

# Language Models

Boris Zubarev

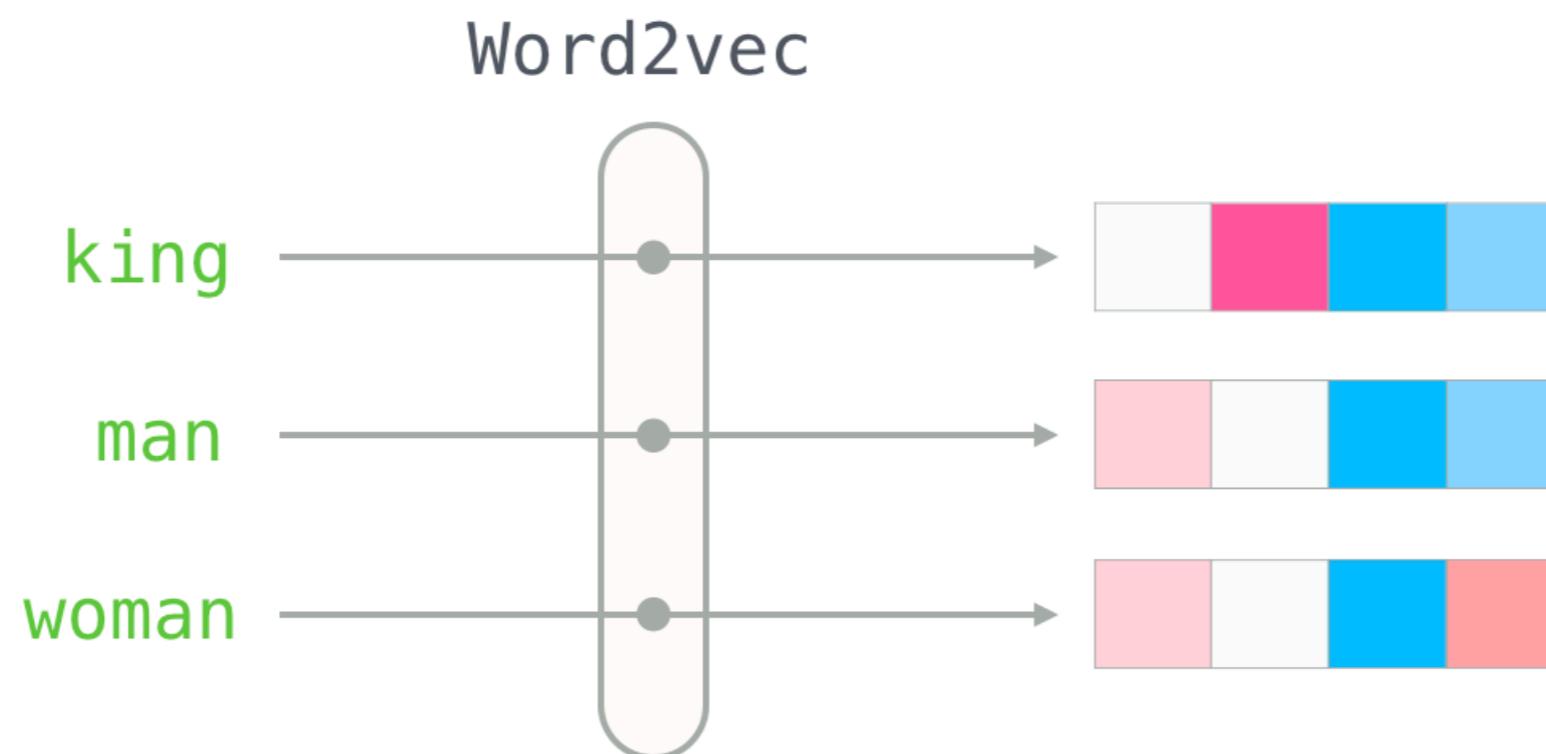


@bobazooba

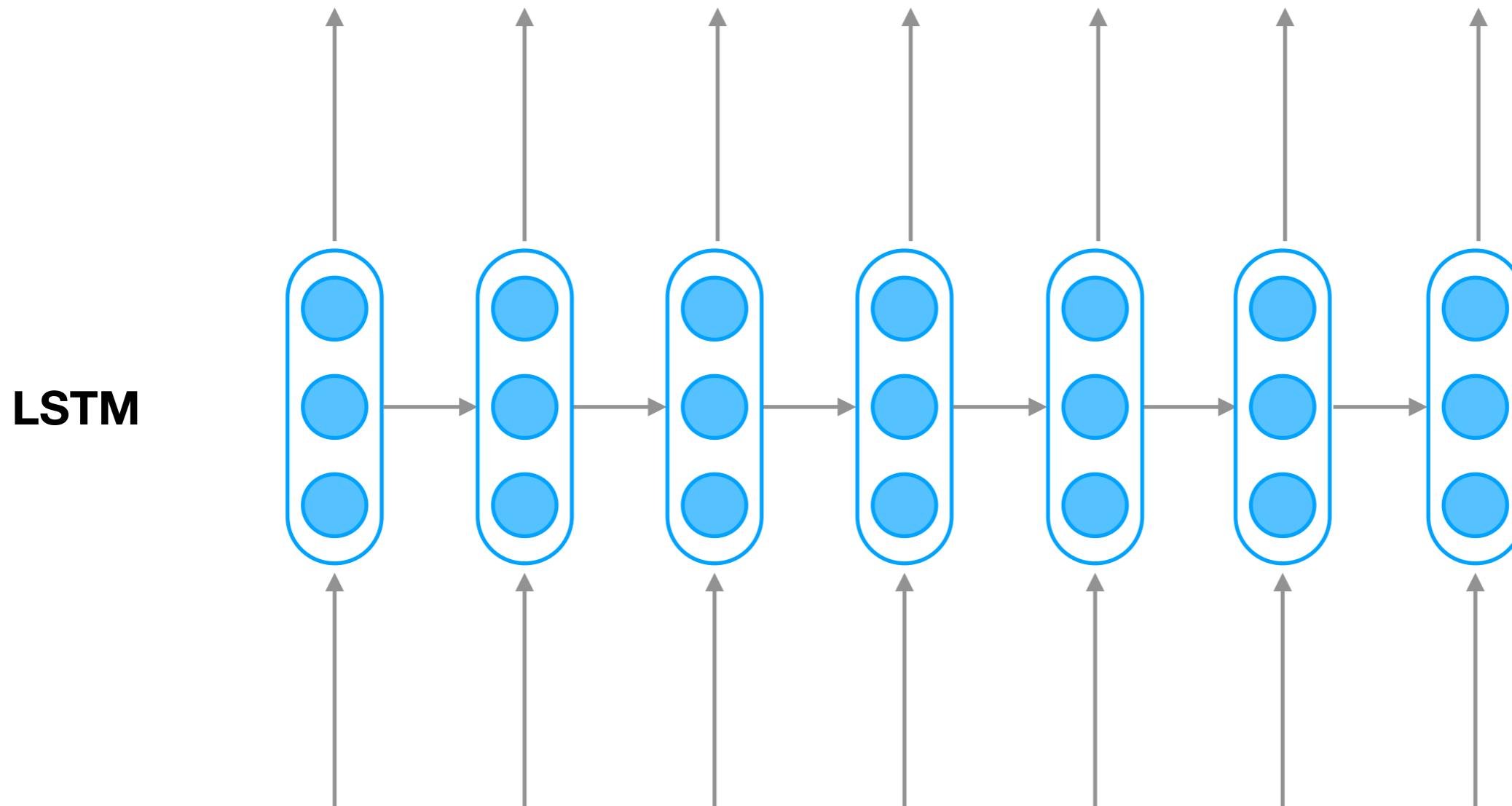
# Word Embeddings

w2v, glove, etc

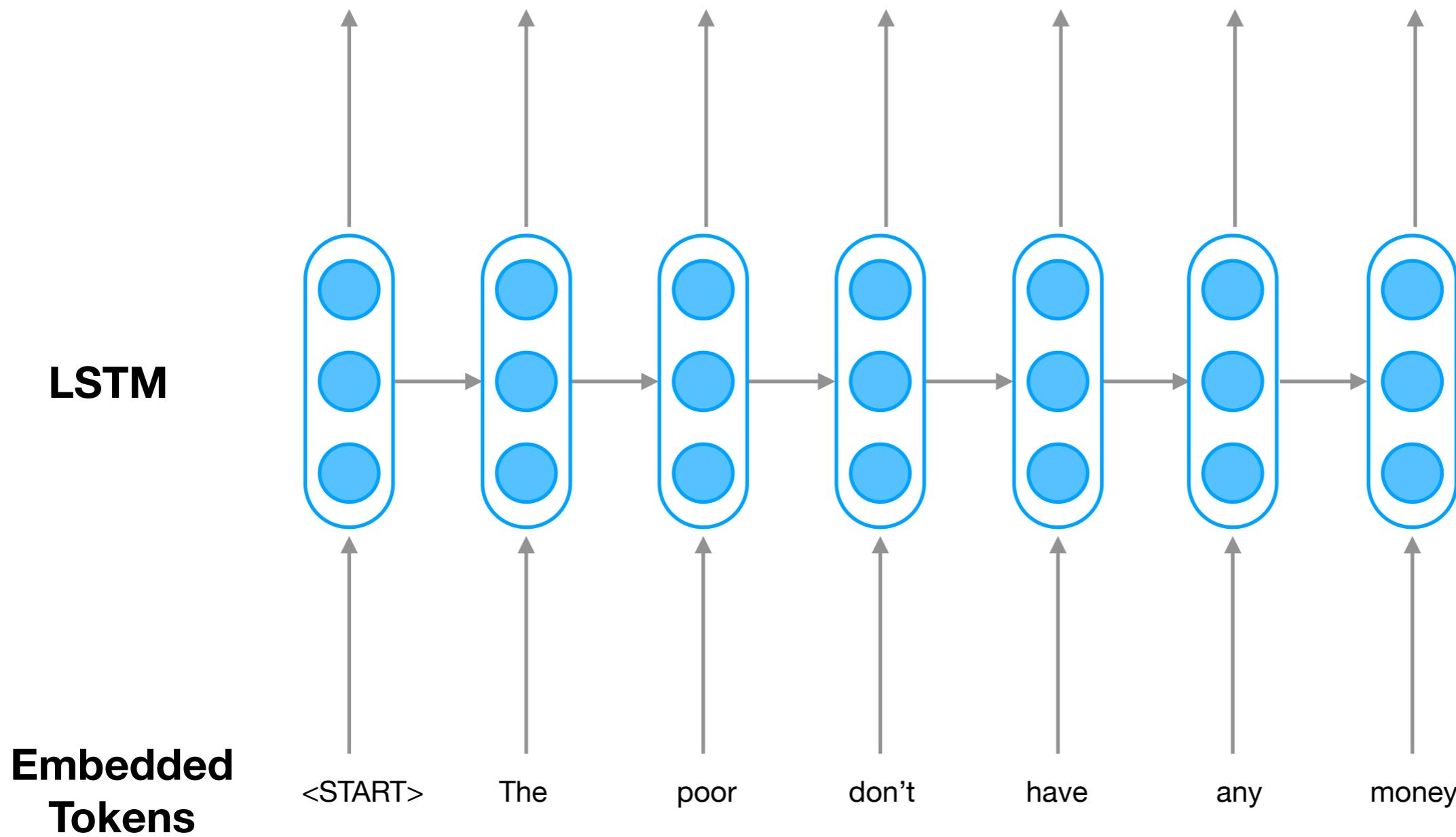
- Just key—value storage at inference
- Don't change from relationships with other words in the current text



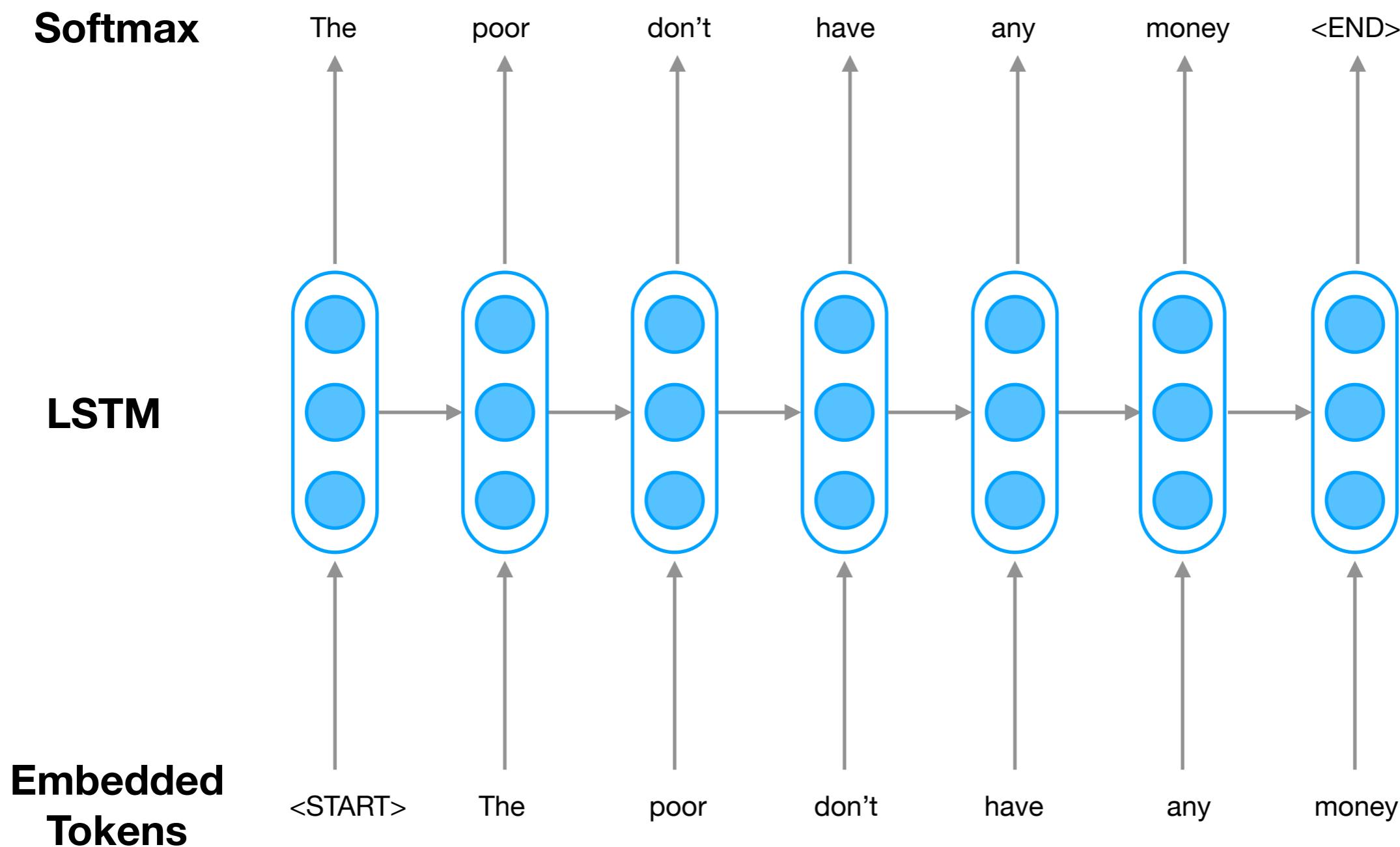
# Language Model



# Language Model



# Language Model



# Language Model

## Inference

LSTM

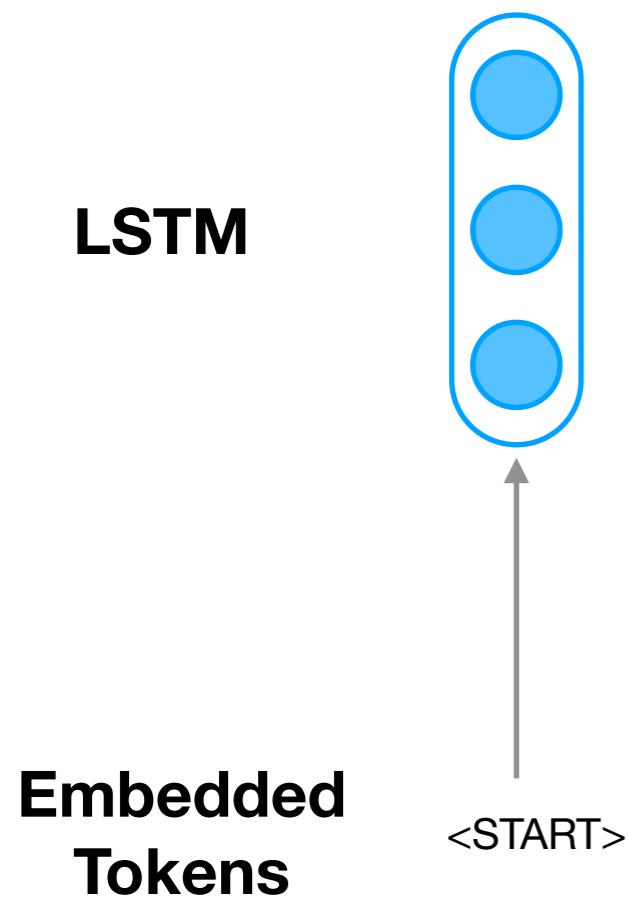
Embedded  
Tokens

<START>

# Language Model

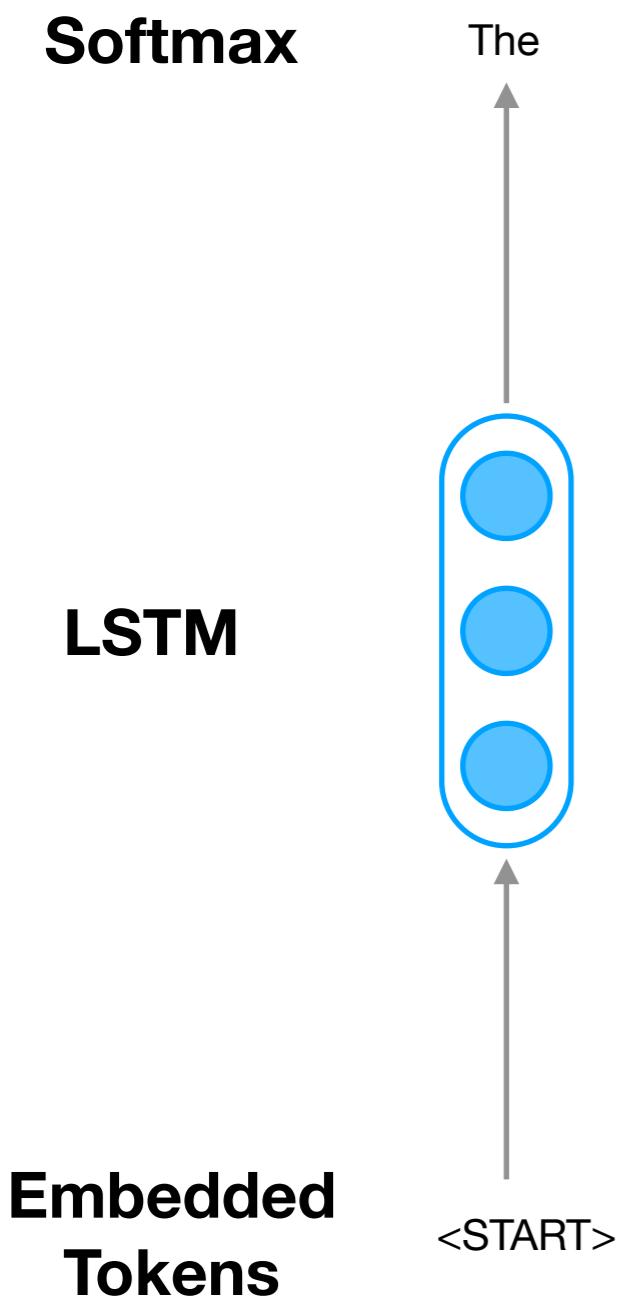
## Inference

Softmax



# Language Model

## Inference



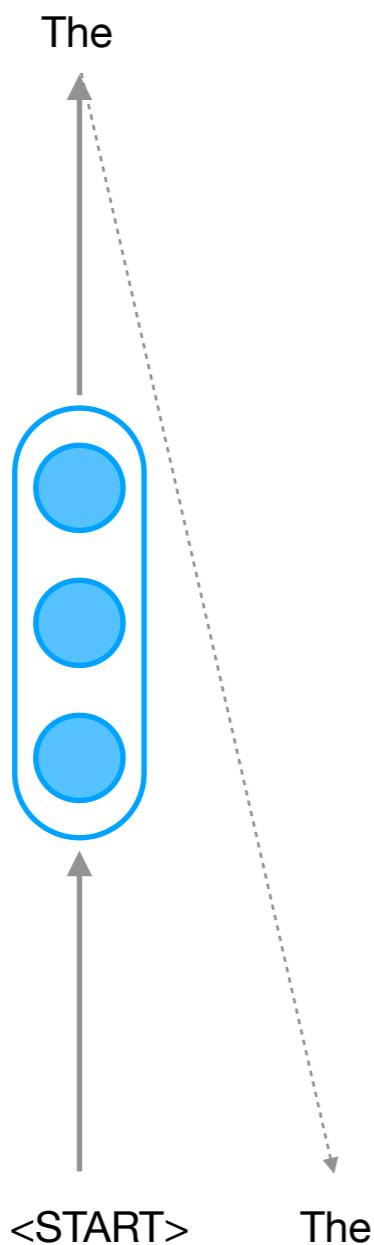
# Language Model

## Inference

**Softmax**

**LSTM**

**Embedded  
Tokens**



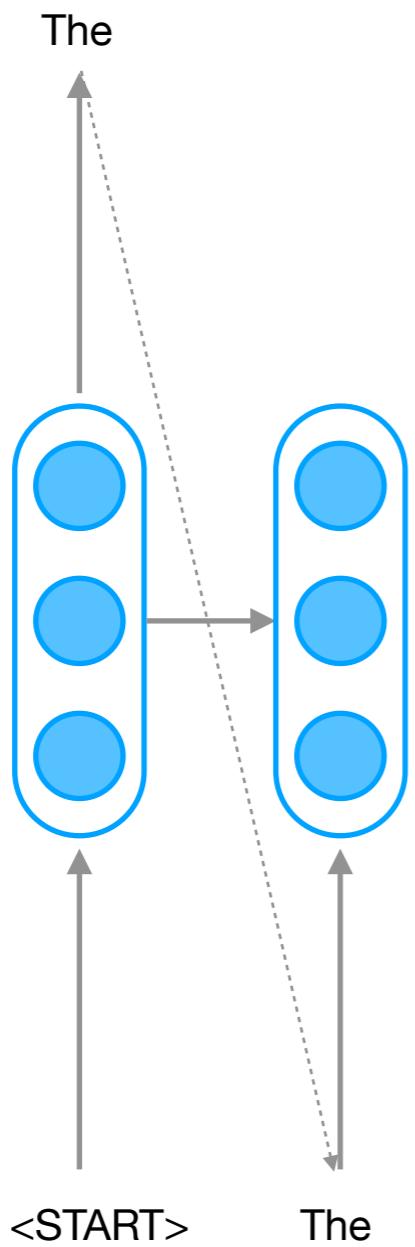
# Language Model

## Inference

**Softmax**

**LSTM**

**Embedded  
Tokens**



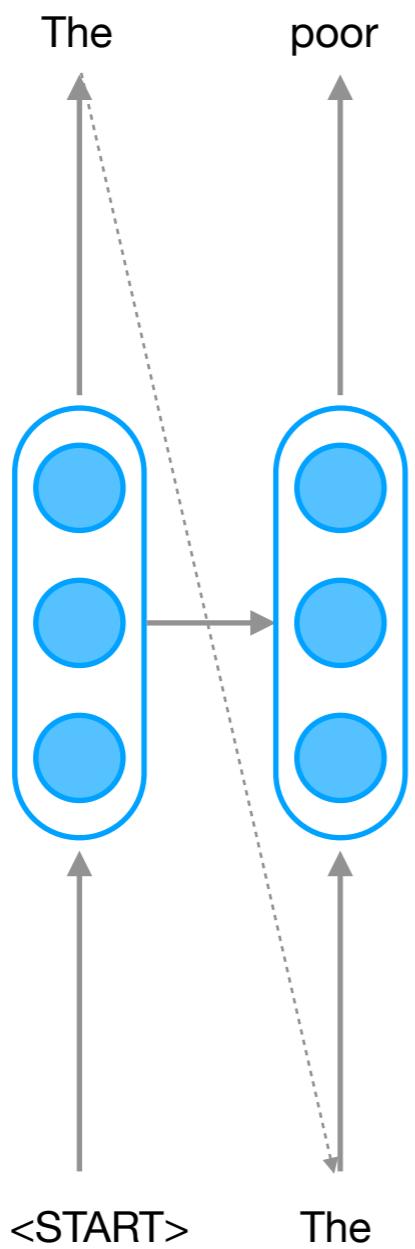
# Language Model

## Inference

**Softmax**

**LSTM**

**Embedded  
Tokens**



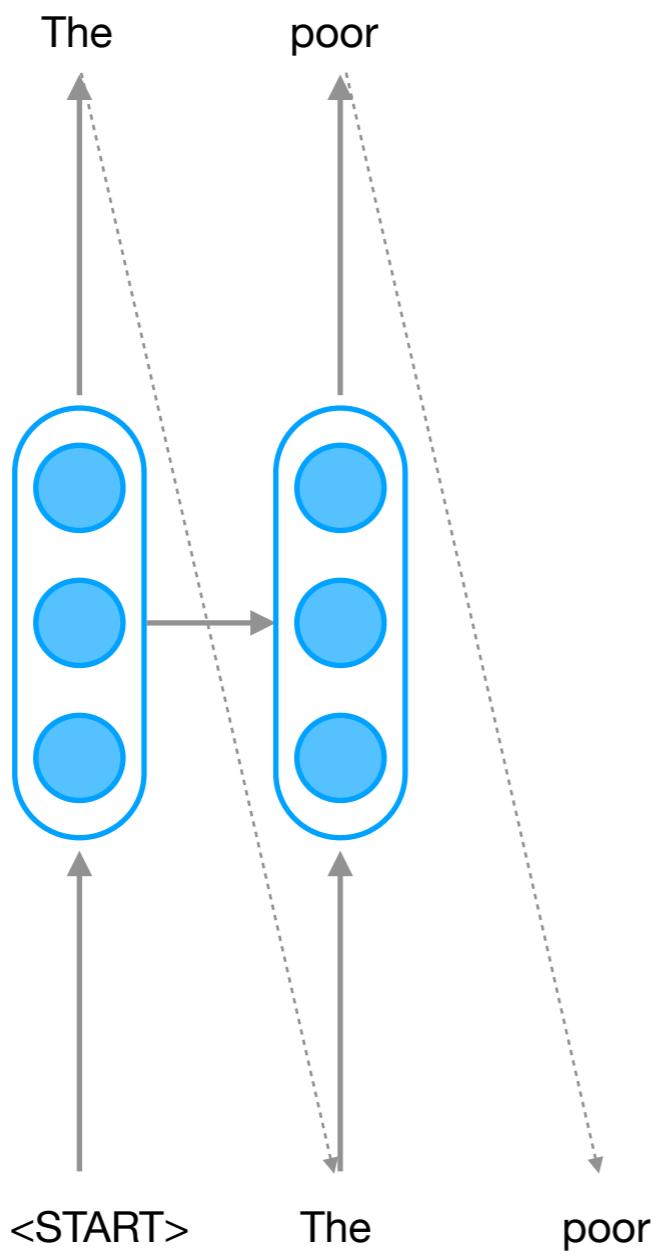
# Language Model

## Inference

**Softmax**

**LSTM**

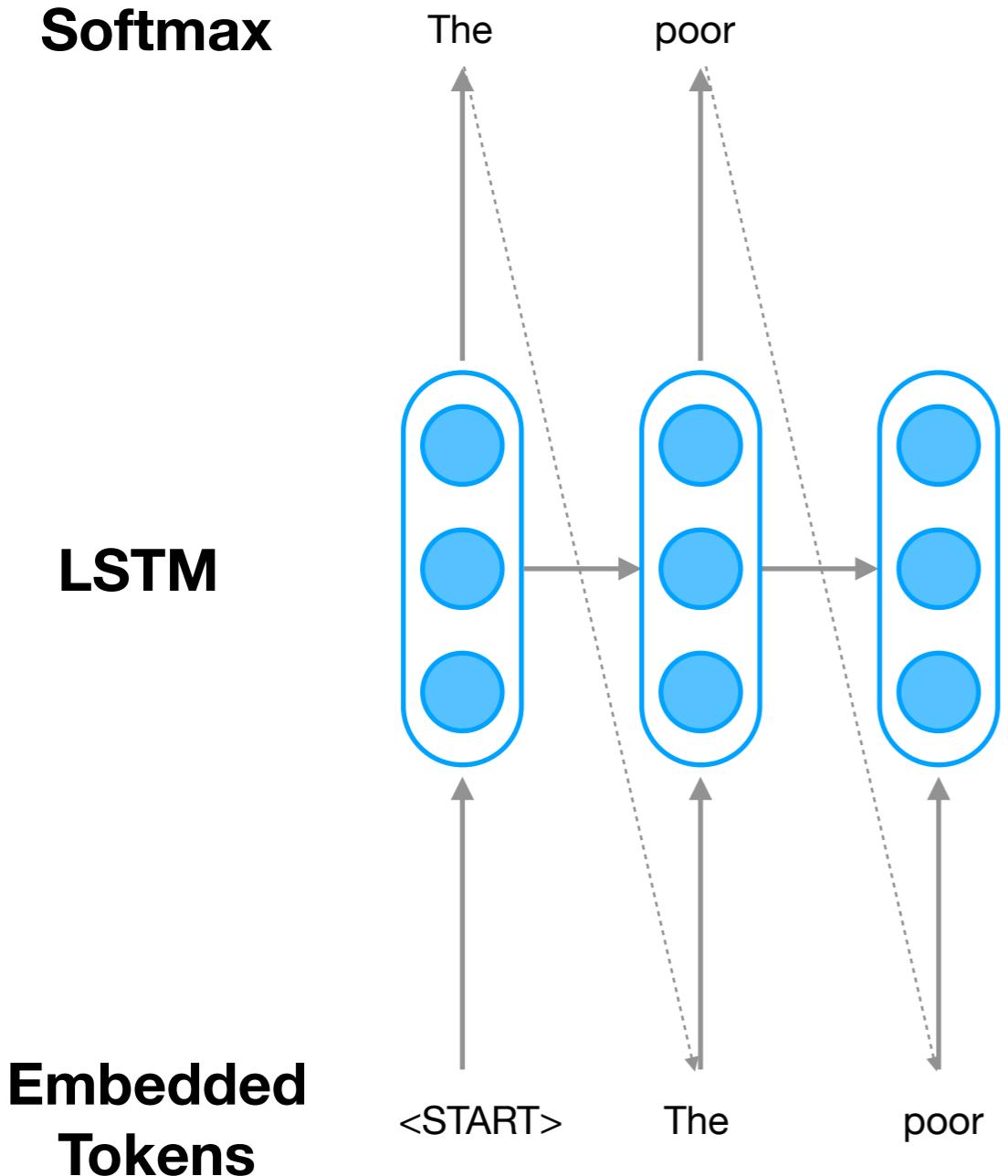
**Embedded  
Tokens**



# Language Model

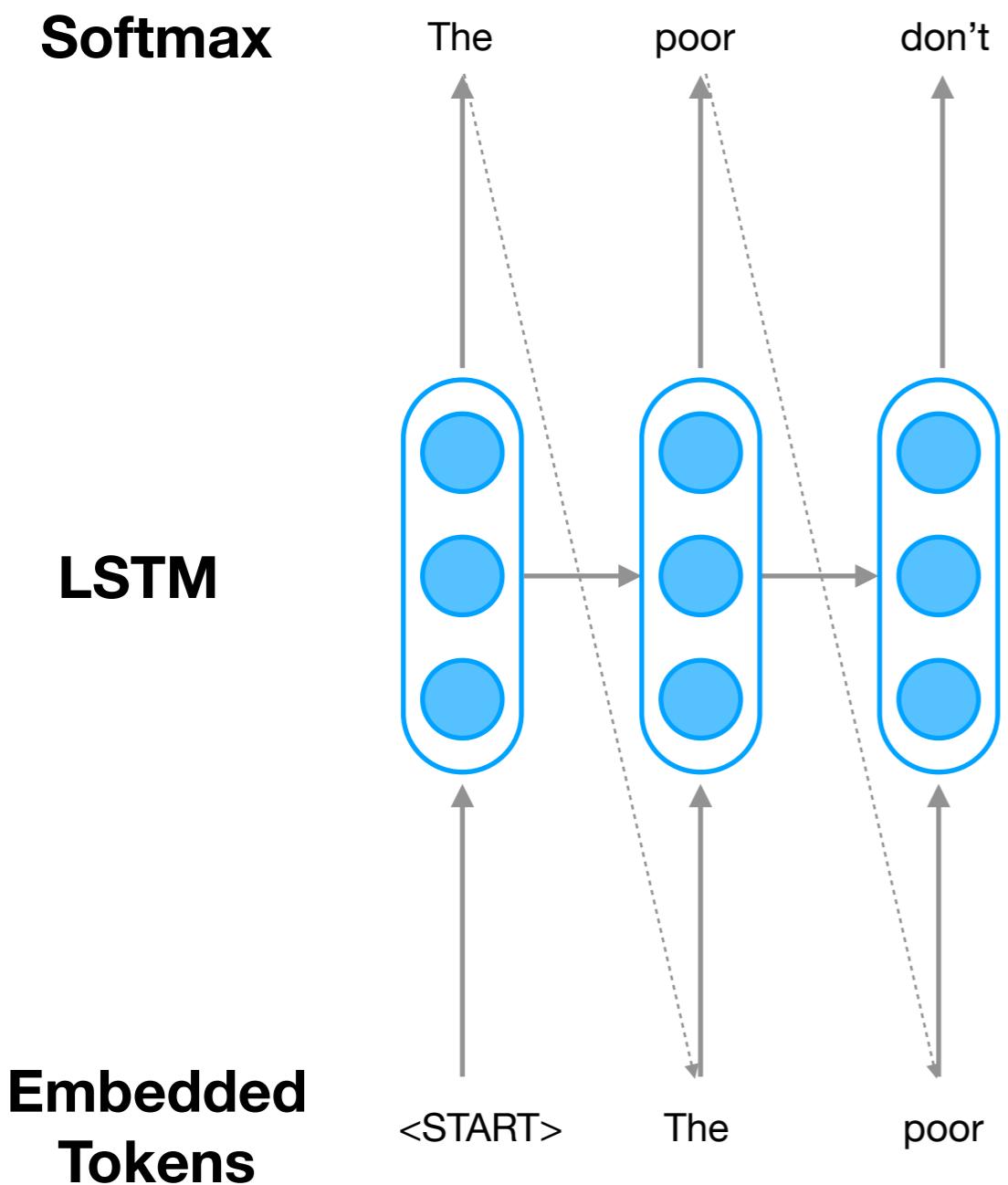
## Inference

**Softmax**



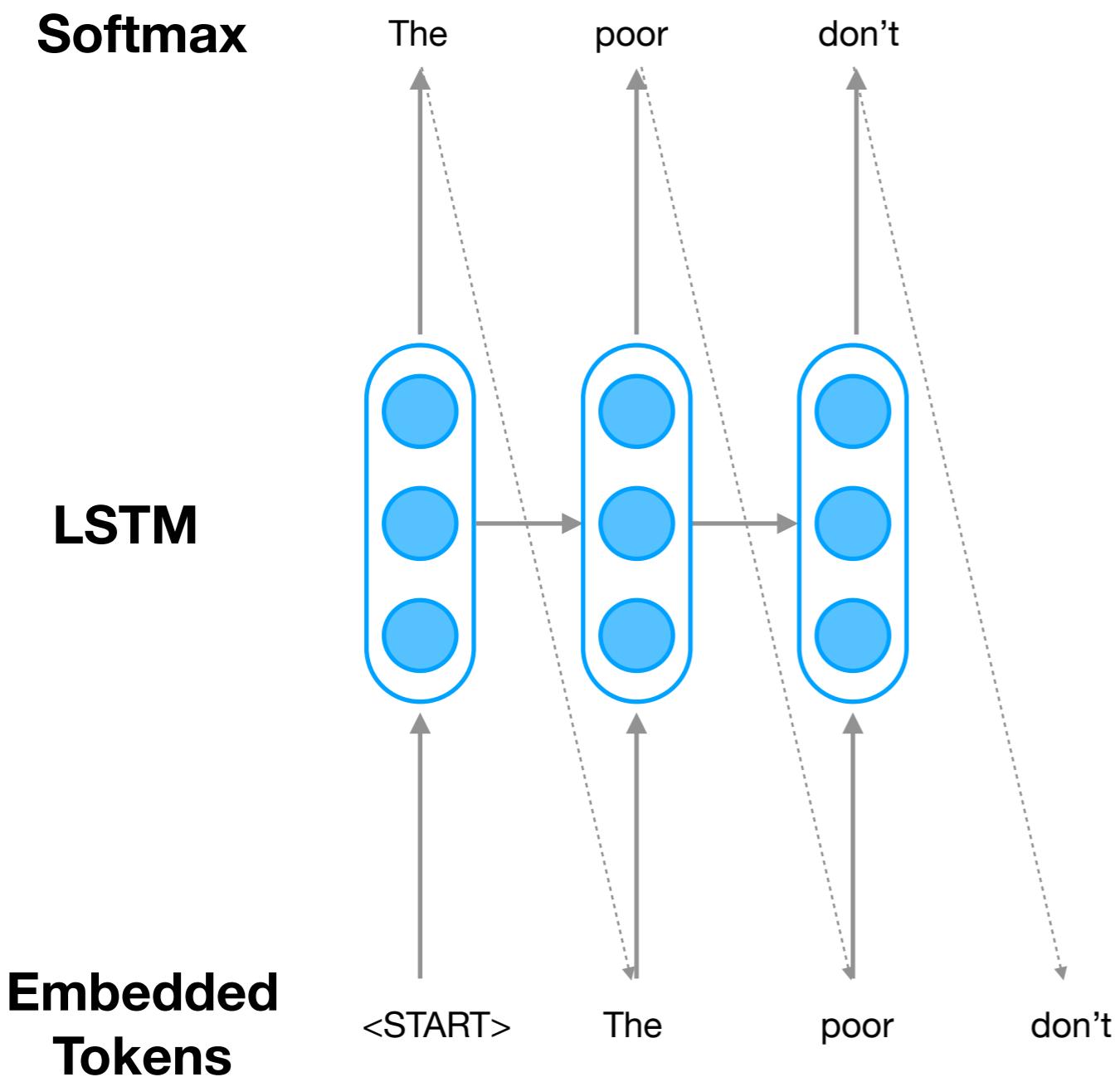
# Language Model

## Inference



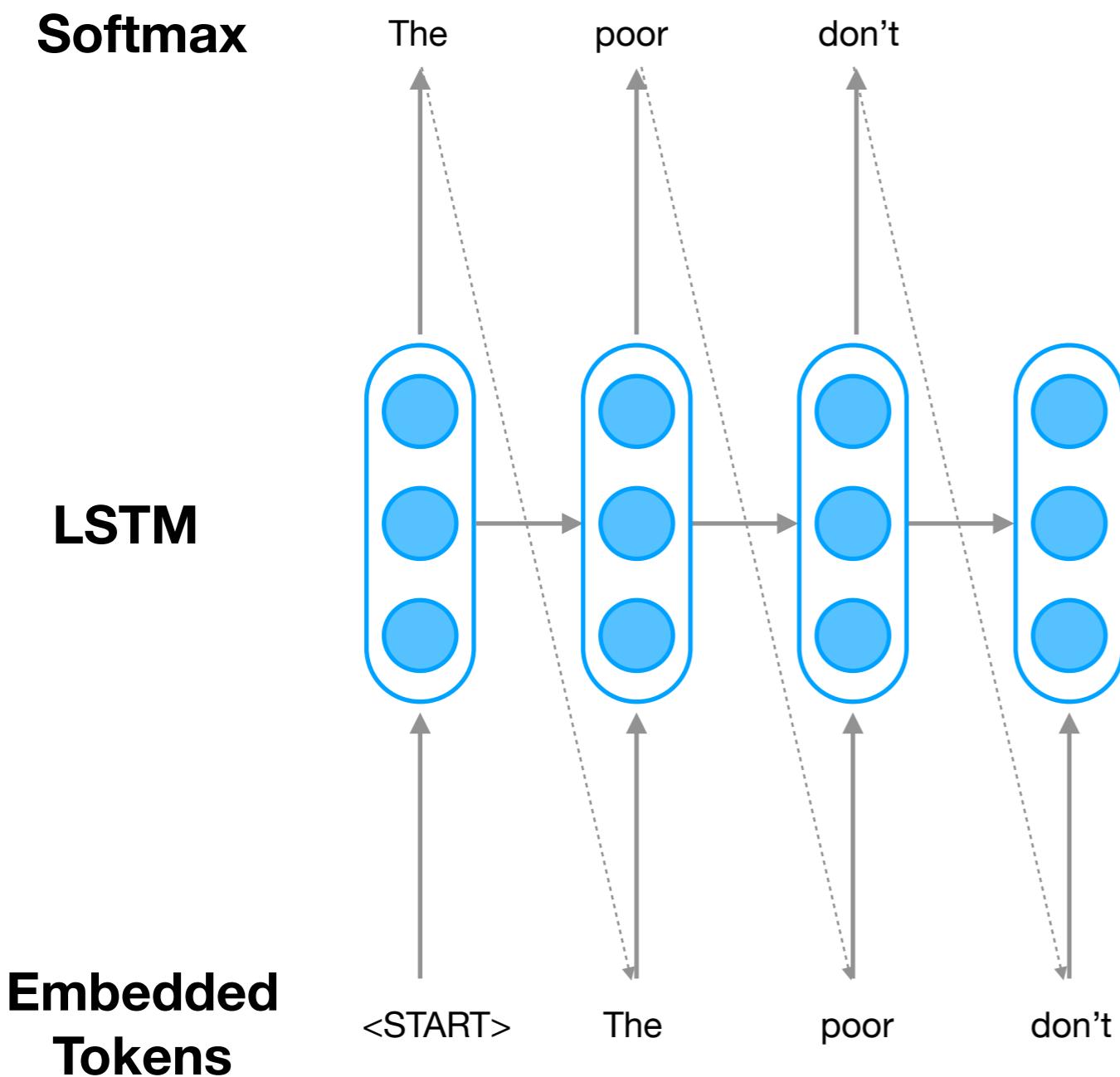
# Language Model

## Inference



# Language Model

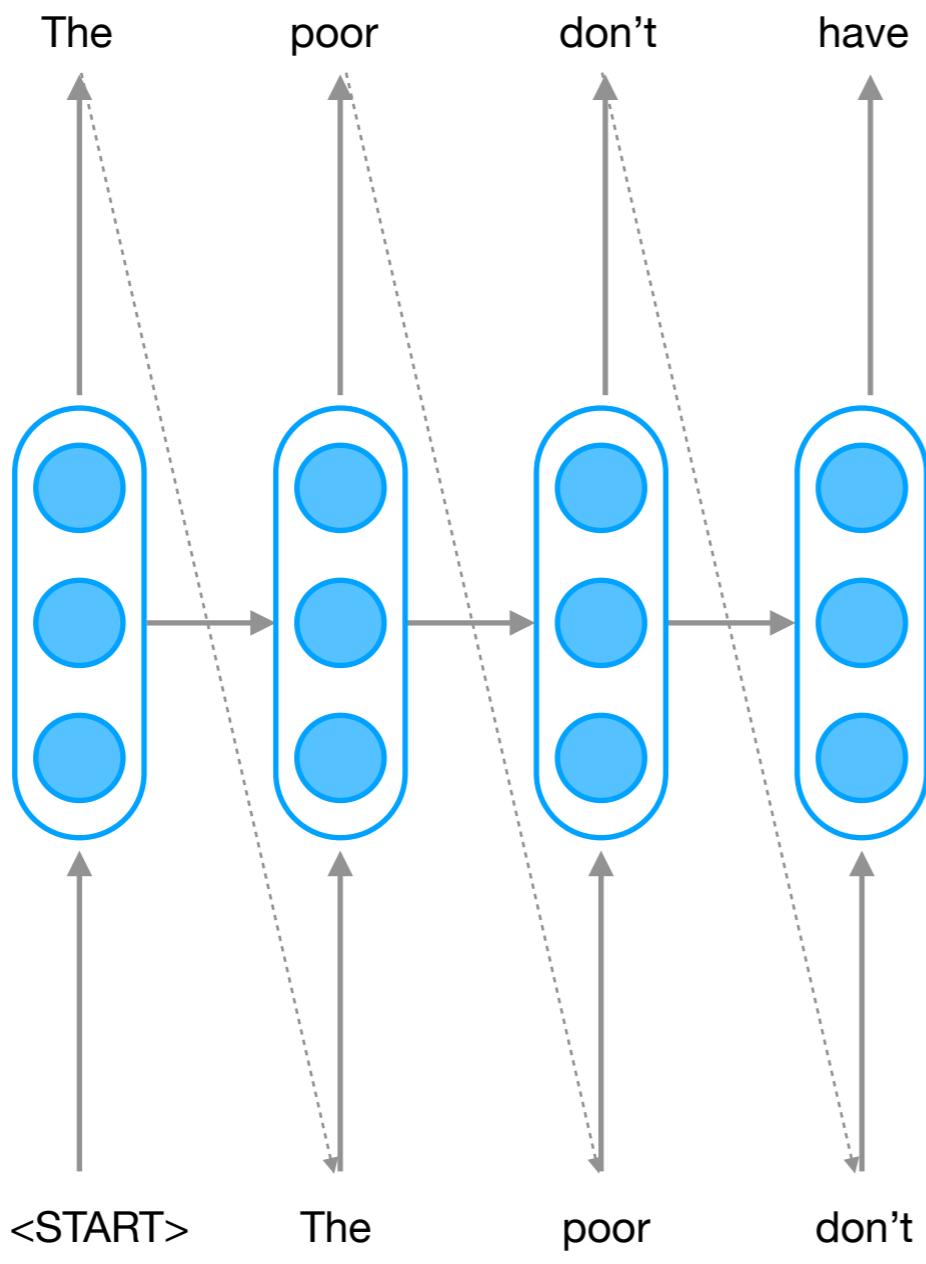
## Inference



# Language Model

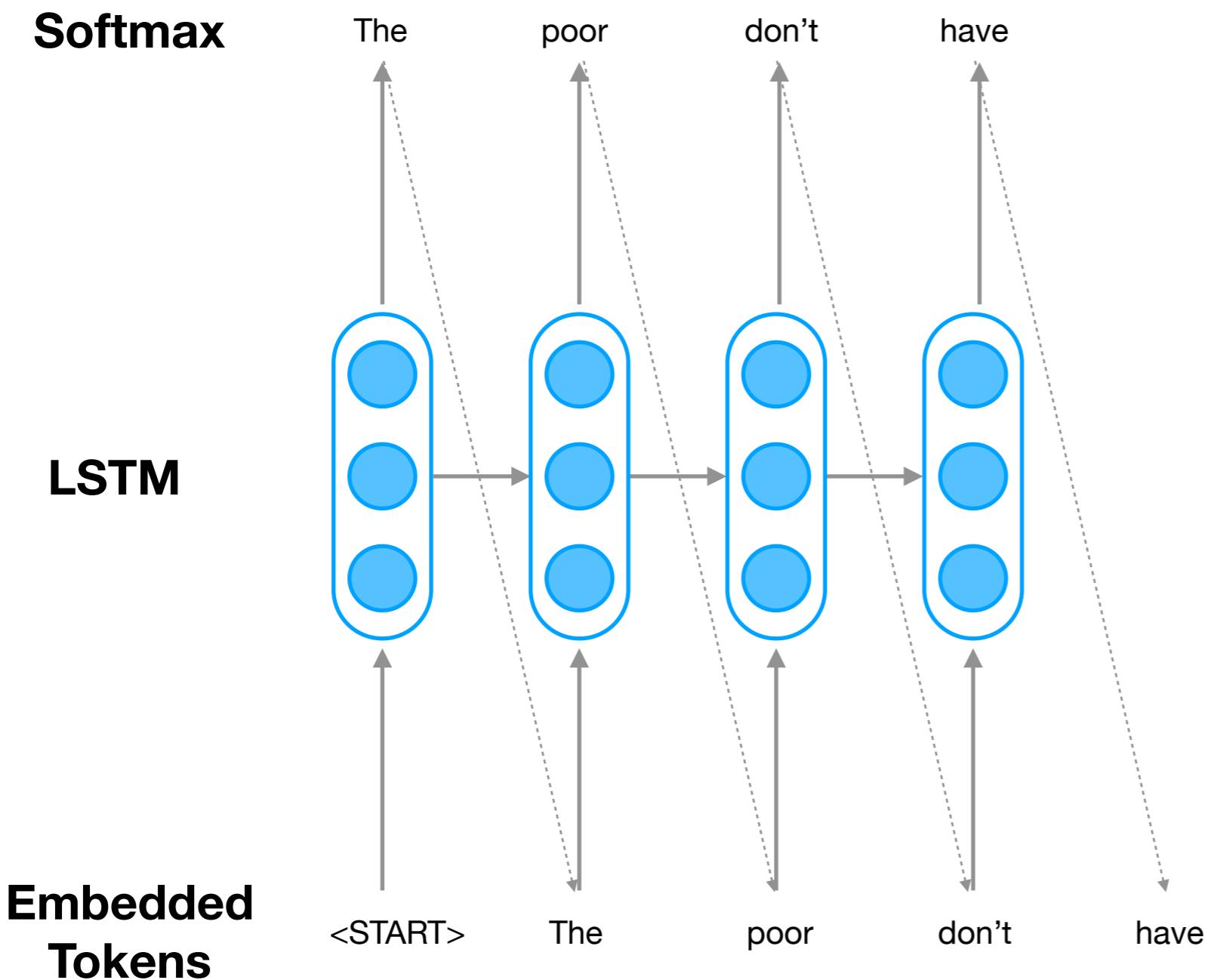
## Inference

**Softmax**



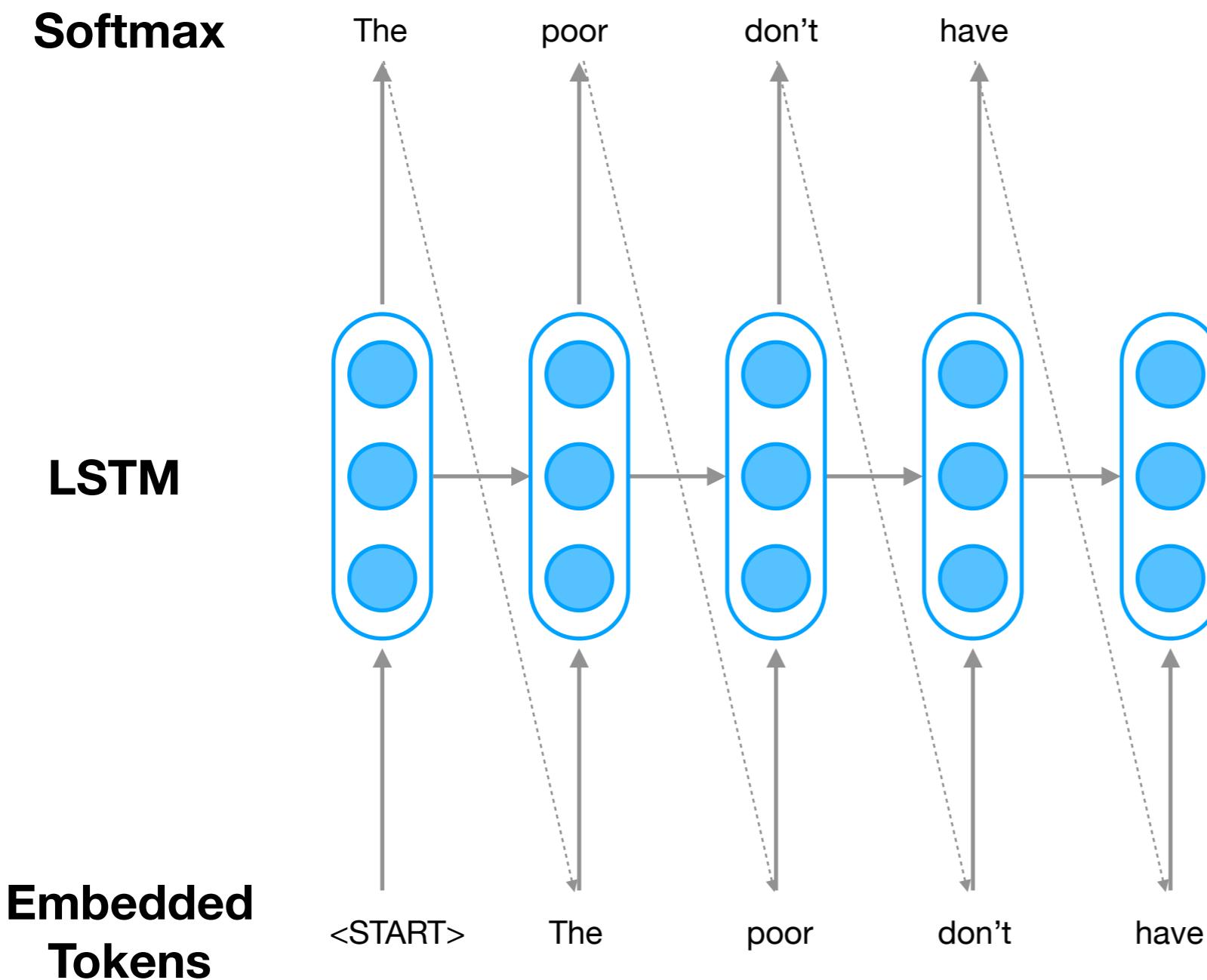
# Language Model

## Inference



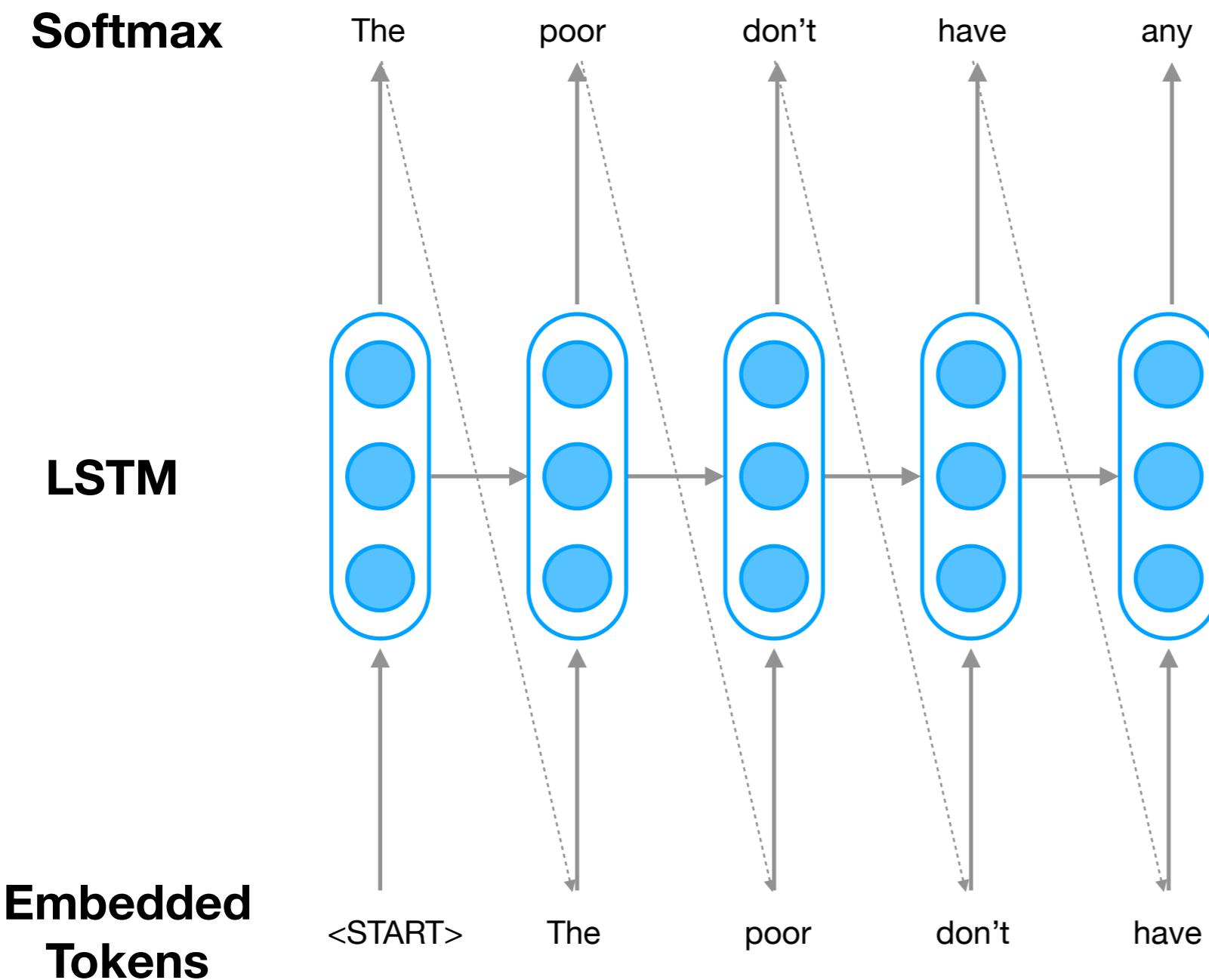
# Language Model

## Inference



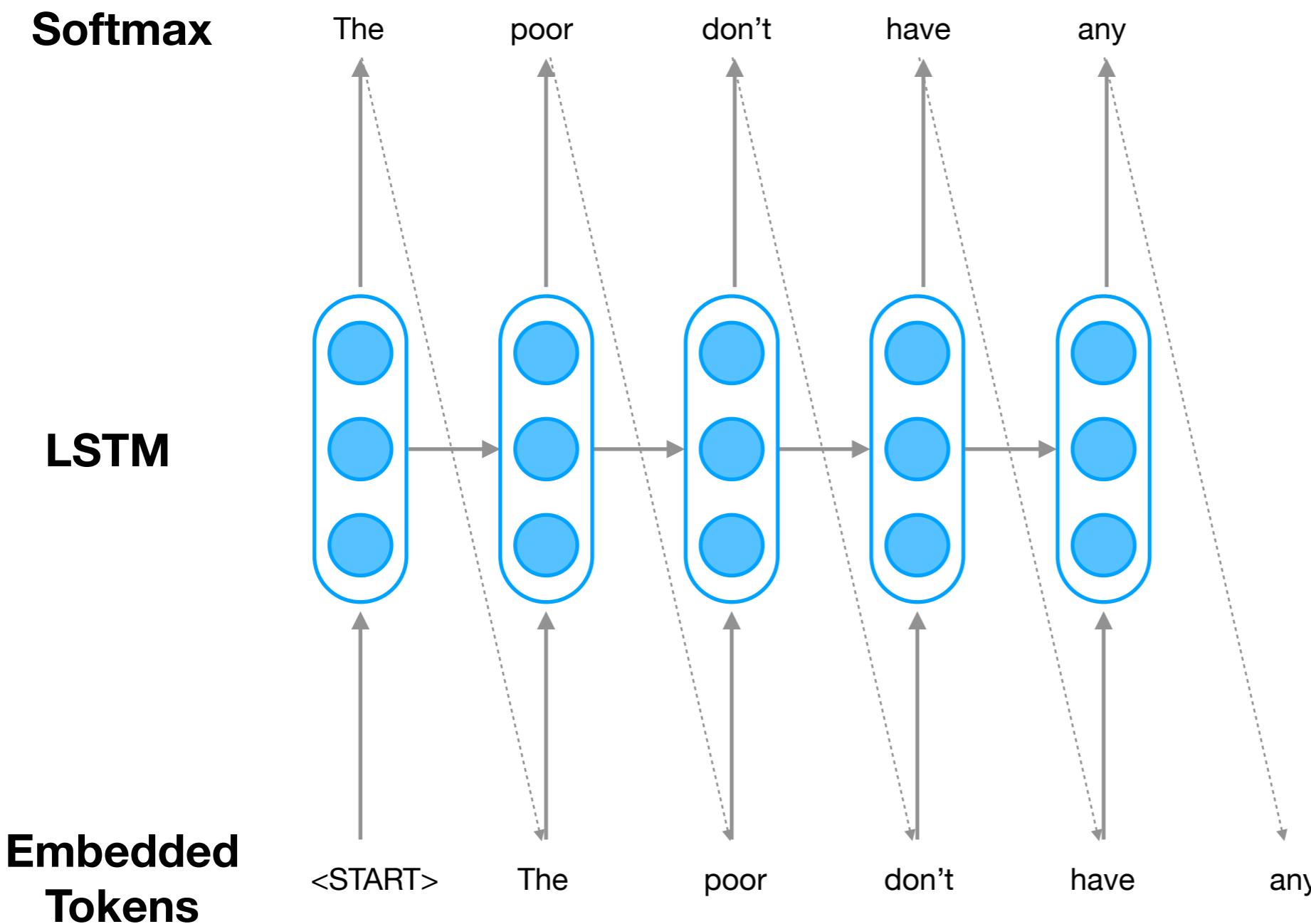
# Language Model

## Inference



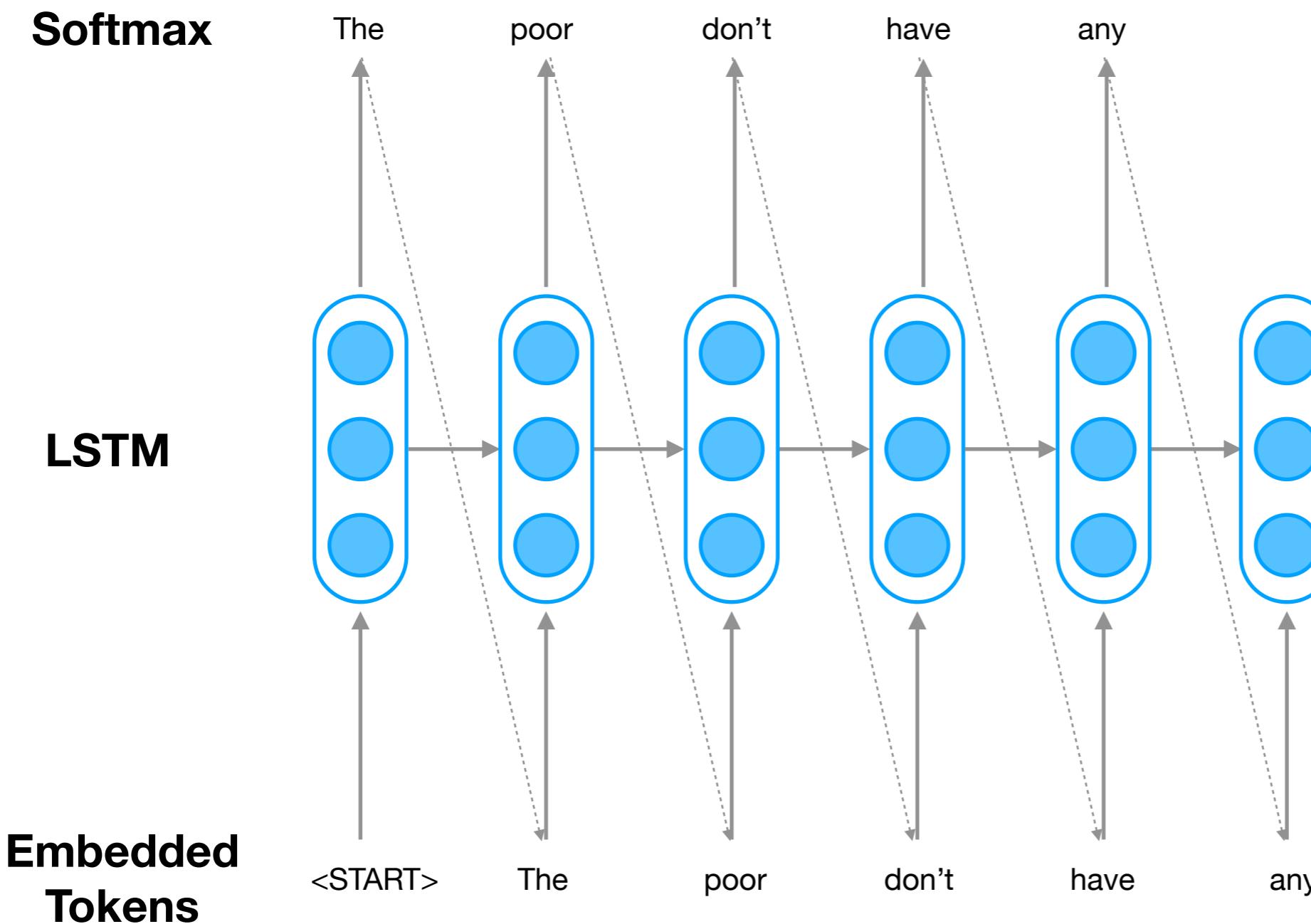
# Language Model

## Inference



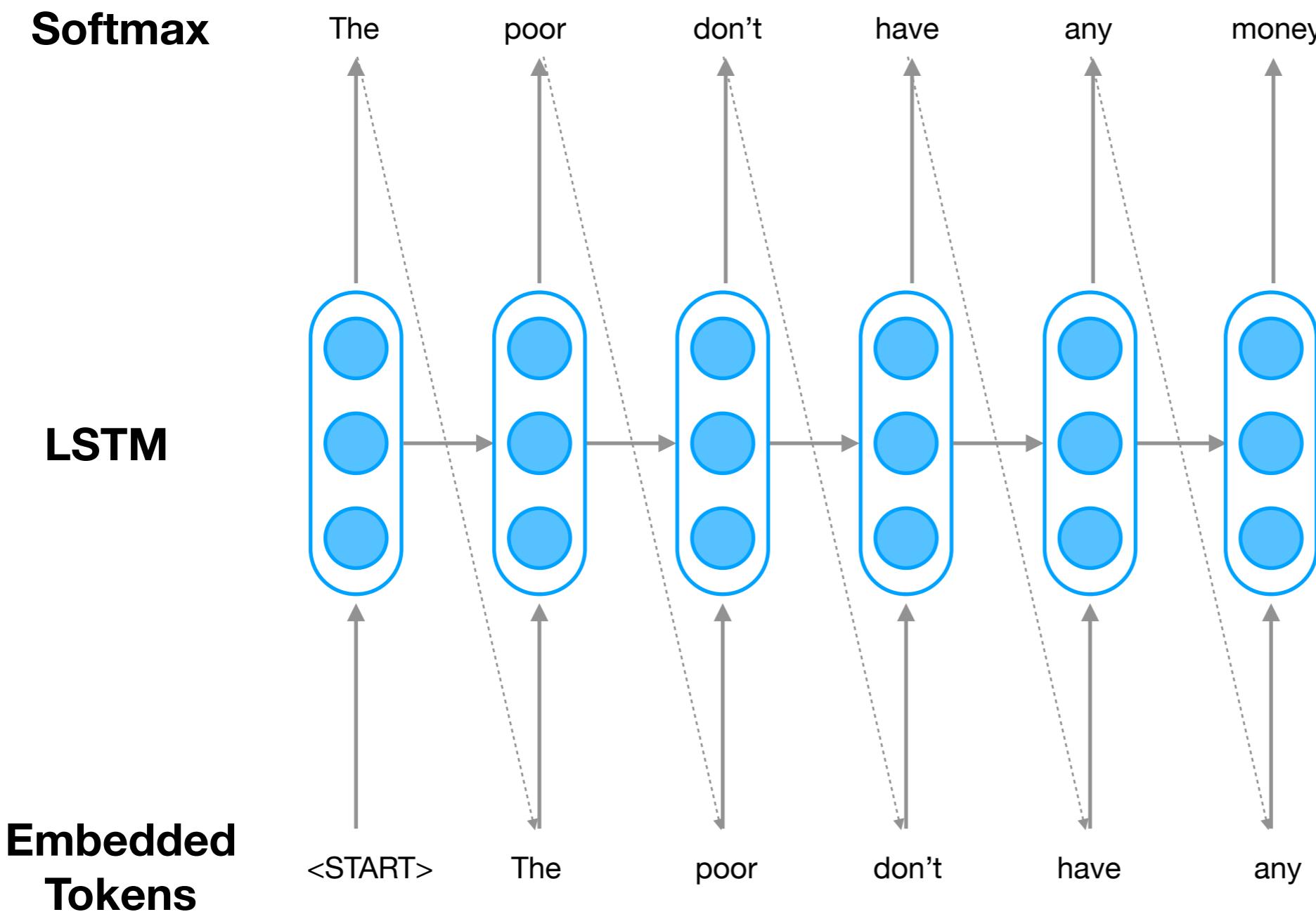
# Language Model

## Inference



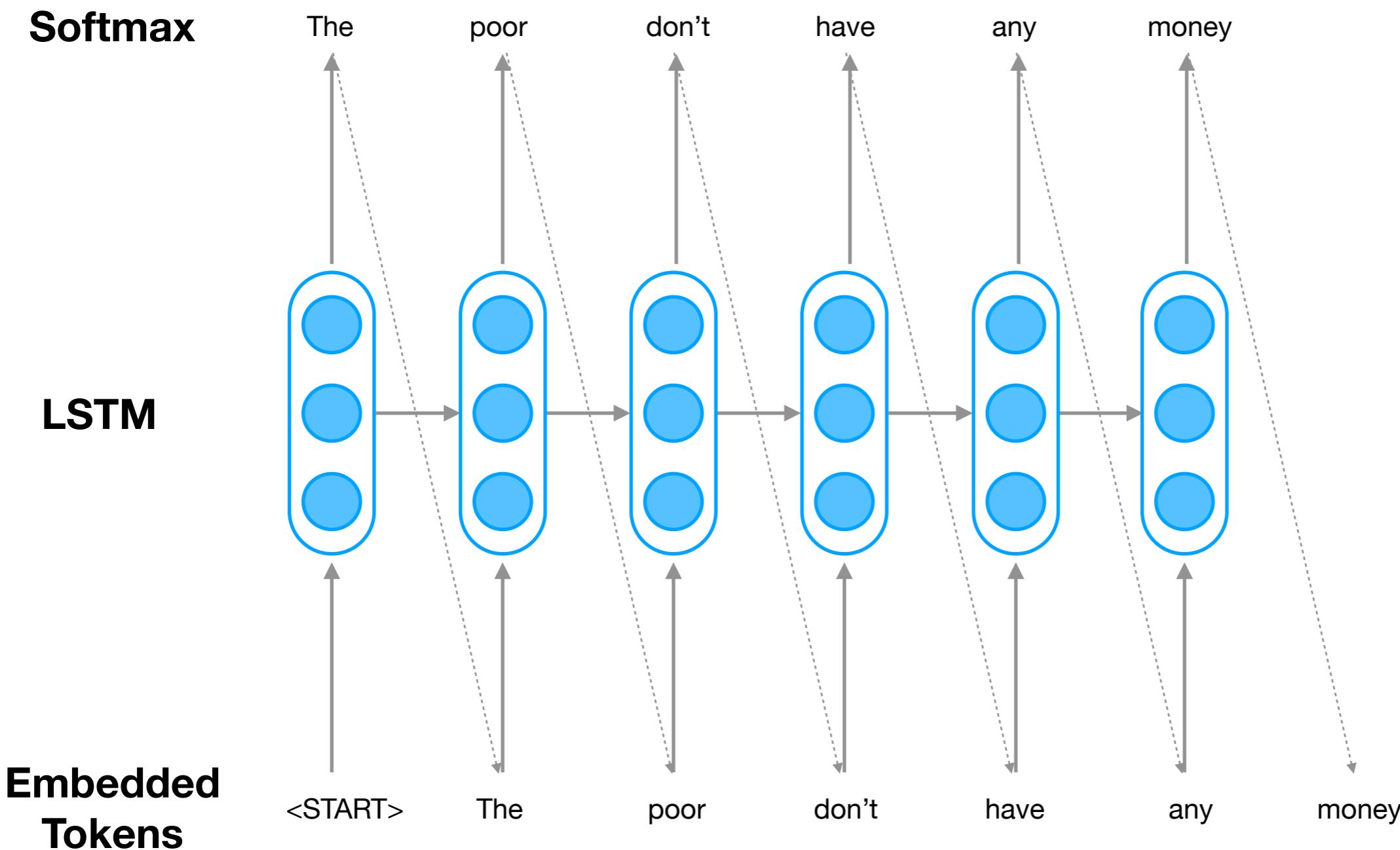
# Language Model

## Inference



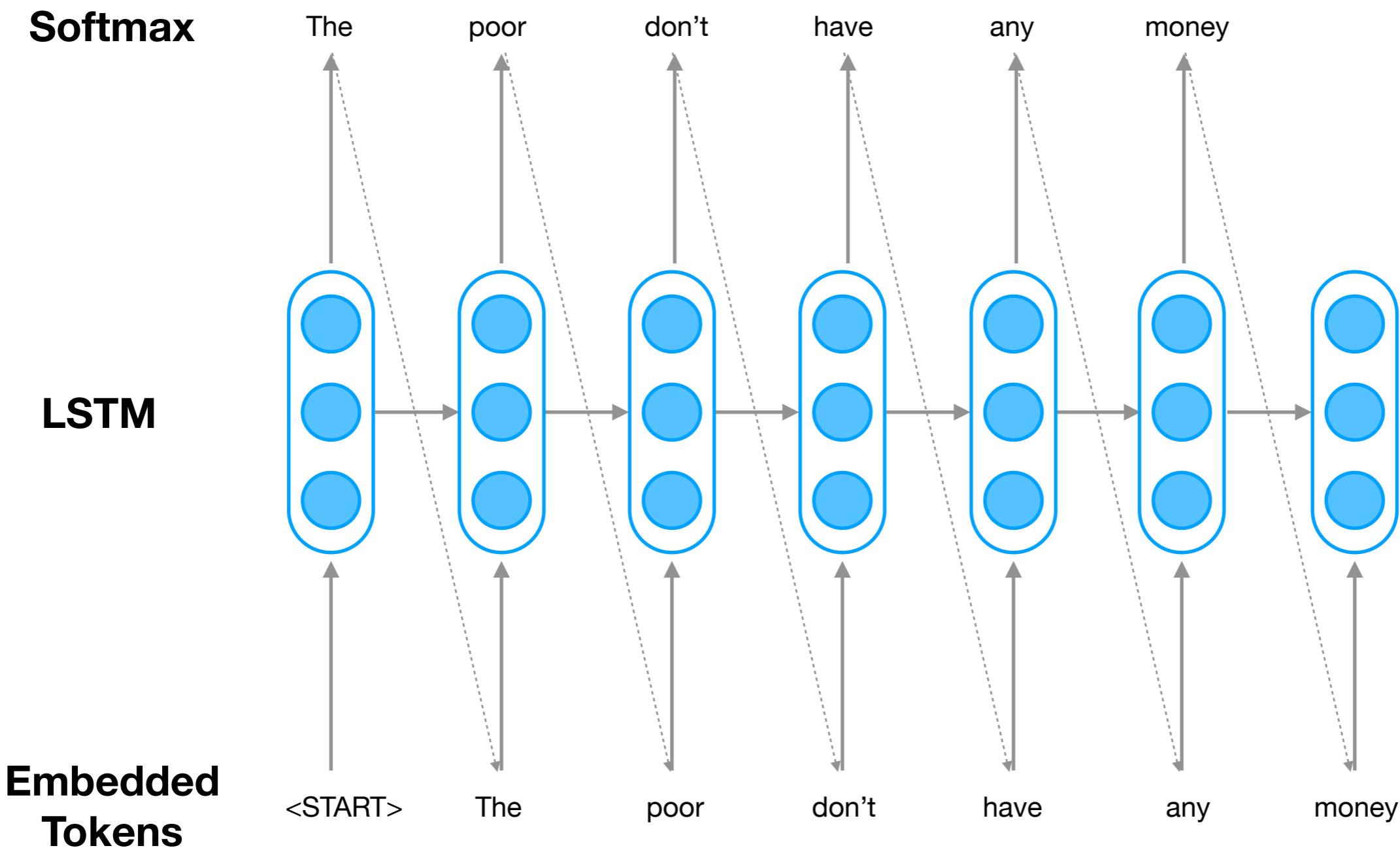
# Language Model

## Inference



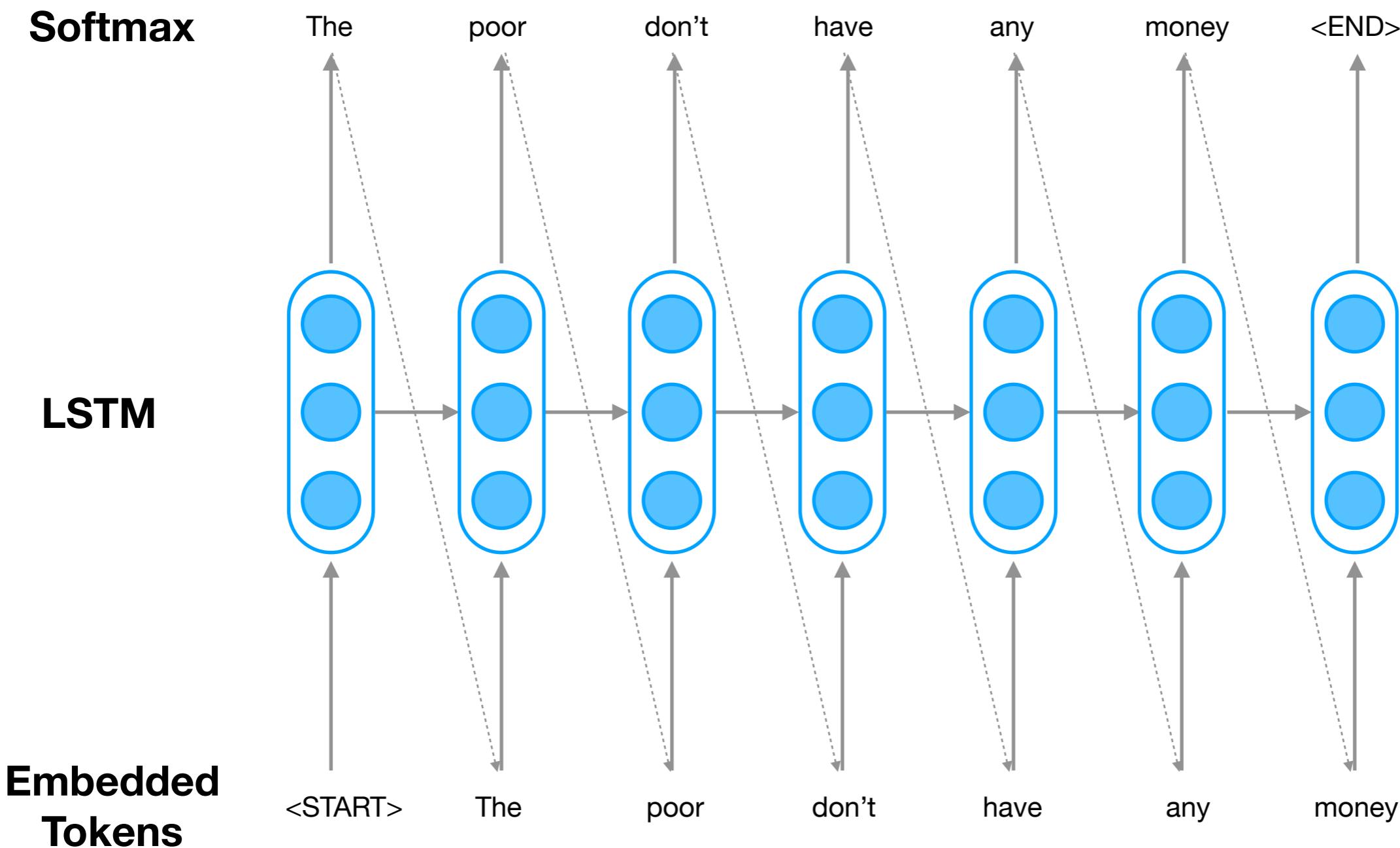
# Language Model

## Inference



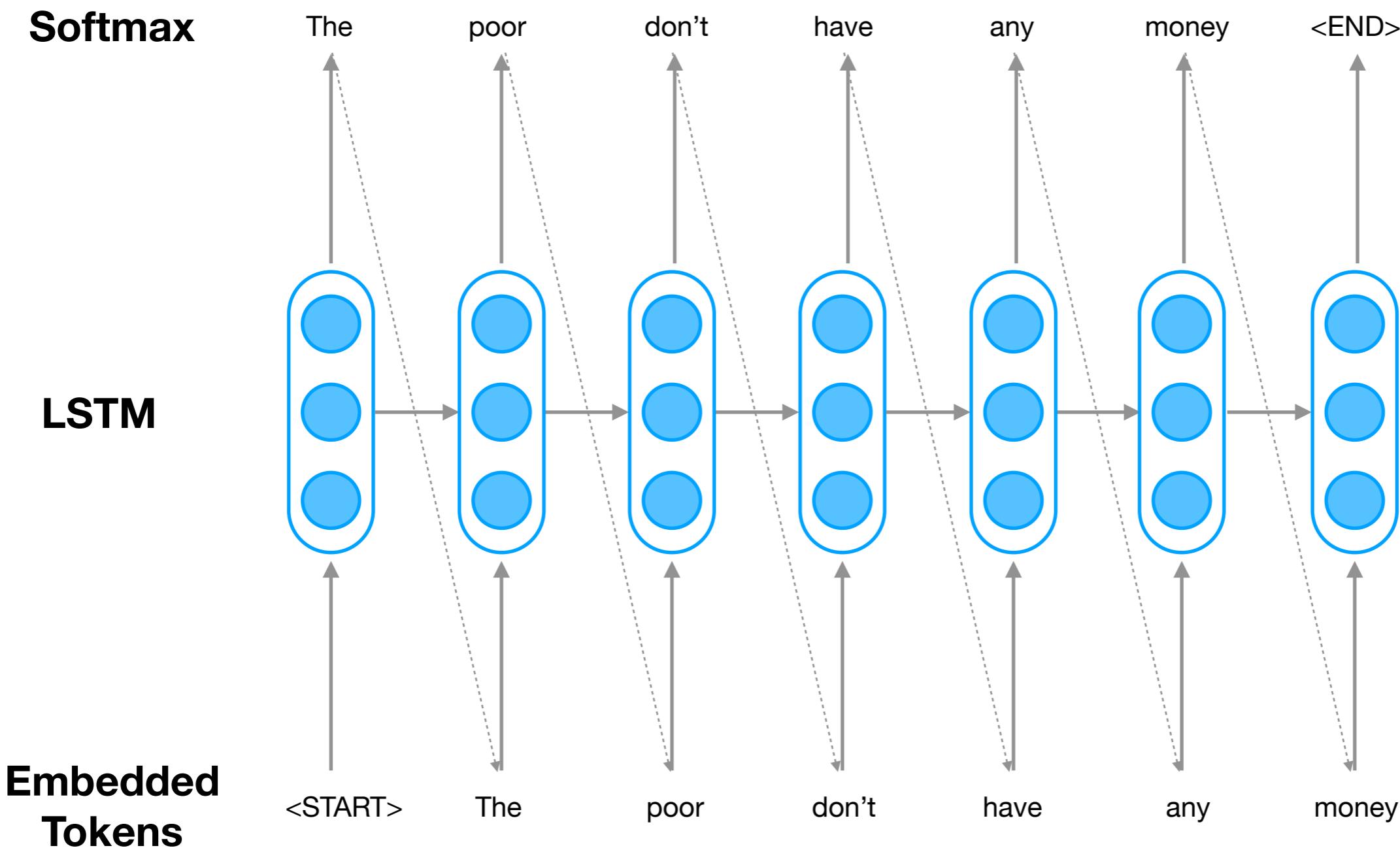
# Language Model

## Inference



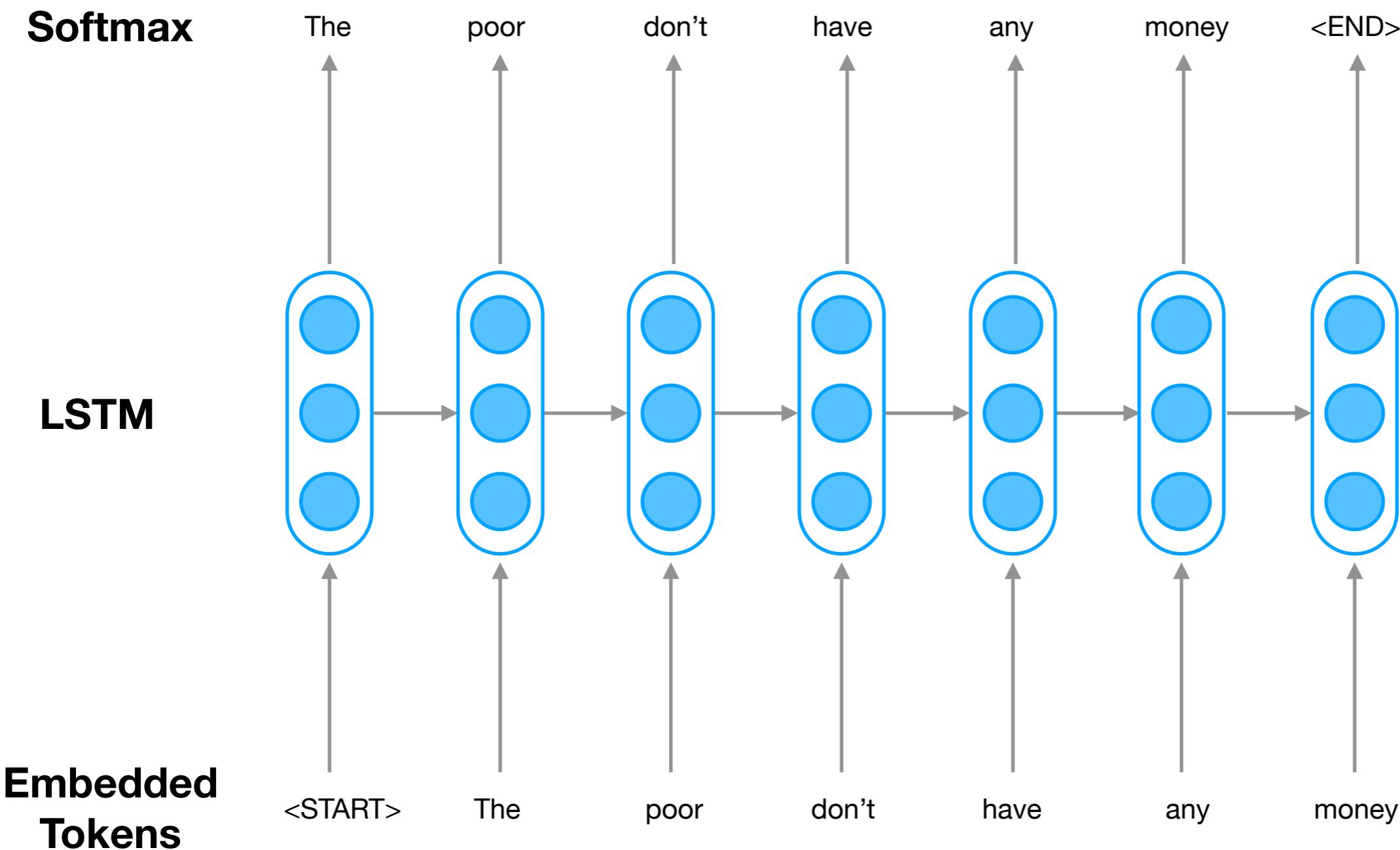
# Language Model

$$p(\mathbf{x}) = \prod_i p(x|x_{<i}) = p(x_0)p(x_1|x_0)p(x_2|x_0, x_1)\dots$$

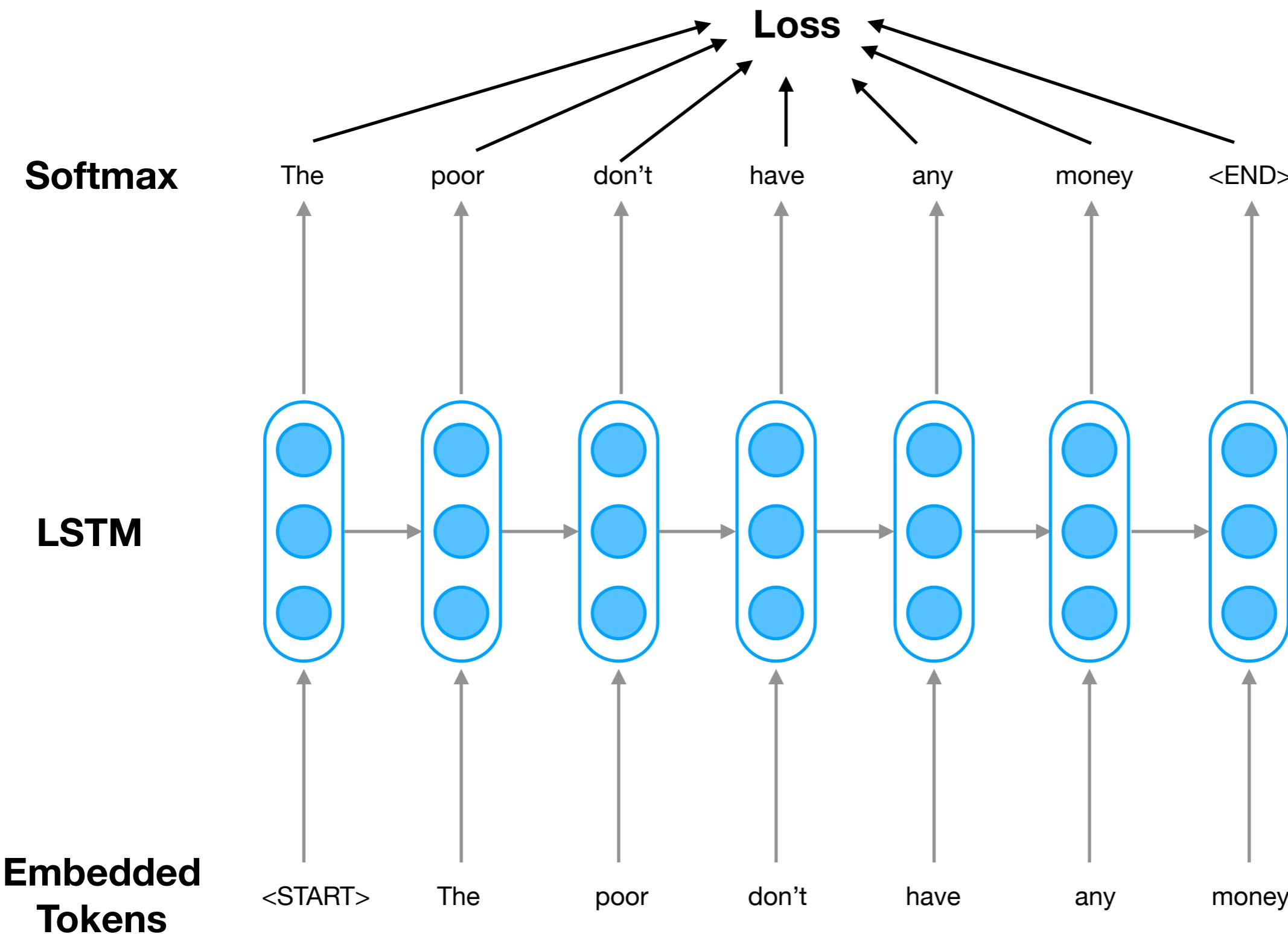


# Language Model

## Training

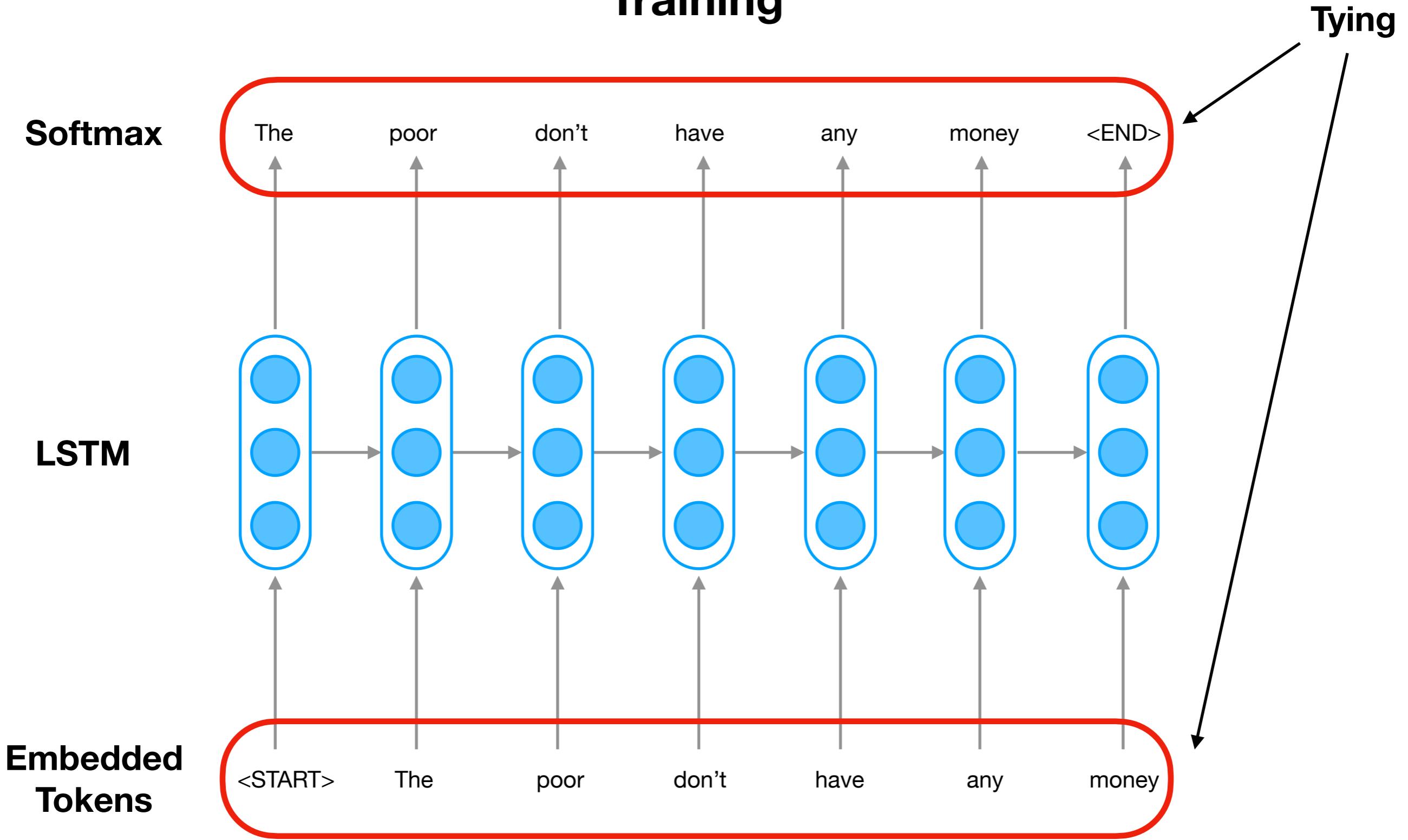


# Language Model

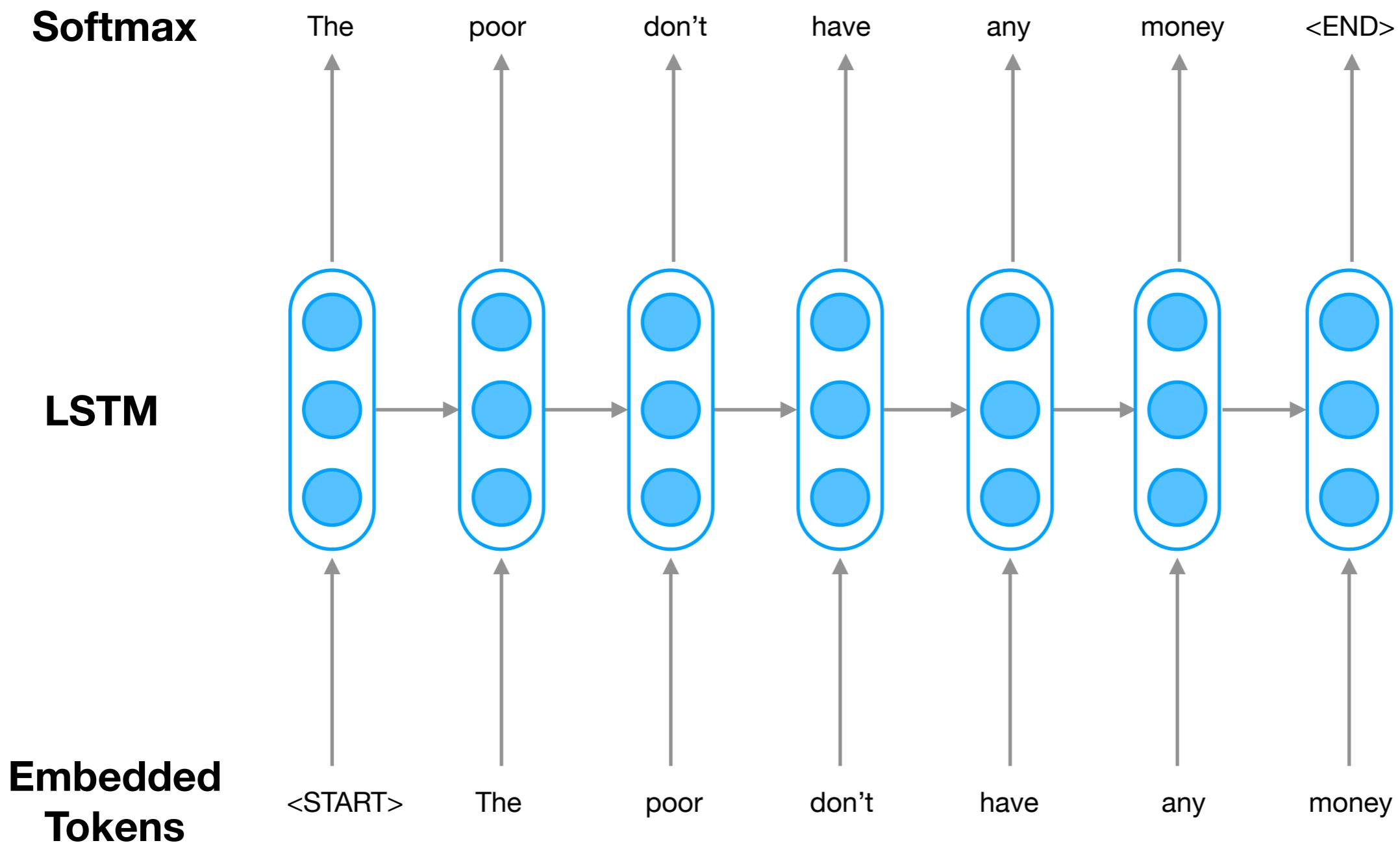


# Language Model

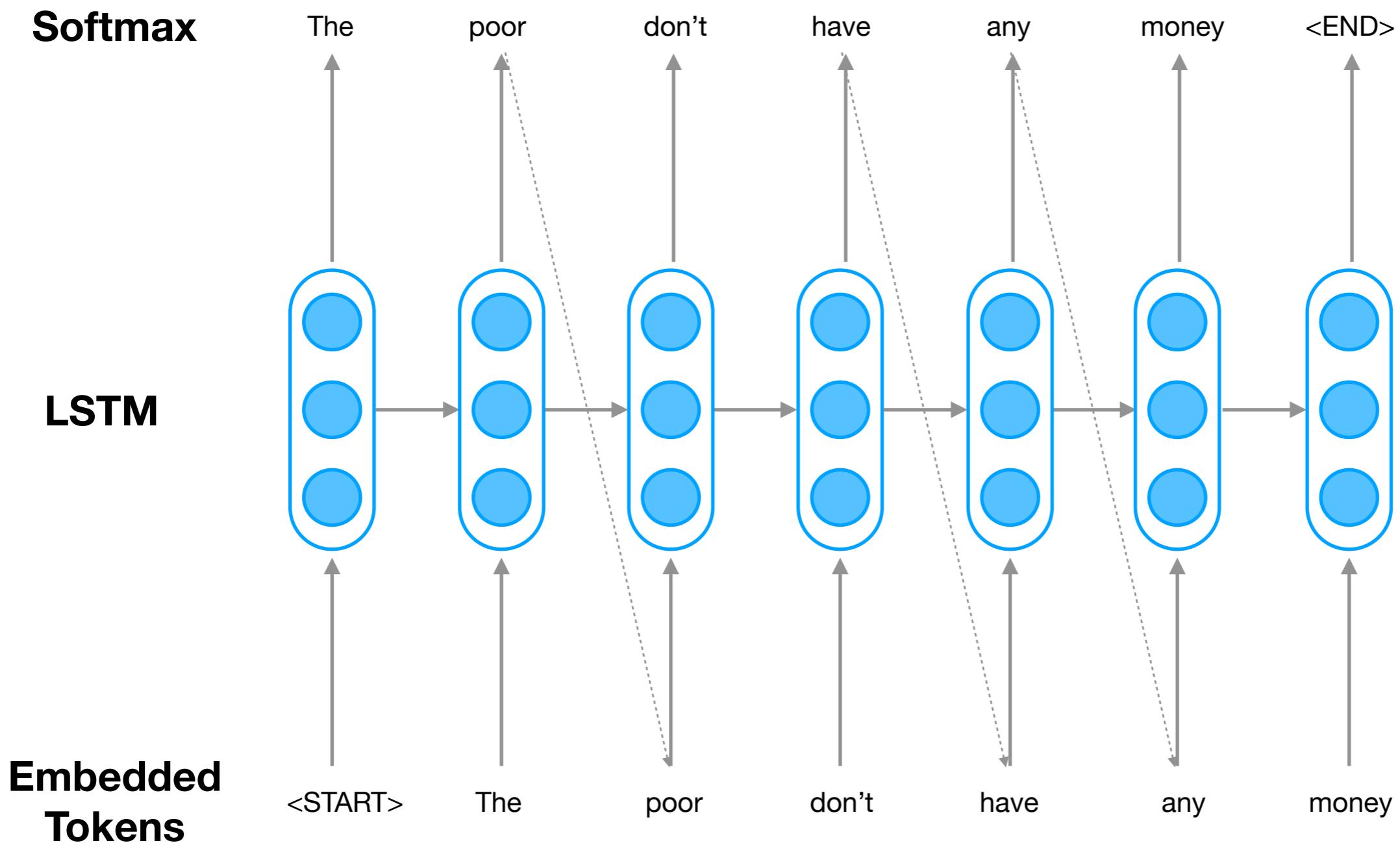
Training



# Teacher Forcing



# Teacher Forcing



# Language Model

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

# Language Model

X

<START>

Y

The

# Language Model

X      <START>      The  
Y      poor

# Language Model

**X**

<START>      The      poor

**Y**

don't

# Language Model

X

<START>      The      poor      don't

Y

have

# Language Model

X

<START>      The      poor      don't      have

Y

any

# Language Model

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

# Language Model

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

Y <END>

# Language Model

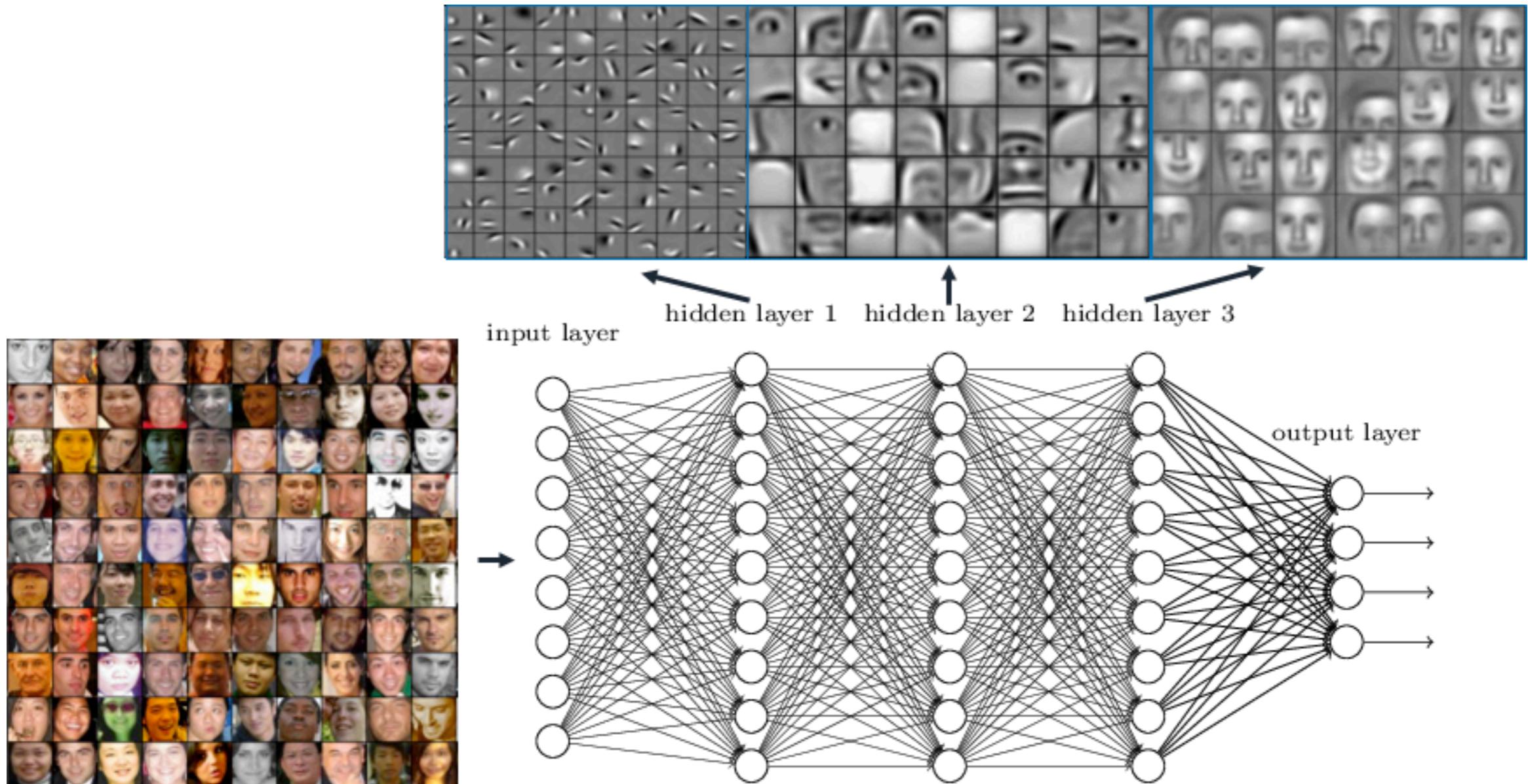
X <START> The poor don't have any money

Y <END>

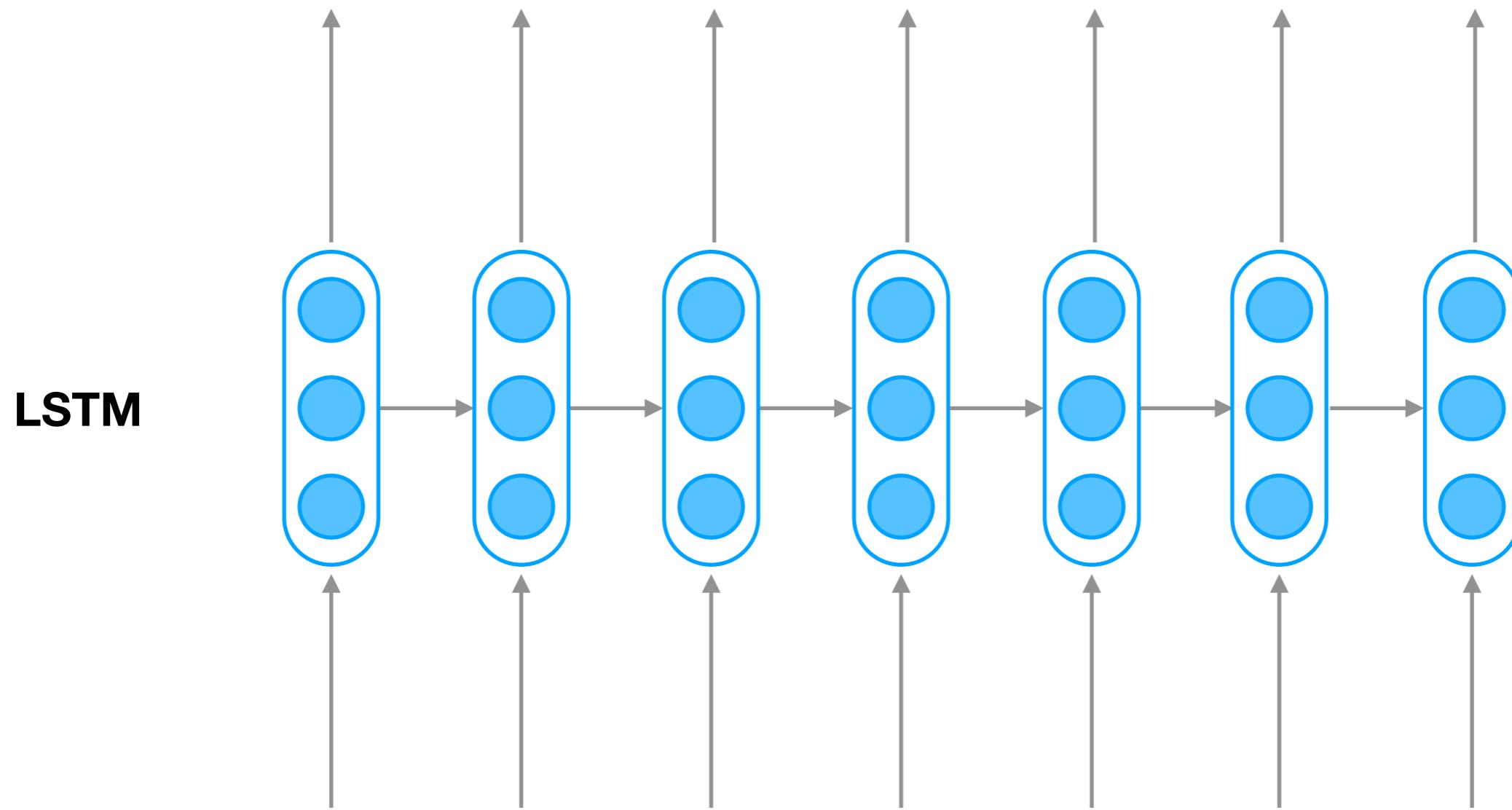
**Sum of Losses**

# Transfer Learning

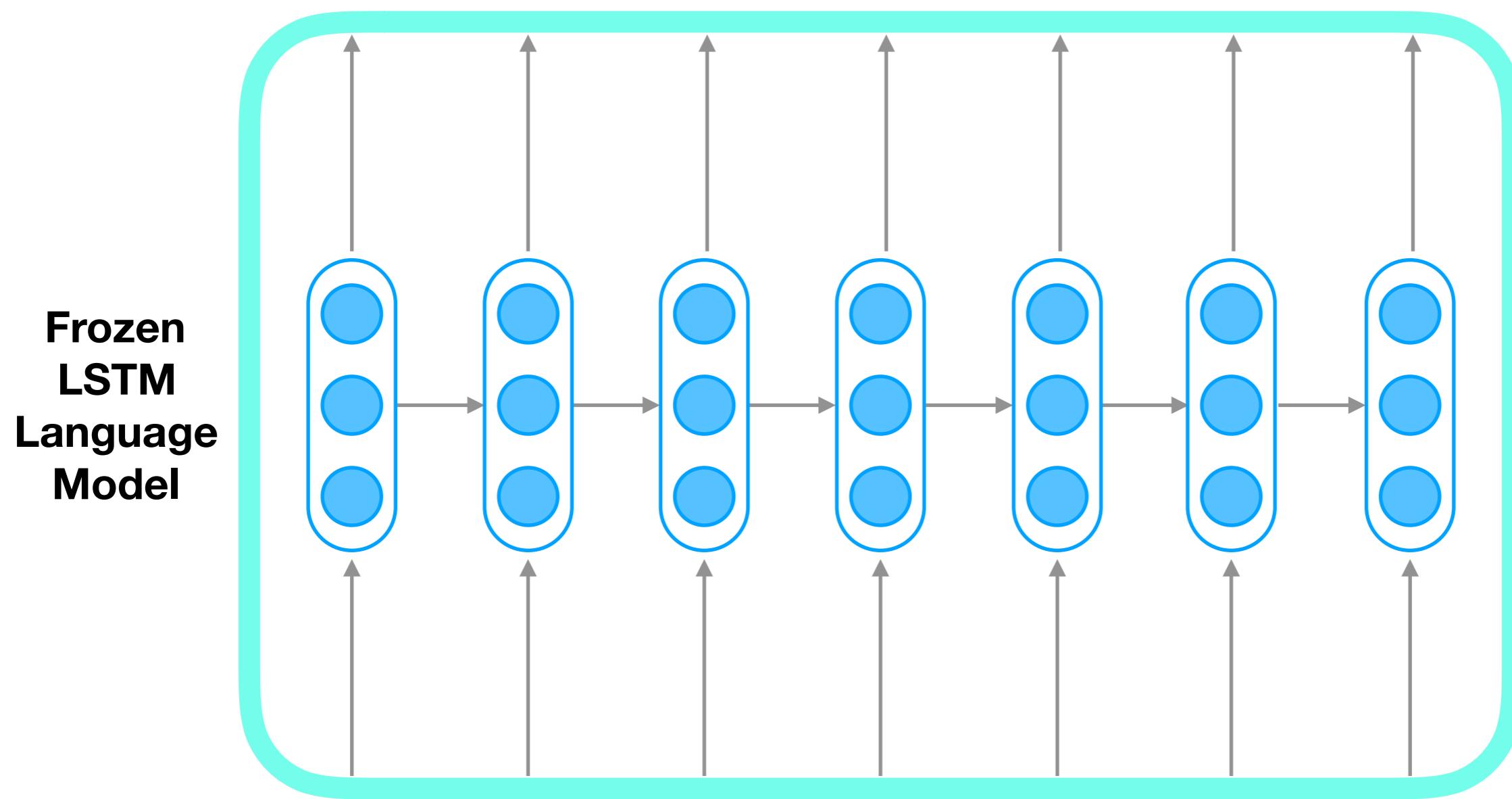
## CNN Intuition



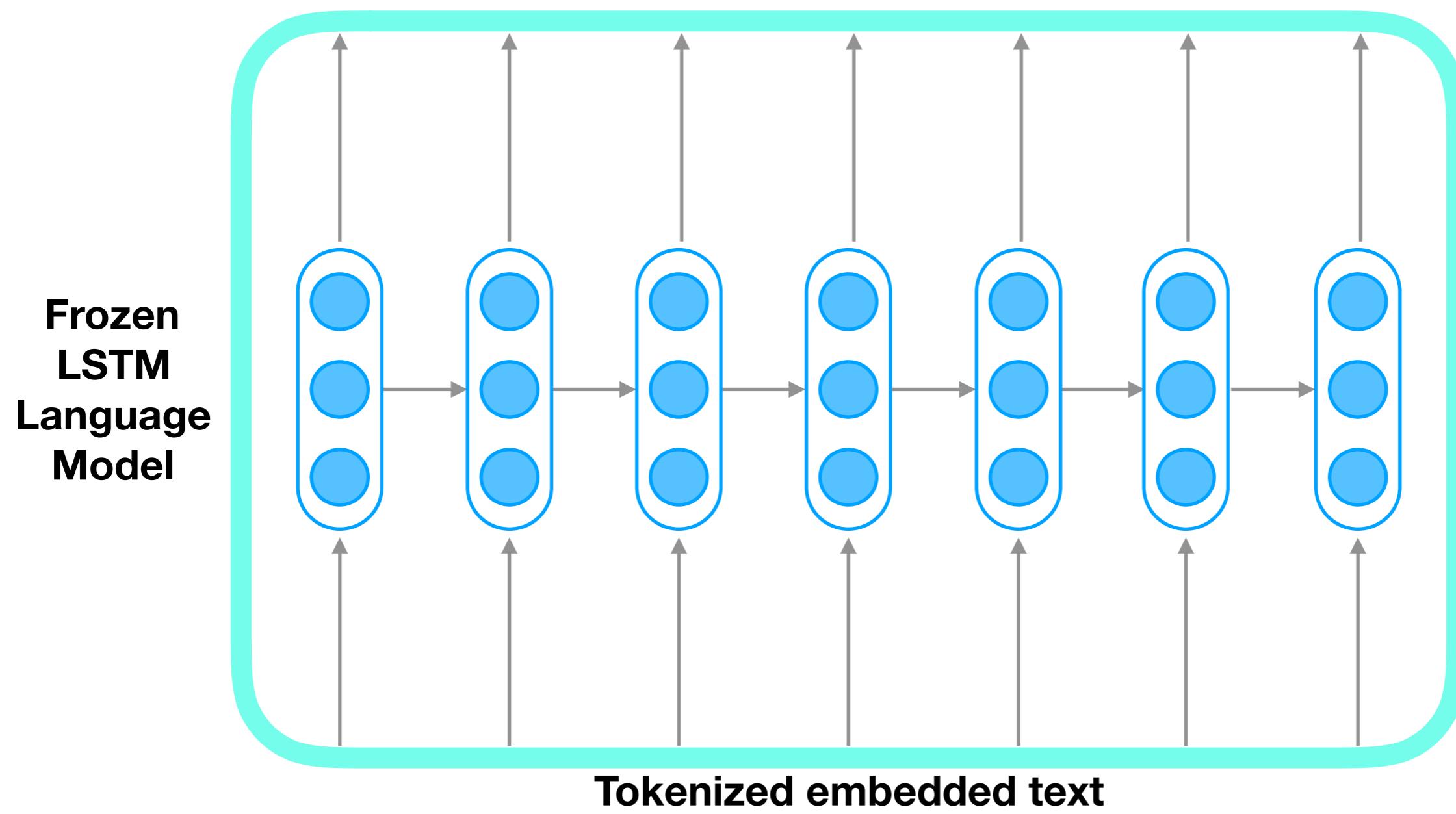
# NLP LM Transfer Learning



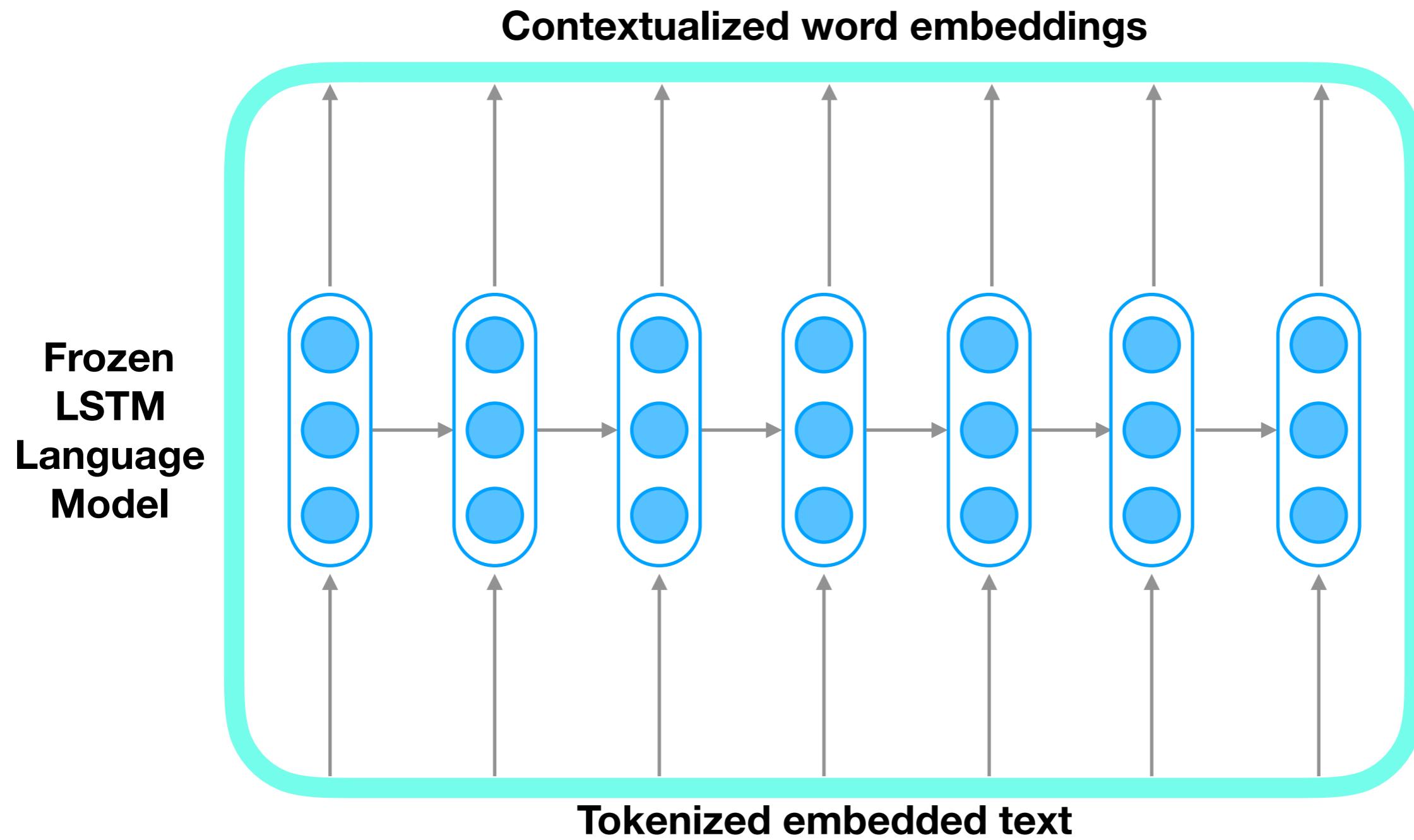
# NLP LM Transfer Learning



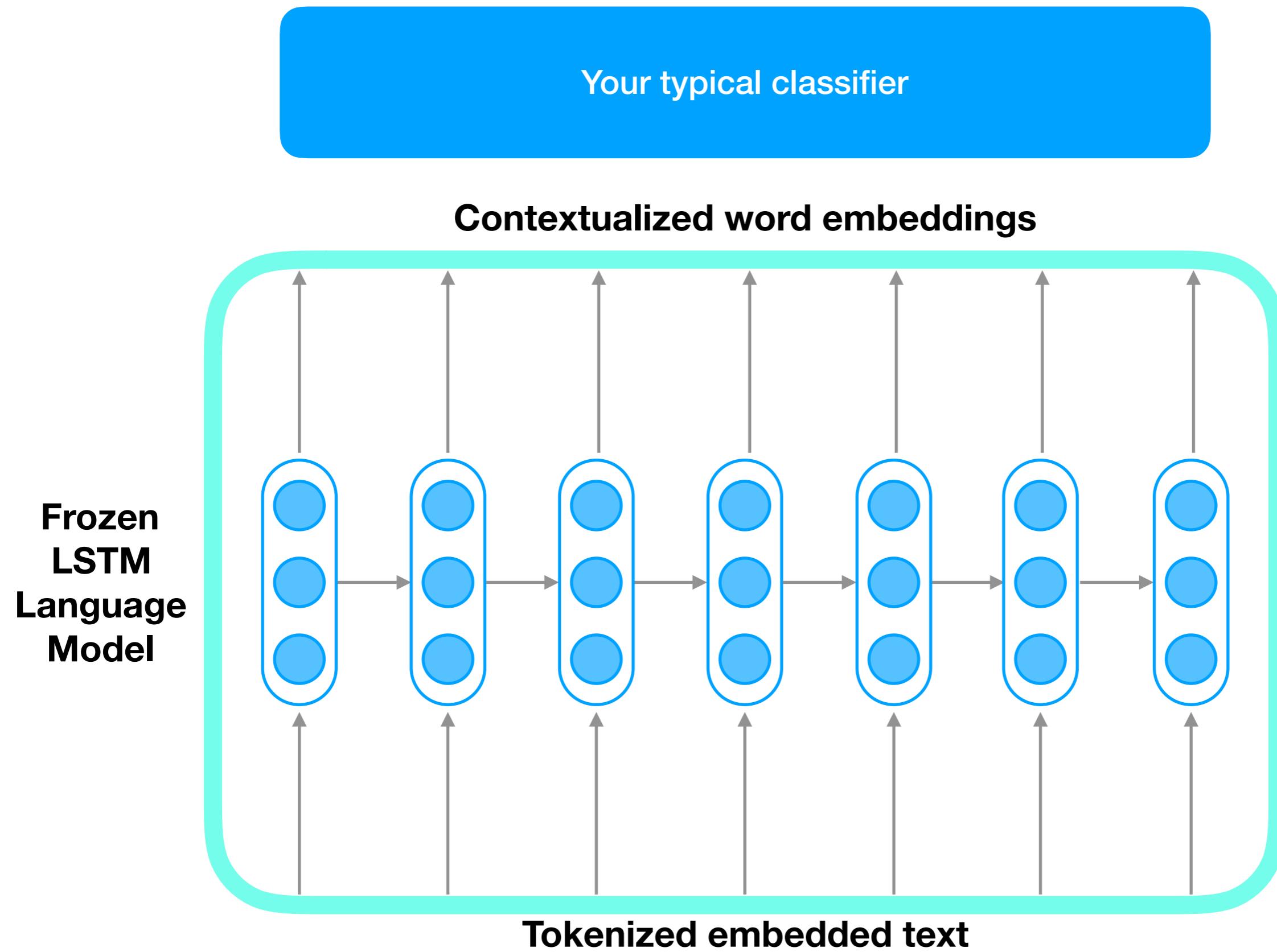
# NLP LM Transfer Learning



# NLP LM Transfer Learning



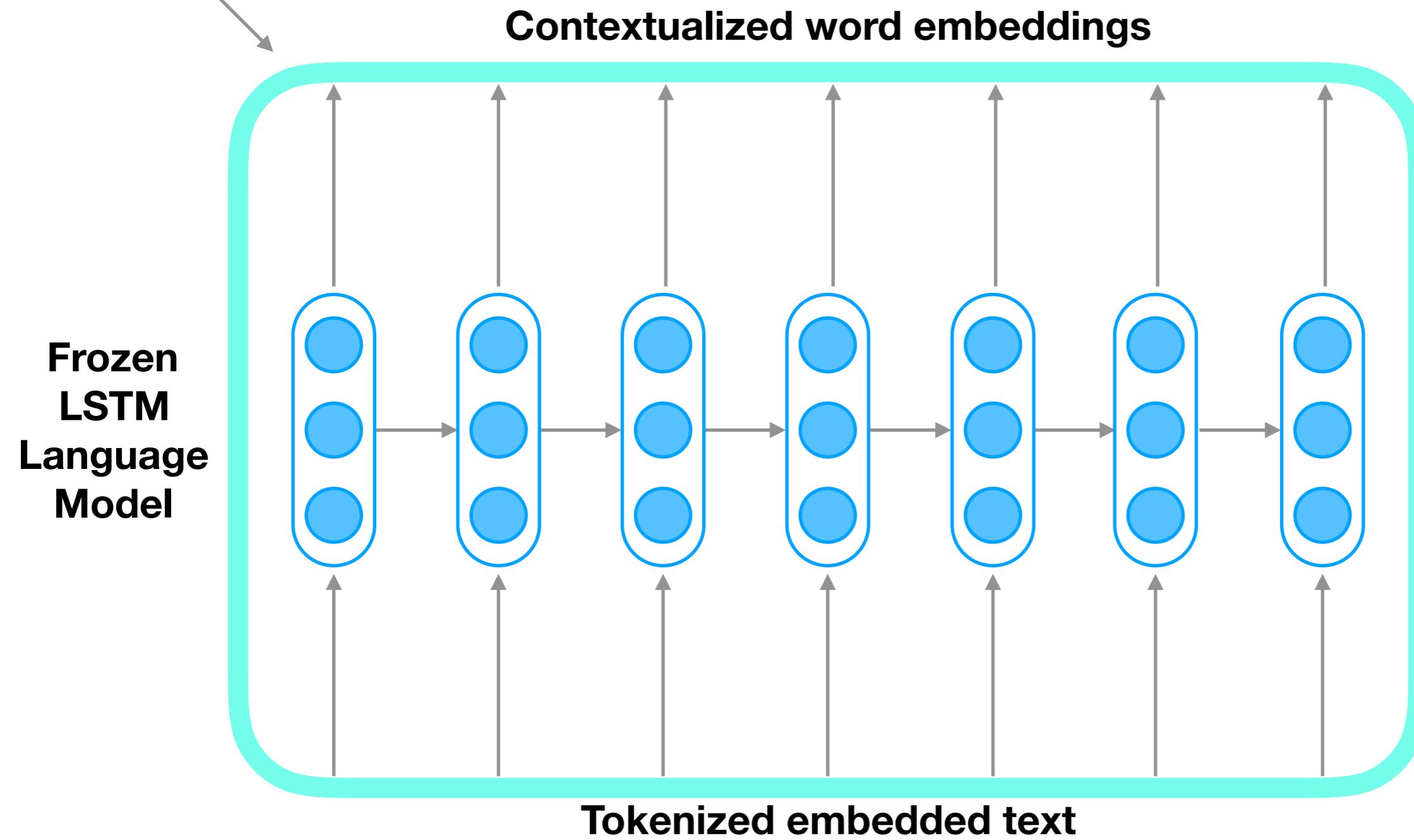
# NLP LM Transfer Learning



# NLP LM Transfer Learning

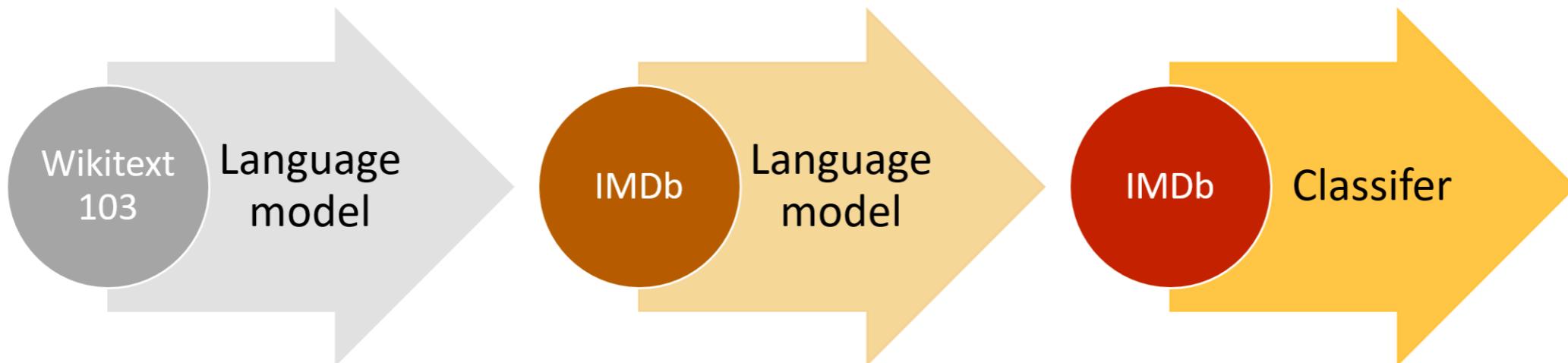
Featuring  
Model

Your typical classifier

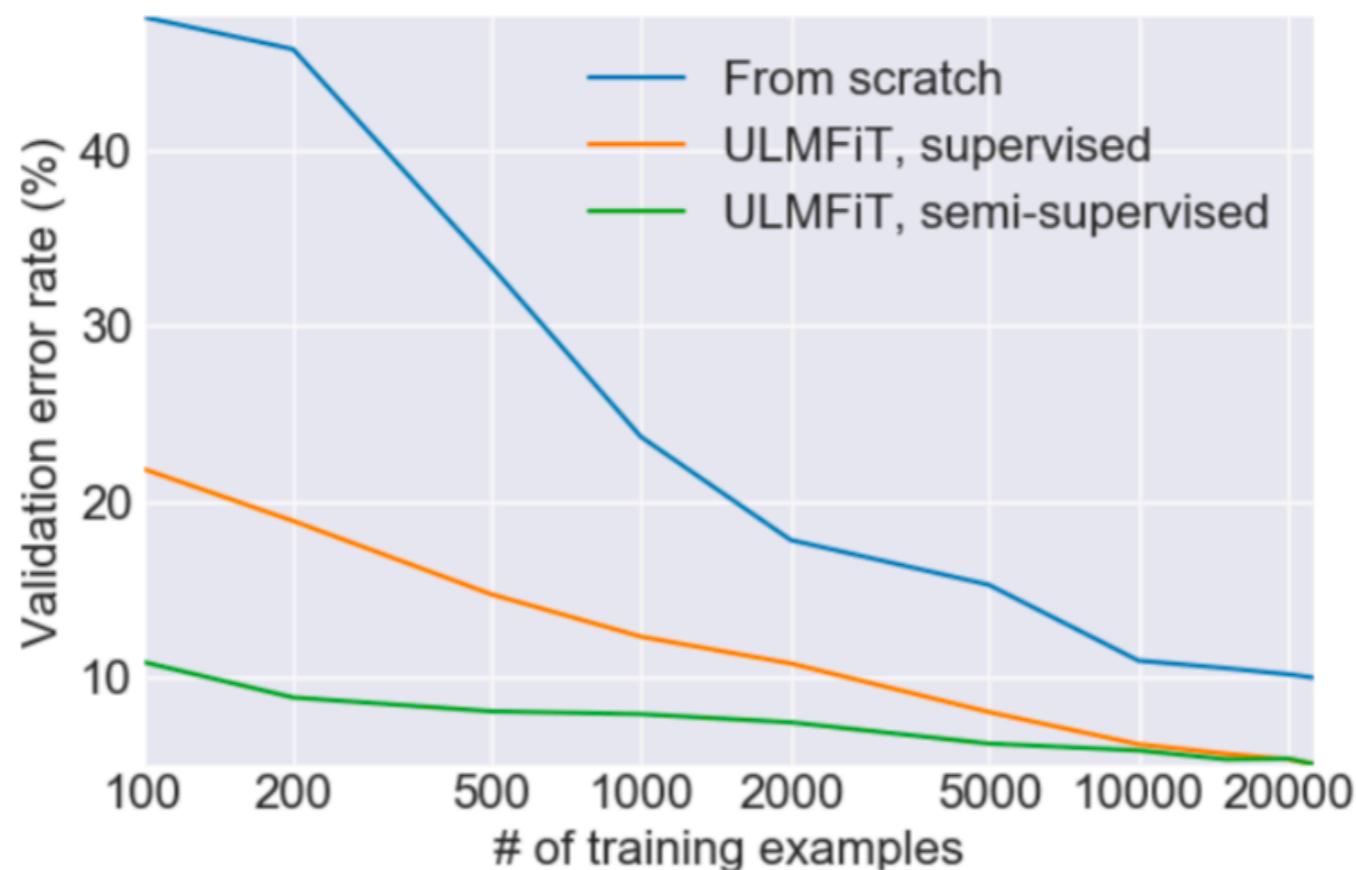
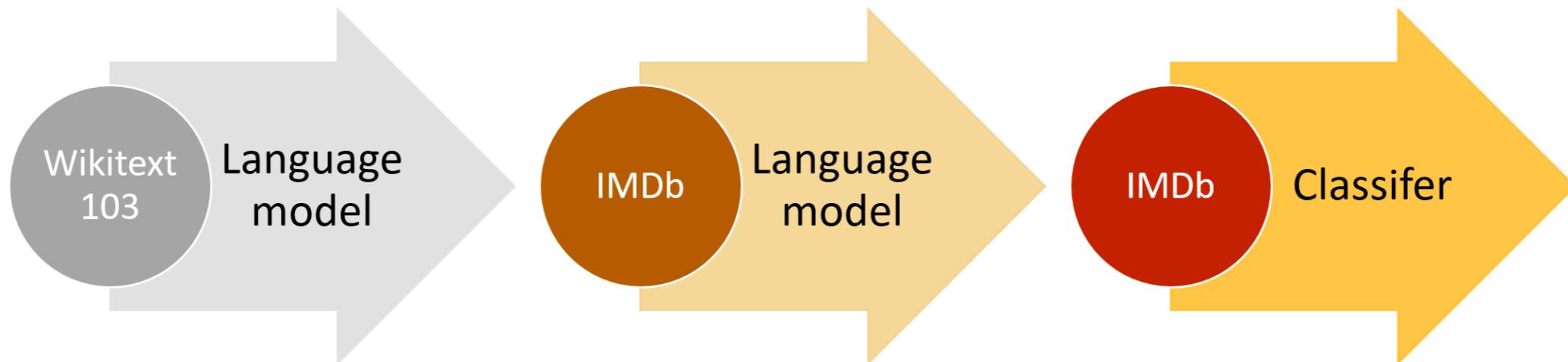


**ULMFIT**

# ULMFiT



# ULMFiT



# ULMFiT

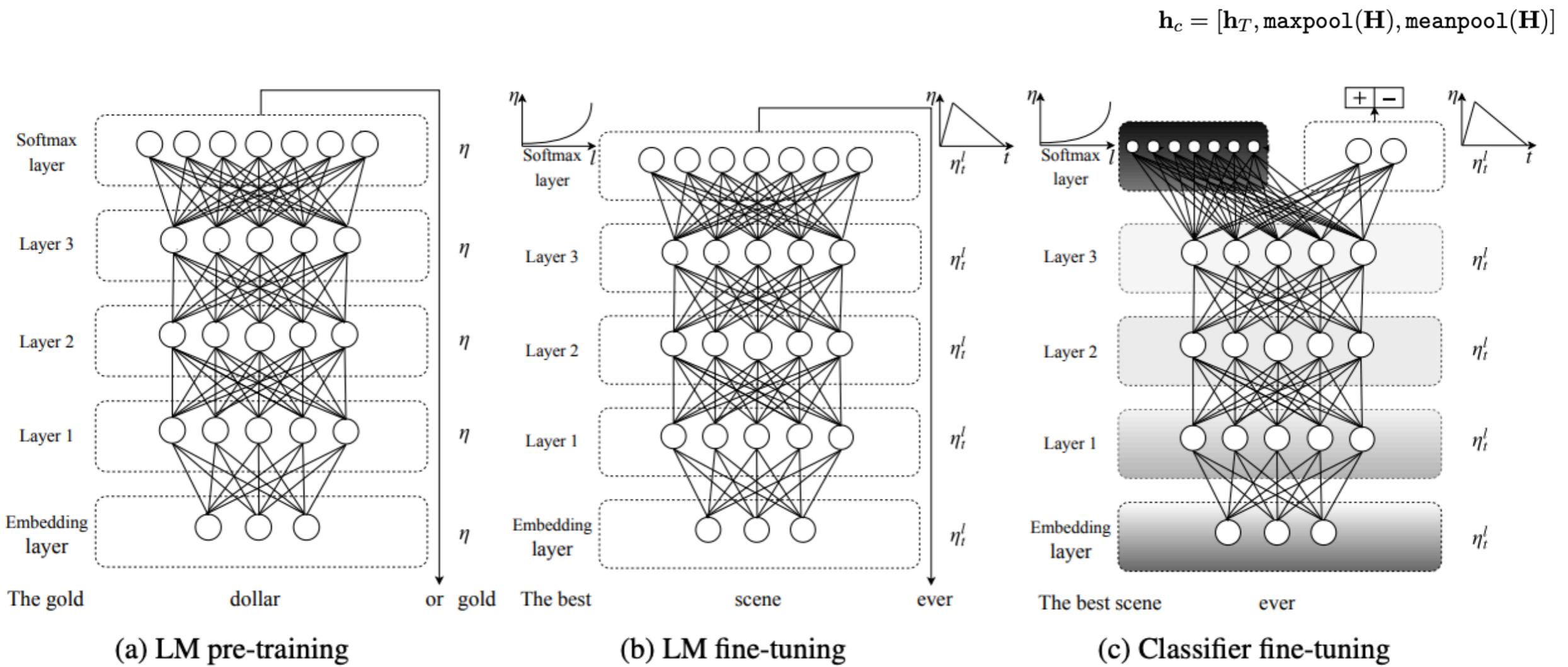


Figure 1: ULMFiT consists of three stages: a) The LM is trained on a general-domain corpus to capture general features of the language in different layers. b) The full LM is fine-tuned on target task data using discriminative fine-tuning ('*Discr*') and slanted triangular learning rates (STLR) to learn task-specific features. c) The classifier is fine-tuned on the target task using gradual unfreezing, '*Discr*', and STLR to preserve low-level representations and adapt high-level ones (shaded: unfreezing stages; black: frozen).

# Scheduling

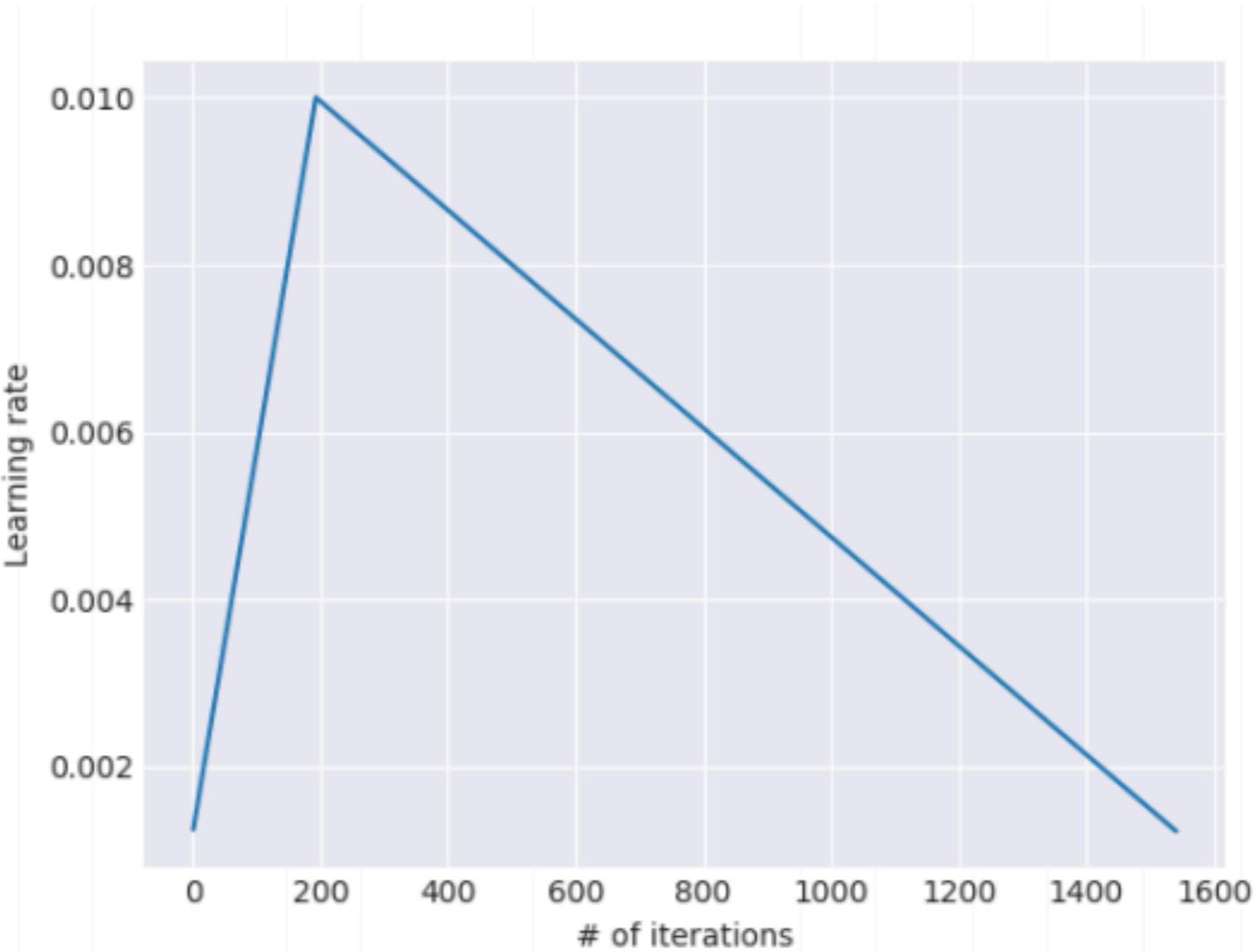
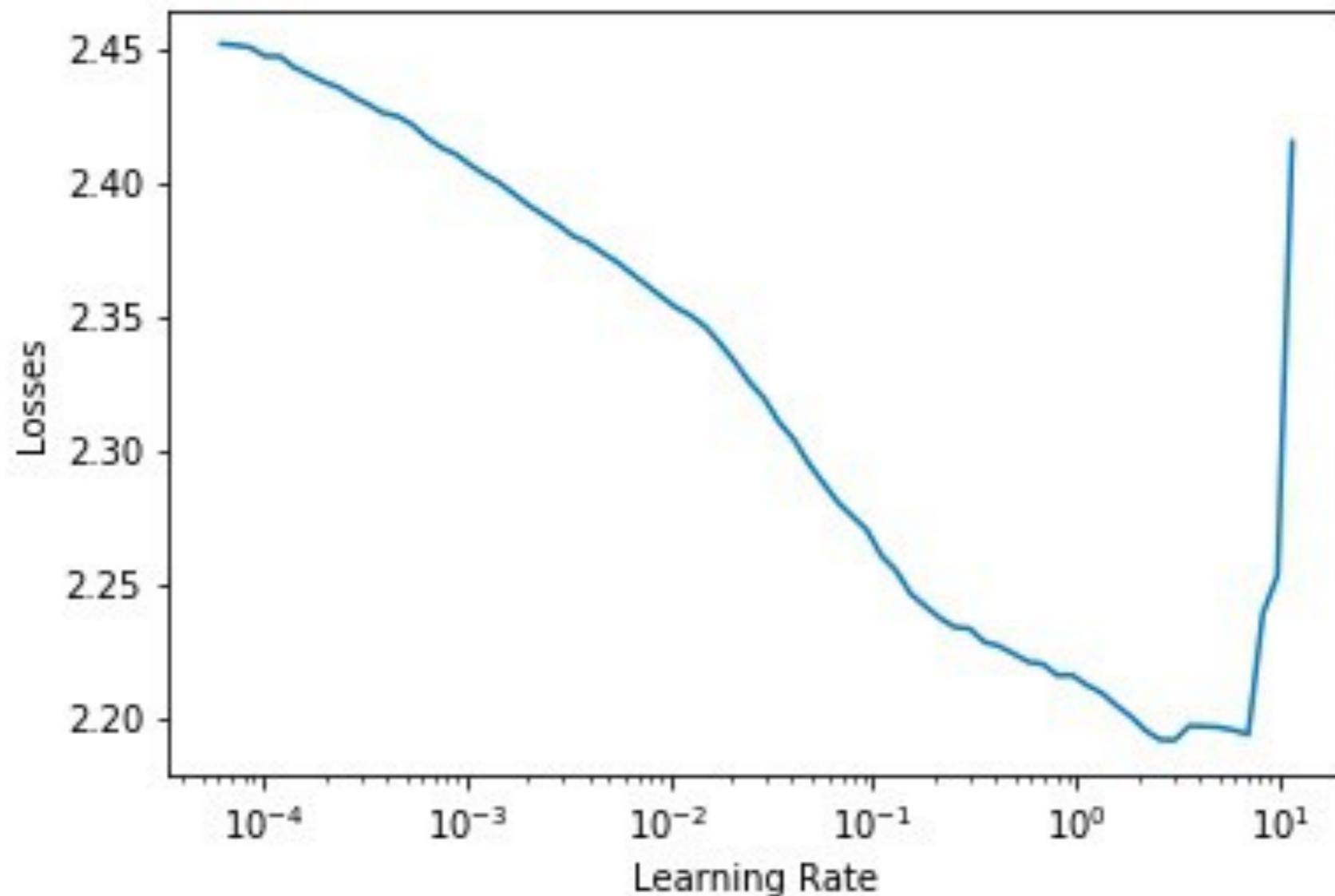


Figure 2: The slanted triangular learning rate schedule used for ULMFiT as a function of the number of training iterations.

# LR Finder



# ULMFiT

## Gradual Unfreezing (Catastrophic Forgetting)

$$h_c = [h_T, \text{maxpool}(\mathbf{H}), \text{meanpool}(\mathbf{H})]$$

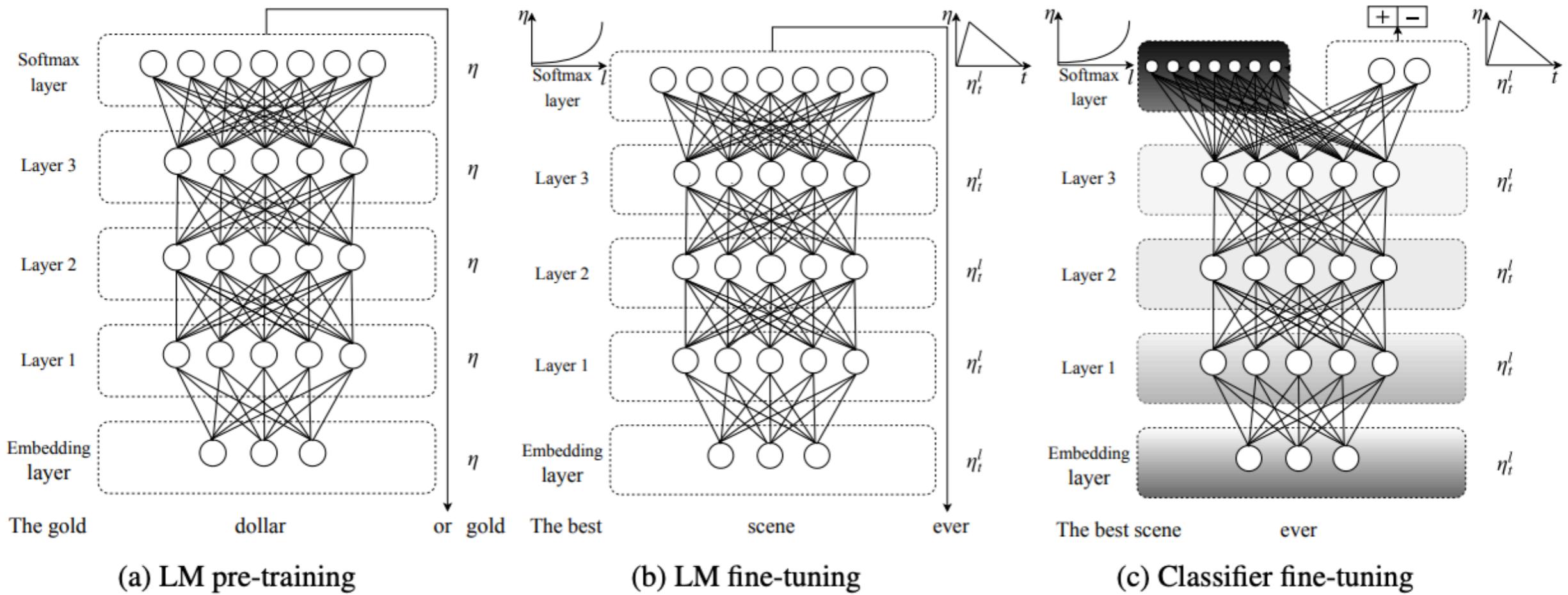


Figure 1: ULMFiT consists of three stages: a) The LM is trained on a general-domain corpus to capture general features of the language in different layers. b) The full LM is fine-tuned on target task data using discriminative fine-tuning ('*Discr*') and slanted triangular learning rates (STLR) to learn task-specific features. c) The classifier is fine-tuned on the target task using gradual unfreezing, '*Discr*', and STLR to preserve low-level representations and adapt high-level ones (shaded: unfreezing stages; black: frozen).

# ULMFiT

## Classifier

$$\mathbf{h}_c = [\mathbf{h}_T, \text{maxpool}(\mathbf{H}), \text{meanpool}(\mathbf{H})]$$

```
(1): PoolingLinearClassifier(  
    (layers): Sequential(  
        (0): BatchNorm1d(900, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (1): Dropout(p=0.4)  
        (2): Linear(in_features=900, out_features=50, bias=True)  
        (3): ReLU(inplace)  
        (4): BatchNorm1d(50, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (5): Dropout(p=0.1)  
        (6): Linear(in_features=50, out_features=2, bias=True)  
    )  
)
```

# AWD-LSTM

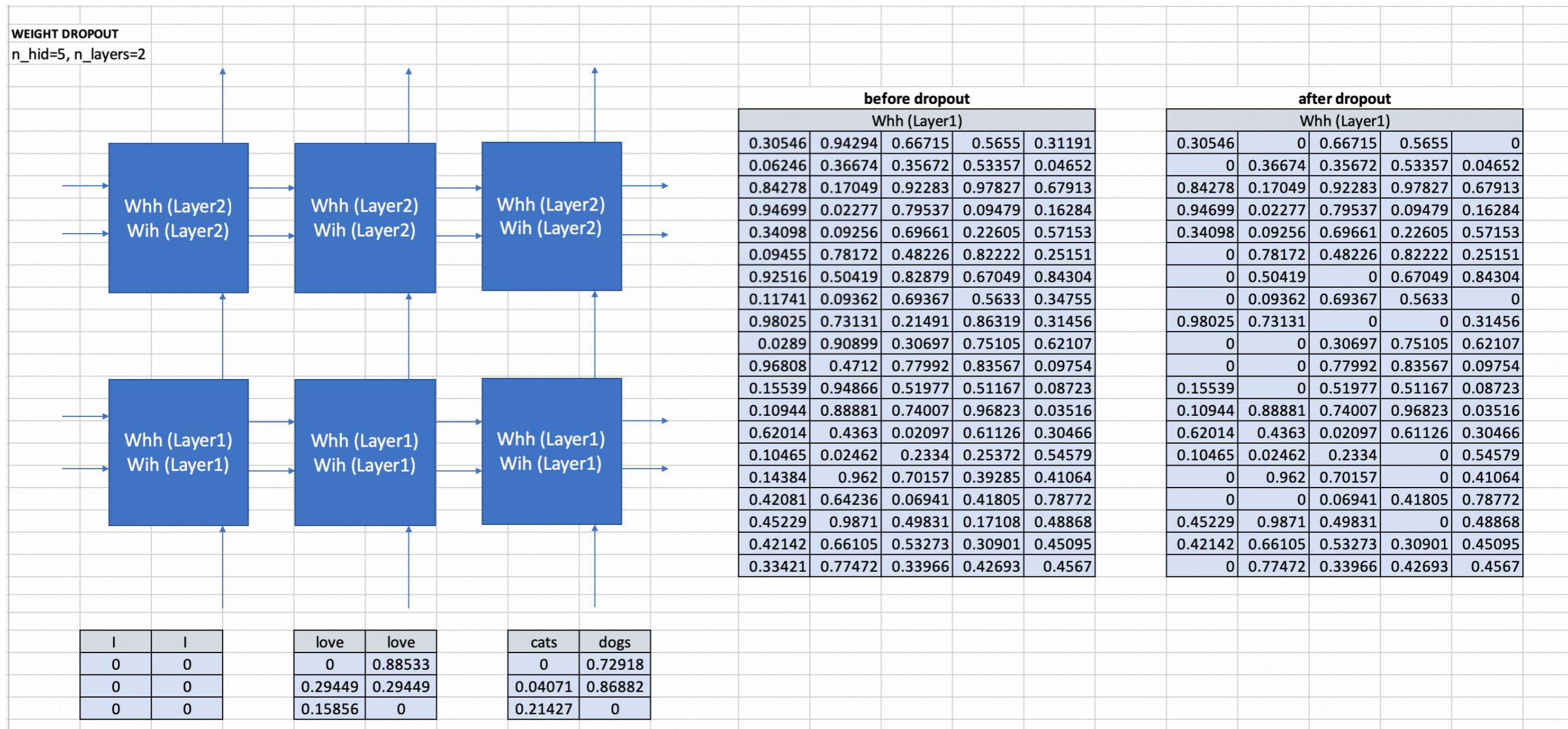
ENCODER DROPOUT					before dropout				after dropout			
	token	d1	d2	...		token	d1	d2	...			
	I	0.399	0.75379	0.62616		I	0	0	0			
	love	0.88533	0.29449	0.15856		love	0.88533	0.29449	0.15856			
	cats	0.48927	0.04071	0.21427		cats	0.48927	0.04071	0.21427			
	dogs	0.72918	0.86882	0.77136		dogs	0.72918	0.86882	0.77136			

# AWD-LSTM

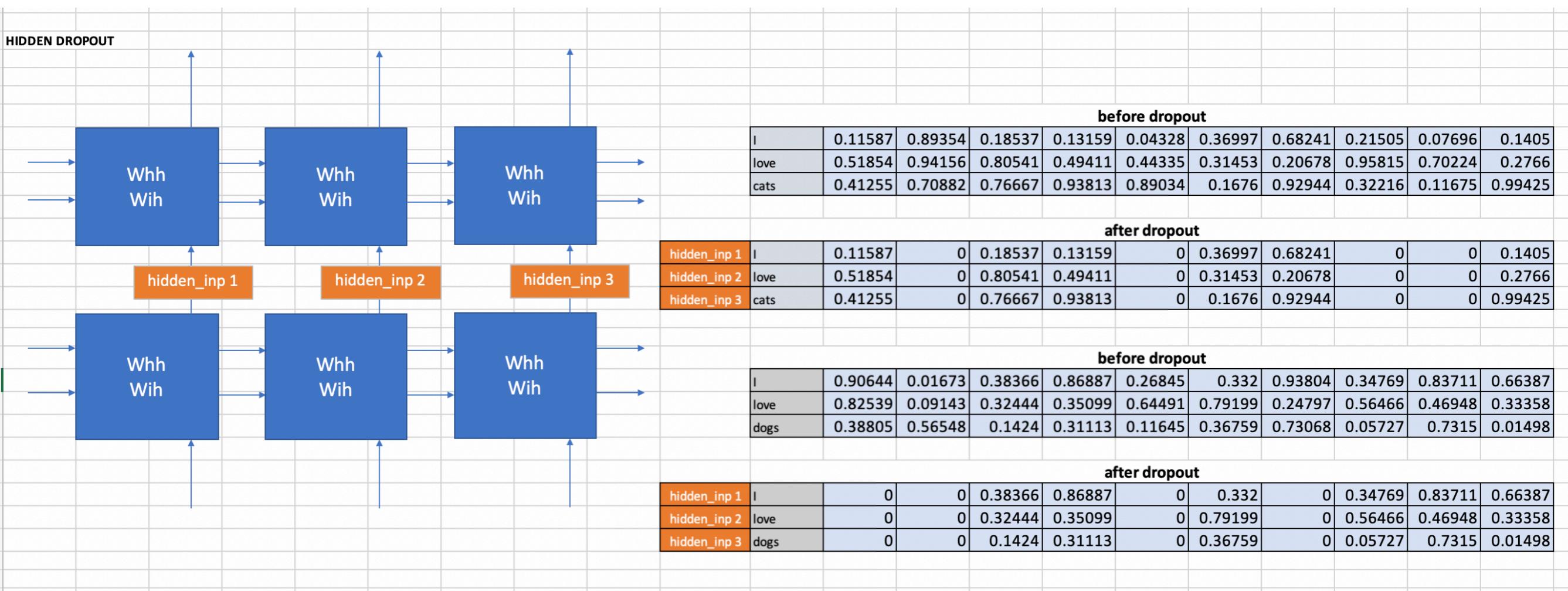
ENCODER DROPOUT				
	before dropout			after dropout
token	d1	d2	...	token
I	0.399	0.75379	0.62616	I
love	0.88533	0.29449	0.15856	love
cats	0.48927	0.04071	0.21427	cats
dogs	0.72918	0.86882	0.77136	dogs

INPUT DROPOUT				
batch: [I love cats, I love dogs]				
	before dropout			after dropout
I	0	0	0	I
love	0.88533	0.29449	0.15856	love
cats	0.48927	0.04071	0.21427	cats
	before dropout			after dropout
I	0	0	0	I
love	0.88533	0.29449	0.15856	love
dogs	0.72918	0.86882	0.77136	dogs

# AWD-LSTM



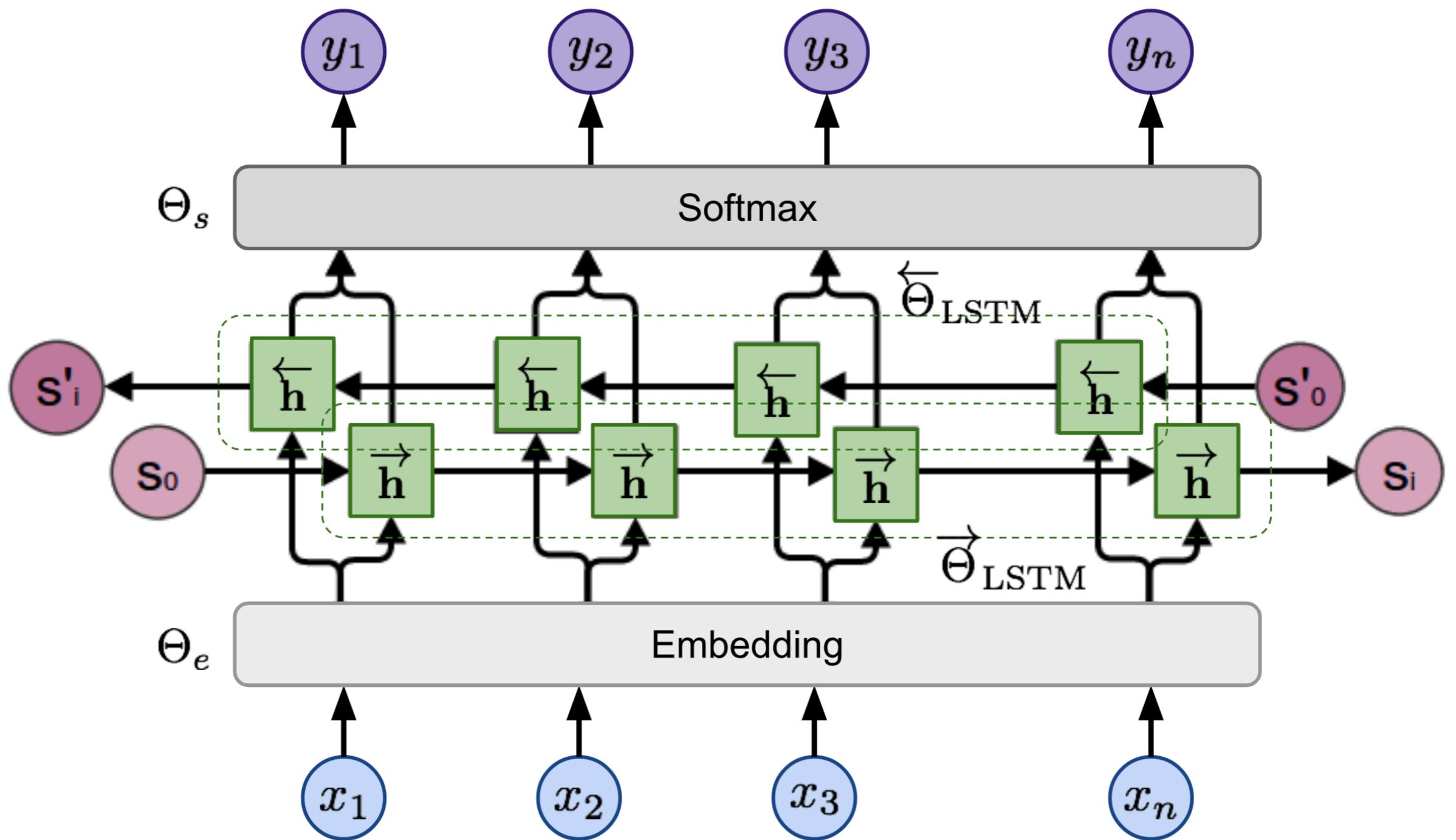
# AWD-LSTM



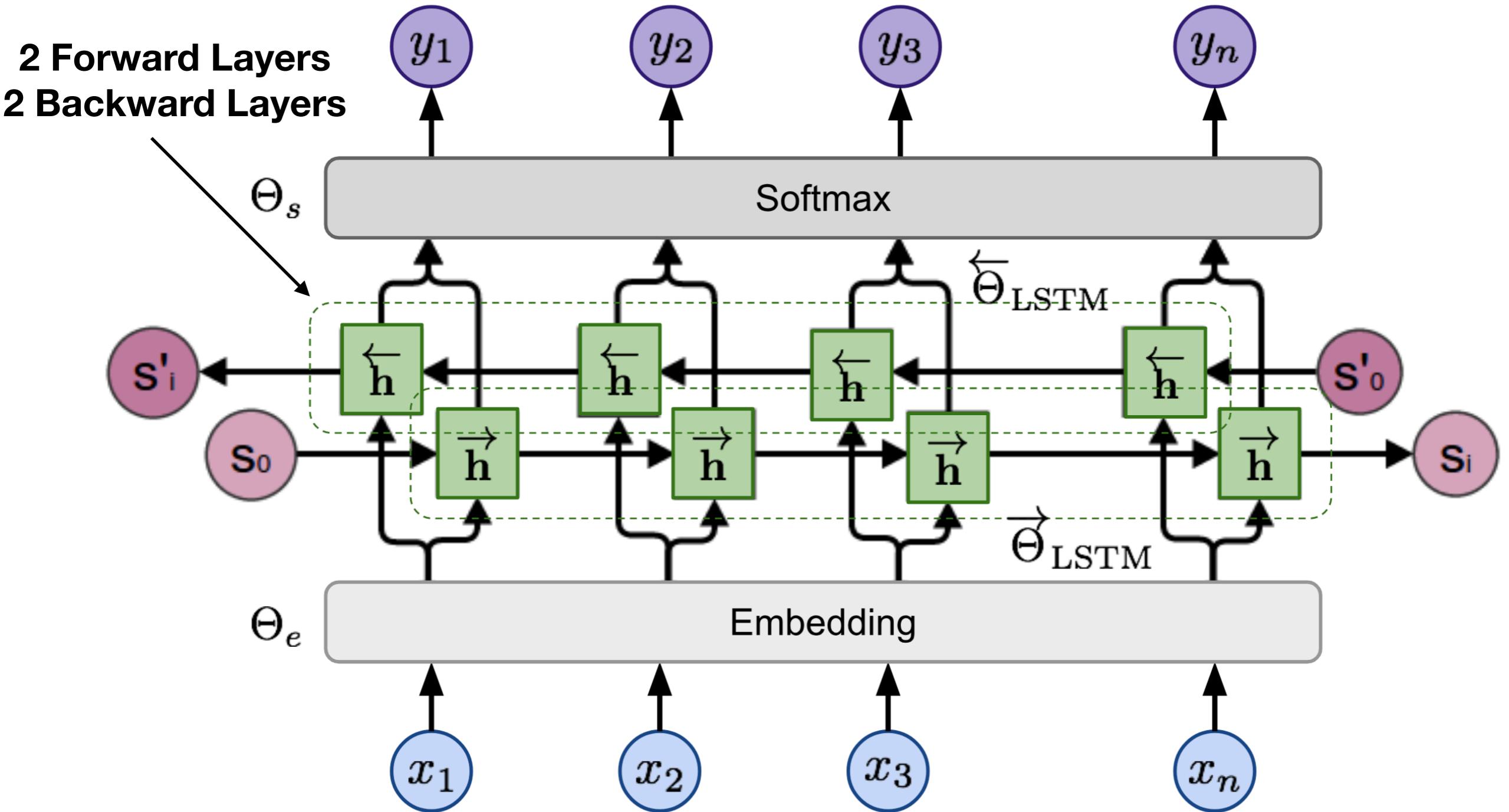
# **ELMo**



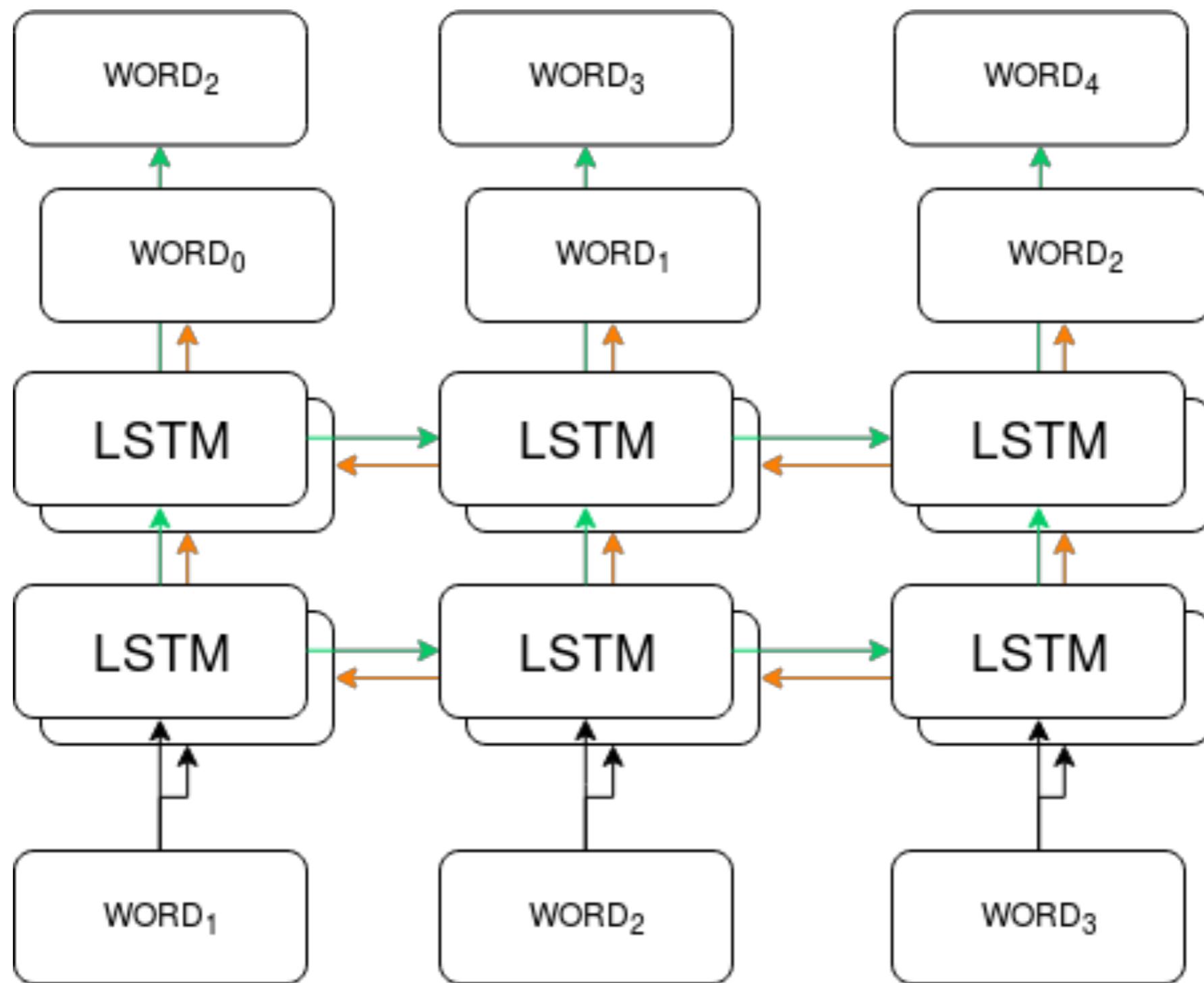
# ELMo



# ELMo



# ELMo



# ELMo

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

# ELMo

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

**Forward**

**Backward**

# ELMo

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

**Forward** <START>

**Backward** <END>

# ELMo

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

**Forward** <START> The

**Backward** <END> money

# ELMo

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

**Forward** <START> The poor

**Backward** <END> 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

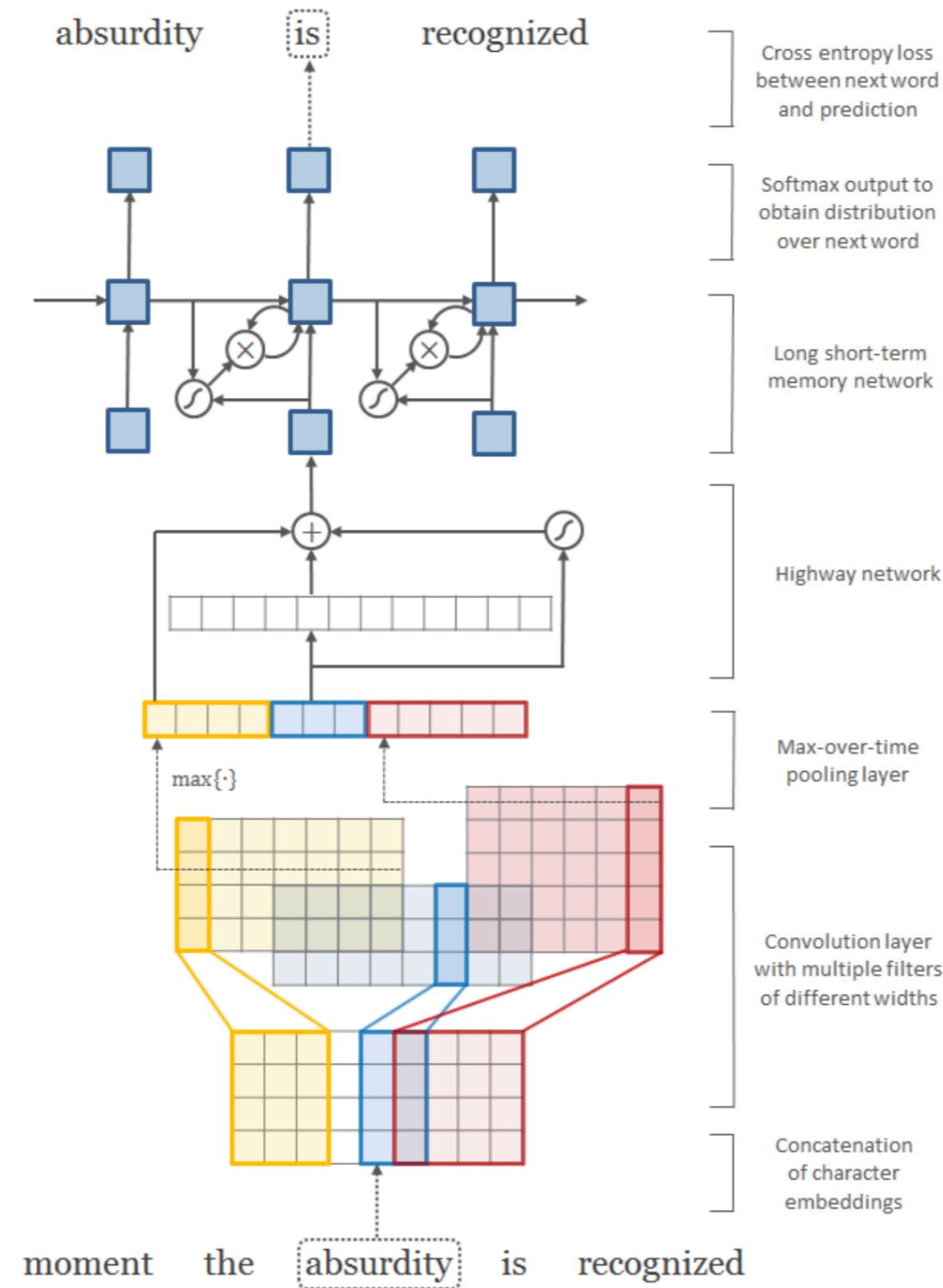
# ELMo

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

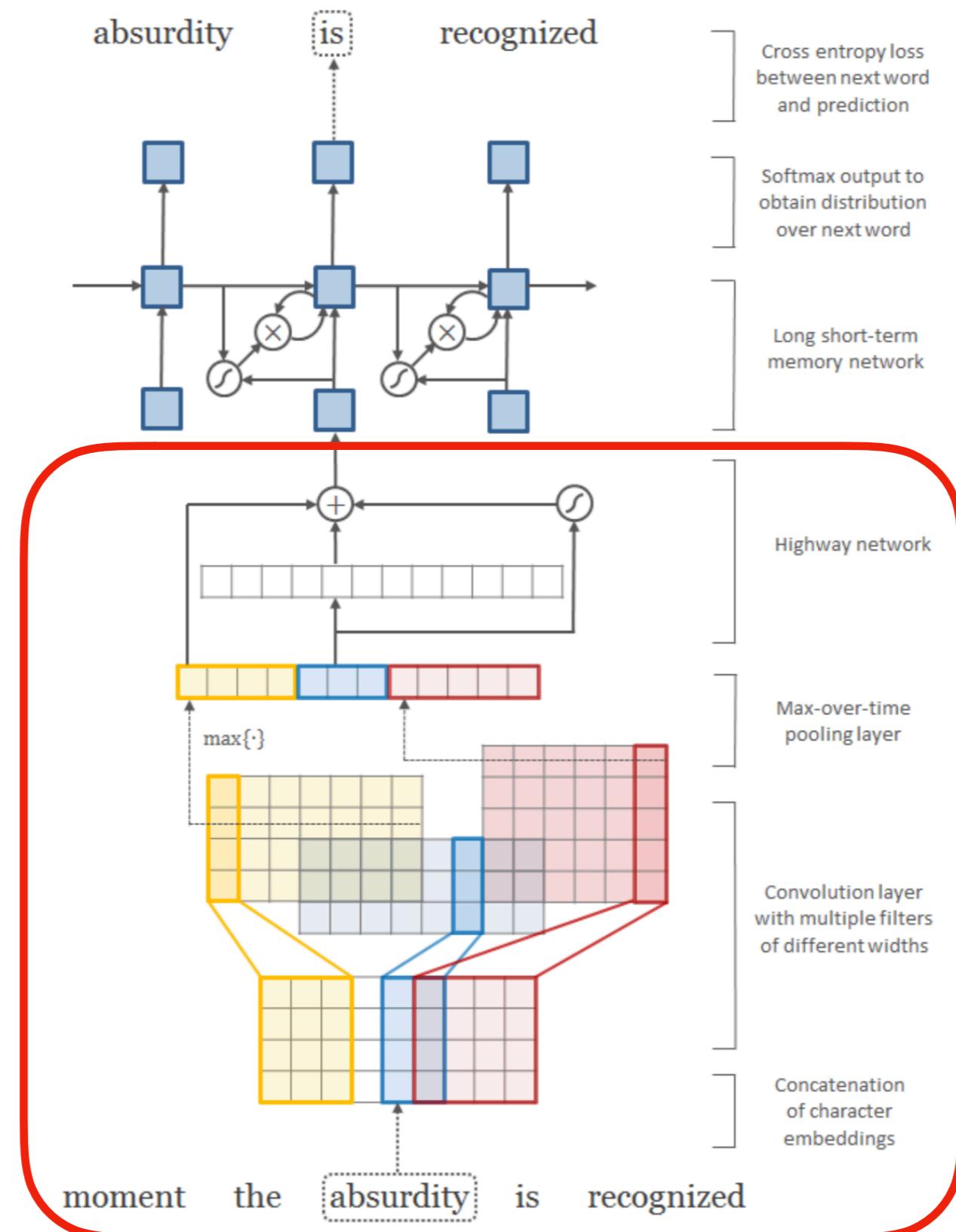
**Forward** <START> The poor don't have any money <END>

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

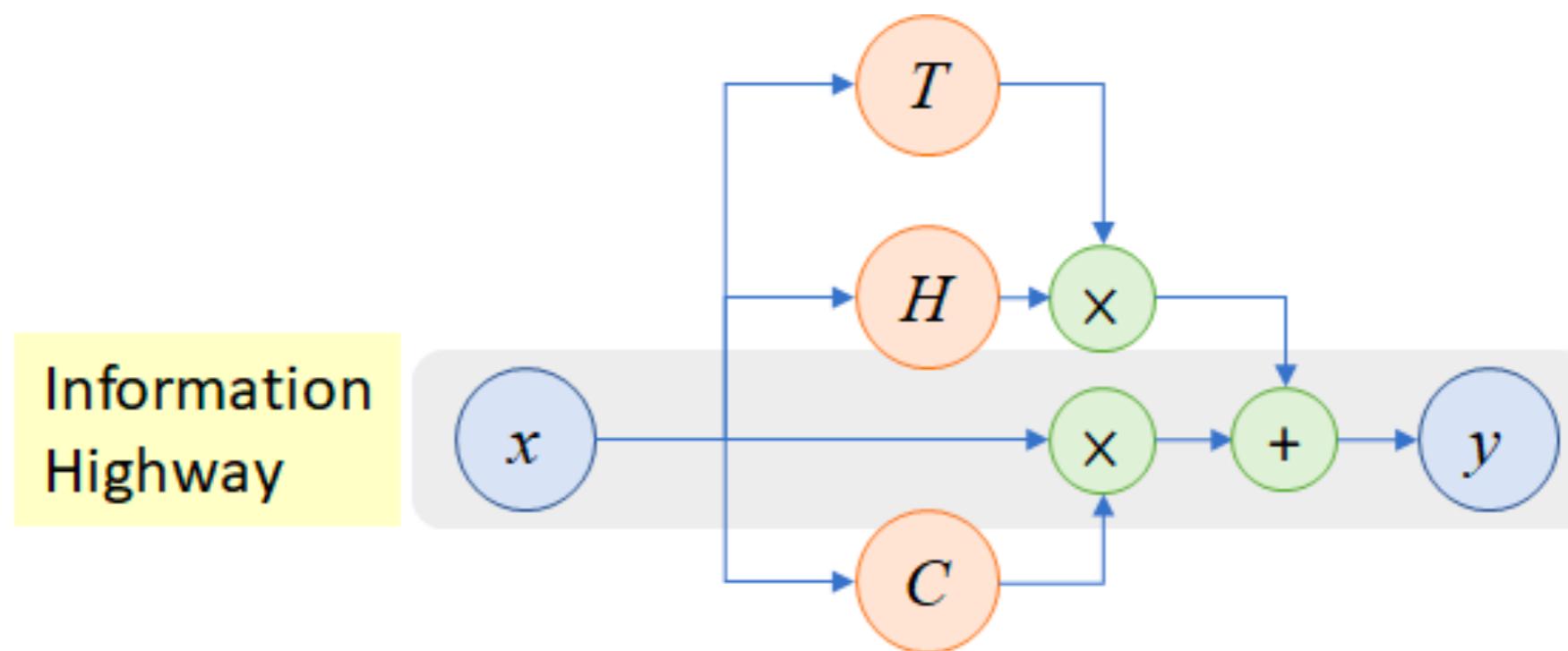
# Character-Aware Neural Language Models



# Character-Aware Neural Language Models

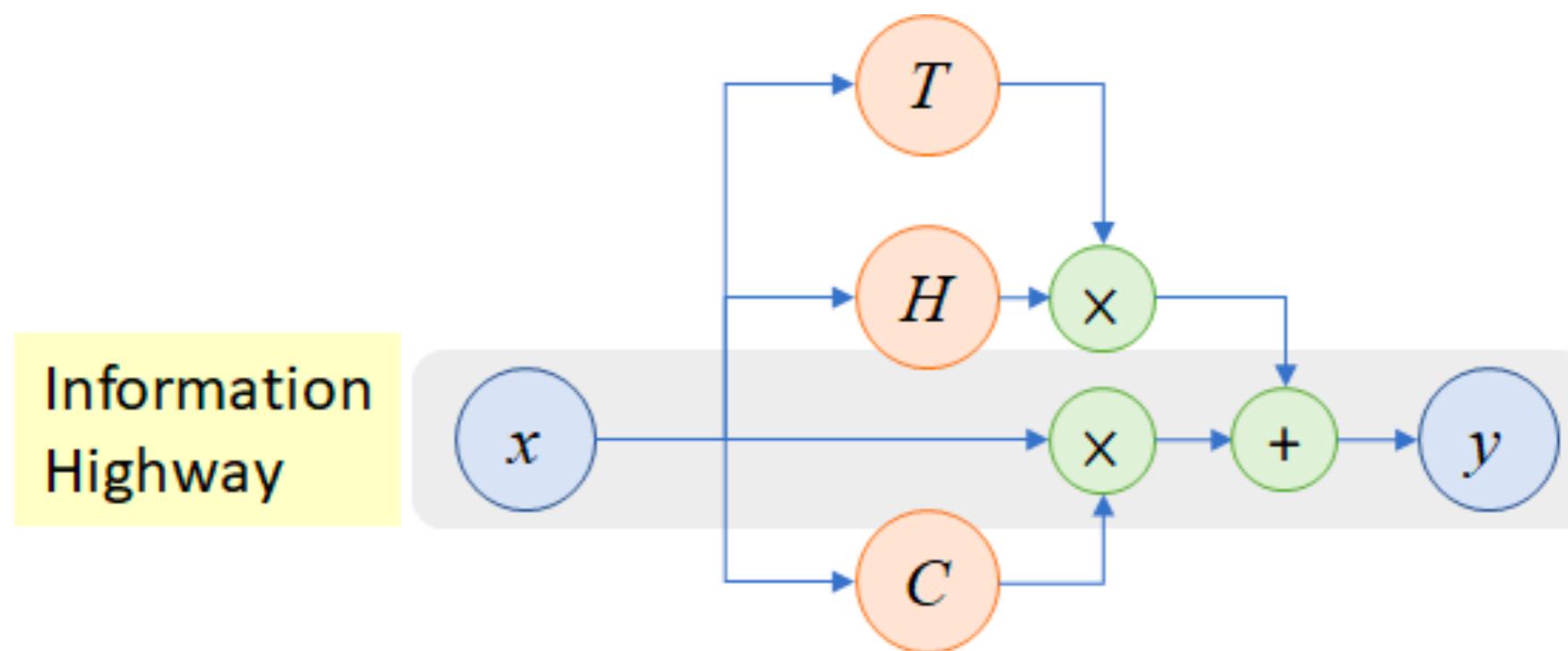


# Highway Network



$$\mathbf{y} = H(\mathbf{x}, \mathbf{W}_H) \cdot T(\mathbf{x}, \mathbf{W}_T) + \mathbf{x} \cdot C(\mathbf{x}, \mathbf{W}_C).$$

# Highway Network



$$\mathbf{y} = H(\mathbf{x}, \mathbf{W}_H) \cdot T(\mathbf{x}, \mathbf{W}_T) + \mathbf{x} \cdot C(\mathbf{x}, \mathbf{W}_C).$$

$$\mathbf{y} = H(\mathbf{x}, \mathbf{W}_H) \cdot T(\mathbf{x}, \mathbf{W}_T) + \mathbf{x} \cdot (1 - T(\mathbf{x}, \mathbf{W}_T)).$$

# Highway Network

```
class Highway(BaseModule):

    def __init__(self, in_features, out_features, activation_function=torch.nn.ReLU()):

        super(Highway, self).__init__()

        self.nonlinear = torch.nn.Linear(in_features, out_features)
        self.linear = torch.nn.Linear(in_features, out_features)
        self.gate = torch.nn.Linear(in_features, out_features)

        self.gate_function = torch.nn.Sigmoid()
        self.activation_function = activation_function

    def forward(self, x):
        """
        :param x: tensor with shape of [batch_size, size]
        :return: tensor with shape of [batch_size, size]
        applies  $\sigma(x) \odot (f(G(x))) + (1 - \sigma(x)) \odot (Q(x))$  transformation / G and Q is affine transformation,
        f is non-linear transformation,  $\sigma(x)$  is affine transformation with sigmoid non-linearity
        and  $\odot$  is element-wise multiplication
        """

        gate = self.gate(x)
        gate = self.gate_function(gate)

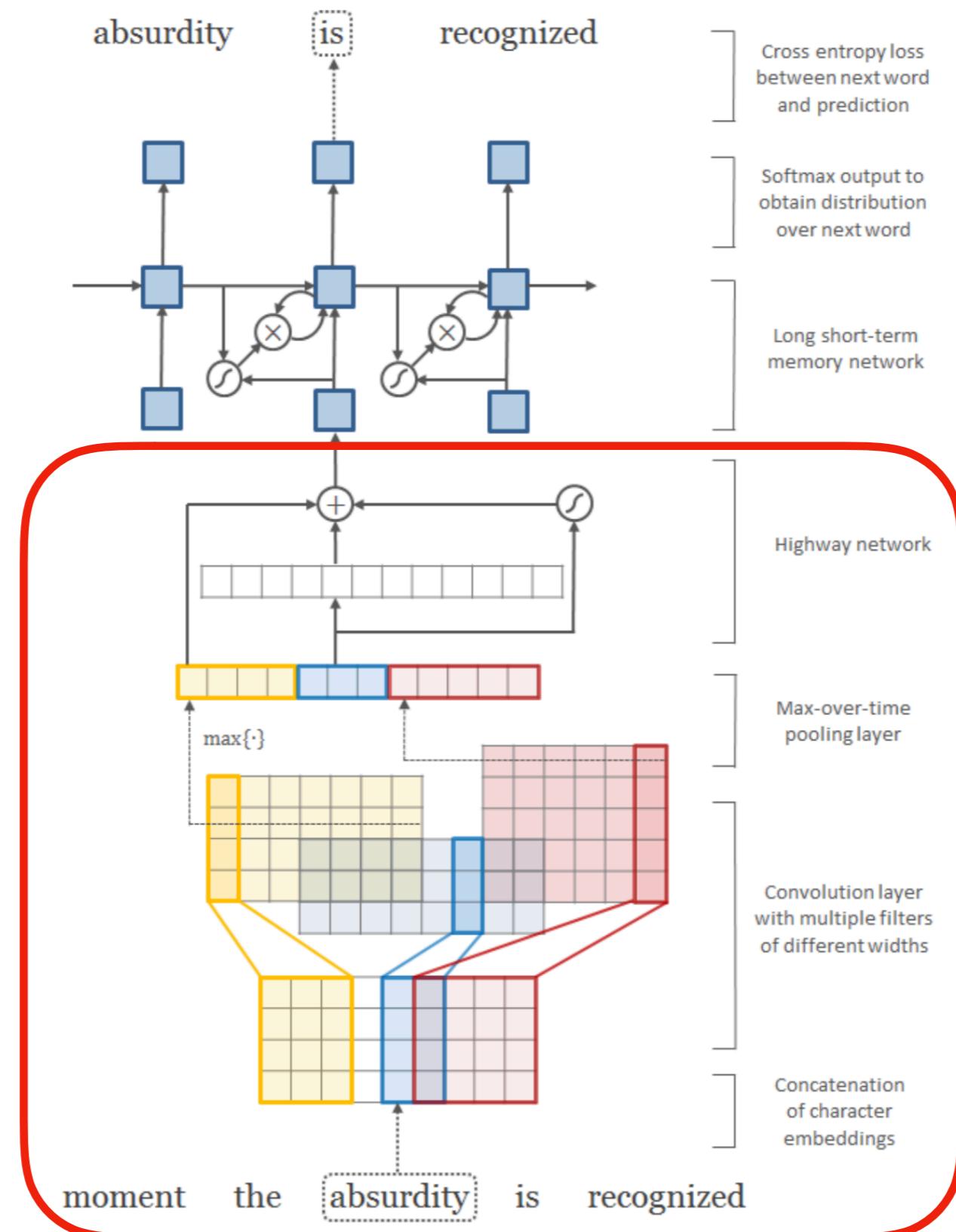
        nonlinear = self.nonlinear(x)
        nonlinear = self.activation_function(nonlinear)

        linear = self.linear(x)

        x = gate * nonlinear + (1 - gate) * linear

        return x
```

# Character-Aware Neural Language Models



# Gated Linear Unit

$$(X^*W+b) \otimes \text{sigmoid}(X^*V+c)$$

```
class GatedLinearUnit(BaseModule):
    """
    Implementation of this https://arxiv.org/pdf/1612.08083.pdf
    """

    def __init__(self, in_features, out_features):
        super().__init__()

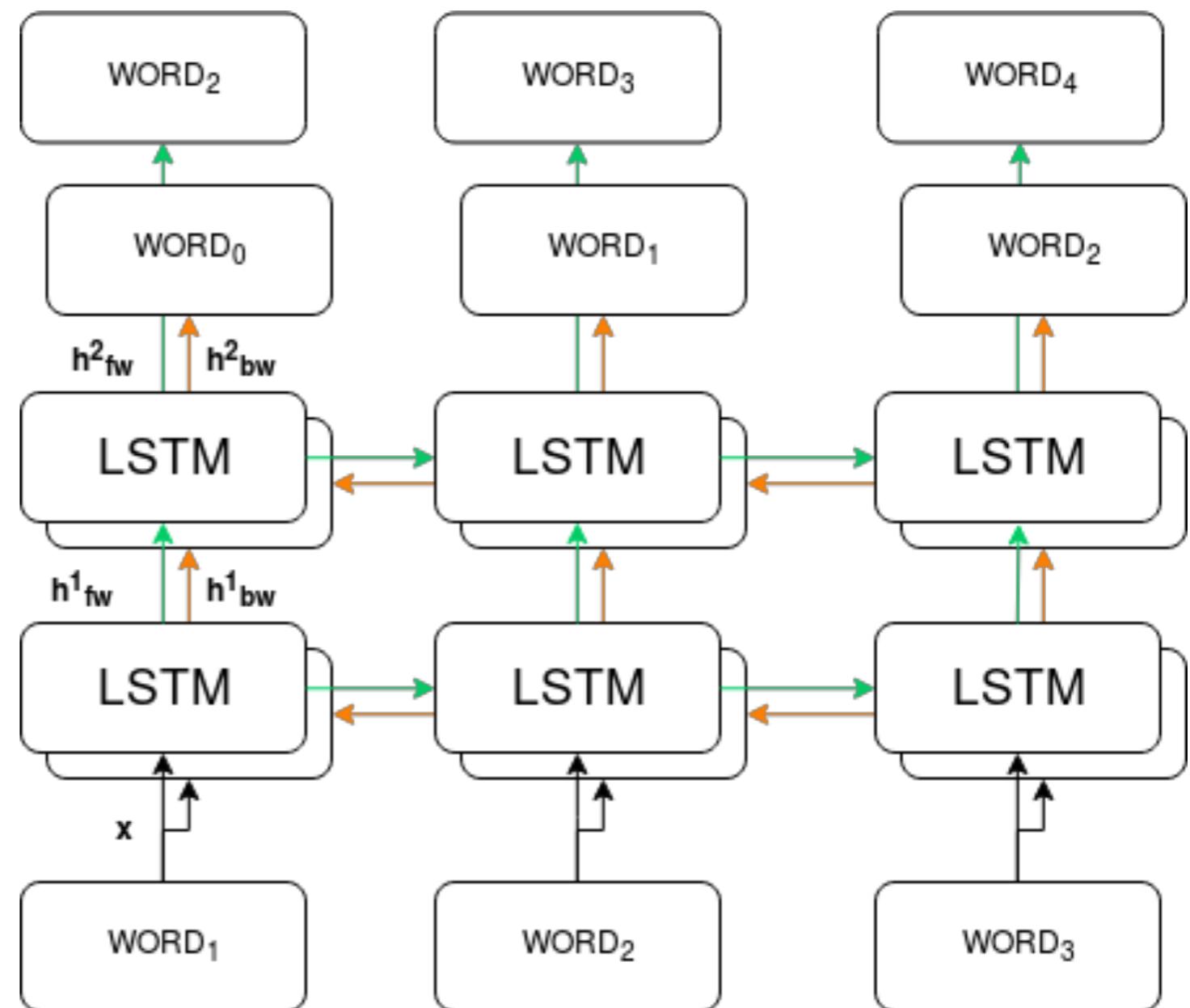
        self.linear = torch.nn.Linear(in_features=in_features, out_features=out_features)
        self.gate = torch.nn.Linear(in_features=in_features, out_features=out_features)

    def forward(self, x):
        x = self.linear(x)
        gate = self.linear(x)

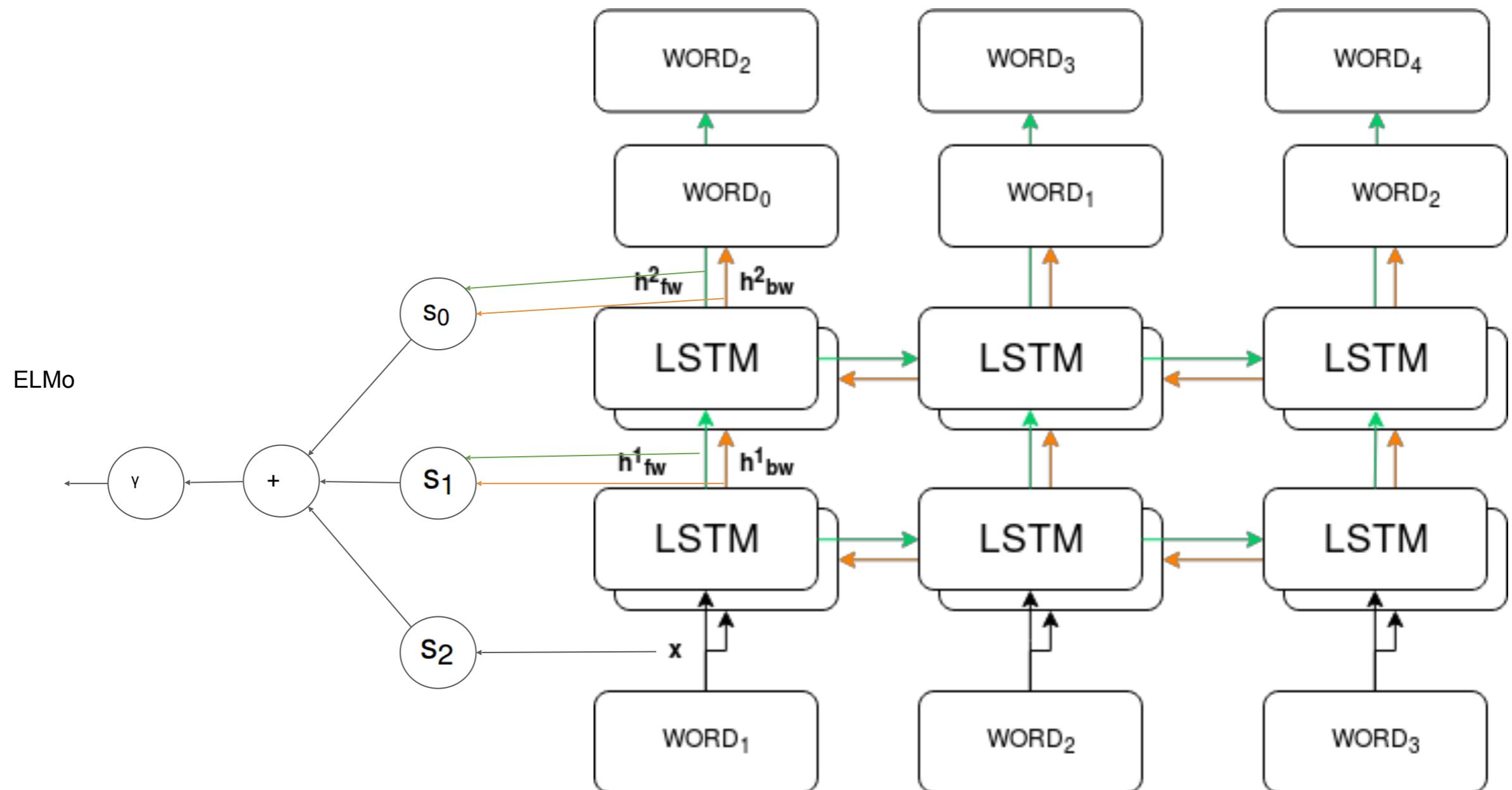
        x = x * torch.sigmoid(gate)

        return x
```

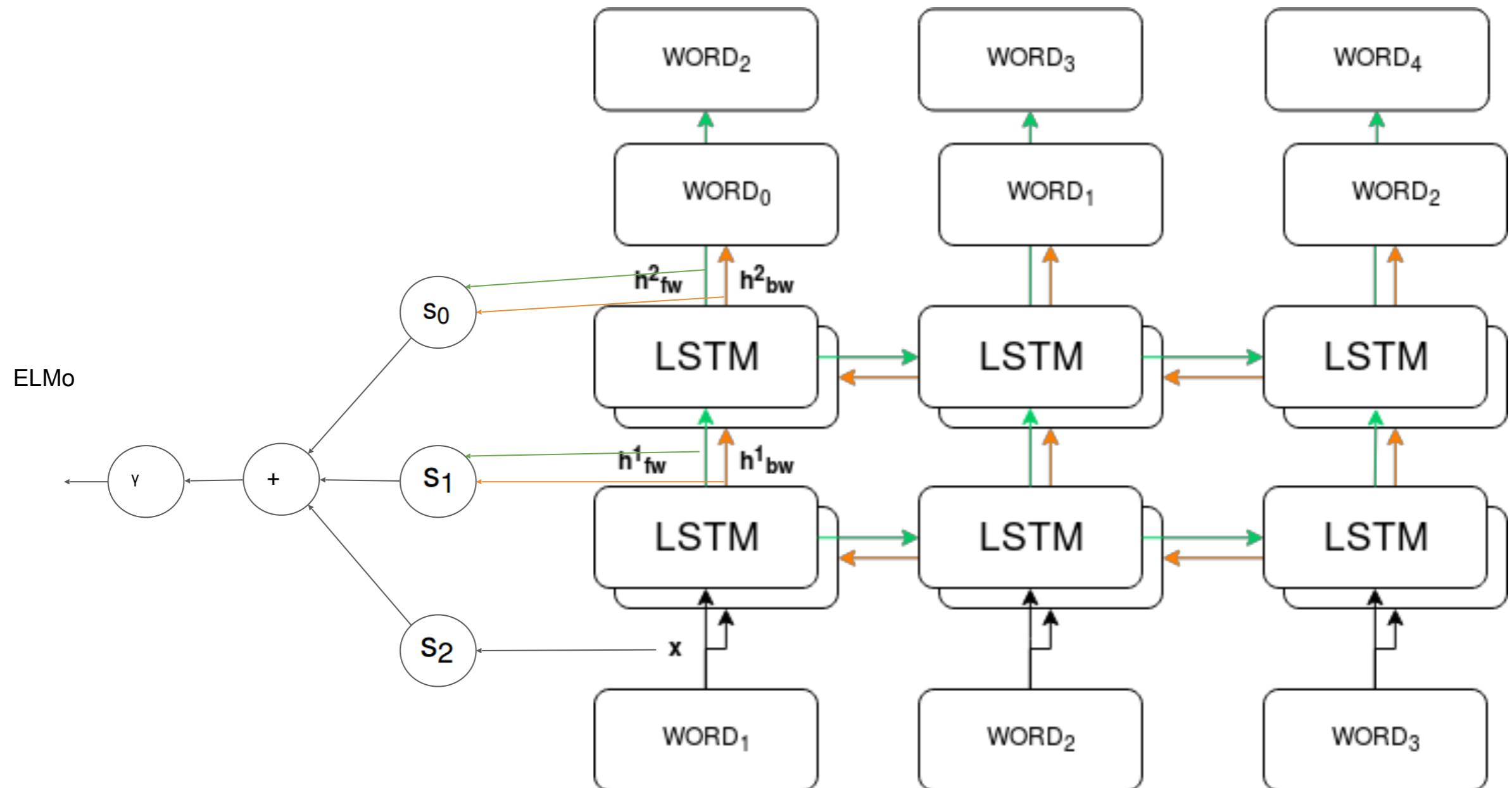
# ELMo



# ELMo



# ELMo



$$\begin{aligned} R_k &= \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\} \\ &= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\}, \end{aligned}$$

$$\text{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM}.$$

# ELMo

## Results

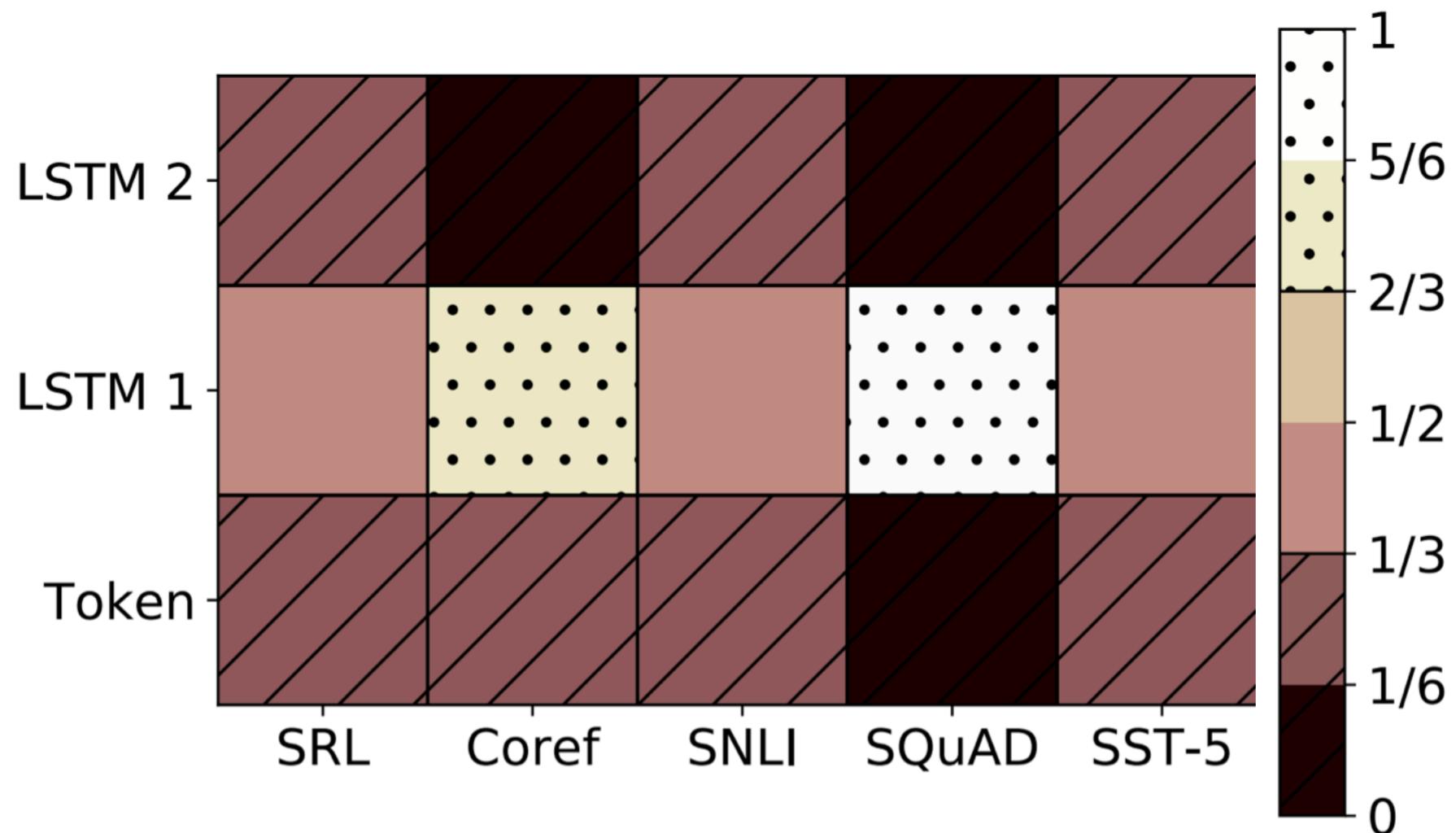
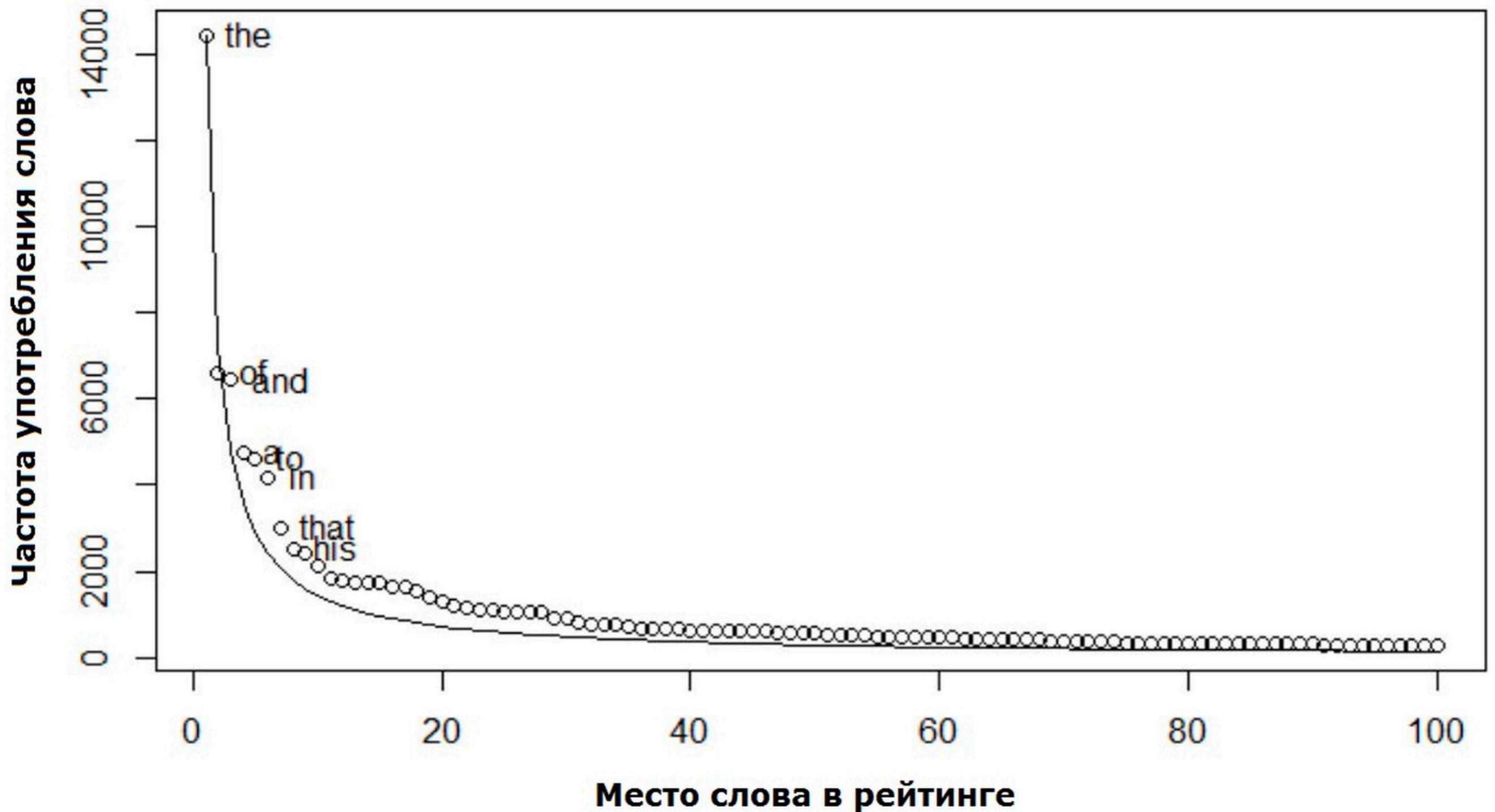


Figure 2: Visualization of softmax normalized biLM layer weights across tasks and ELMo locations. Normalized weights less than  $1/3$  are hatched with horizontal lines and those greater than  $2/3$  are speckled.

# 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']
```

# Tokenization

## Word level

- + Small text length
- Big vocabulary size
- OOV

## Character level

- Long text
- + Small vocabulary size
- + Almost no OOV



# 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

---

```
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
l o w	→	low
e r ·	→	er·

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

# BPE

---

## Algorithm 1 Learn BPE operations

---

```
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 ·

- learning

- word:freq : {low:5, lowest:2, newer:6, wider:3}

- merge & 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>'

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

# BPE

---

## Algorithm 1 Learn BPE operations

---

```
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
l o w	→	low
e r ·	→	er ·

- learning

- word:freq : {low:5, lowest:2, newer:6, wider:3}

- merge & 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

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

# BPE

## SentencePiece

[build failing](#) [build passing](#) [coverage 98%](#) [issues 35 open](#) [code quality A](#) [pypi package 0.1.83](#) [contributions](#) [welcome](#)

License Apache 2.0

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.

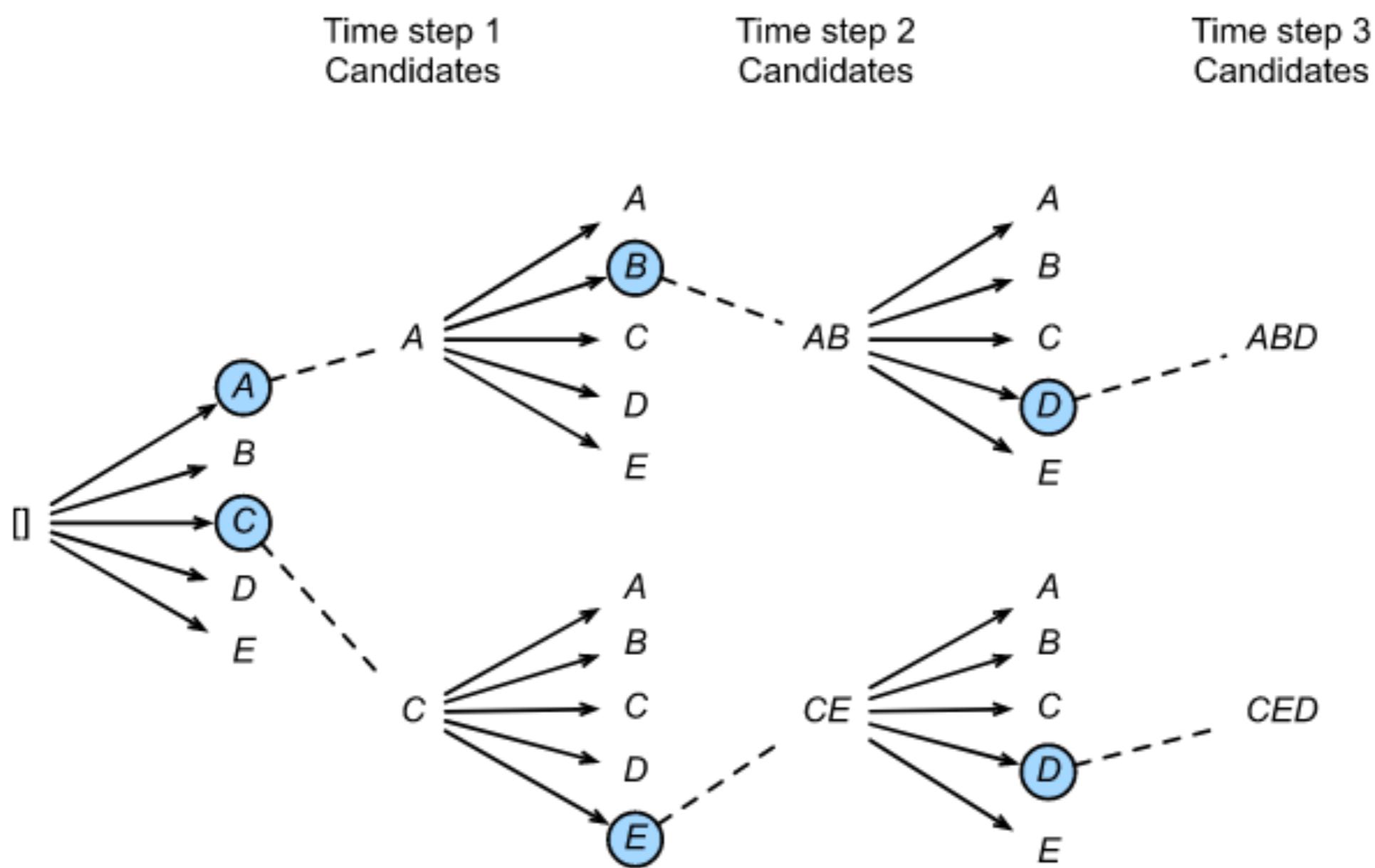
## YouTokenToMe

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

# Beam Search



# Beam Search

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{} | x, y^{<1>}, \dots, y^{$$

# Beam Search

$$p(\mathbf{x}) = \prod_i p(x|x_{<i}) = p(x_0)p(x_1|x_0)p(x_2|x_0, x_1)\dots$$
$$\arg \max_y \prod_{t=1}^{T_y} P(y^{} | x, y^{<1>}, \dots, y^{$$

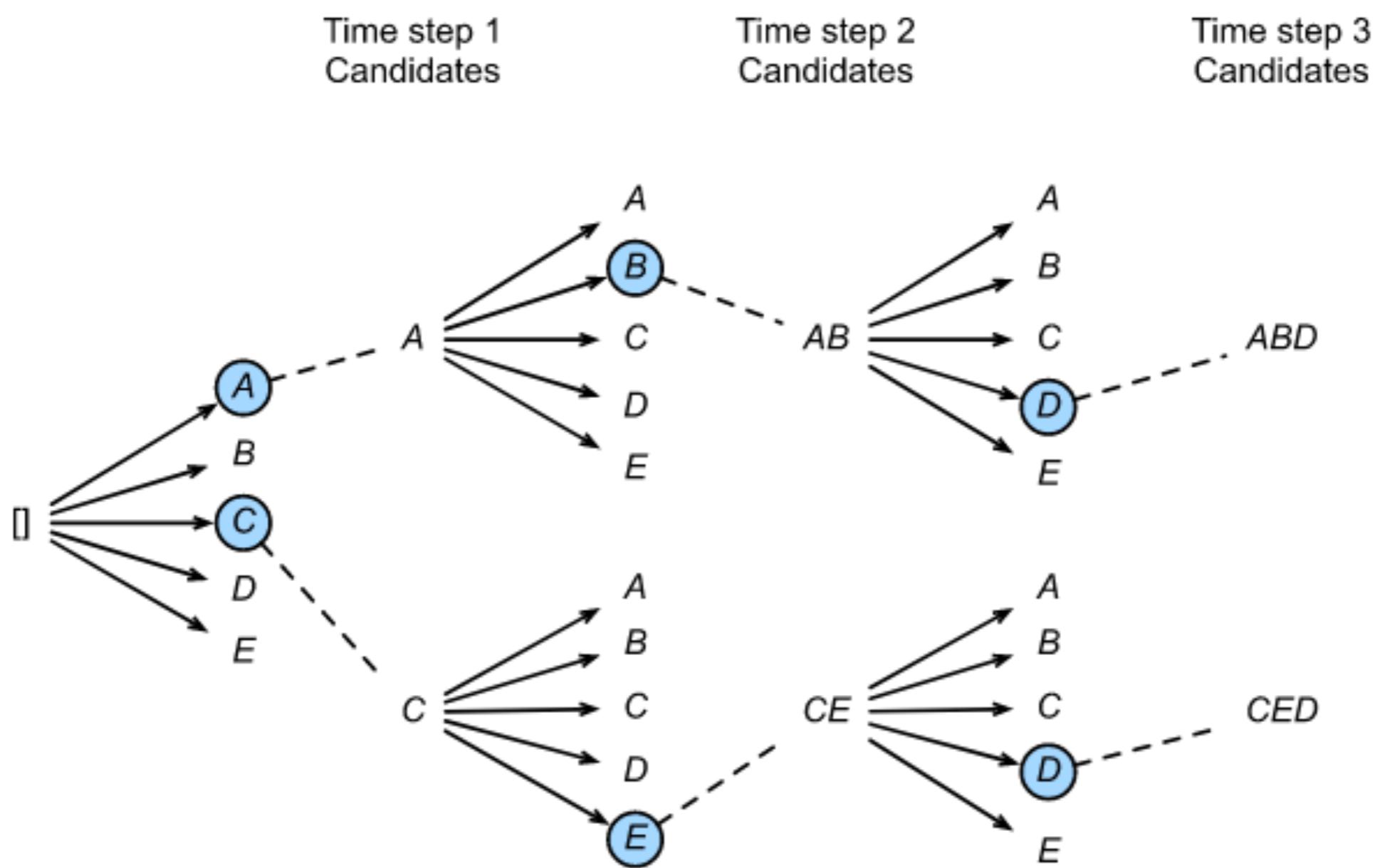
# Beam Search

$$p(\mathbf{x}) = \prod_i p(x|x_{<i}) = p(x_0)p(x_1|x_0)p(x_2|x_0, x_1)\dots$$

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{} | x, y^{<1>}, \dots, y^{})$$

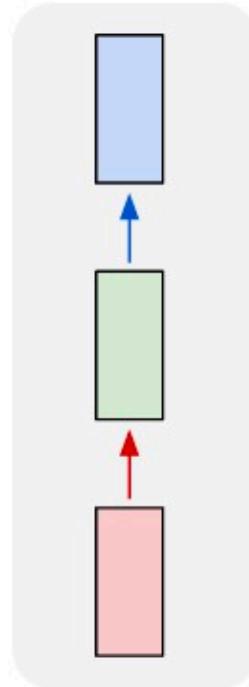
$$\arg \max_y \sum_{y=1}^{T_y} \log P(y^{} | x, y^{<1>}, \dots, y^{})$$

# Beam Search

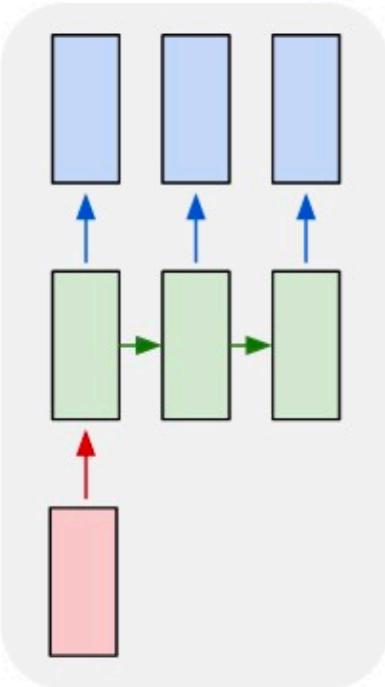


# Sequence to sequence

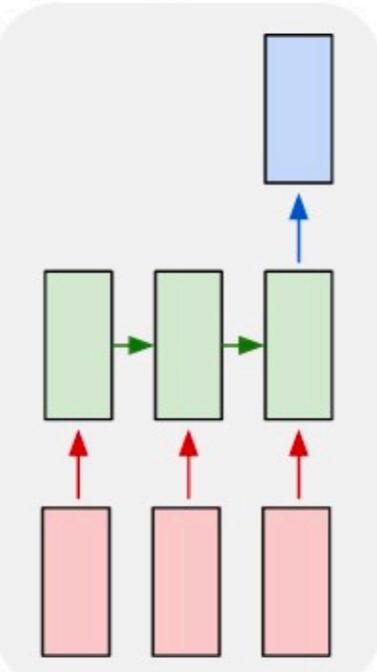
one to one



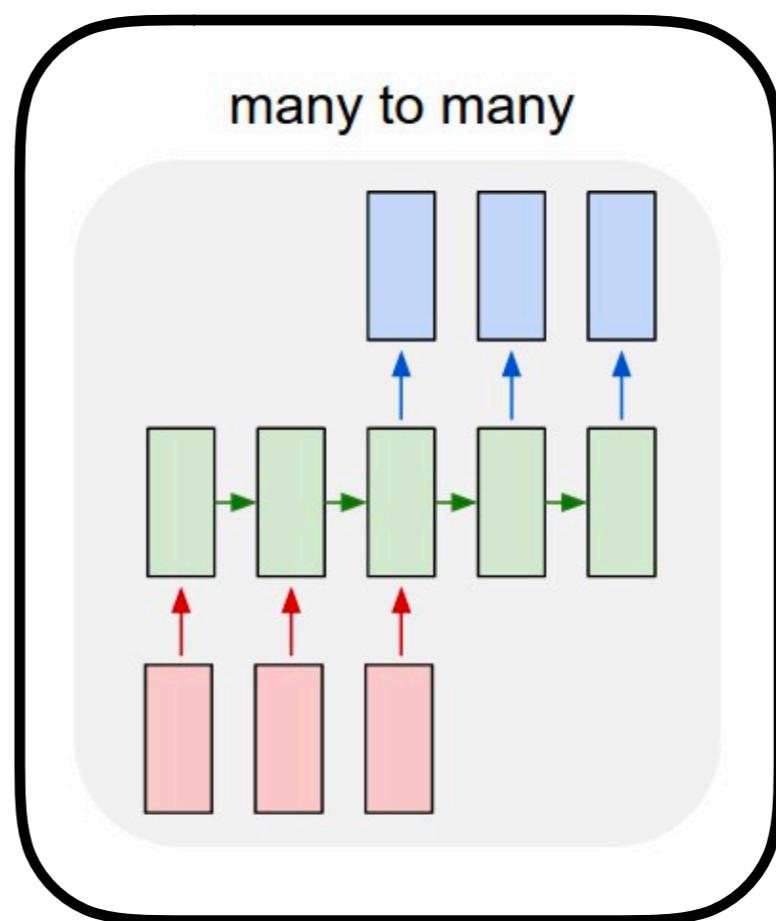
one to many



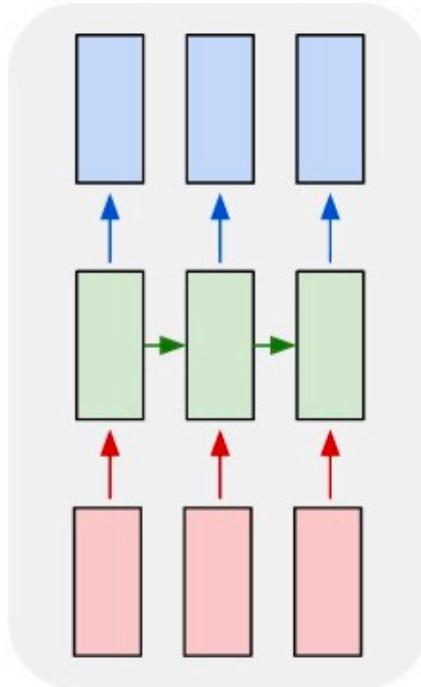
many to one



many to many



many to many



# Sequence to sequence

---

**Source sentence**

# Sequence to sequence

Les pauvres sont démunis

---

**Source sentence**

# Sequence to sequence

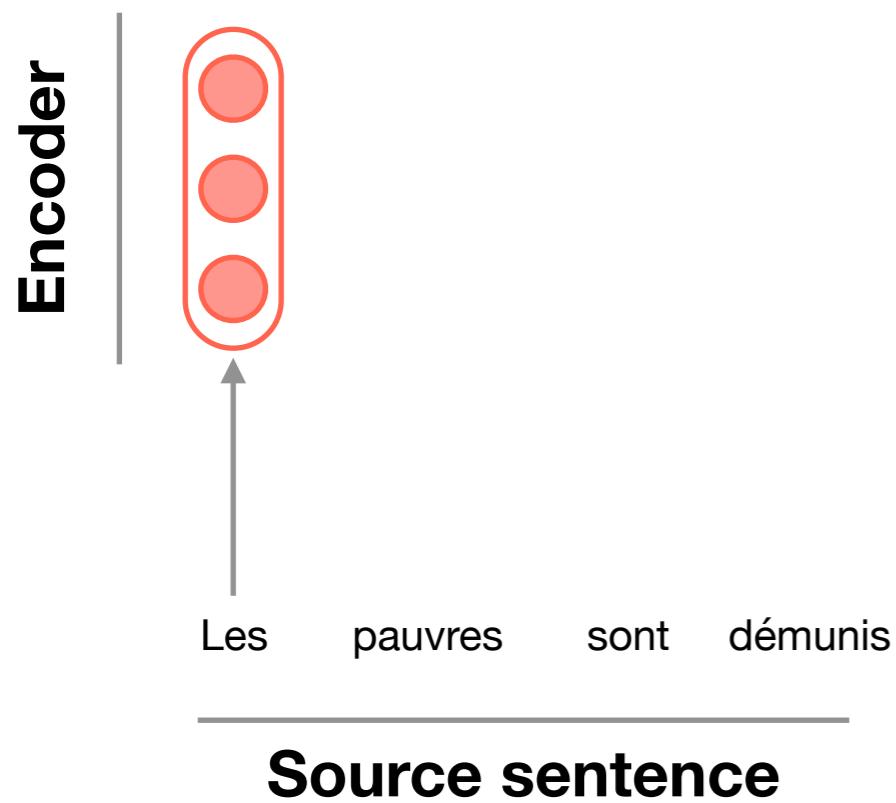
Encoder

Les pauvres sont démunis

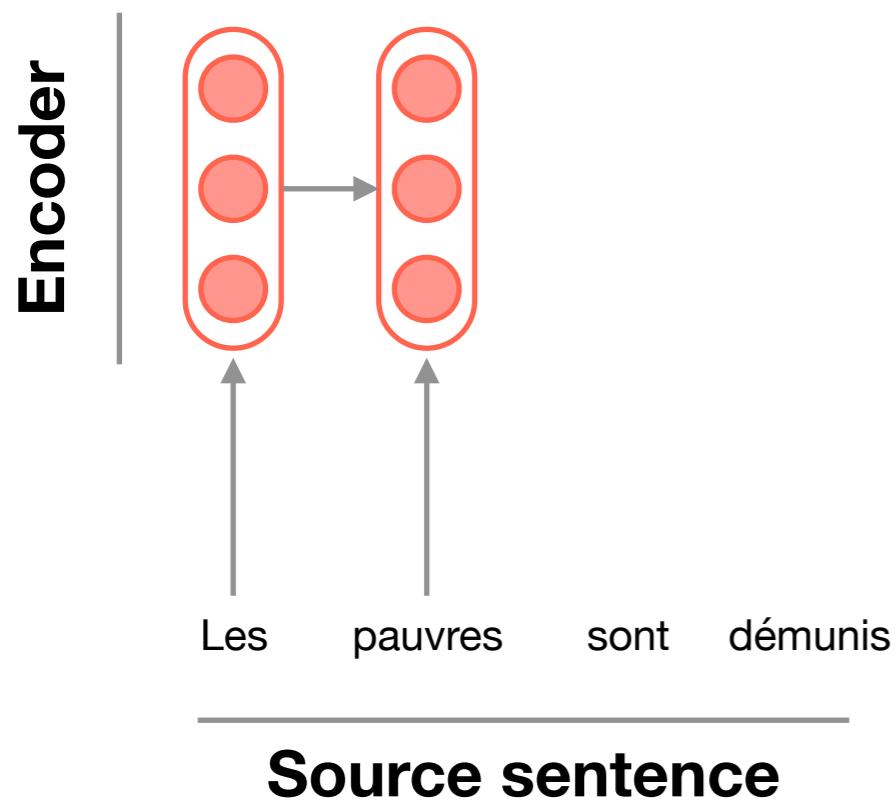
---

**Source sentence**

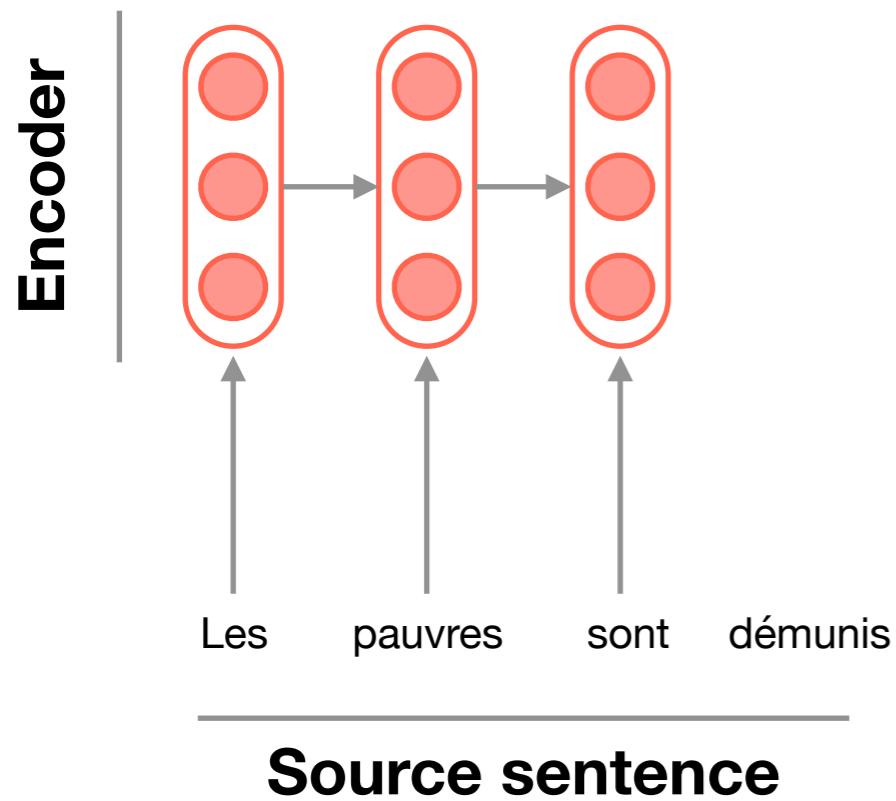
# Sequence to sequence



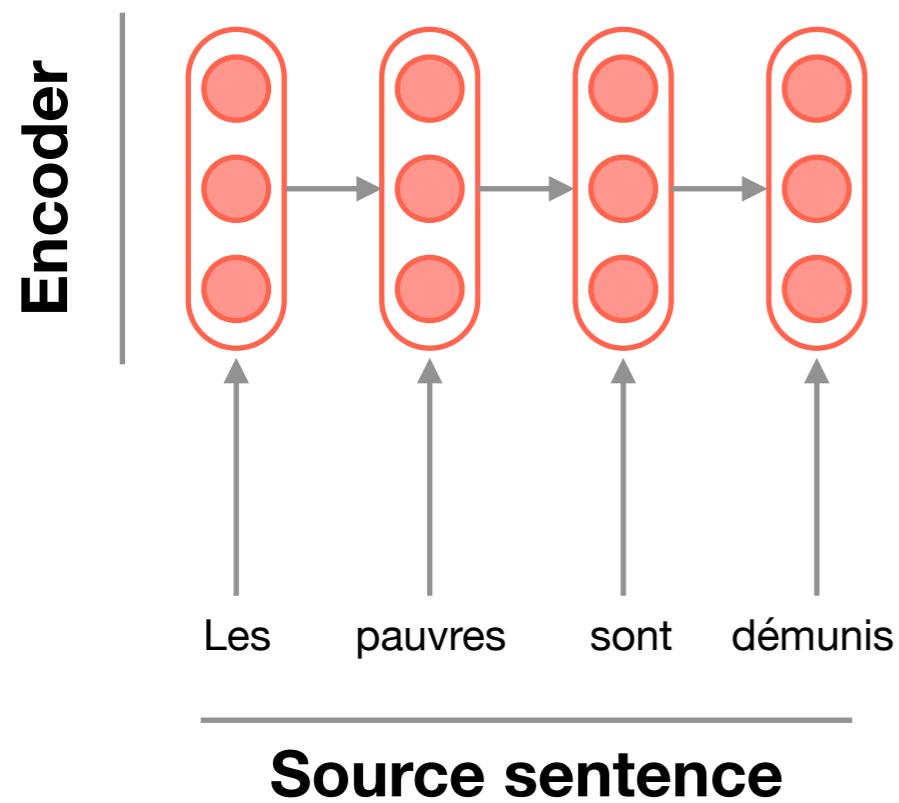
# Sequence to sequence



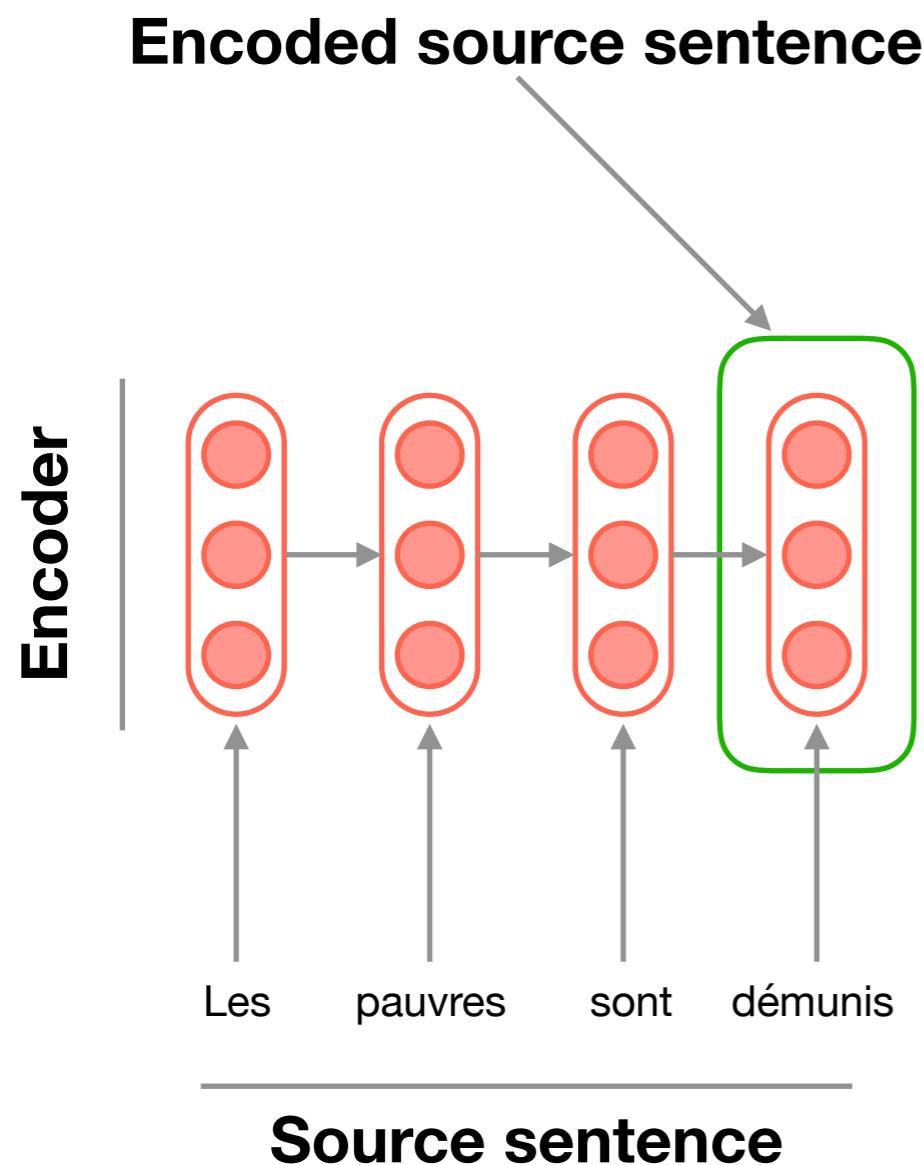
# Sequence to sequence



# Sequence to sequence

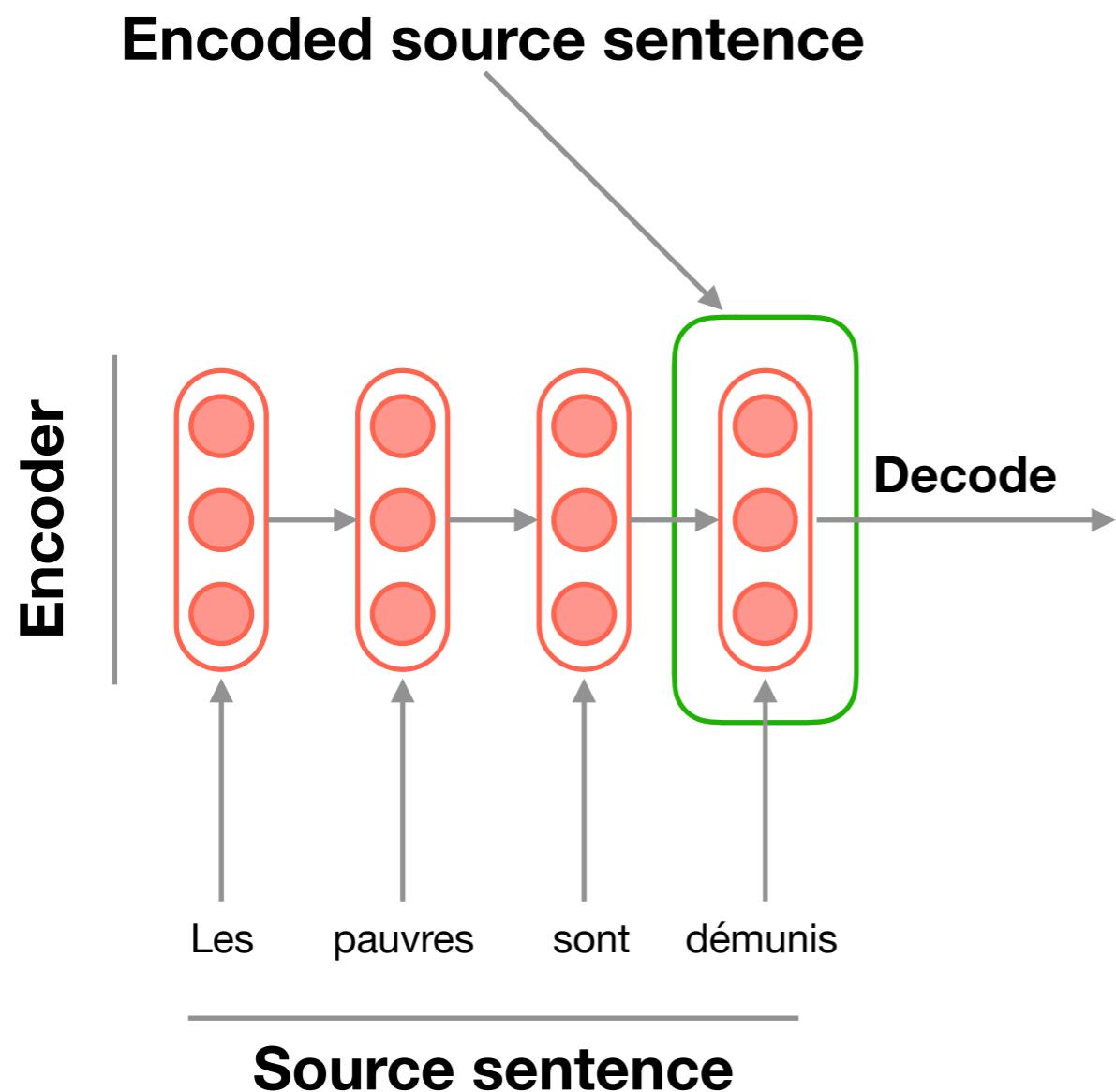


# Sequence to sequence



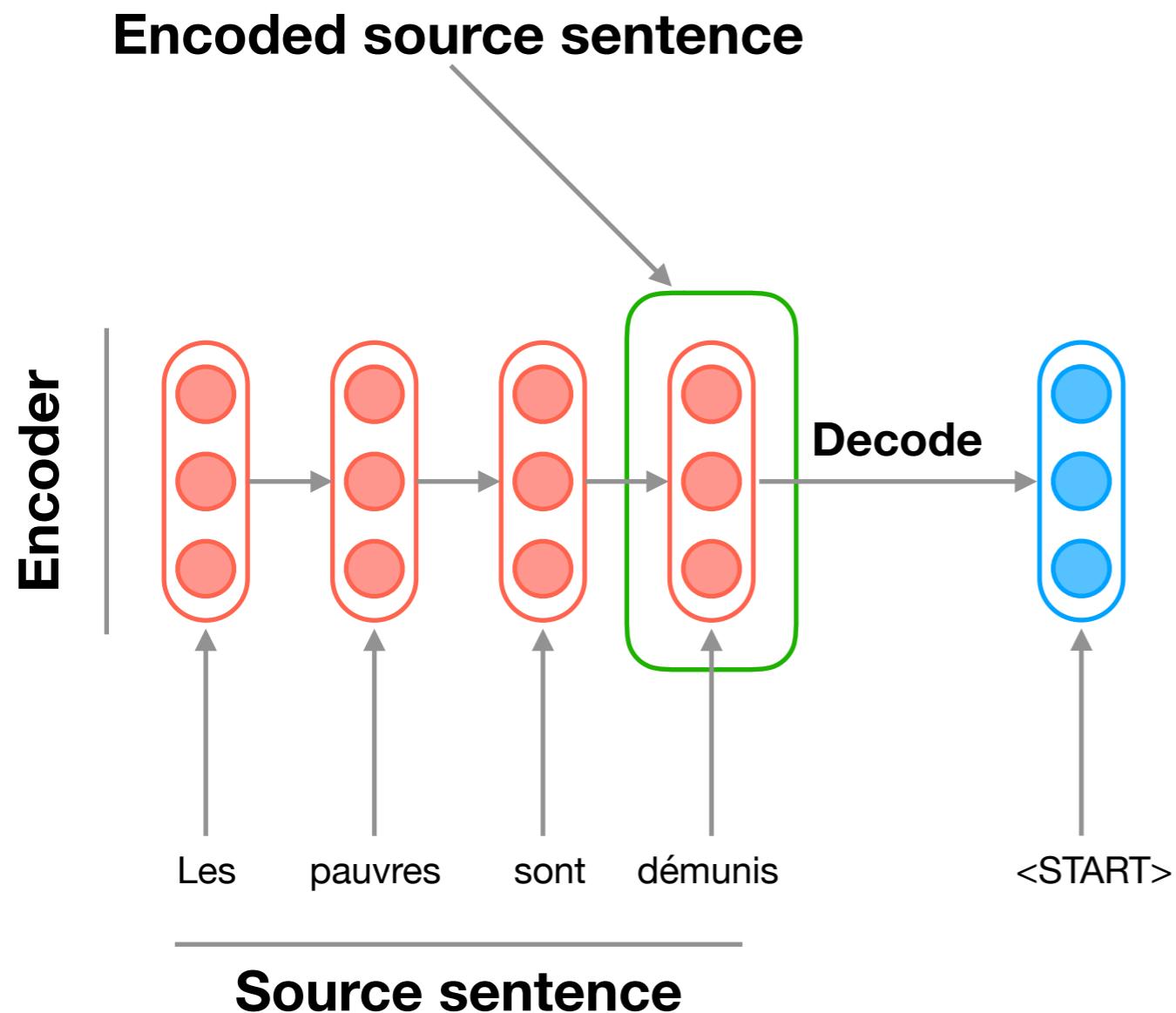
# Sequence to sequence

## Inference



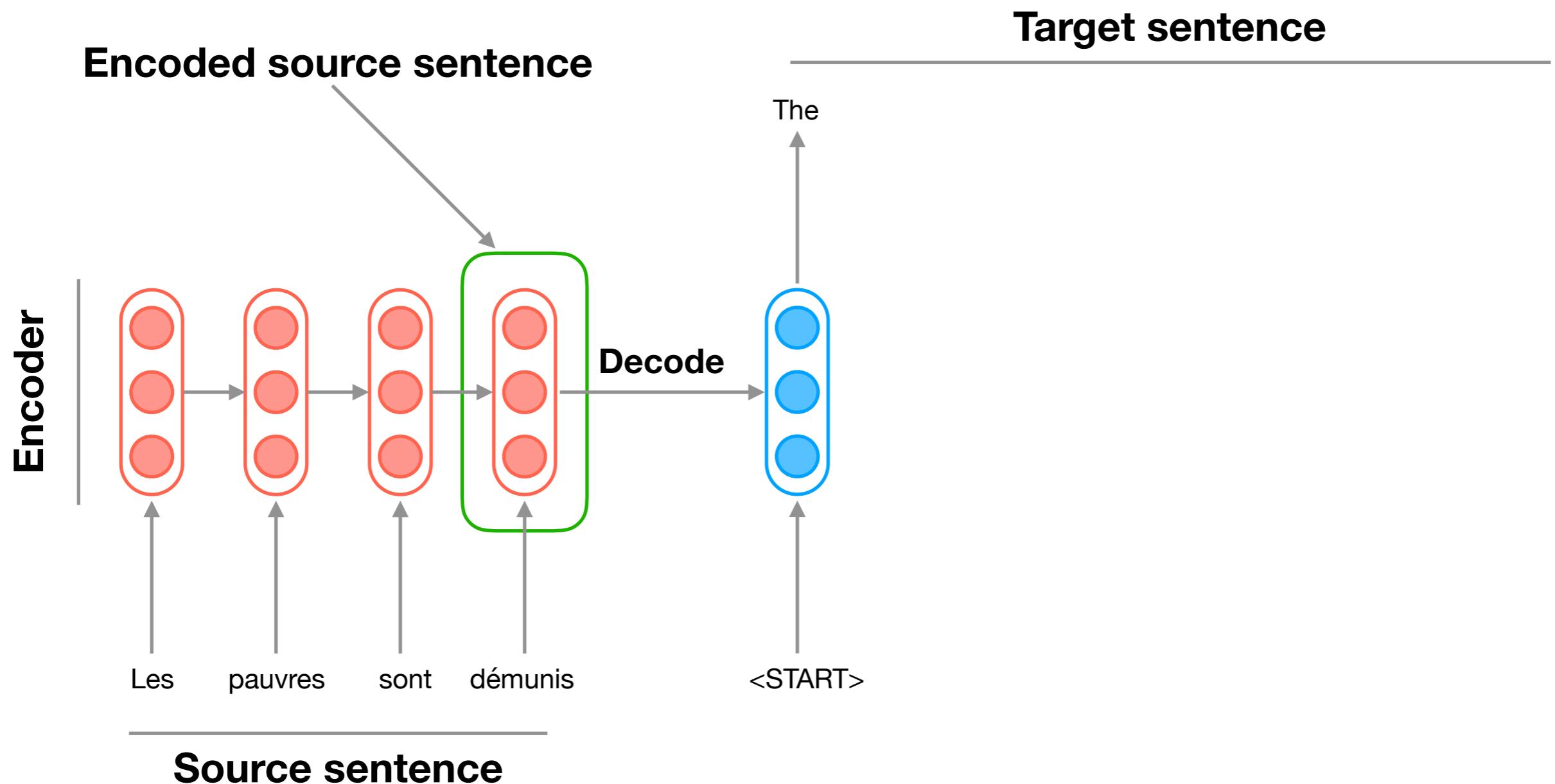
# Sequence to sequence

## Inference



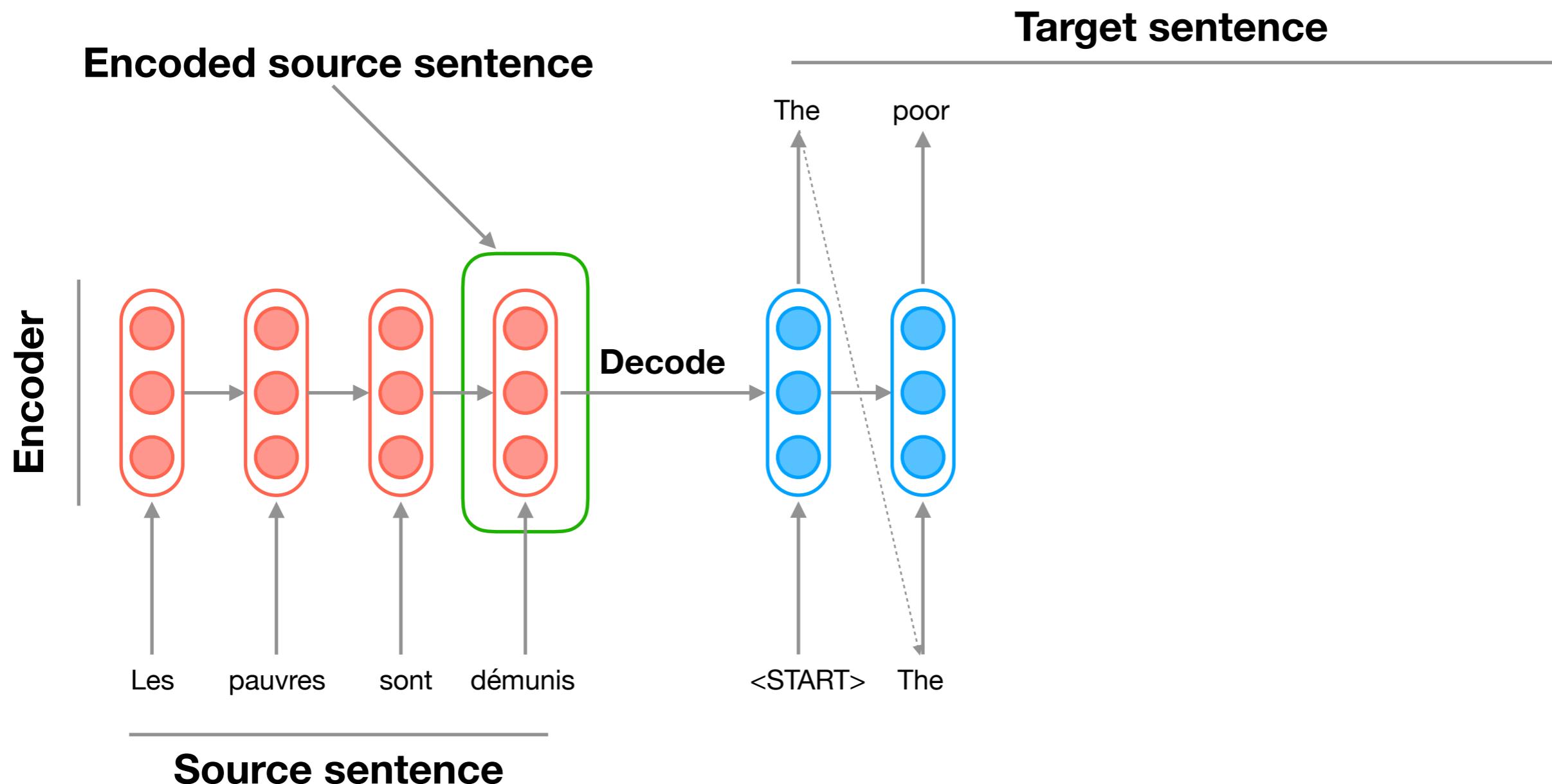
# Sequence to sequence

## Inference



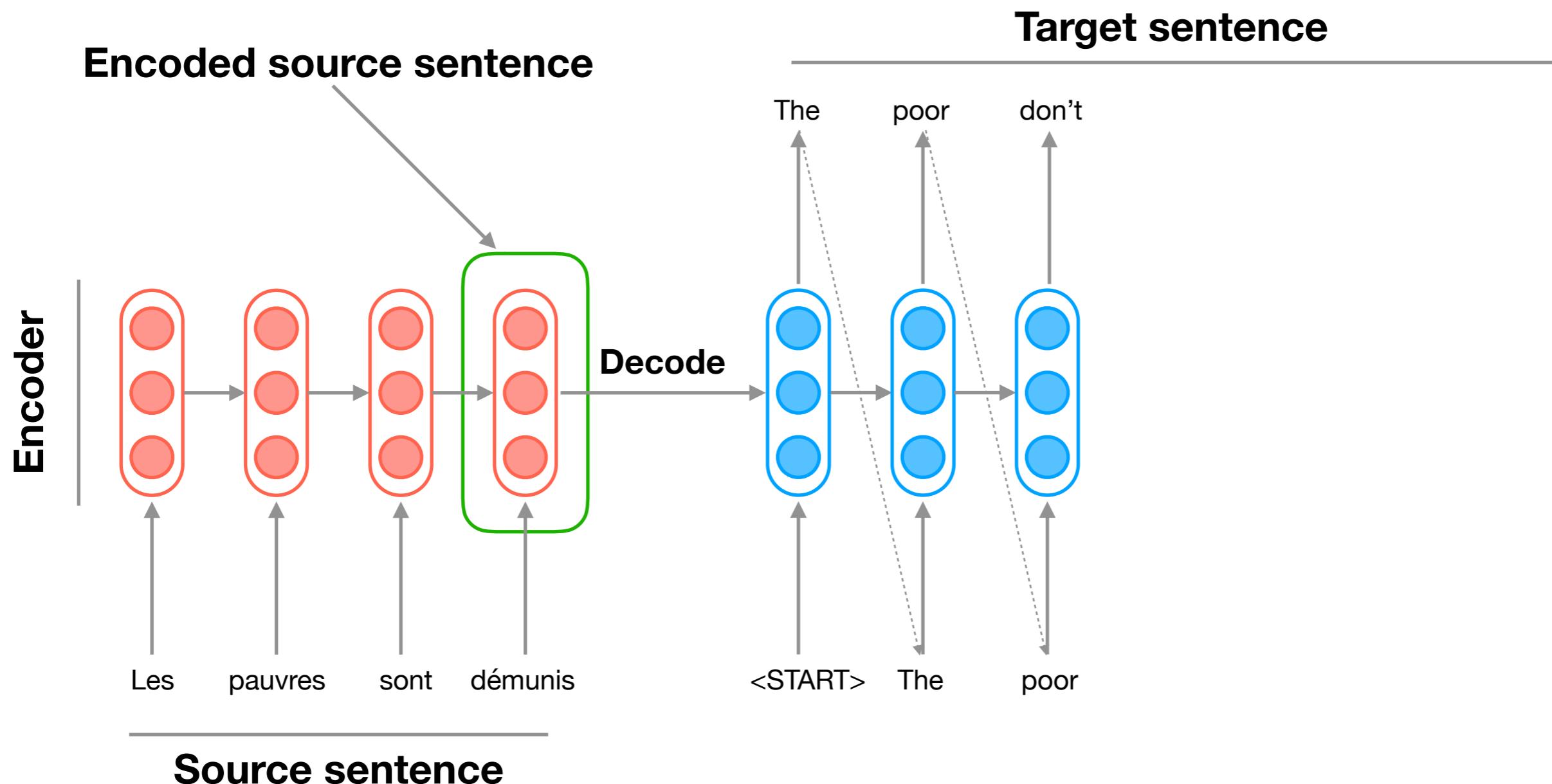
# Sequence to sequence

## Inference



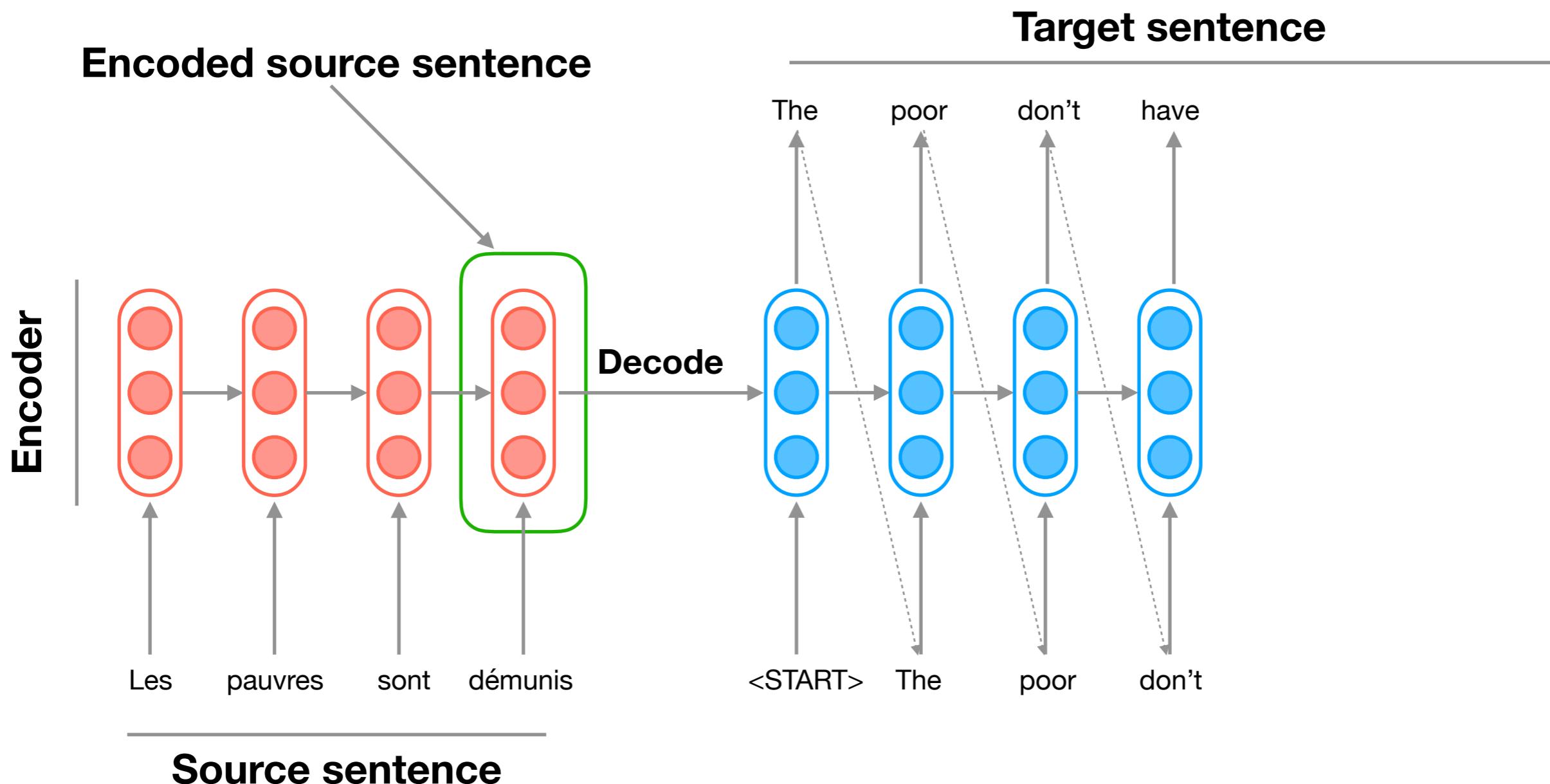
# Sequence to sequence

## Inference



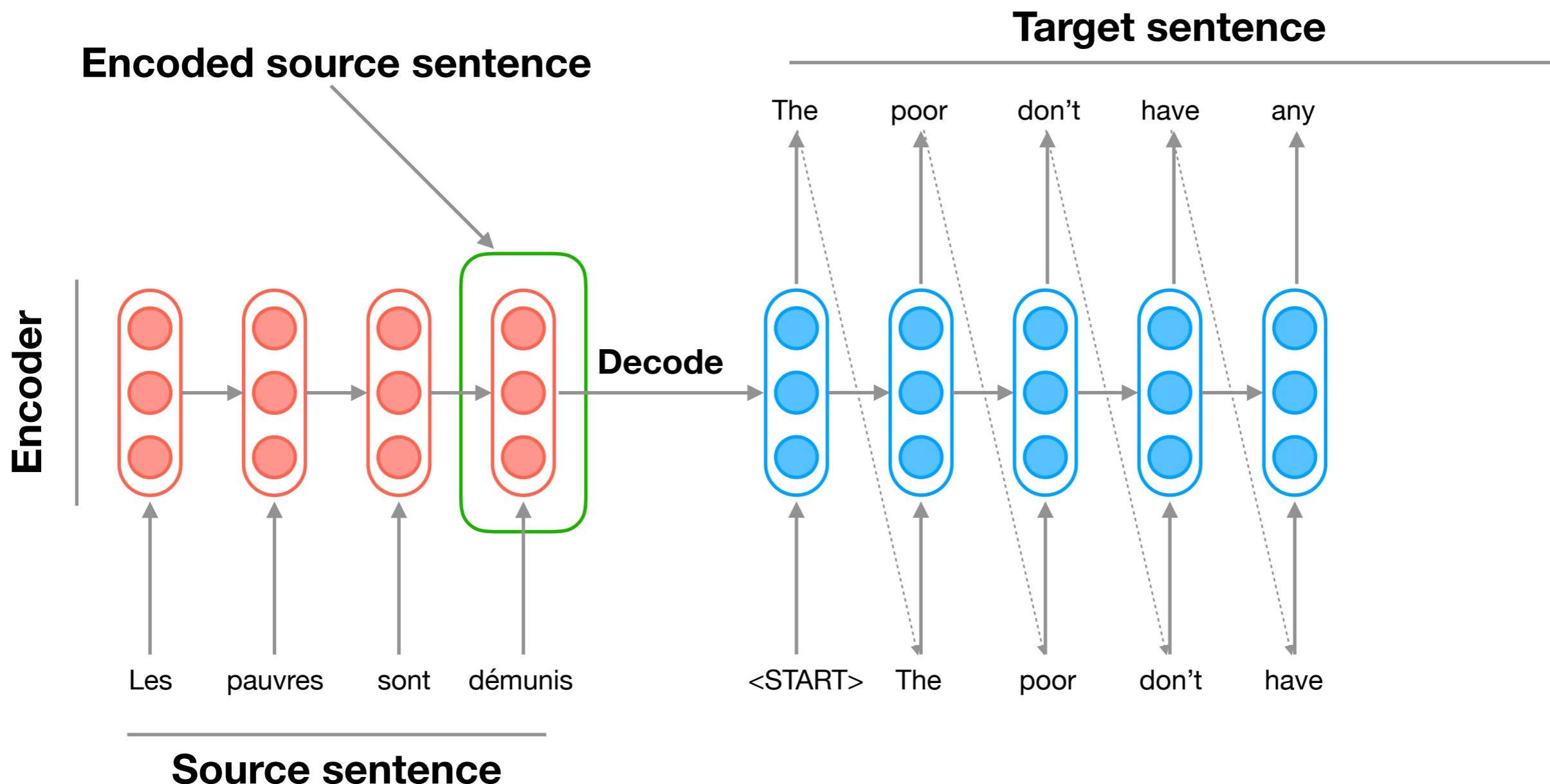
# Sequence to sequence

## Inference



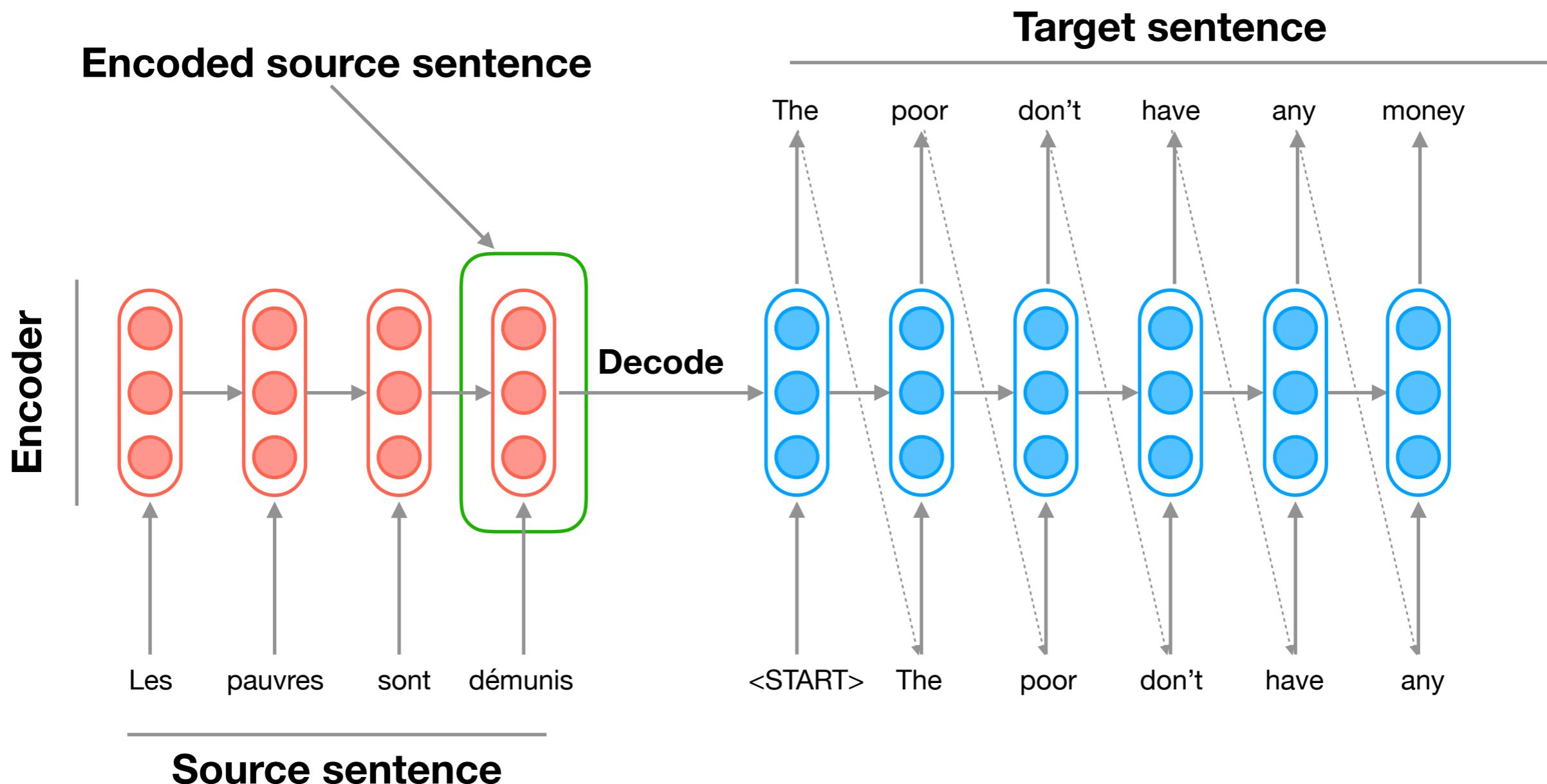
# Sequence to sequence

## Inference



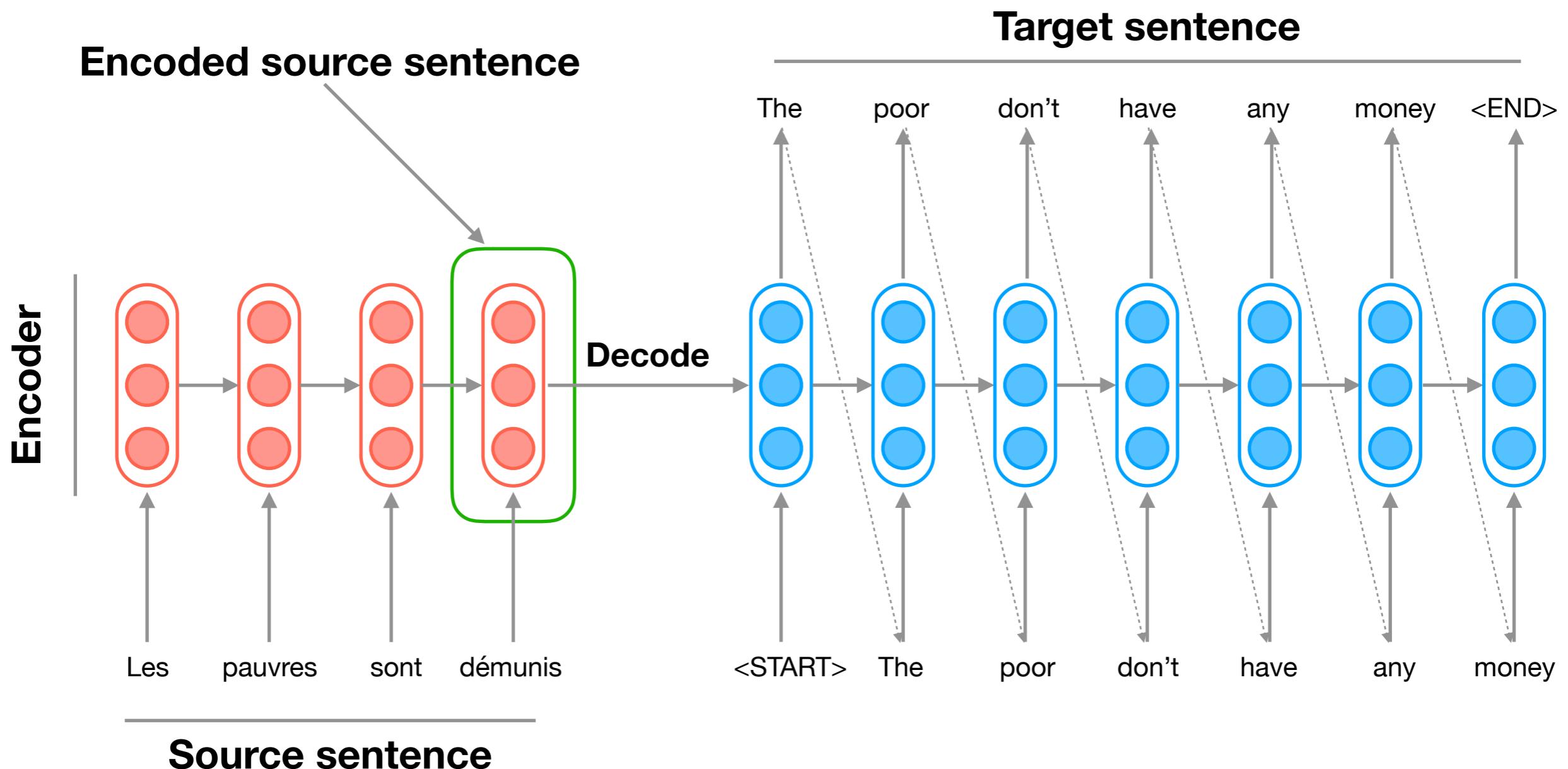
# Sequence to sequence

## Inference



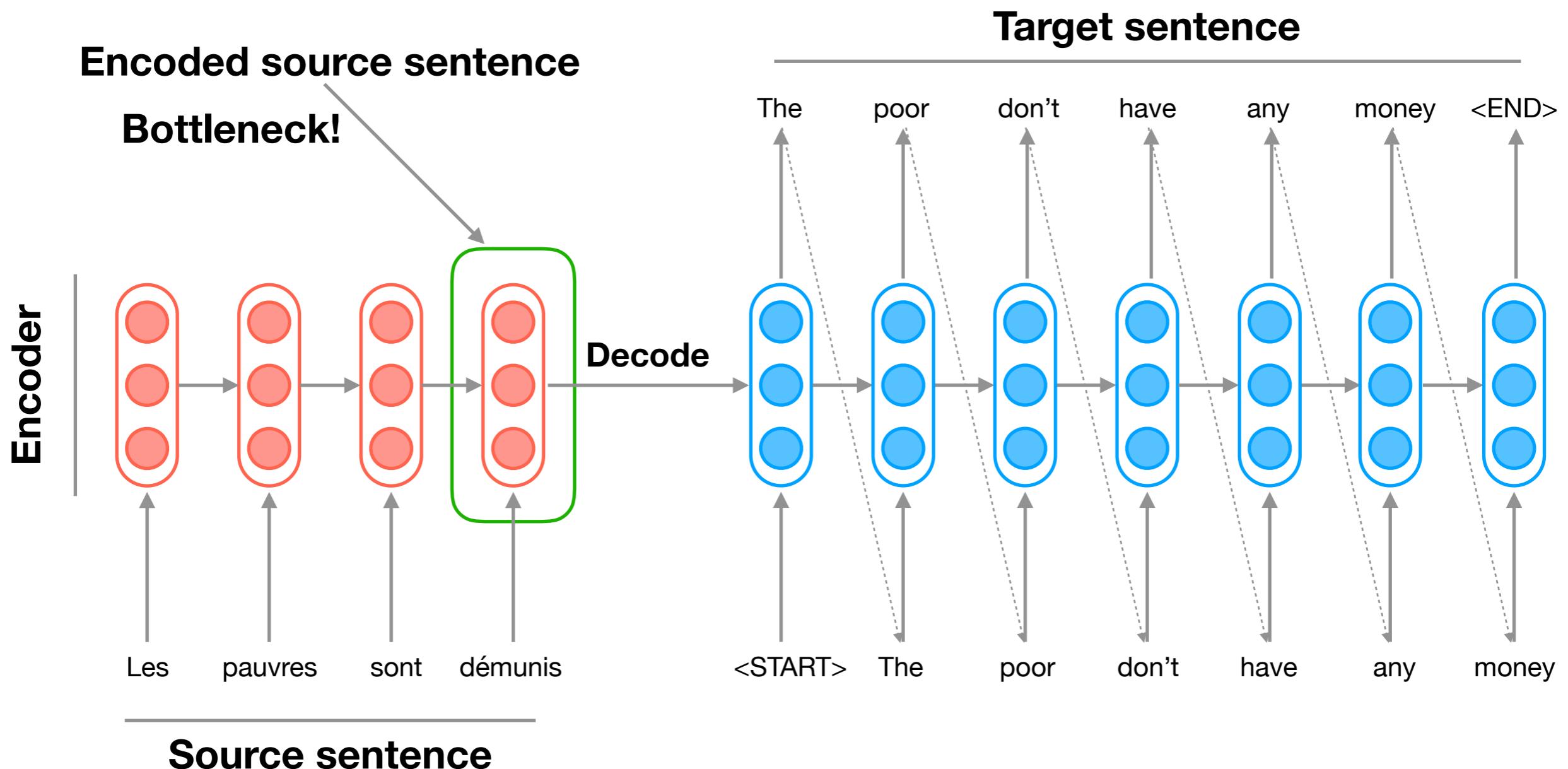
# Sequence to sequence

## Inference



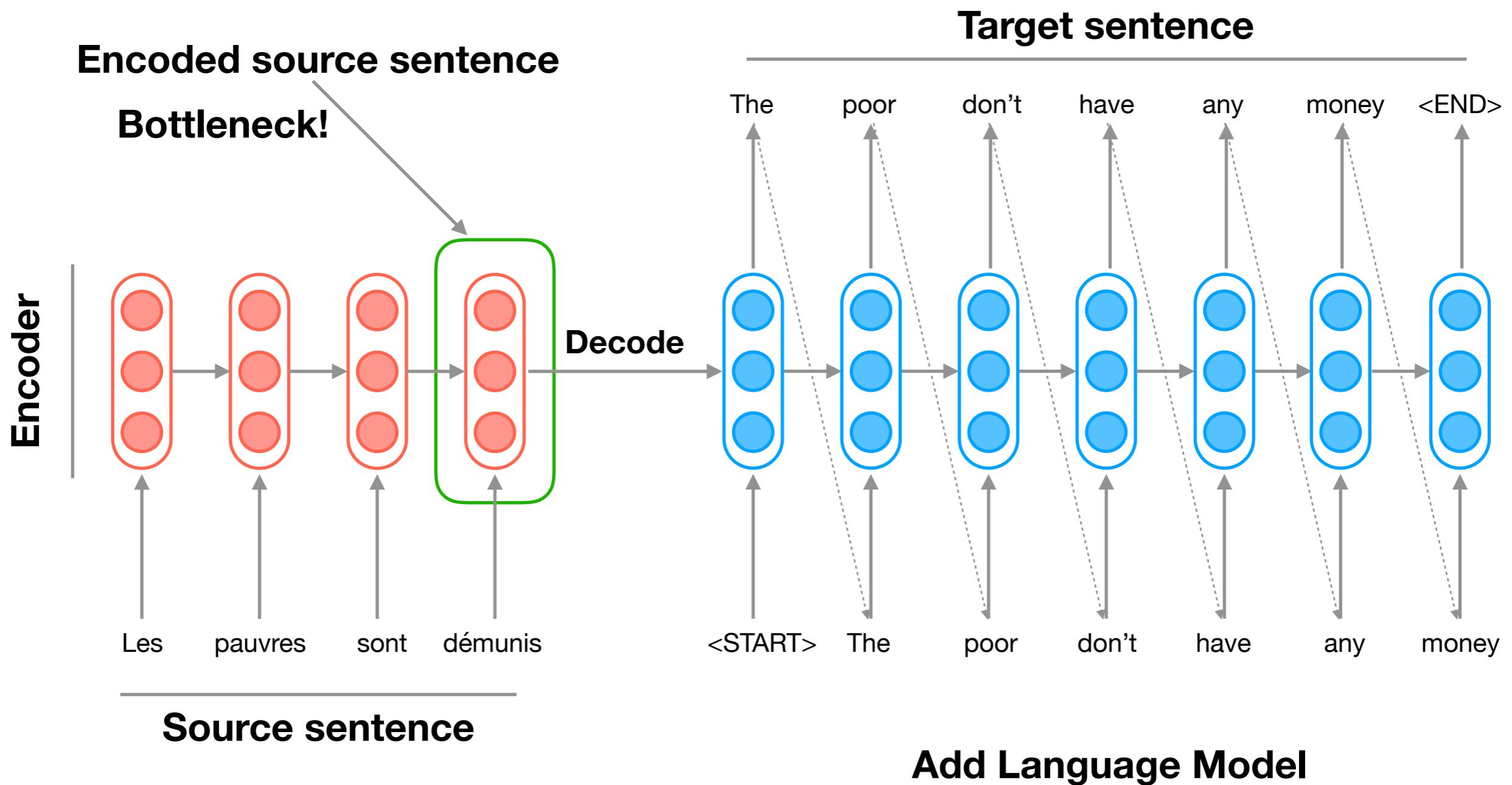
# Sequence to sequence

## Inference



# Sequence to sequence

## Inference



# Contextualized Word Vectors

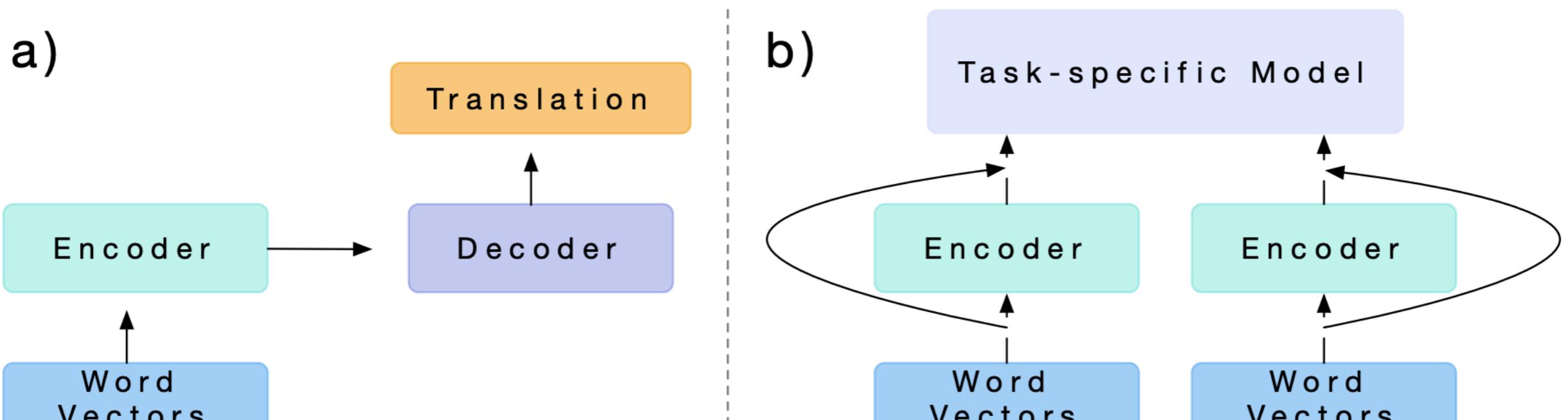


Figure 1: We a) train a two-layer, bidirectional LSTM as the encoder of an attentional sequence-to-sequence model for machine translation and b) use it to provide context for other NLP models.

# Thanks for your Attention!

Boris Zubarev



@bobazooba