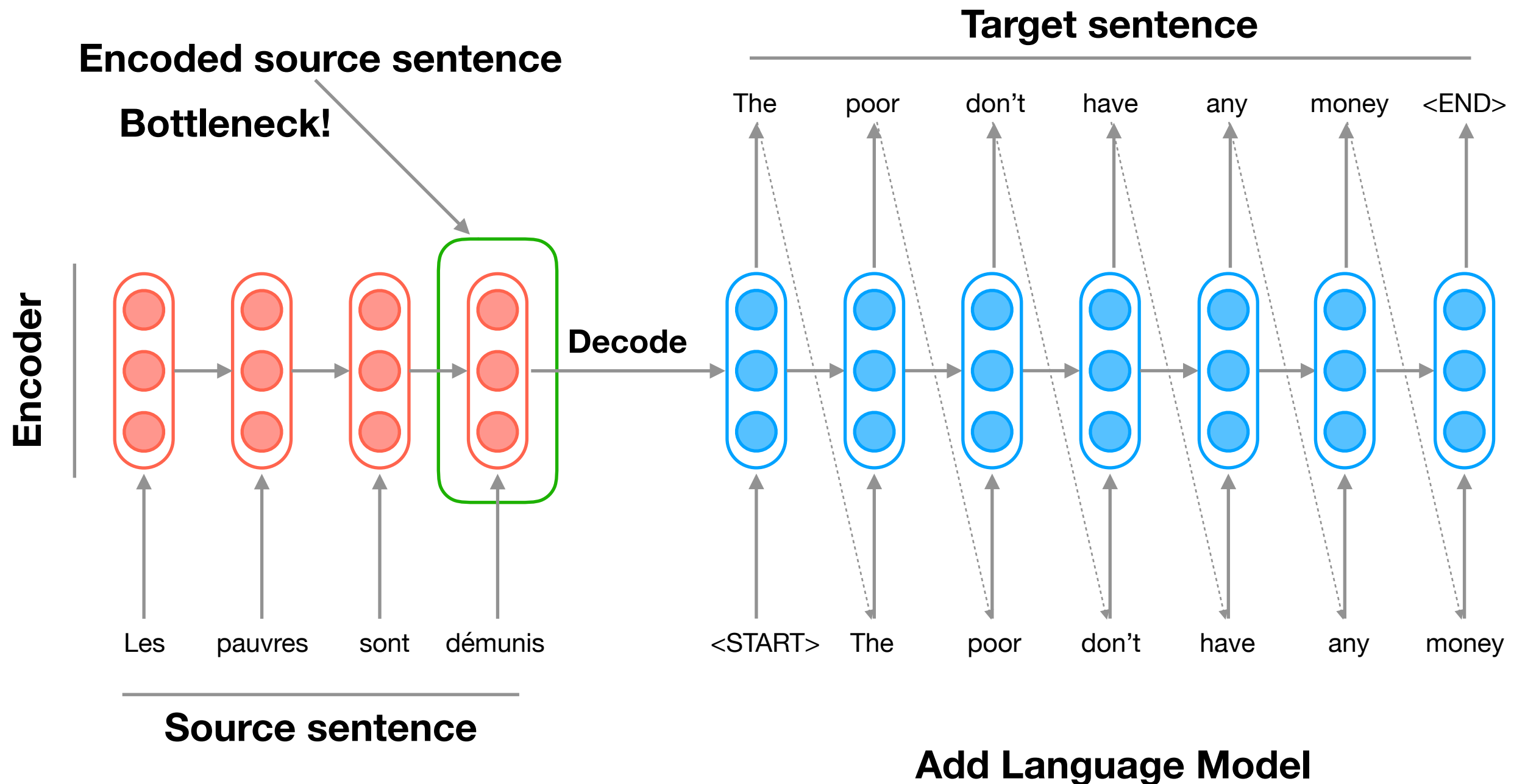
# Sequence to sequence

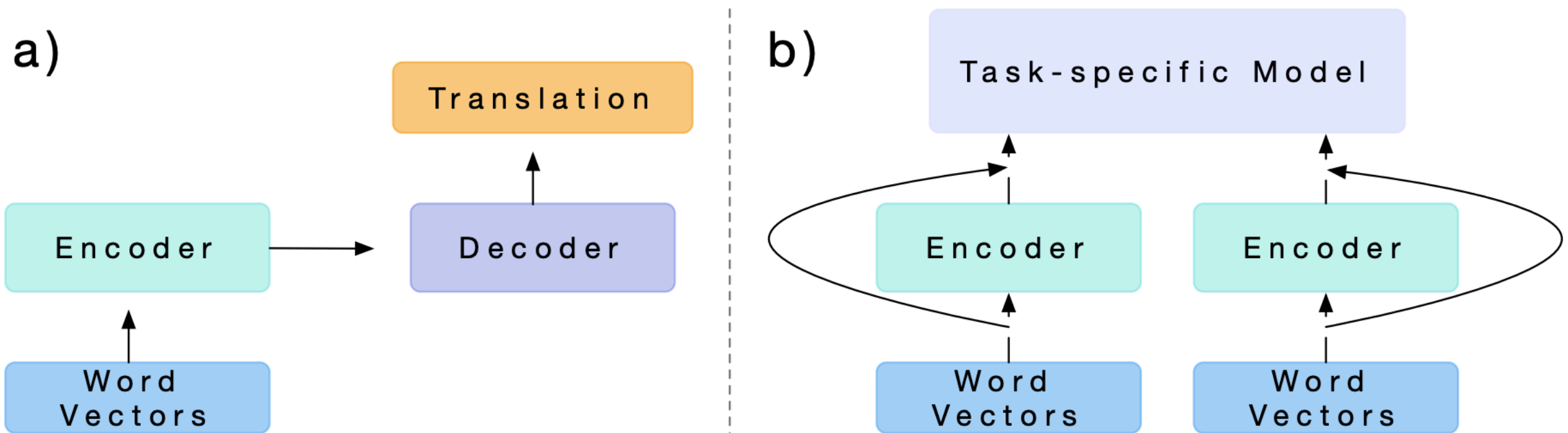## Inference

# Contextualized Word Vectors



Figure 1: We a) train a two-layer, bidirectional LSTM as the encoder of an attentional sequence-to-sequence model for machine translation and b) use it to provide context for other NLP models.

# Thanks for your Attention!

Boris Zubarev

@bobazooba