



Pre-training method in NLP-ELMo

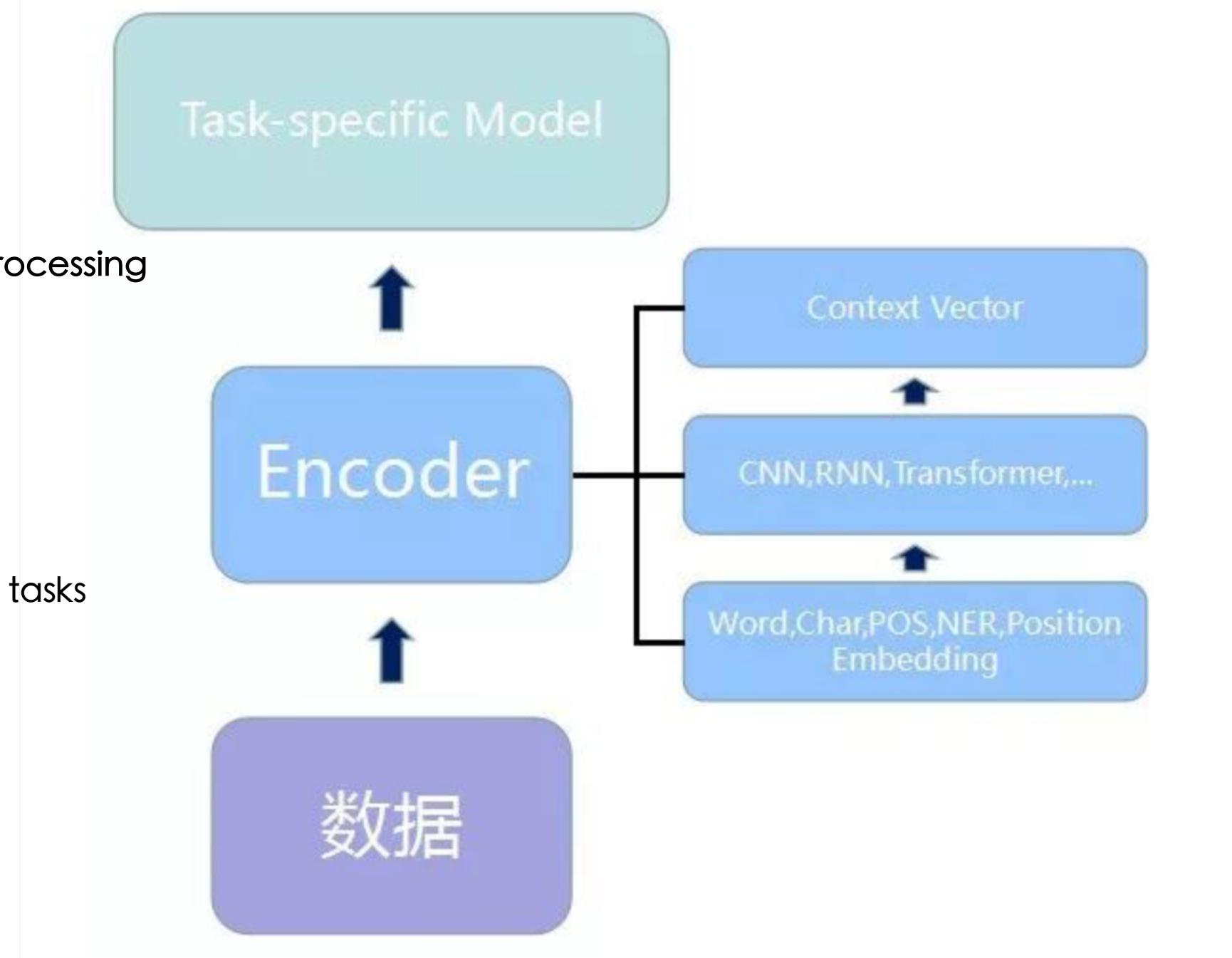
ZhangChenBin

Statistical Machine Intelligence and Learning (SMILE) Lab
University of Electronic Science and Technology of China

outline

- Overview
 - Language model
 - Neural Net Language Model
 - Word2vec
- Analysis of ELMo
 - Character-CNN structure
 - Bidirectional language models
 - EMLO
 - Analysis

- Text data collection and preprocessing
- Encode and characterize text
- Design model to solve specific tasks



Language model

Sentence 1:美联储主席本·伯南克昨天告诉媒体7000亿美元的救助资金

Sentence 2:美主席联储本·伯南克**告诉昨天**媒体7000亿美元的资金救

Sentence 3:美主车席联储本·克告诉昨天公司媒体7000伯南亿美行元



哪个句子更像一个合理的句子？如何量化评估？



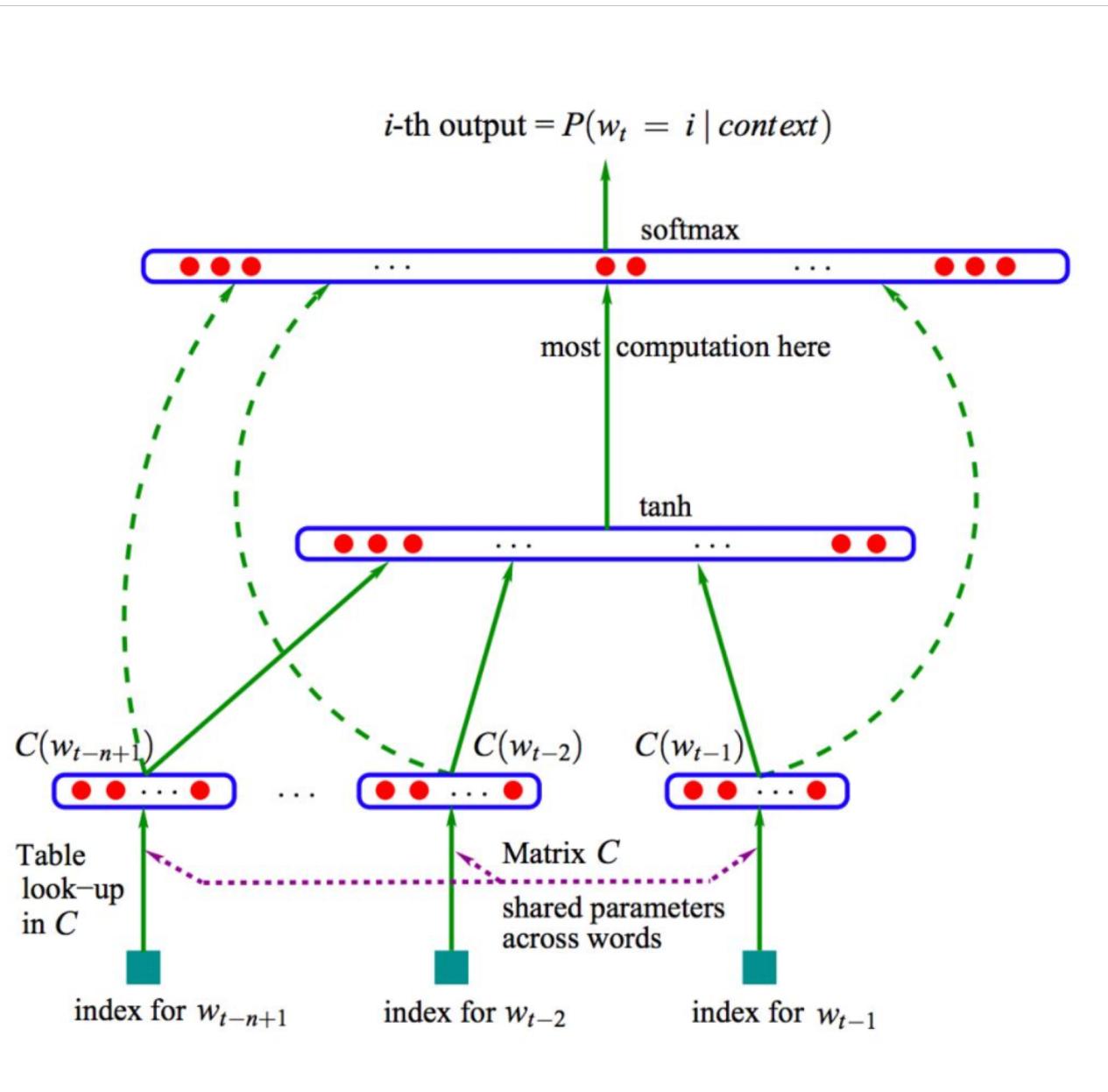
$$P(S) = P(w_1, w_2, \dots, w_n)$$

$$P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \cdots P(w_n|w_1, w_2, \dots, w_{n-1})$$



语言模型： $L = \sum_{w \in C} \log P(w | context(w))$

Neural Net Language Model



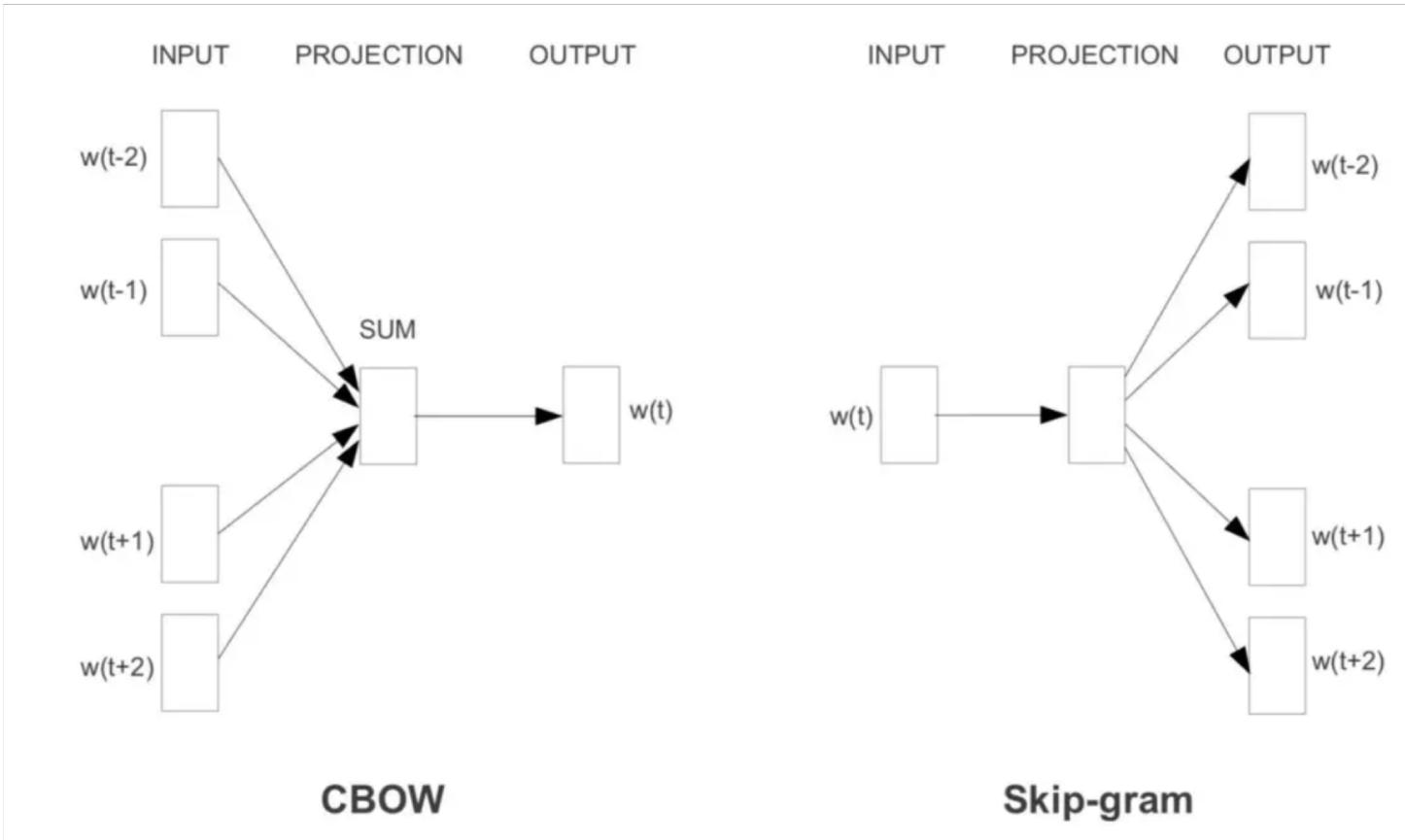
$$L = \frac{1}{T} \sum_t \log P(w_t | w_{t-1}, \dots, w_{t-n+1}; \theta) + R(\theta)$$

$$P(w_t | w_{t-1}, \dots, w_{t-n+1}; \theta) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}}$$

$$y = b + Wx + Utanh(d + Hx)$$

$$x = (C(w_{t-1}), C(w_{t-2}), \dots, C(w_{t-n+1}))$$

Pre-trained word representations -- word2vec



```
In [10]: result = model.most_similar(u"青蛙")
```

```
In [11]: for e in result:  
    print e[0], e[1]
```

```
....:
```

```
老鼠 0.559612870216
```

```
乌龟 0.489831030369
```

```
蜥蜴 0.478990525007
```

```
猫 0.46728849411
```

```
鳄鱼 0.461885392666
```

```
蟾蜍 0.448014199734
```

```
猴子 0.436584025621
```

```
白雪公主 0.434905380011
```

```
蚯蚓 0.433413207531
```

```
螃蟹 0.4314712286
```

```
In [20]: result = model.most_similar(u"清华大学")
```

```
In [21]: for e in result:  
    print e[0], e[1]
```

```
....:
```

```
北京大学 0.763922810555
```

```
化学系 0.724210739136
```

```
物理系 0.694550514221
```

```
数学系 0.684280991554
```

```
中山大学 0.677202701569
```

```
复旦 0.657914161682
```

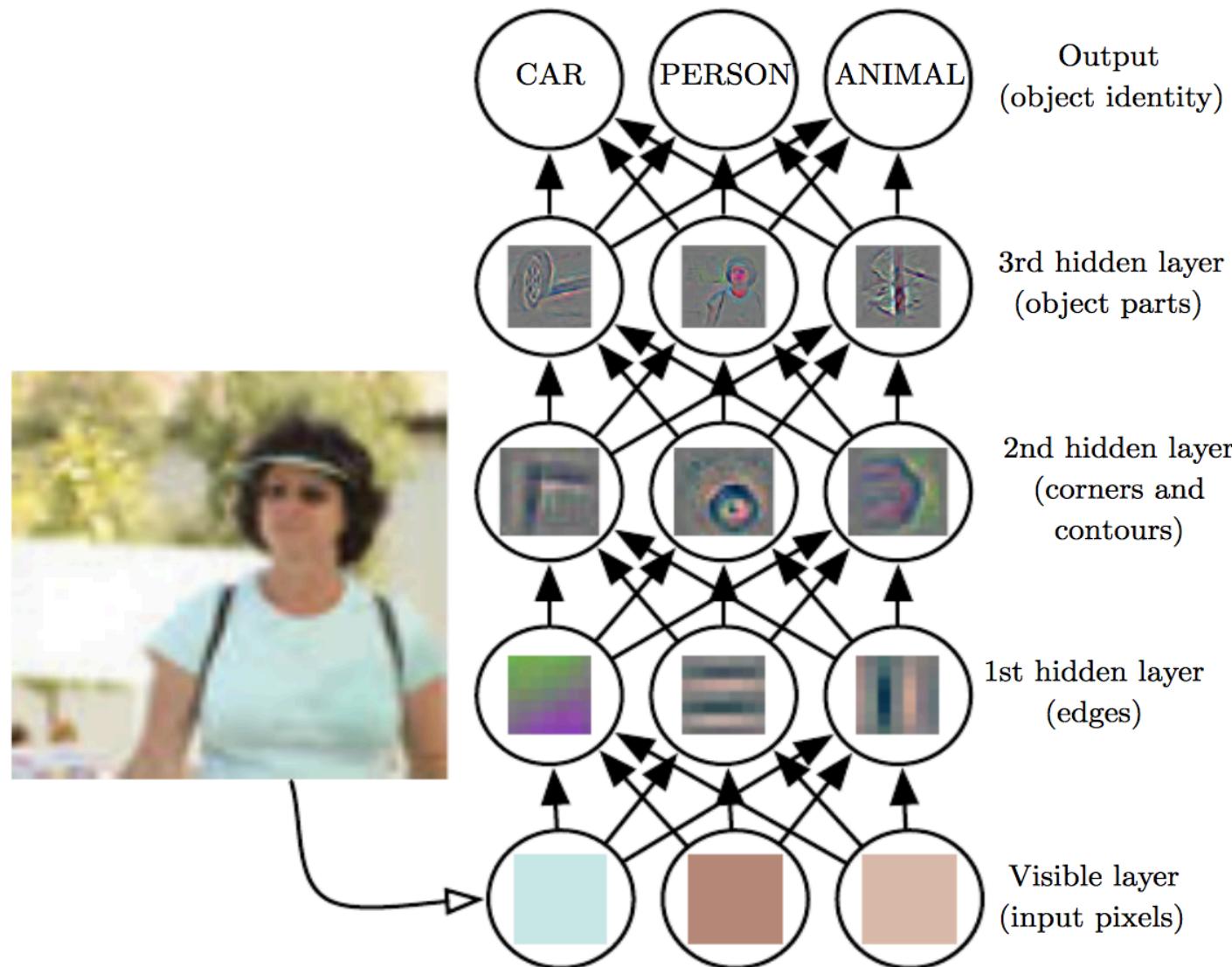
```
师范大学 0.656435549259
```

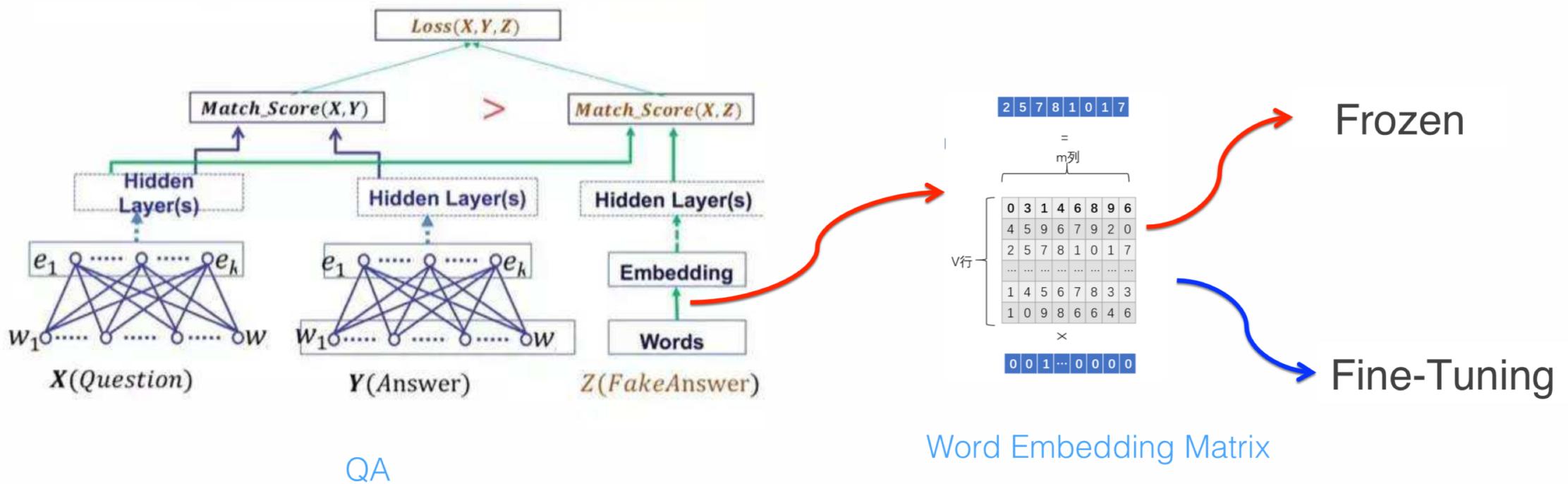
```
哲学系 0.654701948166
```

```
生物系 0.654403865337
```

```
中文系 0.653147578239
```

Pre-training in computer vision





very useful to protect banks or slopes from being washed away by river or raining

the location because it was high, about 100 feet above the bank of river

The bank has plan to branch throughout the country

Do they have the same meaning?

- complex characteristics of word use
 - Syntax
 - Semantics

Question? ?

- how these uses vary across linguistic contexts ?

Deep contextualized word representations

Matthew E. Peters[†], Mark Neumann[†], Mohit Iyyer[†], Matt Gardner[†],
{matthewp, markn, mohiti, mattg}@allenai.org

Christopher Clark^{*}, Kenton Lee^{*}, Luke Zettlemoyer^{†*}
{csquared, kentonl, lsz}@cs.washington.edu

[†]Allen Institute for Artificial Intelligence

^{*}Paul G. Allen School of Computer Science & Engineering, University of Washington

Contextualized Word Representation

- word vectors are **context-independent**. the same word is always the same word vector in different contexts, which obviously leads to the lack of word-disambiguation (WSD) in the word vector model
- People also realized that the **char-level text** also contains some patterns that are difficult to describe by **word-level text**.



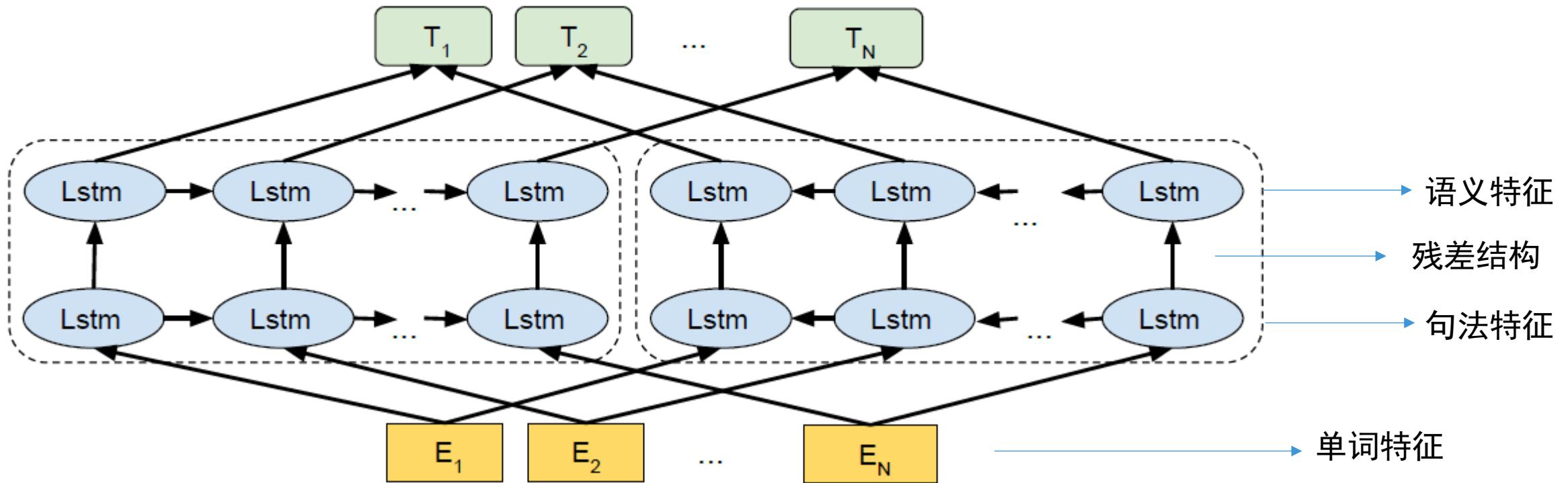
Hey ELMo, what's the embedding
of the word "stick"?

There are multiple possible
embeddings! Use it in a sentence.

Oh, okay. Here:
"Let's stick to improvisation in this
skit"

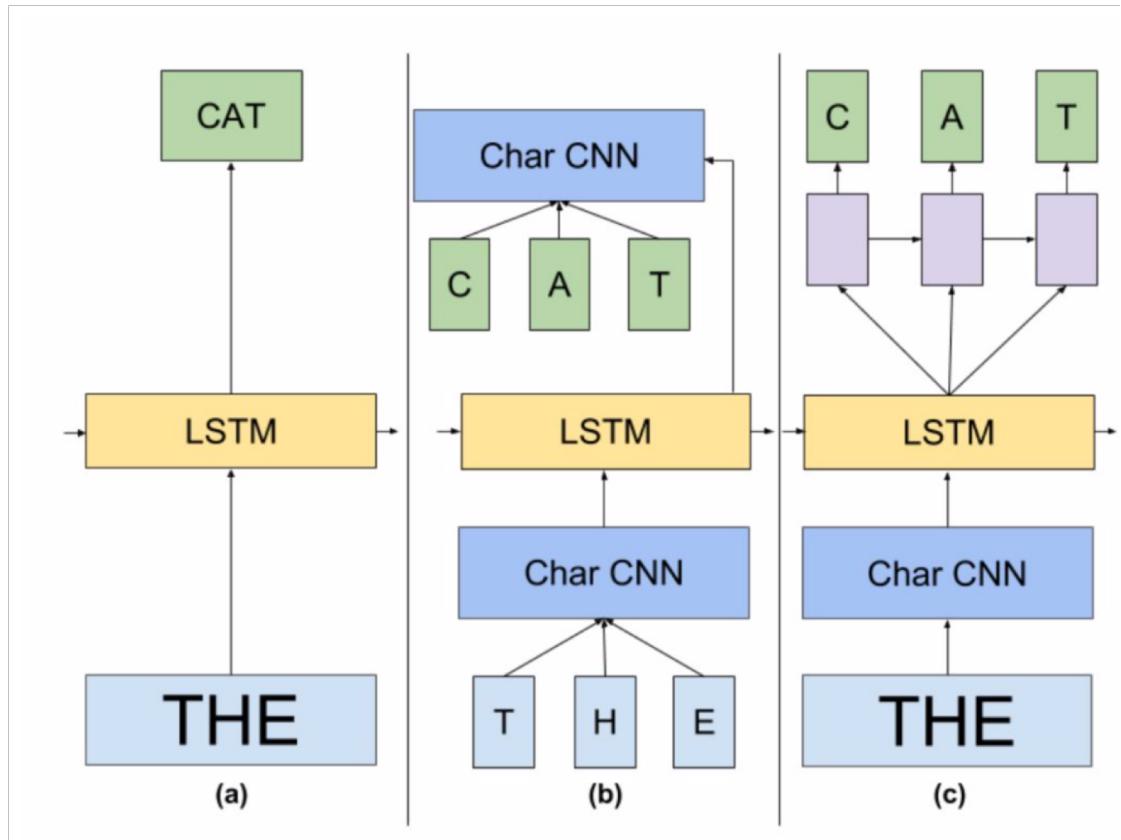
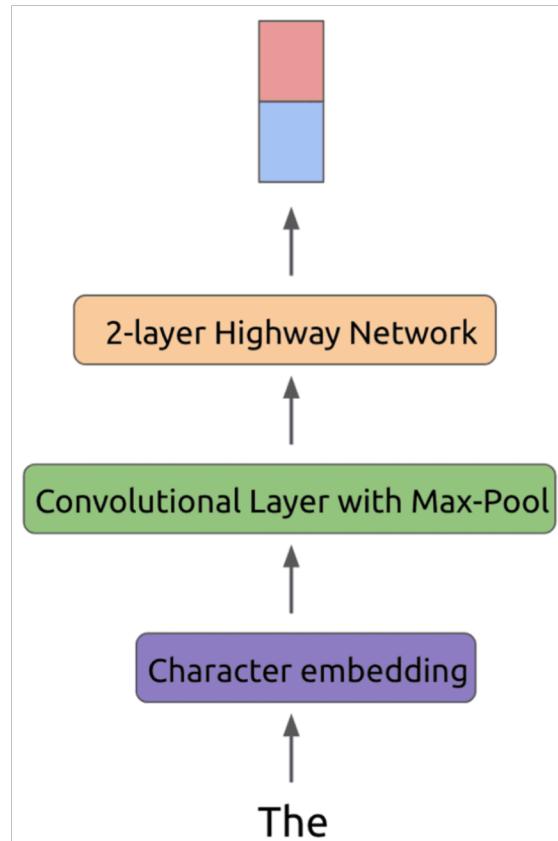
Oh in that case, the embedding is:
-0.02, -0.16, 0.12, -0.1etc

ELMo



- two-layer biLMs with character convolutions
- Different layer represent different information

Encoder: Character-CNN structure



```
1 #将词转化为char_ids
2 def _convert_word_to_char_ids(self, word):
3     code = np.zeros([self.max_word_length], dtype=np.int32)
4     code[:] = self.pad_char
5
6     #将word中每一个字符转化为utf-8编码，然后用数组存起来，例如：
7     #english中， e:101, n:110, g:103, l:108, h:105, s:115, h:104
8     word_encoded = word.encode('utf-8', 'ignore')[:(self.max_word_length-2)]
9     code[0] = self.bow_char
10    #加上词开始和结尾的编码
11    for k, chr_id in enumerate(word_encoded, start=1):
12        code[k] = chr_id
13    code[k + 1] = self.eow_char
14
15    return code
```

```
18 def __init__(self, filename, max_word_length, **kwargs):
19     #调用父类Vocabulary, 生成word和id之间的转换等
20     super(UnicodeCharsVocabulary, self).__init__(filename, **kwargs)
21     #每个词对应最大字符长
22     self._max_word_length = max_word_length
23
24     # char ids 0-255 come from utf-8 encoding bytes
25     # assign 256-300 to special chars
26     self.bos_char = 256 # <begin sentence>
27     self.eos_char = 257 # <end sentence>
28     self.bow_char = 258 # <begin word>
29     self.eow_char = 259 # <end word>
30     self.pad_char = 260 # <padding>
31
32     #单词的个数, 父类中的属性
33     num_words = len(self._id_to_word)
34
35     #每个词都会对应一个char_ids列表
36     self._word_char_ids = np.zeros([num_words, max_word_length],
37                                   dtype=np.int32)
38
39     # the character representation of the begin/end of sentence characters
40     # 对句首或者句尾的token来一个字符的表示
41     def _make_bos_eos(c):
42         r = np.zeros([self.max_word_length], dtype=np.int32)
43         r[:] = self.pad_char
44         #词的开始
45         r[0] = self.bow_char
46         r[1] = c
47         #词的结束
48         r[2] = self.eow_char
49         return r
50     #句子开始对应的char_ids
51     self.bos_chars = _make_bos_eos(self.bos_char)
52     #句子的结尾对应的char_ids
53     self.eos_chars = _make_bos_eos(self.eos_char)
54     #遍历id2word数组, 得到每一个词的char_ids
55     for i, word in enumerate(self._id_to_word):
56         self._word_char_ids[i] = self._convert_word_to_char_ids(word)
57     #将句子开头和结尾当作一个word处理
58     self._word_char_ids[self.bos] = self.bos_chars
59     self._word_char_ids[self.eos] = self.eos_chars
60
```

```
1 #返回word对应的char_ids数组
2 def word_to_char_ids(self, word):
3     if word in self._word_to_id:
4         return self._word_char_ids[self._word_to_id[word]]
5     else:
6         return self._convert_word_to_char_ids(word)
7
8 def encode_chars(self, sentence, reverse=False, split=True):
9     ...
10    Encode the sentence as a white space delimited string of tokens.
11    对一整句话进行编码，编码成chars
12    ...
13    if split:                                              #如果切割了句子
14        chars_ids = [self.word_to_char_ids(cur_word)
15                      for cur_word in sentence.split()]
16    else:
17        chars_ids = [self.word_to_char_ids(cur_word)
18                      for cur_word in sentence]
19    if reverse:
20        return np.vstack([self.eos_chars] + chars_ids + [self.bos_chars]) #在每一条句子上都加了<eos>和<bos>
21    else:
22        return np.vstack([self.bos_chars] + chars_ids + [self.eos_chars])
```

```

1 def _build_word_char_embeddings(self):
2     # 设置初始的batchsize
3     batch_size = self.options['batch_size']
4
5     unroll_steps = self.options['unroll_steps']
6     #最终维度，投影层维度
7     projection_dim = self.options['lstm']['projection_dim']
8
9     cnn_options = self.options['char_cnn']
10    filters = cnn_options['filters']
11    #一个卷积核抽取了一组特征，总的特征组数
12    n_filters = sum(f[1] for f in filters)
13    #每个词最大字符数
14    max_chars = cnn_options['max_characters_per_token']
15    # 映射成词向量的维度
16    char_embed_dim = cnn_options['embedding']['dim']
17    #字典中字符个数
18    n_chars = cnn_options['n_characters']
19    if n_chars != 261:
20        raise InvalidNumberOfCharacters(
21            "Set n_characters=261 for training see the README.md"
22        )
23    if cnn_options['activation'] == 'tanh':
24        activation = tf.nn.tanh
25    elif cnn_options['activation'] == 'relu':
26        activation = tf.nn.relu

```

```

28     #字符的输入
29     self.tokens_characters = tf.placeholder(DTYPE_INT,
30                                              shape=(batch_size, unroll_steps, max_chars),
31                                              name='tokens_characters')
32
33     # 对字符进行词向量的嵌入
34     with tf.device("/cpu:0"):
35         self.embedding_weights = tf.get_variable(
36             "char_embed", [n_chars, char_embed_dim],
37             dtype=DTYPE,
38             initializer=tf.random_uniform_initializer(-1.0, 1.0)
39         )
40         # shape (batch_size, unroll_steps, max_chars, embed_dim)
41         self.char_embedding = tf.nn.embedding_lookup(self.embedding_weights,
42                                                       self.tokens_characters)
43
44     ## 将字符反向，为下面的双向lstm作准备
45     if self.bidirectional:
46         self.tokens_characters_reverse = tf.placeholder(DTYPE_INT,
47                                                       shape=(batch_size, unroll_steps, max_chars),
48                                                       name='tokens_characters_reverse')
49         self.char_embedding_reverse = tf.nn.embedding_lookup(
50             self.embedding_weights, self.tokens_characters_reverse)

```

```

1 # the convolutions
2 def make_convolution(inp, reuse):
3     with tf.variable_scope('CNN', reuse=reuse) as scope:
4         convolutions = []
5         for i, (width, num) in enumerate(filters):
6             # 优化方法
7             if cnn_options['activation'] == 'relu':
8                 w_init = tf.random_uniform_initializer(
9                     minval=-0.05, maxval=0.05)
10            elif cnn_options['activation'] == 'tanh':
11                w_init = tf.random_normal_initializer(
12                    mean=0.0,
13                    stddev=np.sqrt(1.0 / (width * char_embed_dim)))
14            )
15            w = tf.get_variable(
16                "W_cnn_%s" % i,
17                [1, width, char_embed_dim, num],
18                initializer=w_init,
19                dtype=DTYPE)
20            b = tf.get_variable(
21                "b_cnn_%s" % i, [num], dtype=DTYPE,
22                initializer=tf.constant_initializer(0.0))
23
24            conv = tf.nn.conv2d(
25                inp, w,
26                strides=[1, 1, 1, 1],
27                padding="VALID") + b
28            # now max pool
29            conv = tf.nn.max_pool(
30                conv, [1, 1, max_chars-width+1, 1],
31                [1, 1, 1, 1], 'VALID')
32
33            # activation
34            conv = activation(conv)
35            # 将维数为1的维度去掉，保留数值数据。
36            conv = tf.squeeze(conv, squeeze_dims=[2])
37
38            convolutions.append(conv)
39
40    return tf.concat(convolutions, 2)

```

```

1 reuse = tf.get_variable_scope().reuse
2 embedding = make_convolution(self.char_embedding, reuse)
3 self.token_embedding_layers = [embedding]

```

```

1 #highway设计结构
2 def high(x, ww_carry, bb_carry, ww_tr, bb_tr):
3     carry_gate = tf.nn.sigmoid(tf.matmul(x, ww_carry) + bb_carry)
4     transform_gate = tf.nn.relu(tf.matmul(x, ww_tr) + bb_tr)
5     return carry_gate * transform_gate + (1.0 - carry_gate) * x
6
7 if use_highway:
8     highway_dim = n_filters
9
10 for i in range(n_highway):
11     with tf.variable_scope('CNN_high_%s' % i) as scope:
12         W_carry = tf.get_variable(
13             'W_carry', [highway_dim, highway_dim],
14             # glorit init
15             initializer=tf.random_normal_initializer(
16                 mean=0.0, stddev=np.sqrt(1.0 / highway_dim)),
17             dtype=DTYPE)
18         b_carry = tf.get_variable(
19             'b_carry', [highway_dim],
20             initializer=tf.constant_initializer(-2.0),
21             dtype=DTYPE)
22         W_transform = tf.get_variable(
23             'W_transform', [highway_dim, highway_dim],
24             initializer=tf.random_normal_initializer(
25                 mean=0.0, stddev=np.sqrt(1.0 / highway_dim)),
26             dtype=DTYPE)
27         b_transform = tf.get_variable(
28             'b_transform', [highway_dim],
29             initializer=tf.constant_initializer(0.0),
30             dtype=DTYPE)
31
32         embedding = high(embedding, W_carry, b_carry,
33                           W_transform, b_transform)
34     if self.bidirectional:
35         embedding_reverse = high(embedding_reverse,
36                                 W_carry, b_carry,
37                                 W_transform, b_transform)
38     self.token_embedding_layers.append(
39         tf.reshape(embedding,
40                   [batch_size, unroll_steps, highway_dim]))
41

```



```

1 # 进行投影操作
2         if use_proj:
3             embedding = tf.matmul(embedding, W_proj_cnn) + b_proj_cnn
4             if self.bidirectional:
5                 embedding_reverse = tf.matmul(embedding_reverse, W_proj_cnn)
6                     + b_proj_cnn
7             self.token_embedding_layers.append(
8                 tf.reshape(embedding,
9                         [batch_size, unroll_steps, projection_dim]))
10
11
12 # 从newrshape成(batch_size, tokens, dim)
13 if use_highway or use_proj:
14     shp = [batch_size, unroll_steps, projection_dim]
15     embedding = tf.reshape(embedding, shp)
16     if self.bidirectional:
17         embedding_reverse = tf.reshape(embedding_reverse, shp)
18     self.embedding = embedding
19     if self.bidirectional:
20         self.embedding_reverse = embedding_reverse

```

Why? ?

- 1.CNN can solve OOV (Out-of-Vocabulary) problem
- 2.CNN can pre-calculate each word vector and reduce the computational pressure in the inference phase
- 3.This sharing is reflected in the same set of parameters in the CNN used, which greatly reduces the size of the parameters.

$$P(w_t|c_t) = \frac{\exp(e'(w_t)^T x)}{\sum_{i=1}^{|V|} \exp(e'(w_i)^T x)}, x = \sum_{i \in c} e(w_i)$$

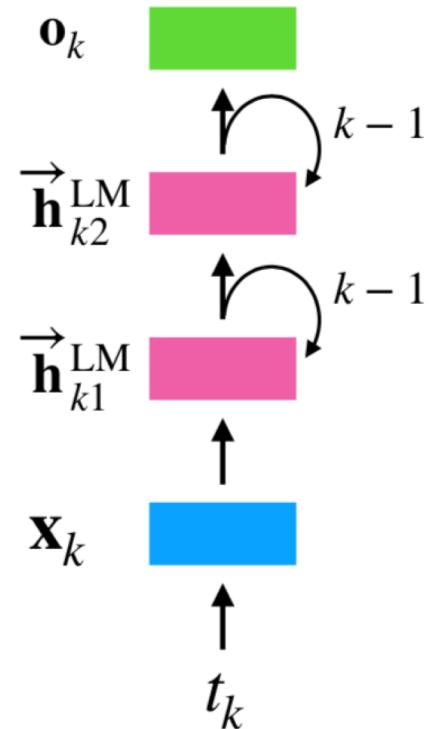
Word2vec Softmax

$$p(t_k|t_1, \dots, t_{k-1}) = \frac{\exp(CNN(t_k)^T h)}{\sum_{i=1}^{|V|} \exp(CNN(t_i)^T h)}, h = LSTM(t_k|t_1, \dots, t_{k-1})$$

CNN Softmax

The forward LM architecture

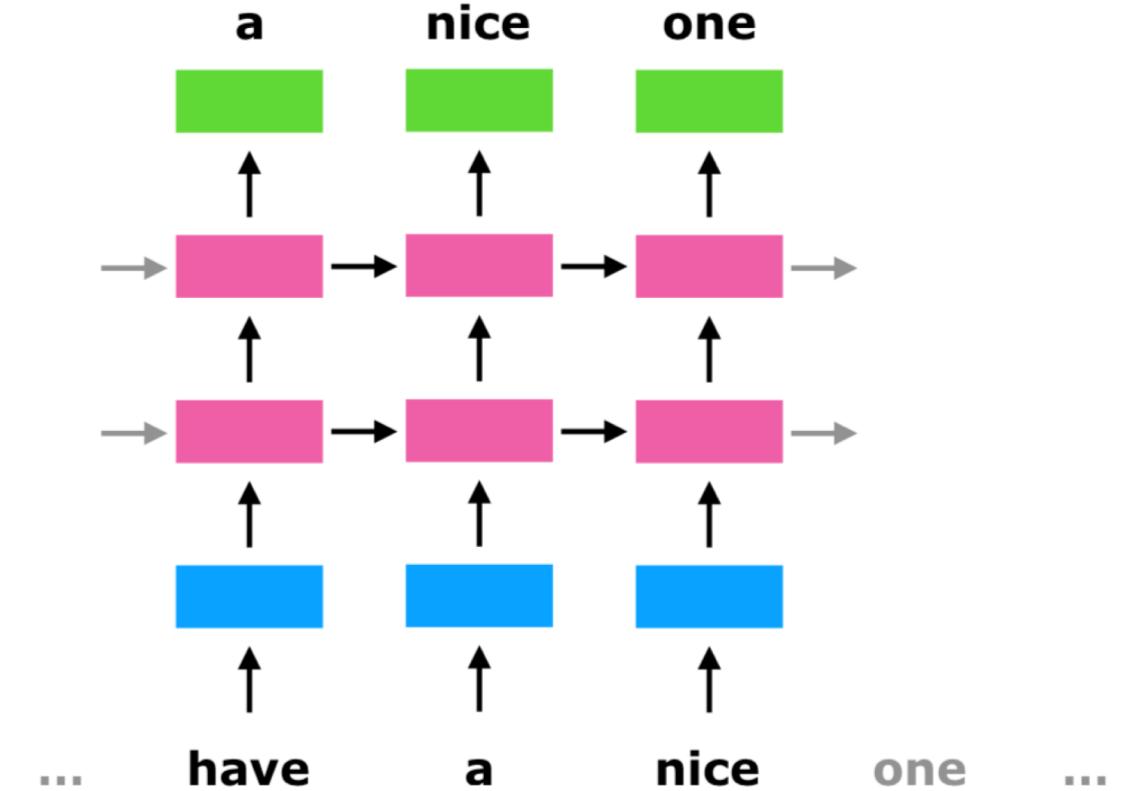
Output layer

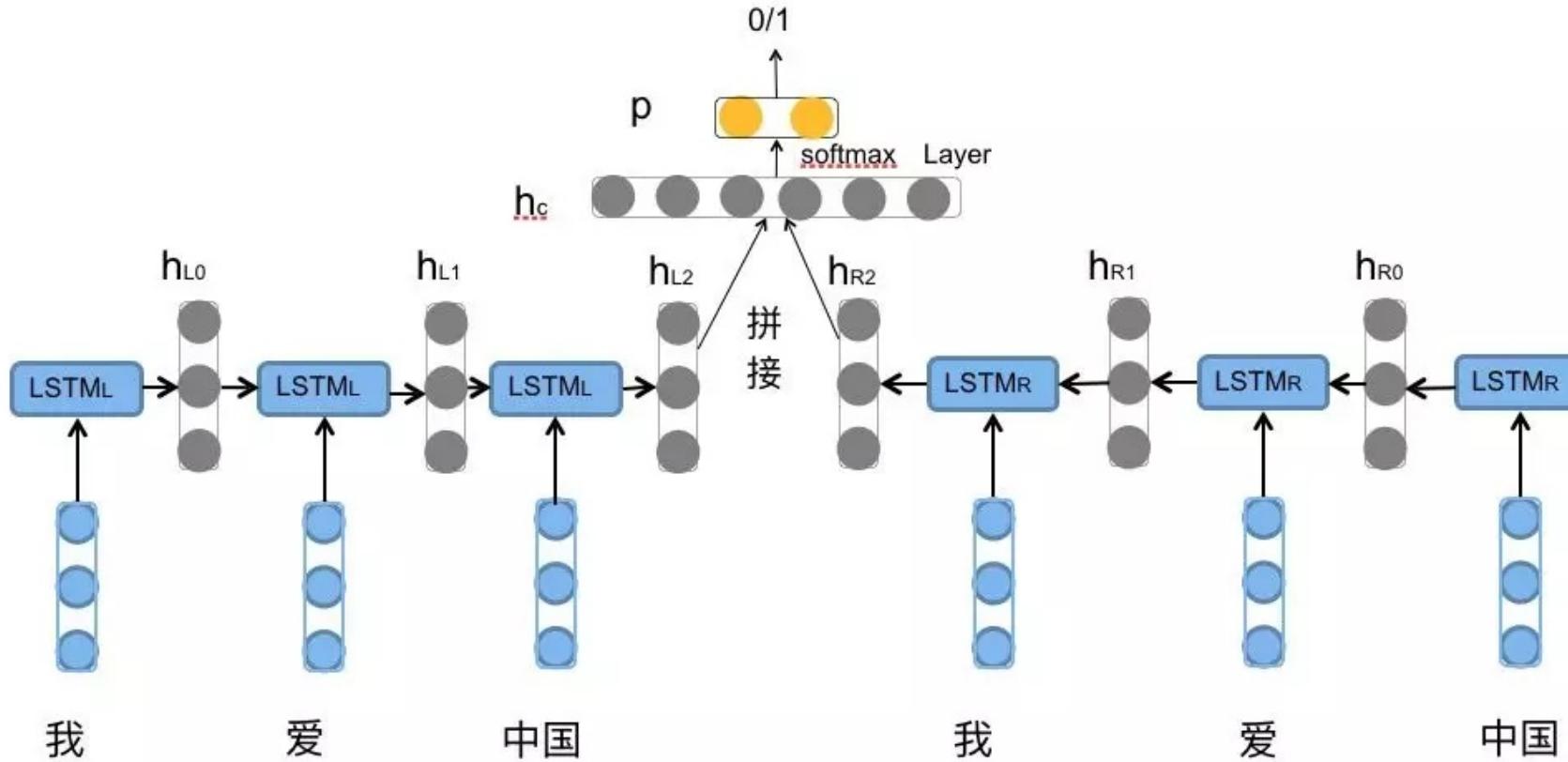


Hidden layers (LSTMs)

Embedding layer

Expanded in the forward direction of k





Why use this structure

- Kim, Yoon, et al. "Character-Aware Neural Language Models." *AAAI*. 2016.
- Jozefowicz R, Vinyals O, Schuster M, et al. "Exploring the limits of language modeling[J]." ICML2016

Different

- Allowing the model to **train two-way parameters**
- Increase **the residual connection** between LSTMs units

Bidirectional language models

Given a sequence of N tokens (t_1, t_2, \dots, t_N):

- Forward language model: Predict token t_k from its history

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k | t_1, t_2, \dots, t_{k-1})$$

Many recent SotA neural language models compute a context-independent representation X_{kj}^{LM} , then pass it through L layers of LSTM, producing h_{kj}^{LM} , where $j = 1, \dots, L$.

Backward language model: Similar as above, but in reverse

$$p(t_1, t_2, \dots, t_N) = \prod_{k=1}^N p(t_k | t_{k+1}, t_{k+2}, \dots, t_N)$$

A LSTM analogous to the forward model produces representations h_{kj}^{LM}

NOTE: direction

Bidirectional language model (biLM): Combines the forward and backward model.
Learn by maximizing:

$$\sum_{k=1}^N (\log p(t_k | t_1, \dots, t_{k-1}; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s) + \log p(t_k | t_{k+1}, \dots, t_N; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s))$$

where:

- Θ_x : token representation parameters (shared between forward and backwards)
- Θ_s : softmax parameters (shared between forward and backwards)
- $\vec{\Theta}_{LSTM}$, $\vec{\Theta}_{LSTM}$: forward and backward LSTM parameters, respectively

The biLM produces $2L + 1$ intermediate representations:

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

where $\mathbf{h}_{k,0}^{LM} = \mathbf{x}_k^{LM}$ is the token layer and
 $\mathbf{h}_{k,j}^{LM} = [\overrightarrow{\mathbf{h}}_{k,j}^{LM}; \overleftarrow{\mathbf{h}}_{k,j}^{LM}]$, for each biLSTM layer.

ELMo: A task specific combination of these features:

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

$$s_i^{task} = e^{s_i} / \sum_j^N e^{s_j}$$

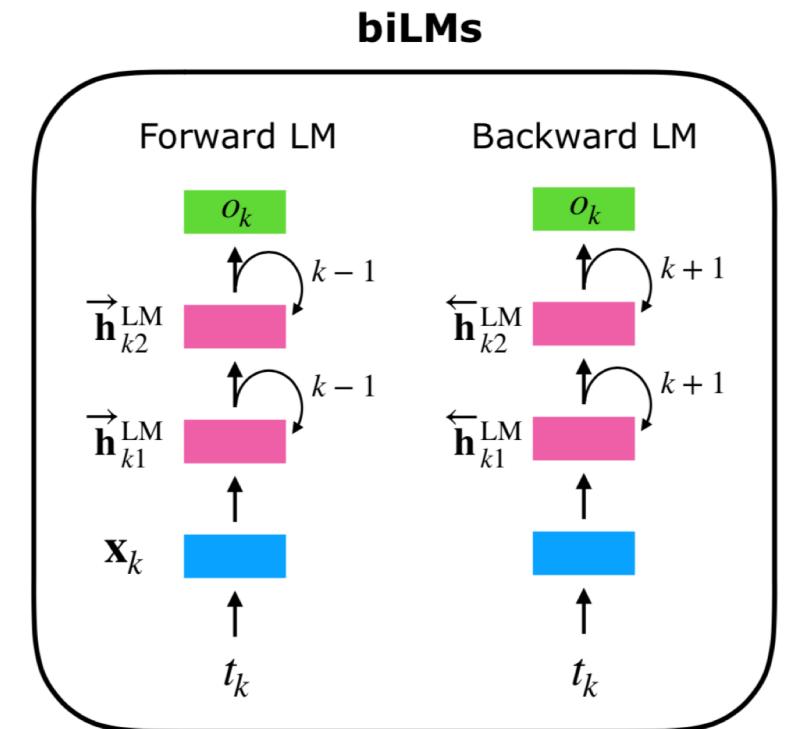
where s_j^{task} are softmax-normalized weights and γ^{task} is a scaling parameter.

ELMo represents a word t_k as a linear combination of corresponding hidden layers

ELMo is a task specific representation. A down-stream task learns weighting parameters

$$\text{ELMo}_k^{\text{task}} = \gamma^{\text{task}} \times \sum \left\{ \begin{array}{l} s_2^{\text{task}} \times \mathbf{h}_{k2}^{\text{LM}} \\ s_1^{\text{task}} \times \mathbf{h}_{k1}^{\text{LM}} \\ s_0^{\text{task}} \times \mathbf{h}_{k0}^{\text{LM}} \\ ([\mathbf{x}_k; \mathbf{x}_k]) \end{array} \right. \xrightarrow{\text{Concatenate hidden layers}} [\vec{\mathbf{h}}_{kj}^{\text{LM}}; \overleftarrow{\mathbf{h}}_{kj}^{\text{LM}}]$$

Unlike usual word embeddings, ELMo is assigned to every *token* instead of a *type*



Intuition

Different layers of **deep Bi-LSTM encode** different types of information.

- Higher-level LSTM states capture context-dependent aspects of word meaning (e.g., they can be used without modification to perform well on supervised word sense disambiguation tasks)
- lower level states model aspects of syntax(e.g., they can be used to do part-of-speech tagging).

The higher layer seemed to learn semantics while the lower layer probably captured syntactic features

Word sense disambiguation

Model	F ₁
WordNet 1st Sense Baseline	65.9
Raganato et al. (2017a)	69.9
Iacobacci et al. (2016)	70.1
CoVe, First Layer	59.4
CoVe, Second Layer	64.7
biLM, First layer	67.4
biLM, Second layer	69.0

Table 5: All-words fine grained WSD F₁. For CoVe and the biLM, we report scores for both the first and second layer biLSTMs.

PoS tagging

Model	Acc.
Collobert et al. (2011)	97.3
Ma and Hovy (2016)	97.6
Ling et al. (2015)	97.8
CoVe, First Layer	93.3
CoVe, Second Layer	92.8
biLM, First Layer	97.3
biLM, Second Layer	96.8

Table 6: Test set POS tagging accuracies for PTB. For CoVe and the biLM, we report scores for both the first and second layer biLSTMs.

The higher layer seemed to learn semantics while the lower layer probably captured syntactic features???

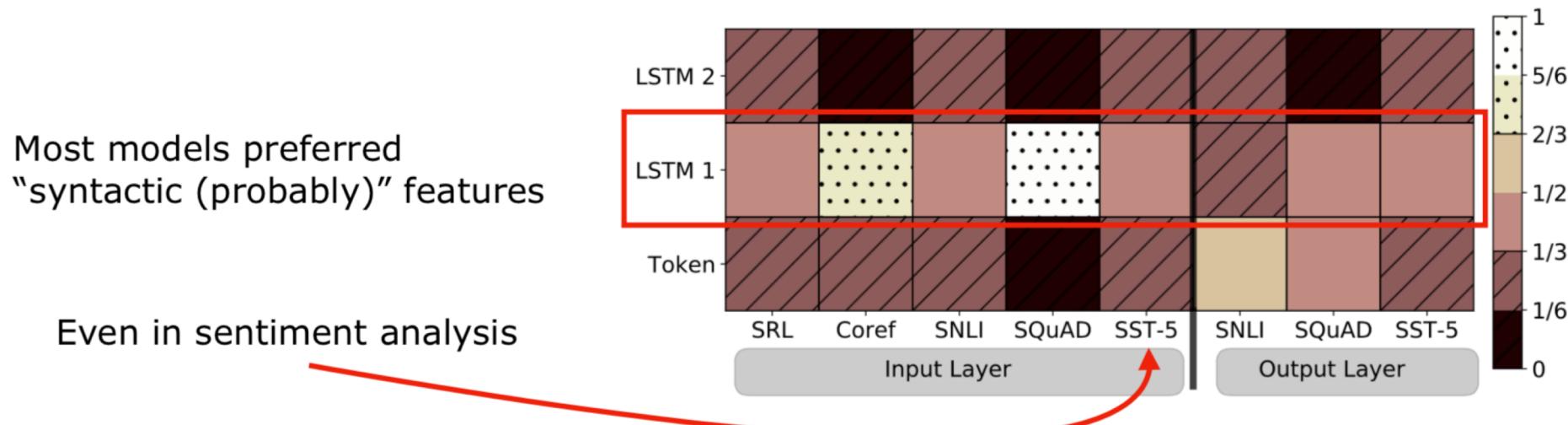


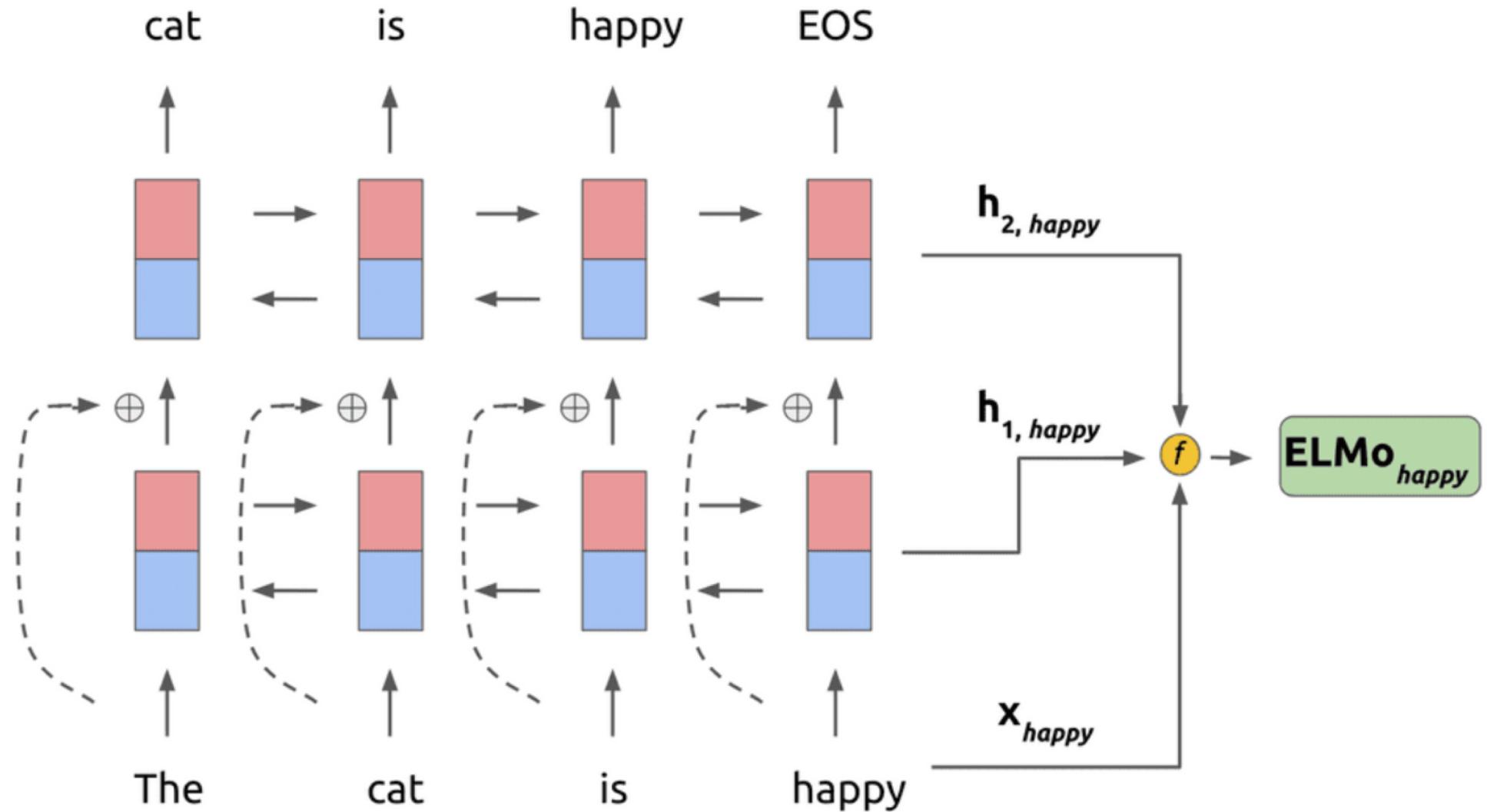
Figure 2: Visualization of softmax normalized biLM layer weights across tasks and ELMo locations. Normalized weights less than $1/3$ are hatched with horizontal lines and those greater than $2/3$ are speckled.

Using biLMs for supervised NLP tasks

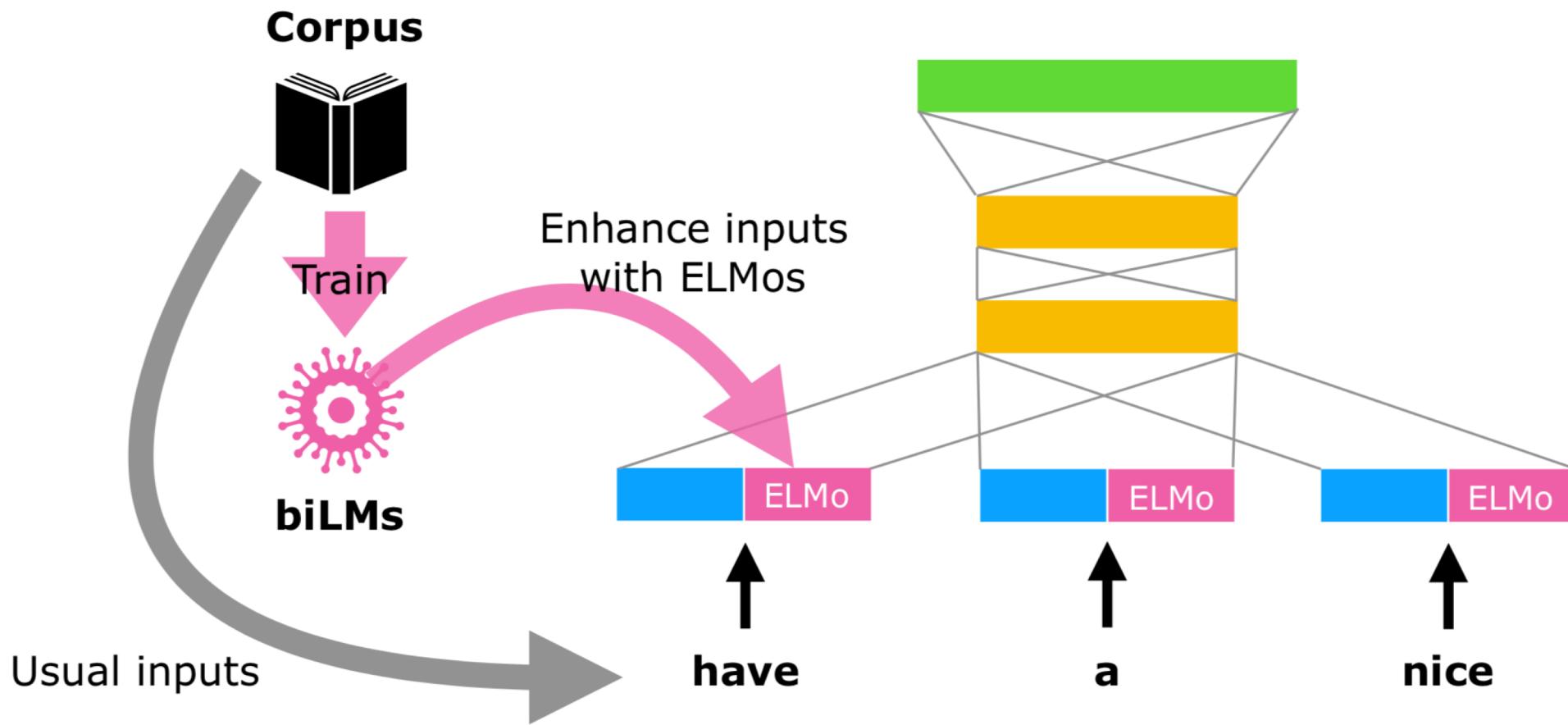
Given a pre-trained biLM, ELMo can easily be dropped into many supervised NLP architectures. Typical NLP architecture:

$$t_k \rightarrow \mathbf{x}_k \rightarrow \mathbf{h}_k$$

- To use ELMo, concatenate ELMo vector to token representation: $[x ; \text{ELMo}^{\text{task}}]$
- Occasionally also helpful to add to the task representation: $[h ; \text{ELMo}^{\text{task}}]$
- Moderate dropout to ELMo beneficial
- Occasionally also helpful to add regularization to the loss: $\lambda ||w||_2$



ELMo can be integrated to almost all neural NLP tasks with simple concatenation to the embedding layer



Many linguistic tasks are improved by using ELMo

TASK	PREVIOUS SOTA	OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/RELATIVE)
Q&A	SQuAD Liu et al. (2017)	84.4	81.1	85.8
Textual entailment	SNLI Chen et al. (2017)	88.6	88.0	88.7 ± 0.17
Semantic role labelling	SRL He et al. (2017)	81.7	81.4	84.6
Coreference resolution	Coref Lee et al. (2017)	67.2	67.2	70.4
Named entity recognition	NER Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10
Sentiment analysis	SST-5 McCann et al. (2017)	53.7	51.4	54.7 ± 0.5

Table 1: Test set comparison of ELMo enhanced neural models with state-of-the-art single model baselines across six benchmark NLP tasks. The performance metric varies across tasks – accuracy for SNLI and SST-5; F_1 for SQuAD, SRL and NER; average F_1 for Coref. Due to the small test sizes for NER and SST-5, we report the mean and standard deviation across five runs with different random seeds. The “increase” column lists both the absolute and relative improvements over our baseline.

ELMo-enhanced models can make use of small datasets more efficiently

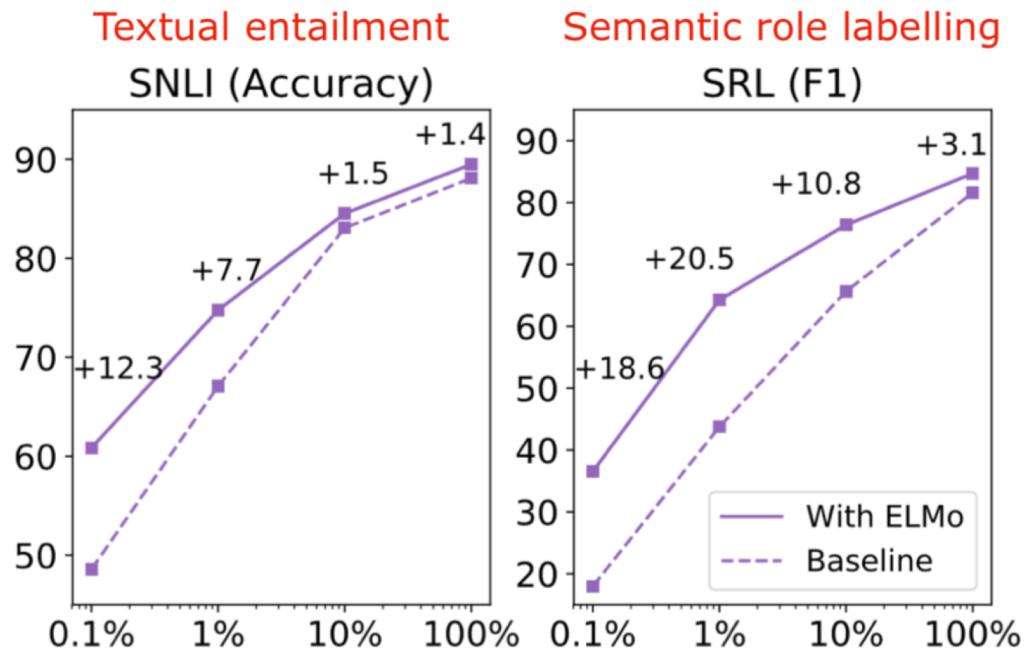


Figure 1: Comparison of baseline vs. ELMo performance for SNLI and SRL as the training set size is varied from 0.1% to 100%.

Continue.....

ELMO/ULMFit ->GPT -> BERT

- 1.Peters, Matthew E., et al. "Deep contextualized word representations." *arXiv preprint arXiv:1802.05365* (2018).
2. Howard, Jeremy, and Sebastian Ruder. "Universal language model fine-tuning for text classification." *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vol. 1. 2018.
3. Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training. URL https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/language-unsupervised/language_understanding_paper.pdf.
4. Devlin J, Chang M W, Lee K, et al. Bert: Pre-training of deep bidirectional transformers for language understanding[J]. arXiv preprint arXiv:1810.04805, 2018.

Thank you !
Q&A