# NLP LM Transfer Learning

**Featuring Model**

Your typical classifier

**Contextualized word embeddings**

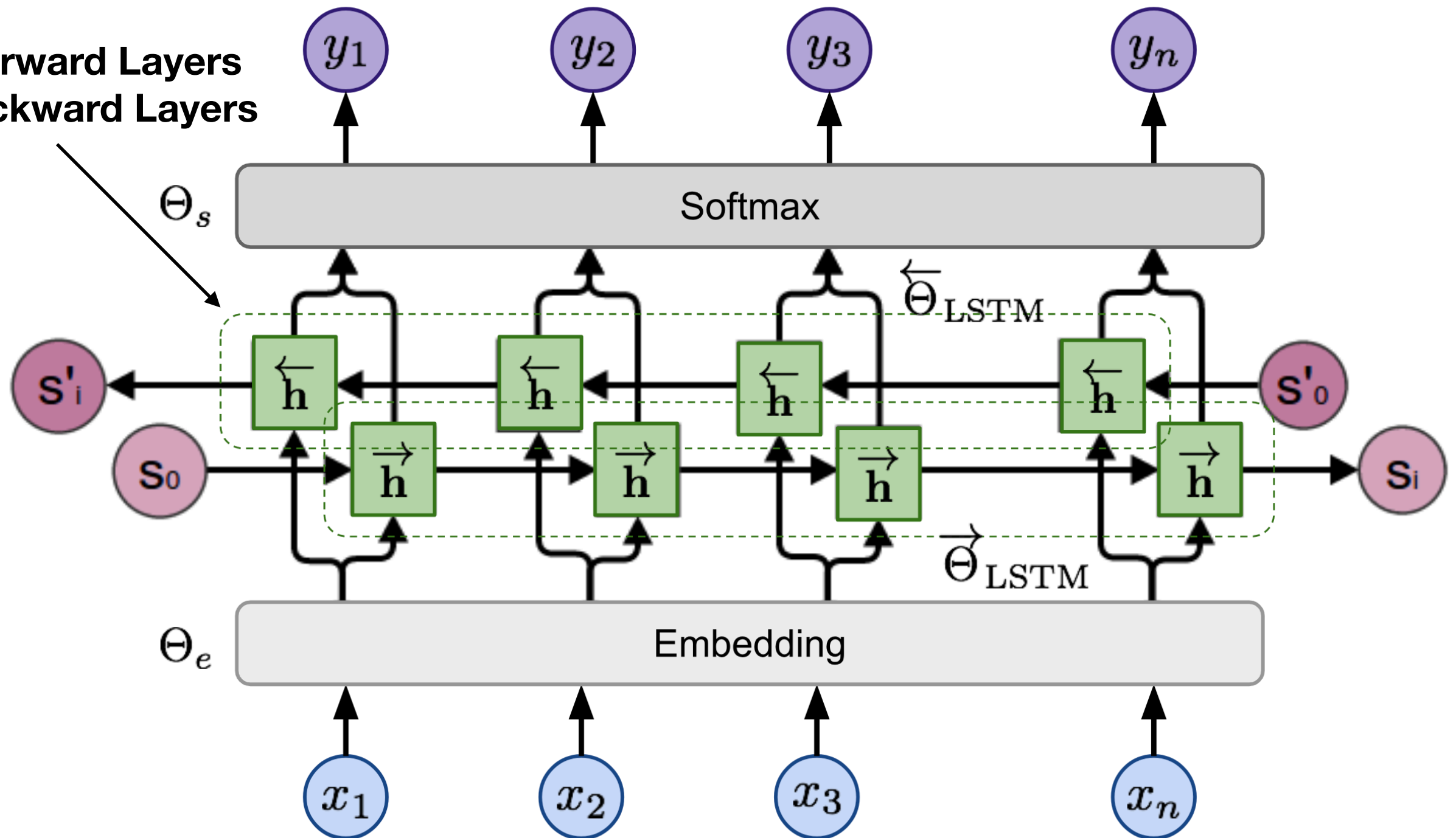**Frozen LSTM Language Model**

**Tokenized embedded text**

# ELMo

# ELMo

# ELMo

# ELMo

# ELMo

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Source** | <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> |

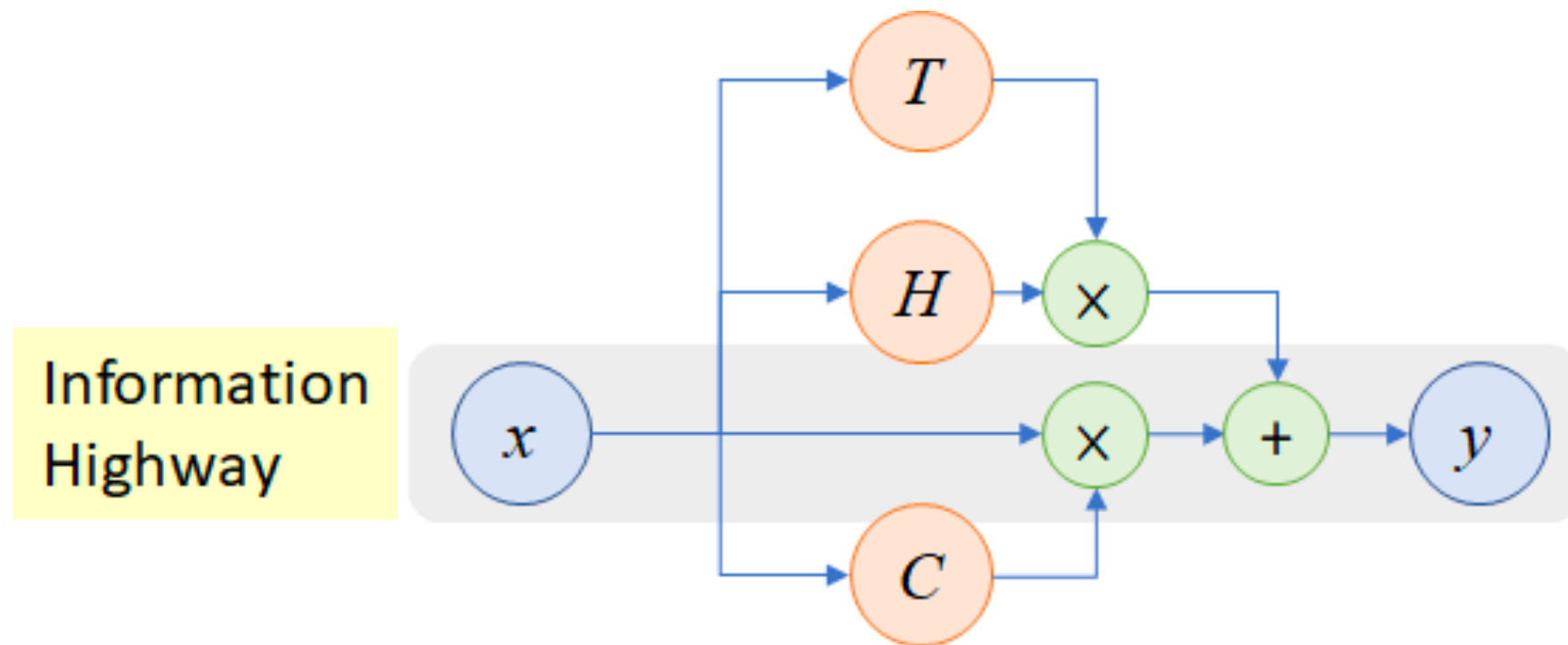# Character-Aware Neural Language Models
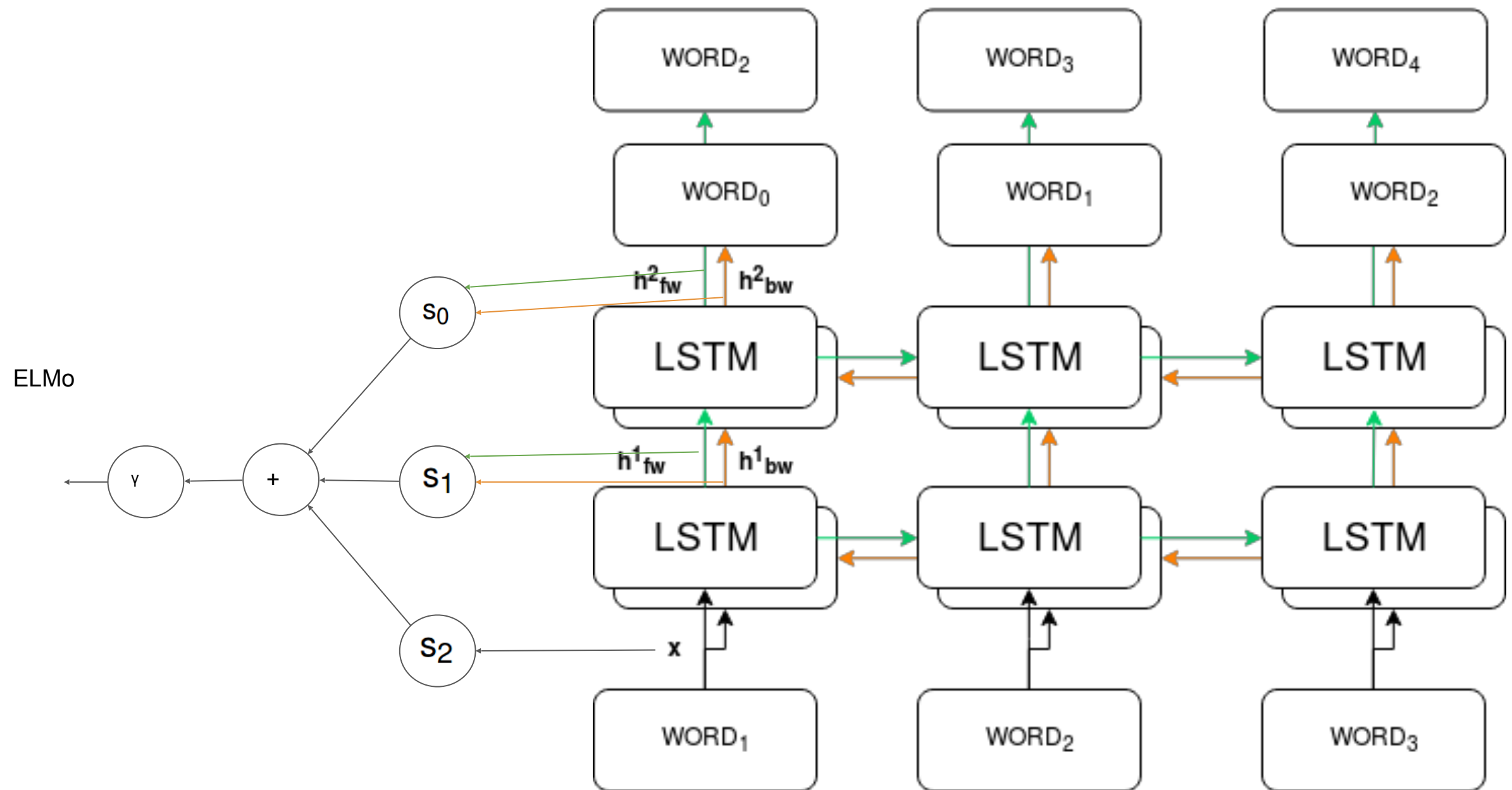
# Character-Aware Neural Language Models
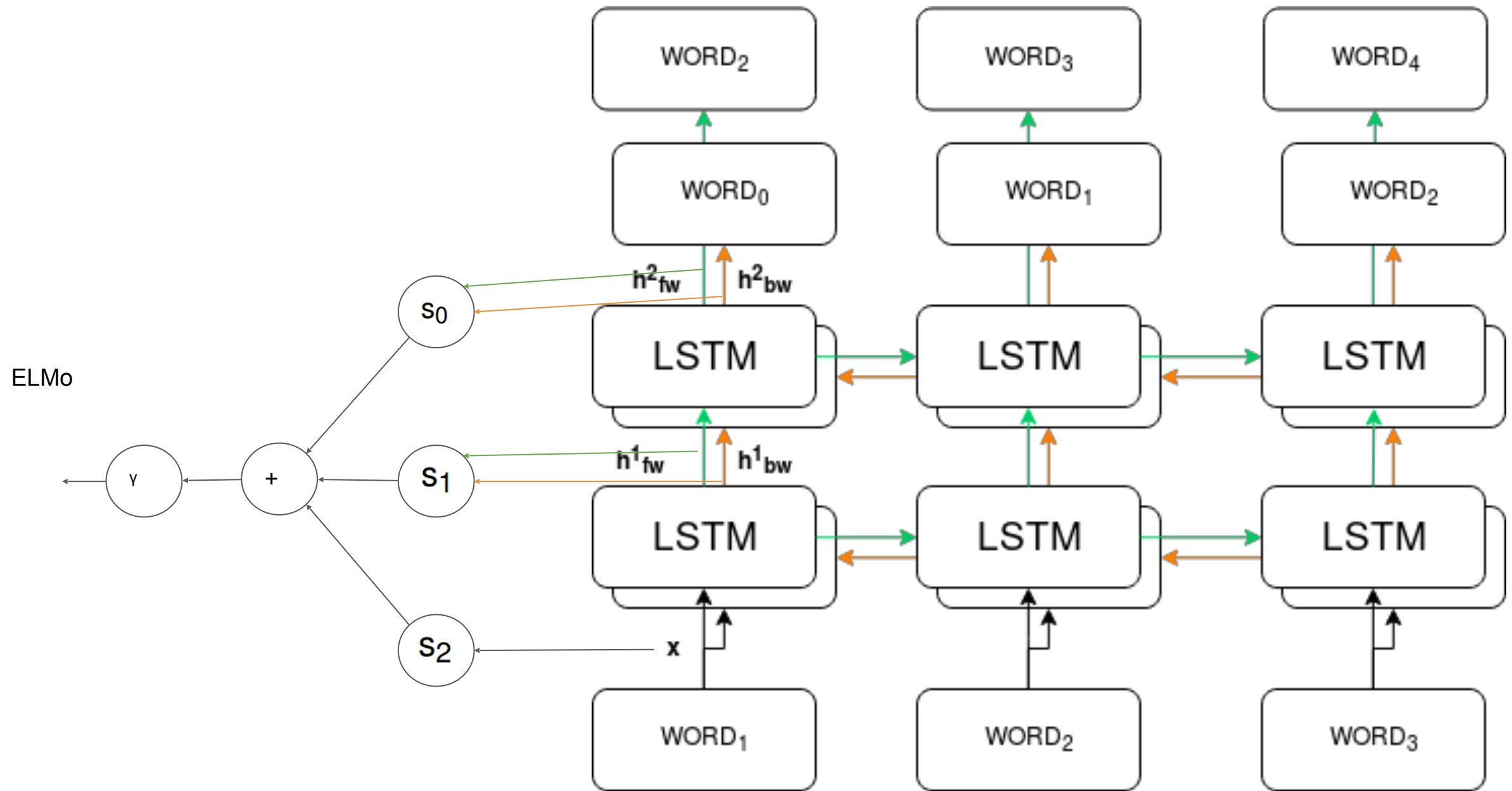
# Highway Network



$$\mathbf{y} = H(\mathbf{x}, \mathbf{W_H}) \cdot T(\mathbf{x}, \mathbf{W_T}) + \mathbf{x} \cdot C(\mathbf{x}, \mathbf{W_C}).$$

$$\mathbf{y} = H(\mathbf{x}, \mathbf{W_H}) \cdot T(\mathbf{x}, \mathbf{W_T}) + \mathbf{x} \cdot (1 - T(\mathbf{x}, \mathbf{W_T})).$$

# ELMo

# ELMo



$$R_k = \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \ldots, L\}$$
$$= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \ldots, L\},$$

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

# ELMo

**Results**



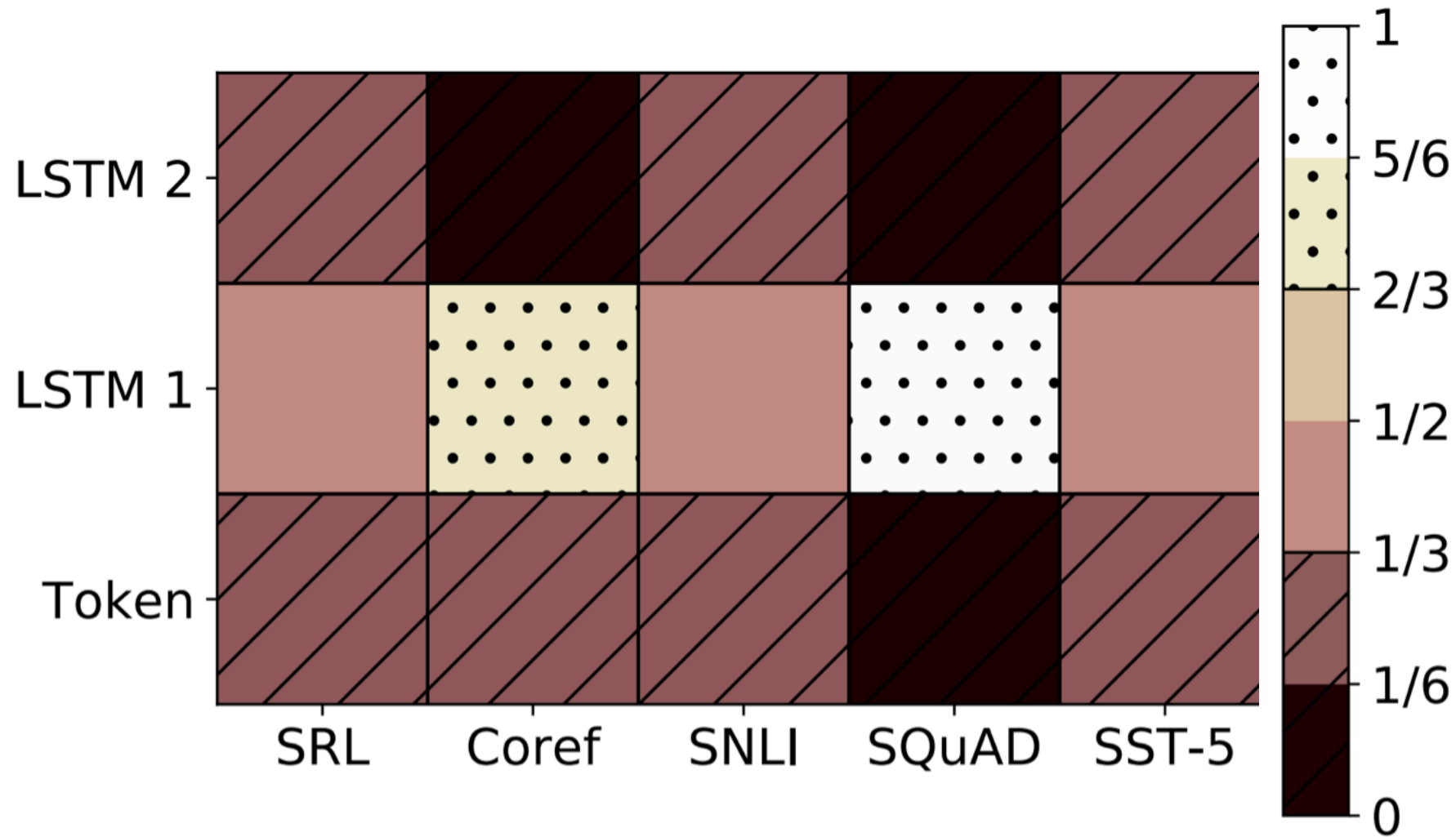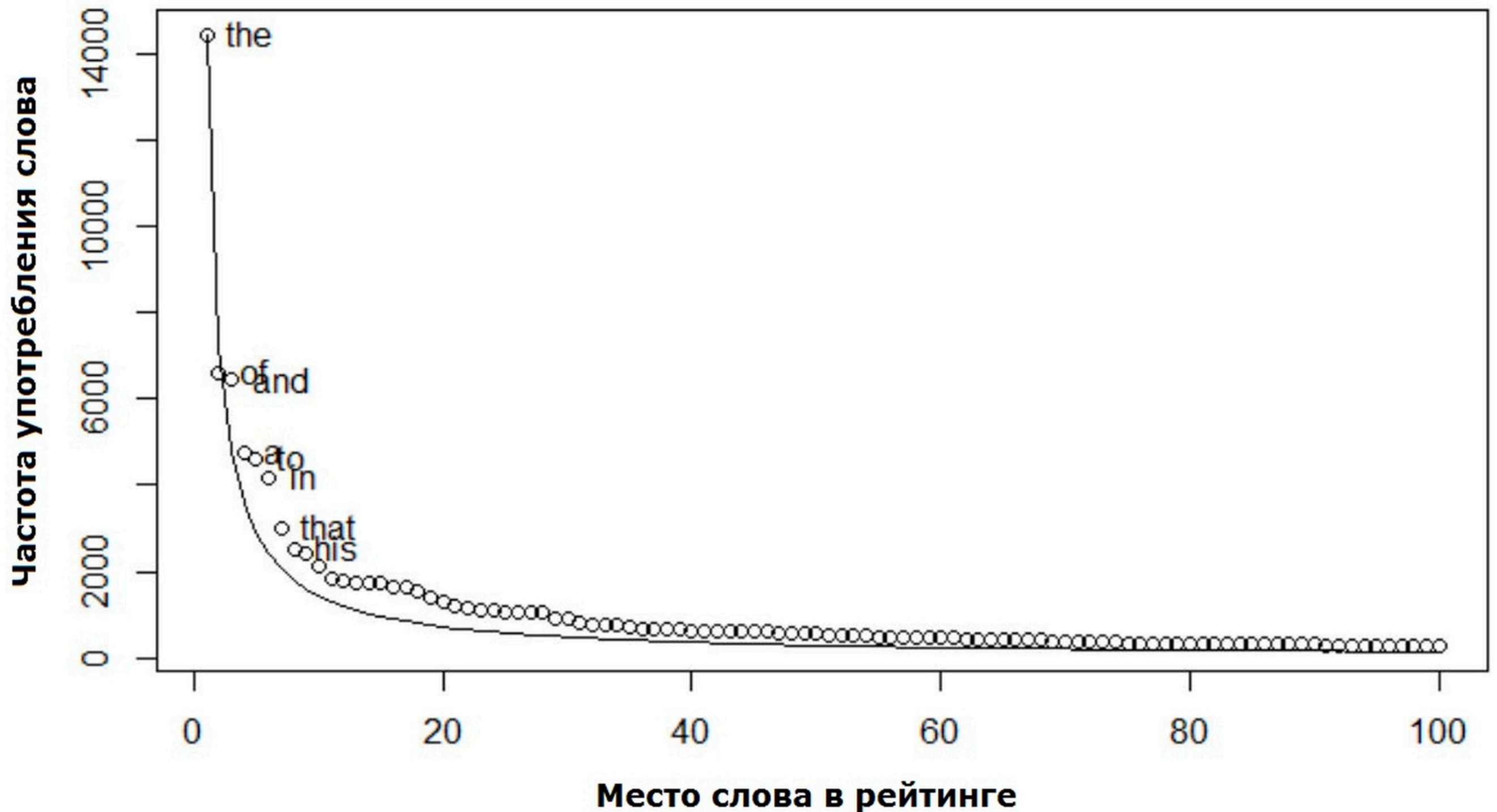Figure 2: Visualization of softmax normalized biLM layer weights across tasks and ELMo locations. Normalized weights less then 1/3 are hatched with horizontal lines and those greater then 2/3 are speckled.

# 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']

# 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

```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'}.

- 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.

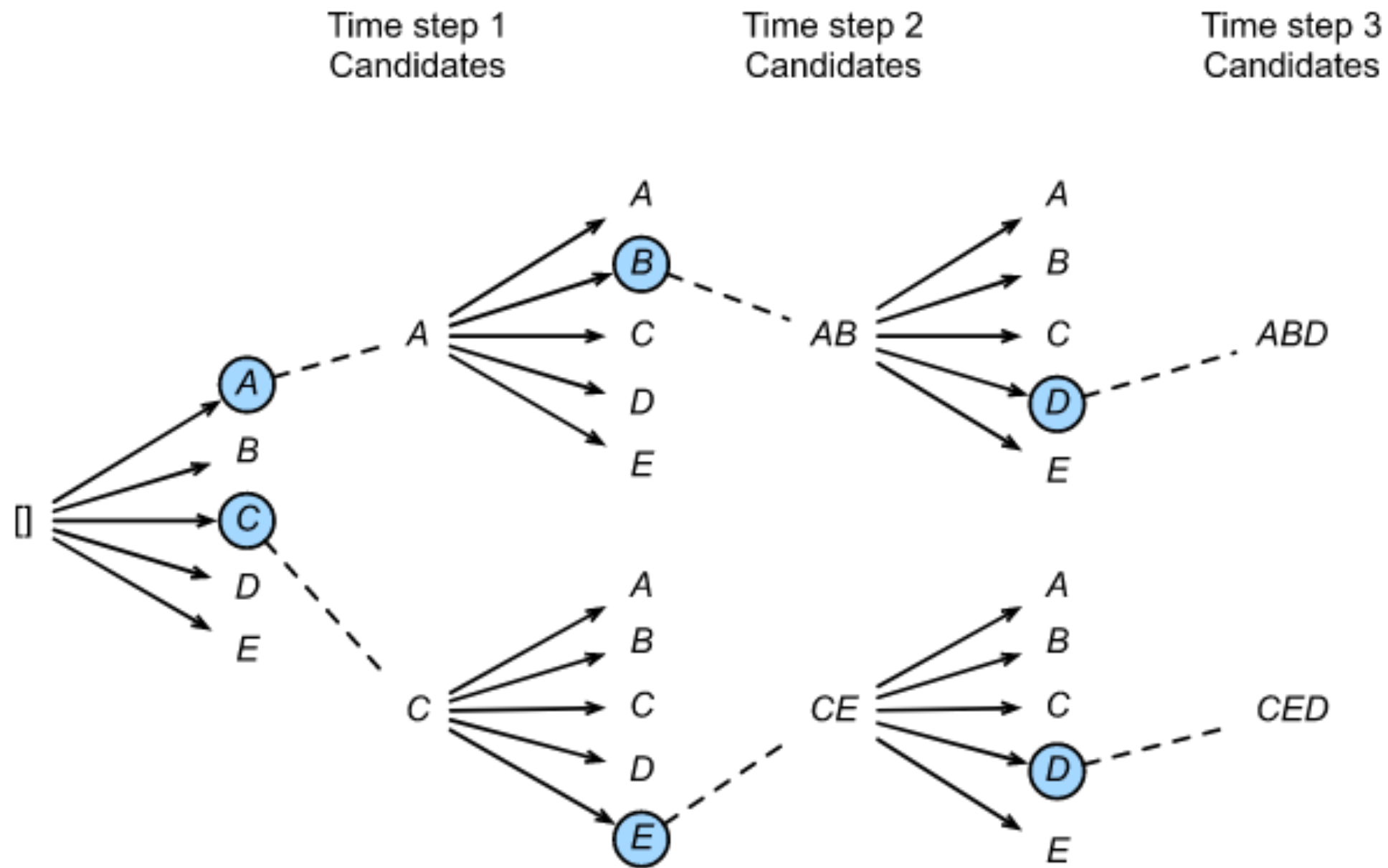**This is not an official Google product.**

## 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++
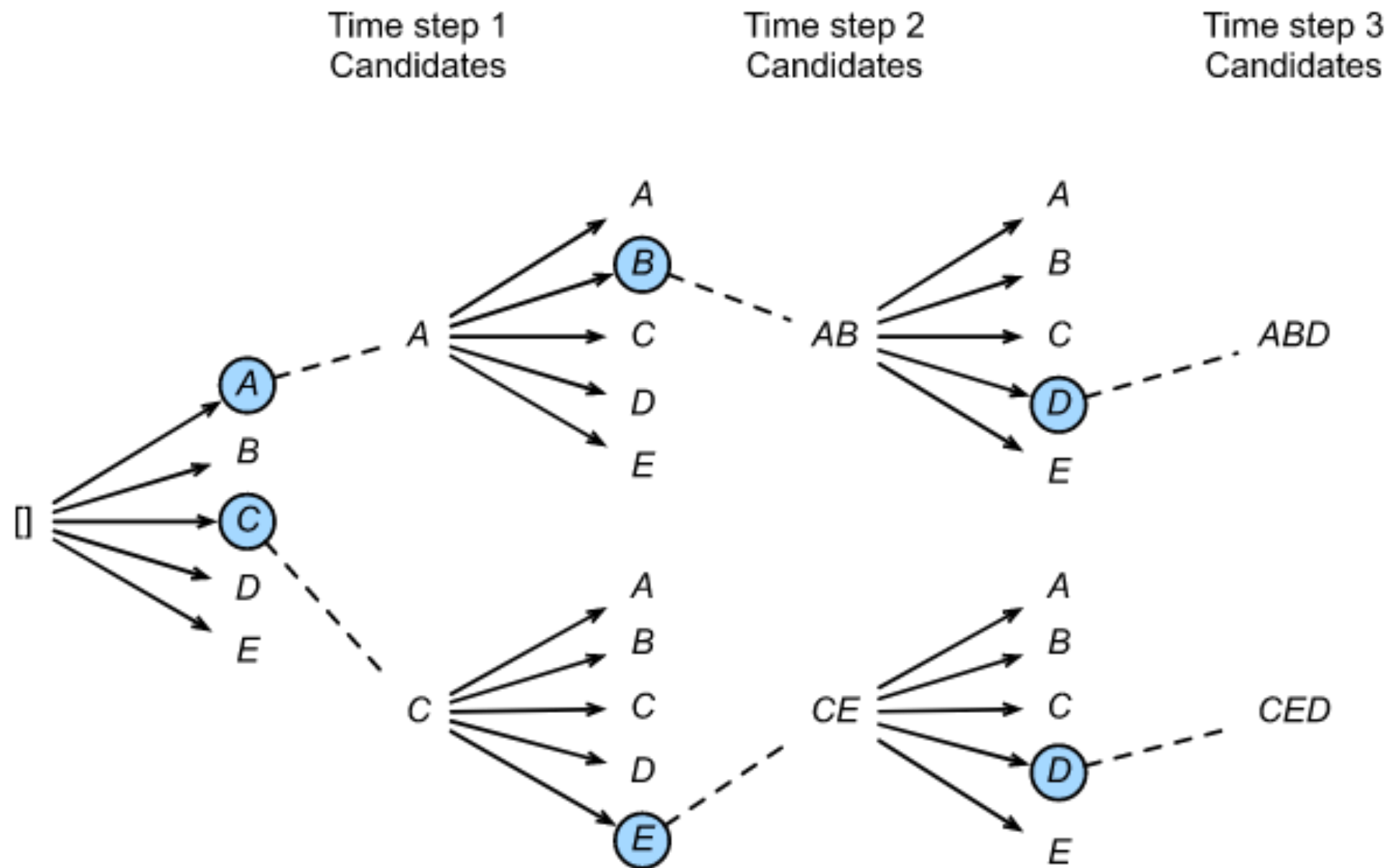- Python wrapper and command-line interface

# Beam Search

# Beam Search

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

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, ..., y^{<t-1>})$$
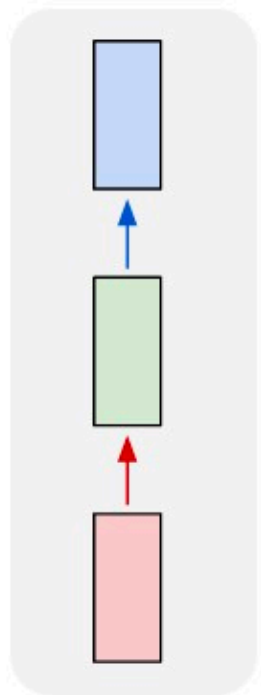
$$\arg \max_y \sum_{y=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, ..., y^{<t-1>})$$
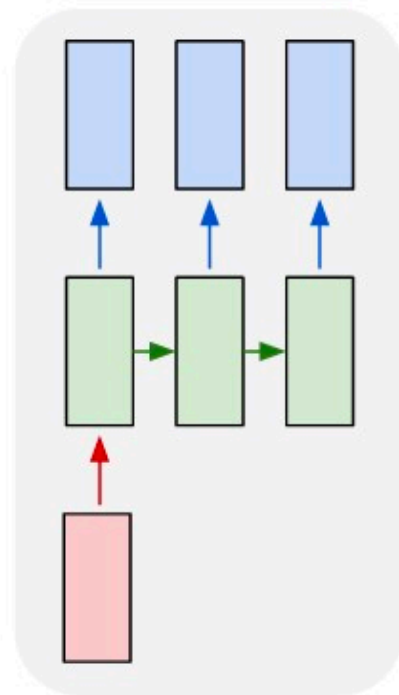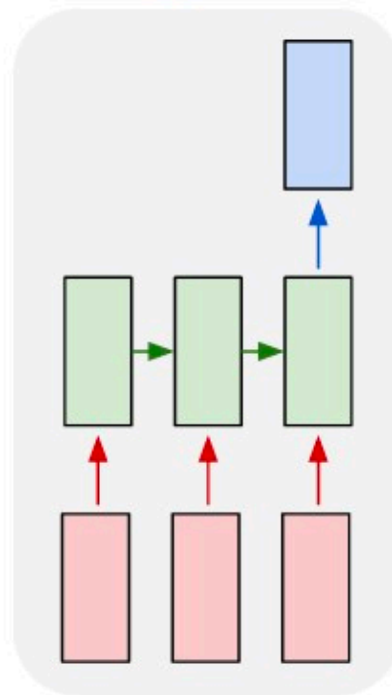
# Beam Search

# Sequence to sequence

# Sequence to sequence

## Inference



Encoded source sentence

Bottleneck!

Target sentence

The  poor  don't  have  any  money  <END>

Encoder

Decode

Les  pauvres  sont  démunis

Source sentence

<START>  The  poor  don't  have  any  money