

# Neural Networks in NLP: The Curse of Indifferentiability

Lili Mou

[doublepower.mou@gmail.com](mailto:doublepower.mou@gmail.com)

<http://sei.pku.edu.cn/~moull12>

# Outline

- **Preliminary**
  - **Word embeddings**
  - Sequence-to-sequence generation
- Indifferentiability, solutions, and applications
- A case study in semantic parsing

# Language Modeling

- One of the most fundamental problems in NLP.
- Given a corpus  $\mathbf{w}=w_1w_2\dots w_t$ , the goal is to maximize  $p(\mathbf{w})$

# Language Modeling

- One of the most fundamental problems in NLP.
- Given a corpus  $\mathbf{w}=w_1w_2\dots w_t$ , the goal is to maximize  $p(\mathbf{w})$
- **Philosophical discussion:** Does “probability of a corpus/sentence” make sense?
  - Recognize speech
  - Wreck a nice beach
- All in all, NLP (especially publishing in NLP) is pragmatic.

# Decomposition of the Joint probability

- $p(\mathbf{w}) = p(w_1)p(w_2|w_1)p(w_3|w_1w_2) \dots p(w_t|w_1w_2\dots w_{t-1})$

# Decomposition of the Joint probability

- $p(\mathbf{w}) = p(w_1)p(w_2|w_1)p(w_3|w_1w_2) \dots p(w_t|w_1w_2\dots w_{t-1})$

## Minor question:

- Can we decompose any probabilistic distribution into this form? Yes.
- Is it necessary to decompose a probabilistic distribution into this form? No.

# Markov Assumption

- $p(\mathbf{w}) = p(w_1)p(w_2|w_1)p(w_3|w_1w_2) \dots p(w_t|w_1w_2\dots w_{t-1})$   
 $\approx p(w_1)p(w_2|w_1)p(w_3|w_2) \dots p(w_t|w_{t-1})$

- A word is dependent only on its previous  $n-1$  words and independent of its position,

I.e., provided with the previous  $n-1$  words, the current word is independent of other random variables.

$$p(\mathbf{w}) \approx \prod_{t=1}^m p(w_t | \mathbf{w}_{t-n+1}^{t-1})$$

# Multinomial Estimate

- Maximum likelihood estimation for a multinomial distribution is merely counting.

$$p(w_n | \mathbf{w}_1^{n-1}) = \frac{\#w_1^n}{\#w_1^{n-1}}$$

- Problems
  - #para grows exp. w.r.t.  $n$
  - Even for very small  $n$  (e.g., 2 or 3), we come across severe data sparsity because of the Zipf distribution

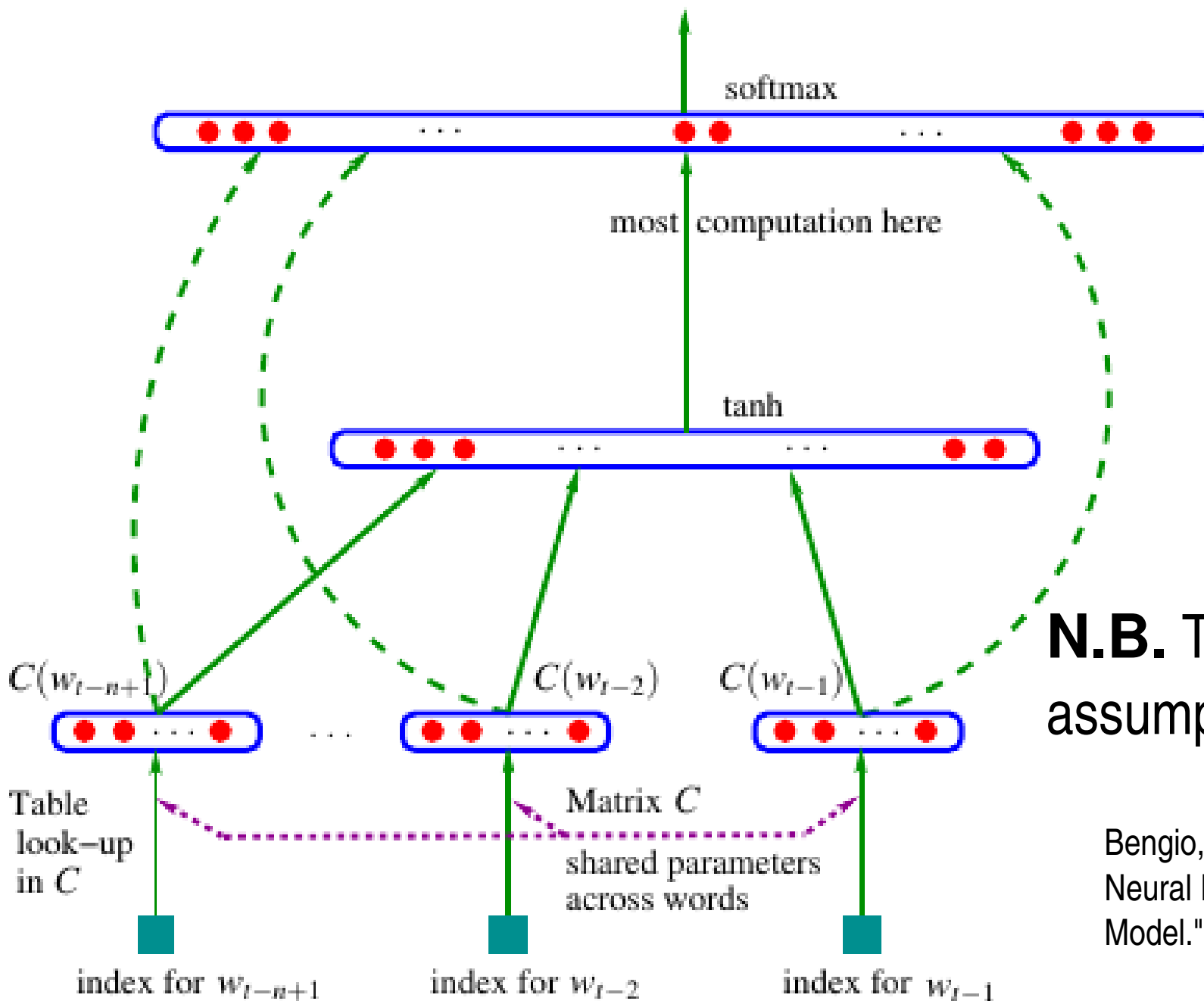


# Parameterizing LMs with Neural Networks

- Each word is mapped to a real-valued vector, called *embeddings*.
- Neural layers capture context information (typically previous words).
- The probability  $p(w| \cdot)$  is predicted by a softmax layer.

# Feed-Forward Language Model

$$i\text{-th output} = P(w_t = i \mid \text{context})$$

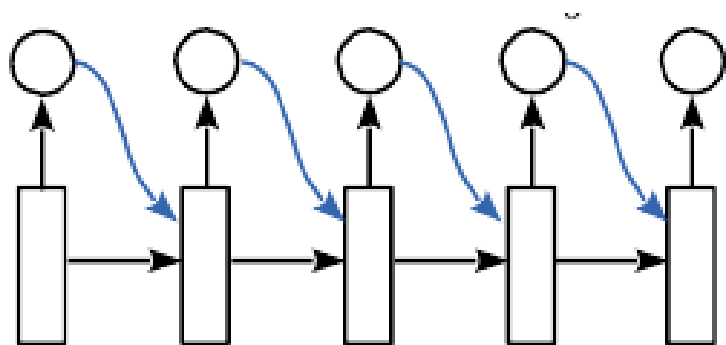


**N.B.** The Markov assumption also holds.

Bengio, Yoshua, et al. "A Neural Probabilistic Language Model." JMLR. 2003.

# Recurrent Neural Language Model

- RNN keeps one or a few hidden states
- The hidden states change at each time step according to the input



$$\mathbf{h}_t = \text{RNN}(\mathbf{x}_t, \mathbf{h}_{t-1})$$

$$= f(W_{\text{in}}\mathbf{x}_t + W_{\text{hid}}\mathbf{h}_{t-1})$$

$$p(w_t | \mathbf{w}_0^{t-1}) \approx \text{softmax}(W_{\text{out}}\mathbf{h}_t)$$

- RNN directly parametrizes  $p(\mathbf{w}) = \prod_{t=1}^m p(w_t | \mathbf{w}_1^{t-1})$

rather than  $p(\mathbf{w}) \approx \prod_{t=1}^m p(w_t | \mathbf{w}_{t-n+1}^{t-1})$

# Complexity Concerns

- Time complexity
  - Hierarchical softmax [1]
  - Negative sampling: Hinge loss [2], Noisy contrastive estimation [3]
- Memory complexity
  - Compressing LM [4]
- Model complexity
  - Shallow neural networks are still too “deep.”
  - CBOW, SkipGram [3]

[1] Mnih A, Hinton GE. A scalable hierarchical distributed language model. NIPS, 2009.

[2] Collobert R, Weston J, Bottou L, Karlen M, Kavukcuoglu K, Kuksa P. Natural language processing (almost) from scratch. JMLR, 2011.

[3] Mikolov T, Chen K, Corrado G, Dean J. Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781. 2013

[4] Yunchuan Chen, Lili Mou, Yan Xu, Ge Li, Zhi Jin. "Compressing neural language models by sparse word representations." In ACL, 2016.

# The Role of Word Embeddings?

- Word embeddings are essentially a connectional weight matrix, whose input is a one-hot vector.
- Implementing by a look-up table is much faster than matrix multiplication.
- Each column of the matrix corresponds to a word.

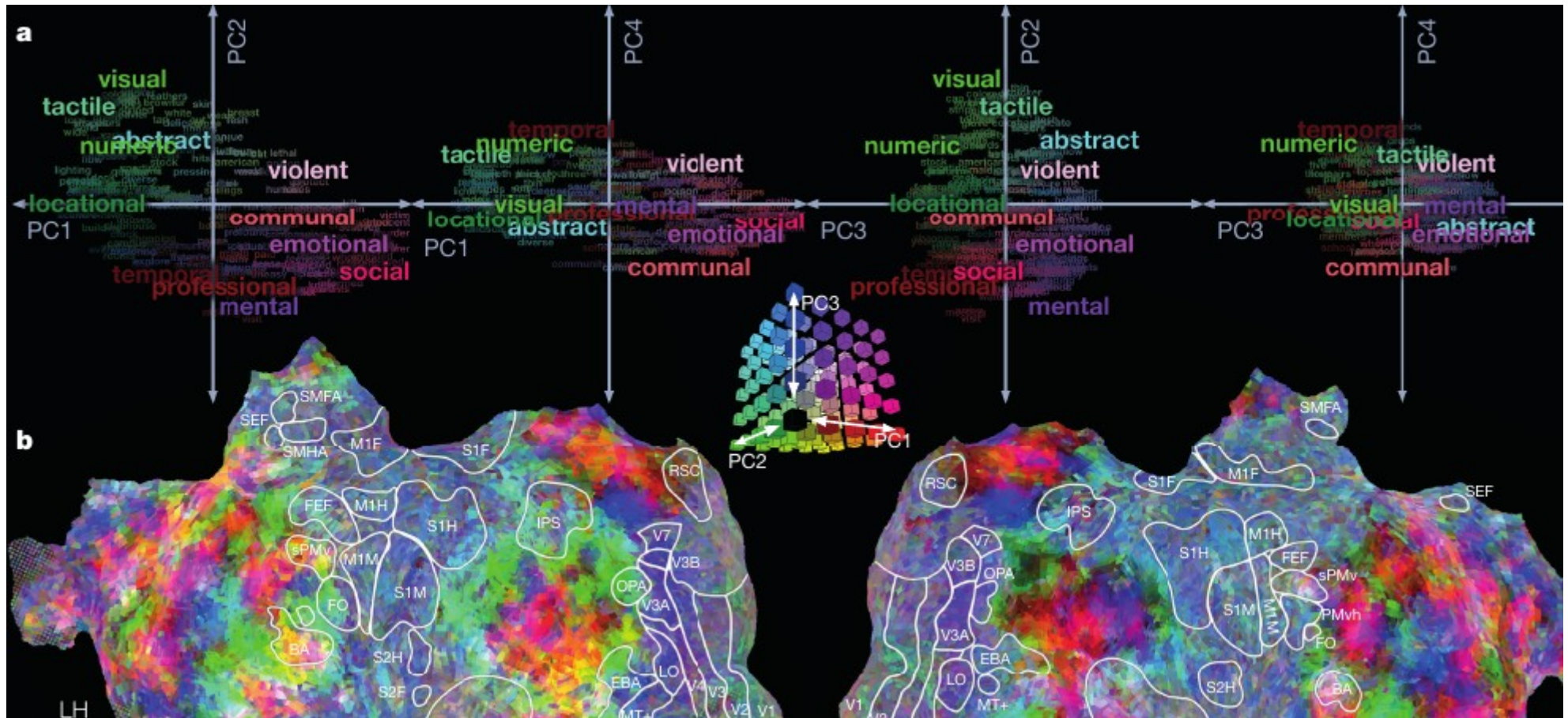
The diagram illustrates the process of retrieving a word embedding. On the left, a large left-facing square bracket contains a vertical blue bar representing a column in a matrix. Below this is the text "Embedding of word  $i$  retrieved by matrix-vector multiplication". To the right of the blue bar is a dot, followed by a large right-facing square bracket containing a vertical list of values: 0, 1, a dot, a dot, and 0. This represents a one-hot vector.

One-hot representation of word  $i$  (sparse)

# How can we use word embeddings?

- Embeddings demonstrate the internal structures of words
  - Relation represented by vector offset
    - “man” – “woman” = “king” – “queen”
  - Word similarity
- Embeddings can serve as the initialization of almost every supervised task
  - A way of pretraining
  - **N.B.:** may not be useful when the training set is large enough

# Word Embeddings in our Brain



Huth, Alexander G., et al. "Natural speech reveals the semantic maps that tile human cerebral cortex." *Nature* 532.7600 (2016): 453-458.

# “Somatotopic Embeddings” in our Brain

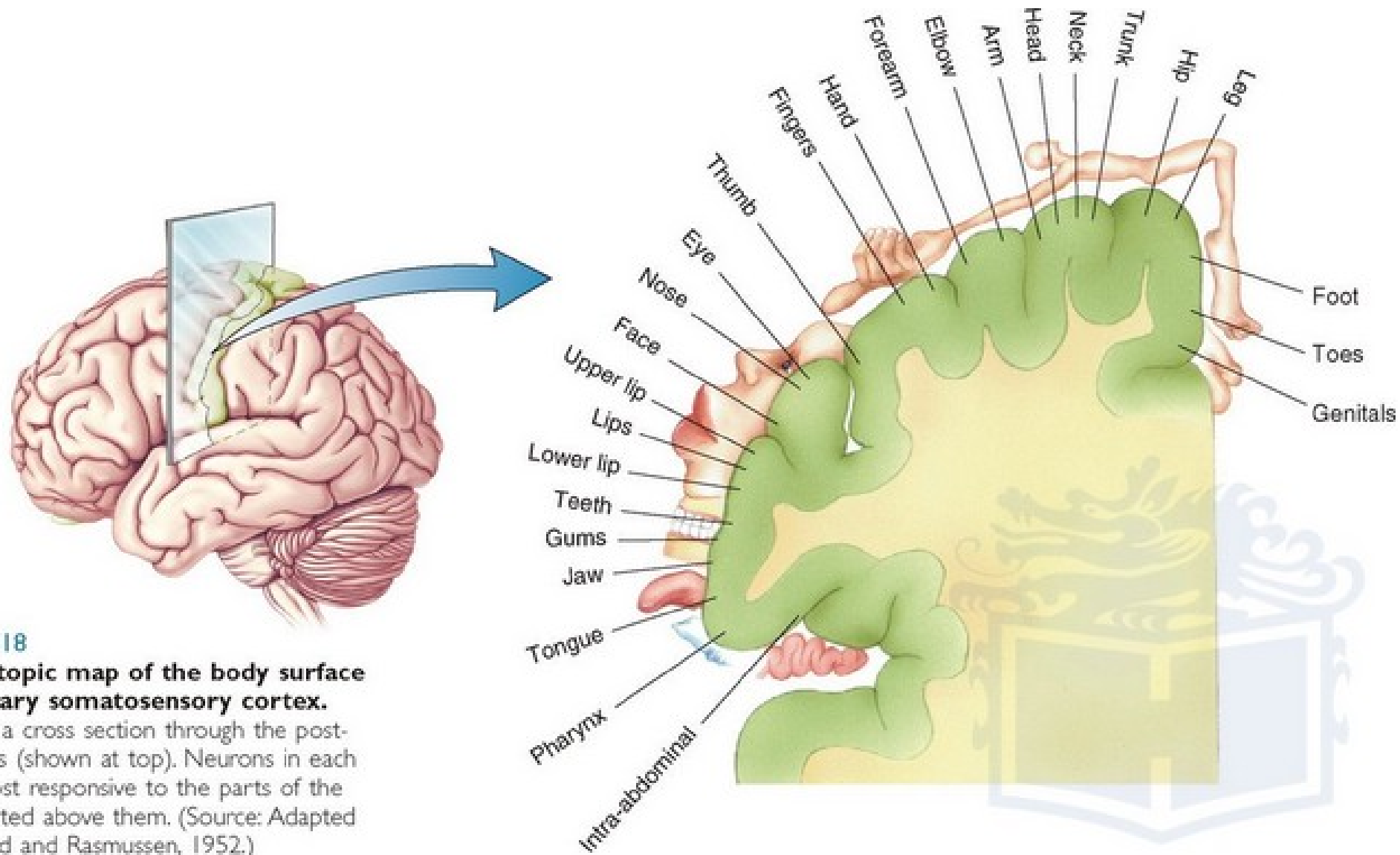
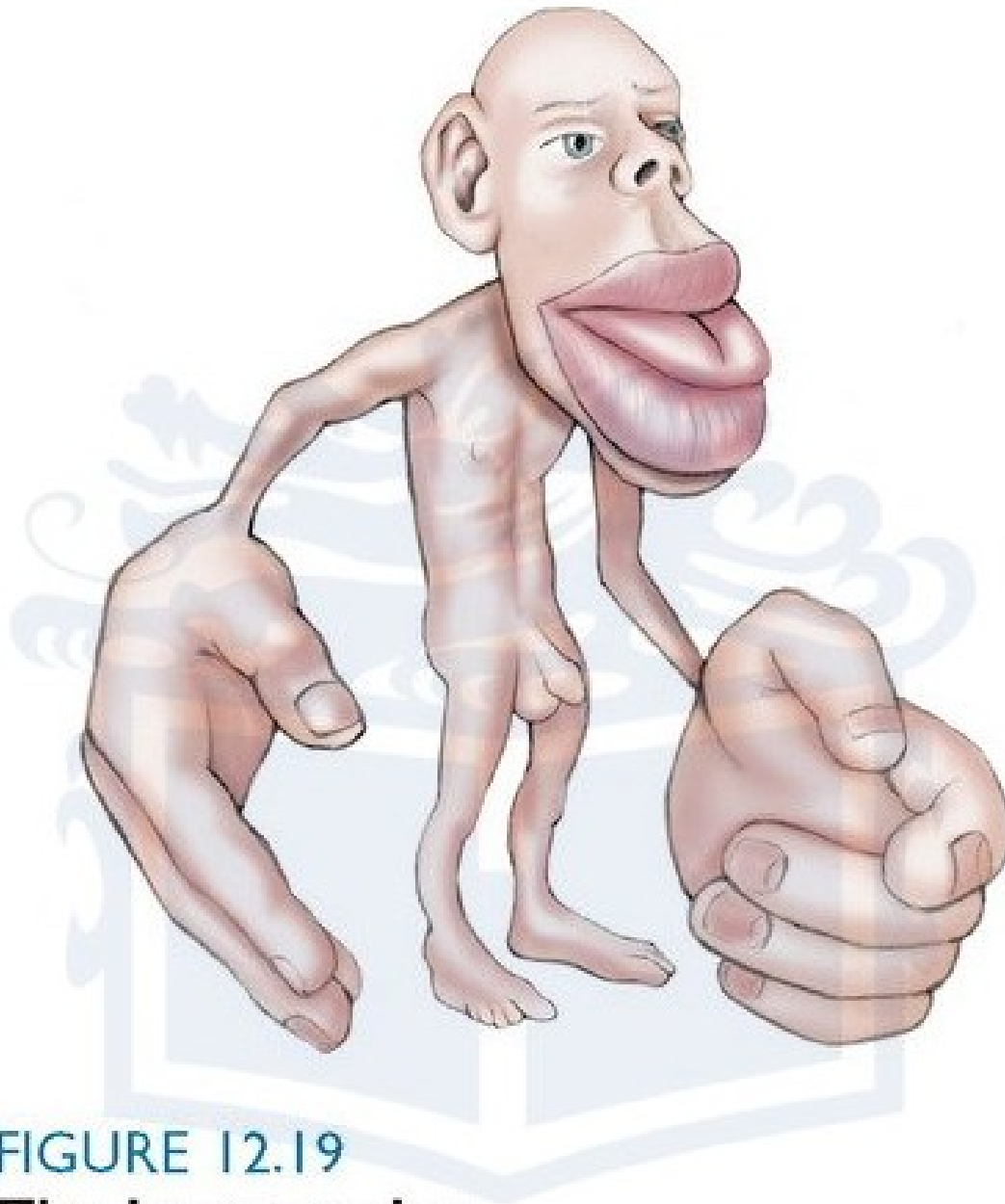


FIGURE 12.18

**A somatotopic map of the body surface onto primary somatosensory cortex.**

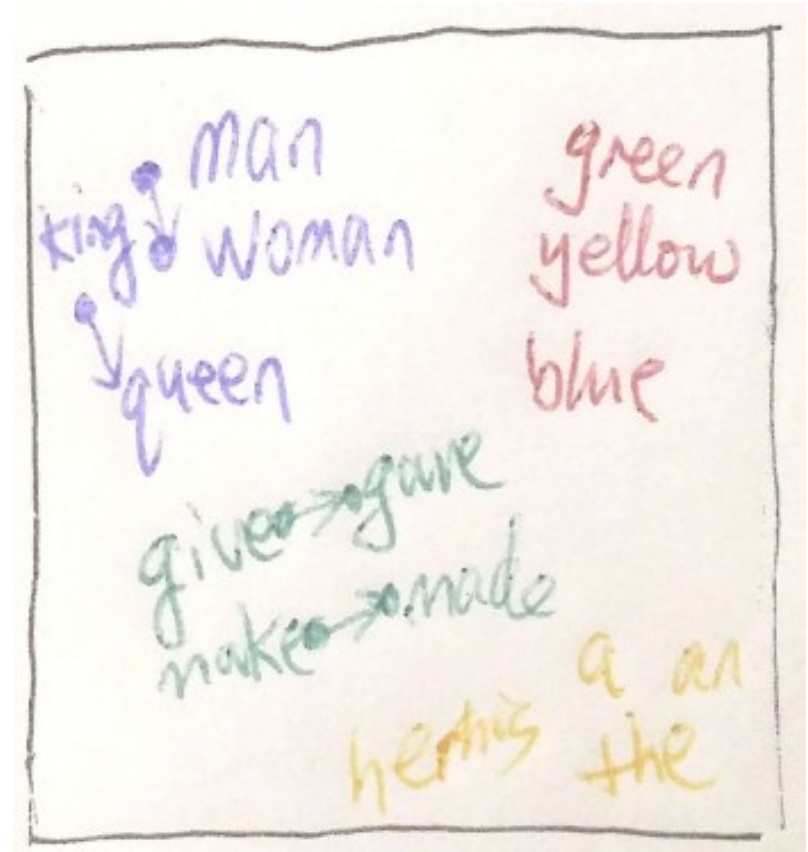
This map is a cross section through the post-central gyrus (shown at top). Neurons in each area are most responsive to the parts of the body illustrated above them. (Source: Adapted from Penfield and Rasmussen, 1952.)



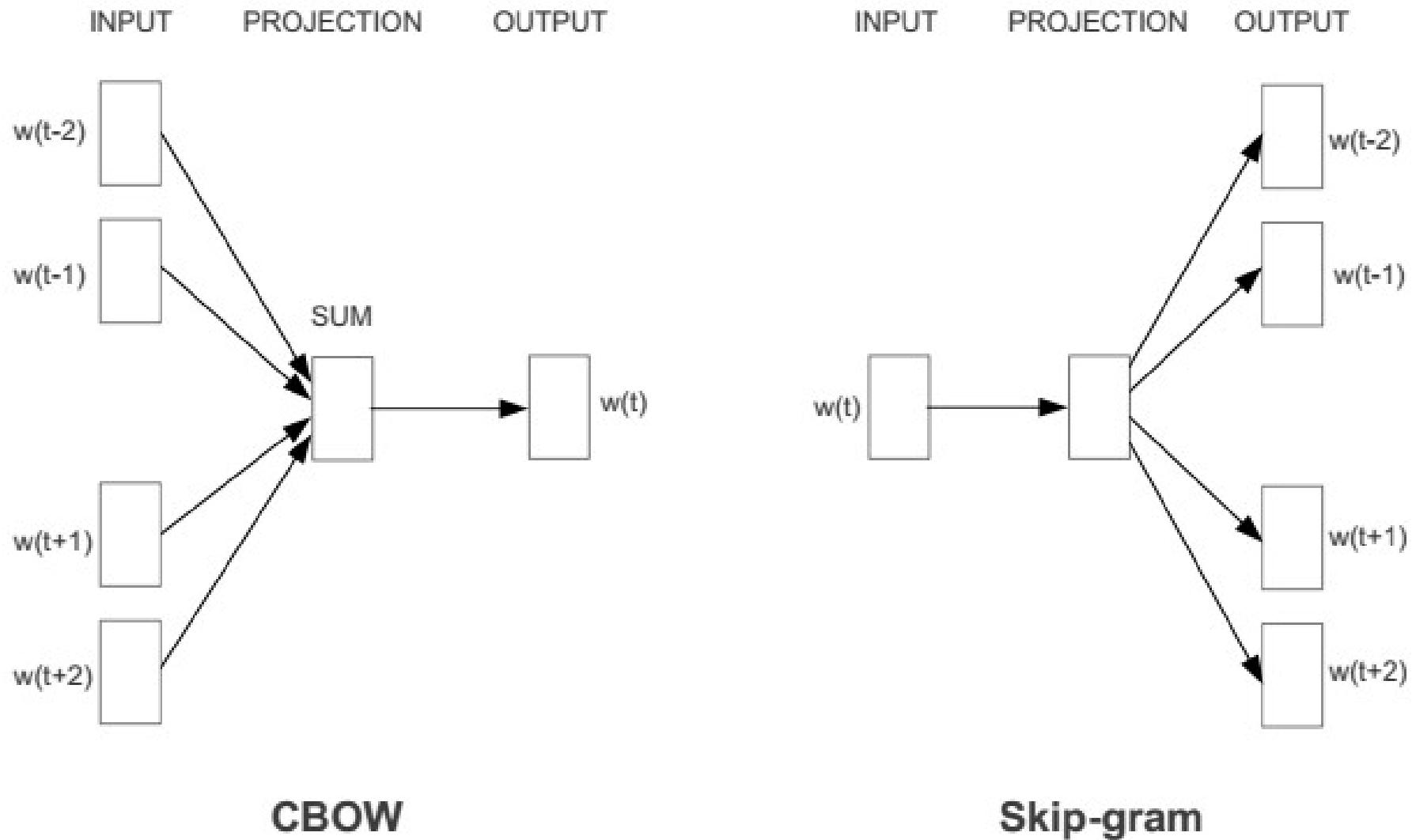


**FIGURE 12.19**  
**The homunculus.**

Deep neural networks:  
To be, or not to be? That is the question.



# CBOW, SkipGram (word2vec)



[6] Mikolov T, Chen K, Corrado G, Dean J. Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781. 2013

# Hierarchical Softmax and Negative Contrastive Estimation

- HS

$$p(w|w_I) = \prod_{j=1}^{L(w)-1} \sigma \left( \mathbb{I}[n(w, j+1) = \text{ch}(n(w, j))] \cdot v'_{n(w, j)}{}^\top v_{w_I} \right)$$

- NCE

$$\log \sigma(v'_{w_O}{}^\top v_{w_I}) + \sum_{i=1}^k \mathbb{E}_{w_i \sim P_n(w)} \left[ \log \sigma(-v'_{w_i}{}^\top v_{w_I}) \right]$$

# Tricks in Training Word Embeddings

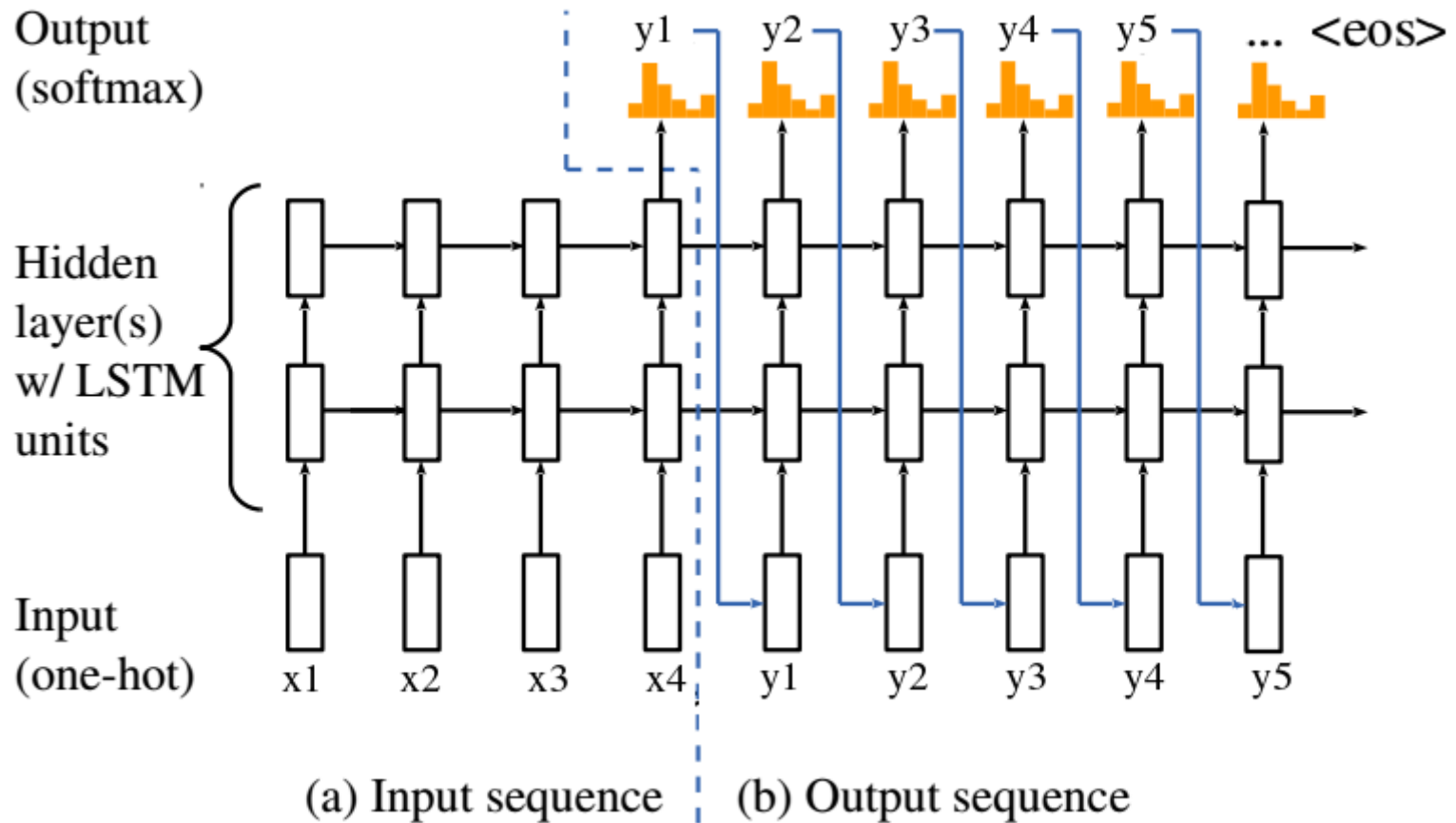
- The # of negative samples?
  - The more, the better.
- The distribution from which negative samples are generated? Should negative samples be close to positive samples?
  - The closer, the better.
- Full softmax vs. NCE vs. HS vs. hinge loss?

# Outline

- **Preliminary**
  - Word embeddings
  - **Sequence-to-sequence generation**
- Indifferentiability, solutions, and applications
- A case study in semantic parsing

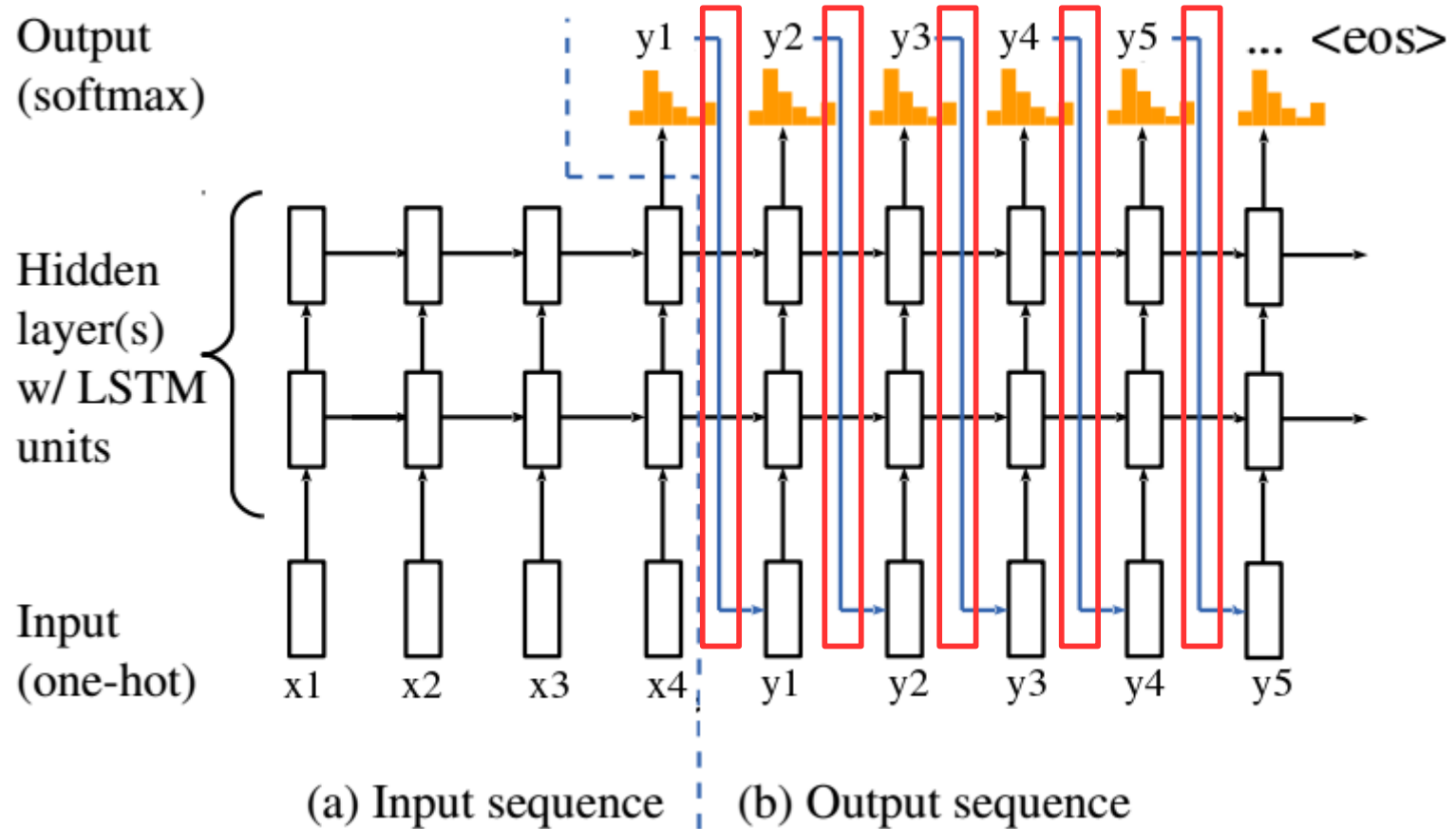
# Seq2Seq Networks

Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. "Sequence to sequence learning with neural networks." NIPS. 2014



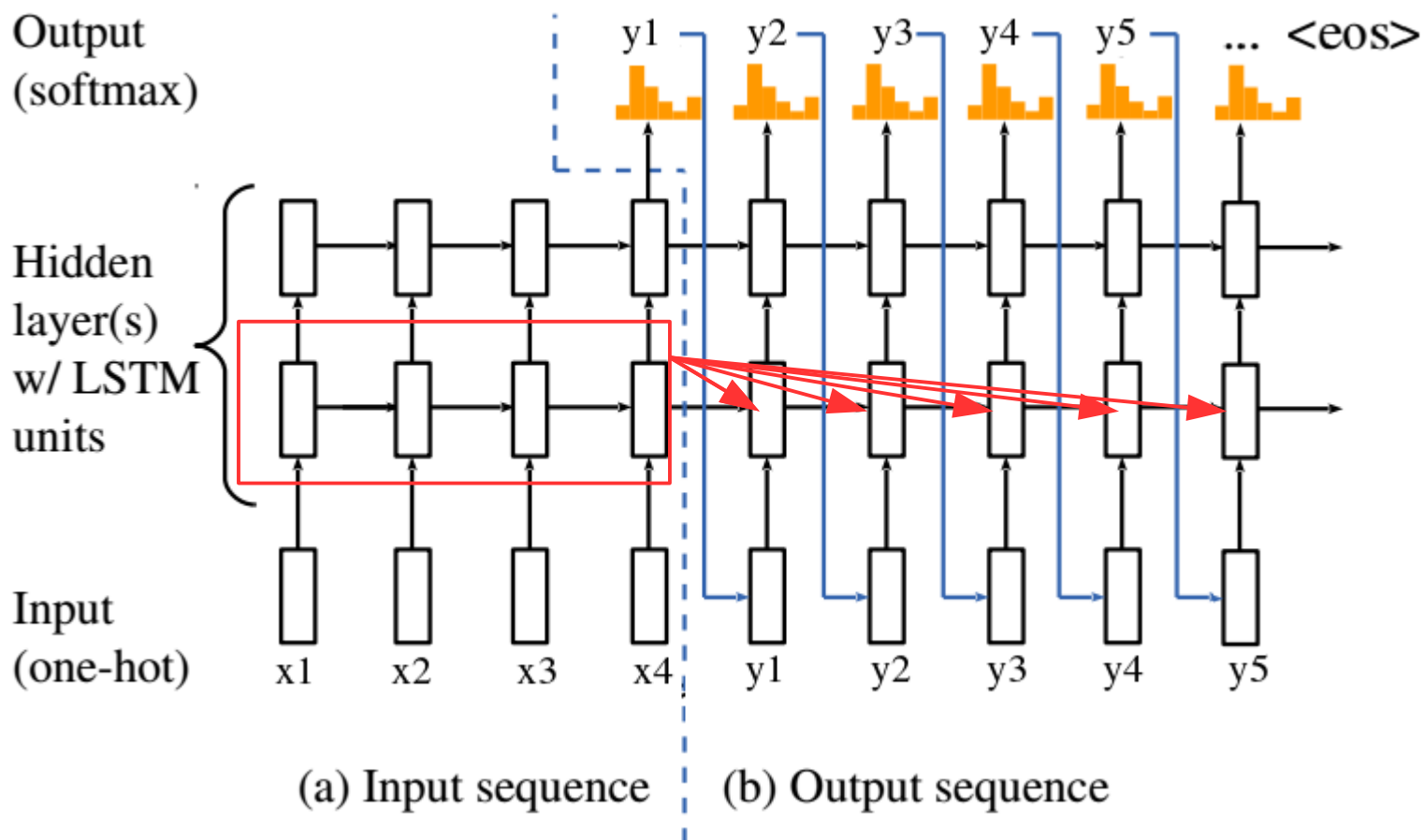
# Feeding back the output

- Mode selection
- Potential chaotic behaviors





# Attention Mechanisms



Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." ICLR, 2014.

# Context Vector

The context vector  $\mathbf{c}$  is a combination of the input sequence's states

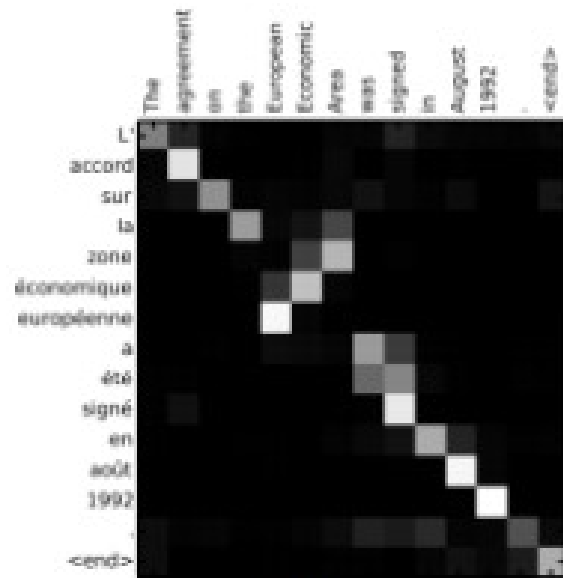
$$\mathbf{c} = \sum_i \alpha_i \mathbf{c}_i$$

where the coefficient  $\alpha_i$  is related to

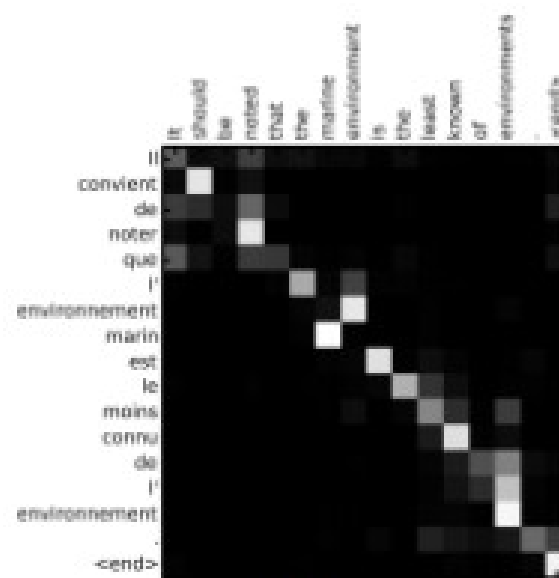
- The local context  $\mathbf{c}_i$ , and
- The last output state  $\mathbf{h}_{t-1}$
- $\alpha_i$  is normalized

$$\alpha_i = \frac{\exp\{\tilde{\alpha}_i\}}{\sum_j \exp\{\tilde{\alpha}_j\}} \quad \tilde{\alpha}_i = W[\mathbf{h}_{t-1}; \mathbf{c}_i]$$

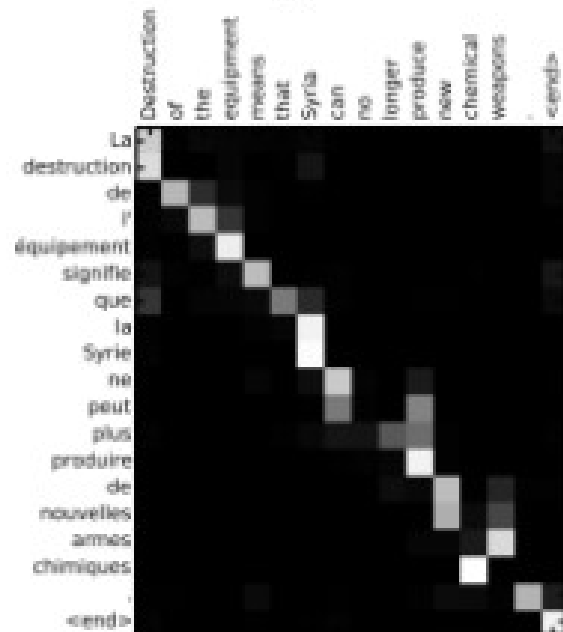
# Attention as Alignment



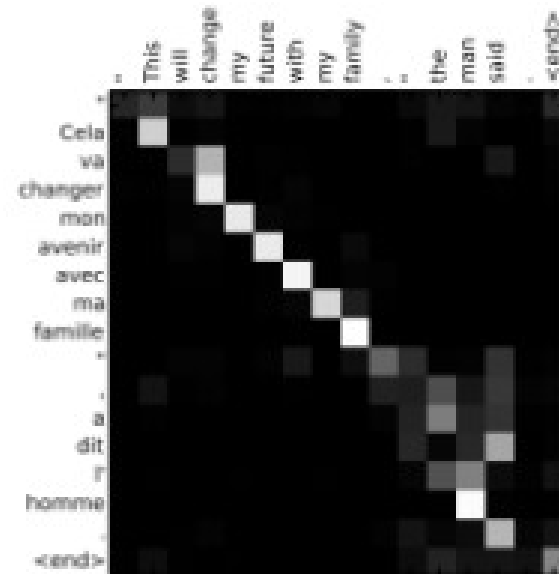
(a)



(b)



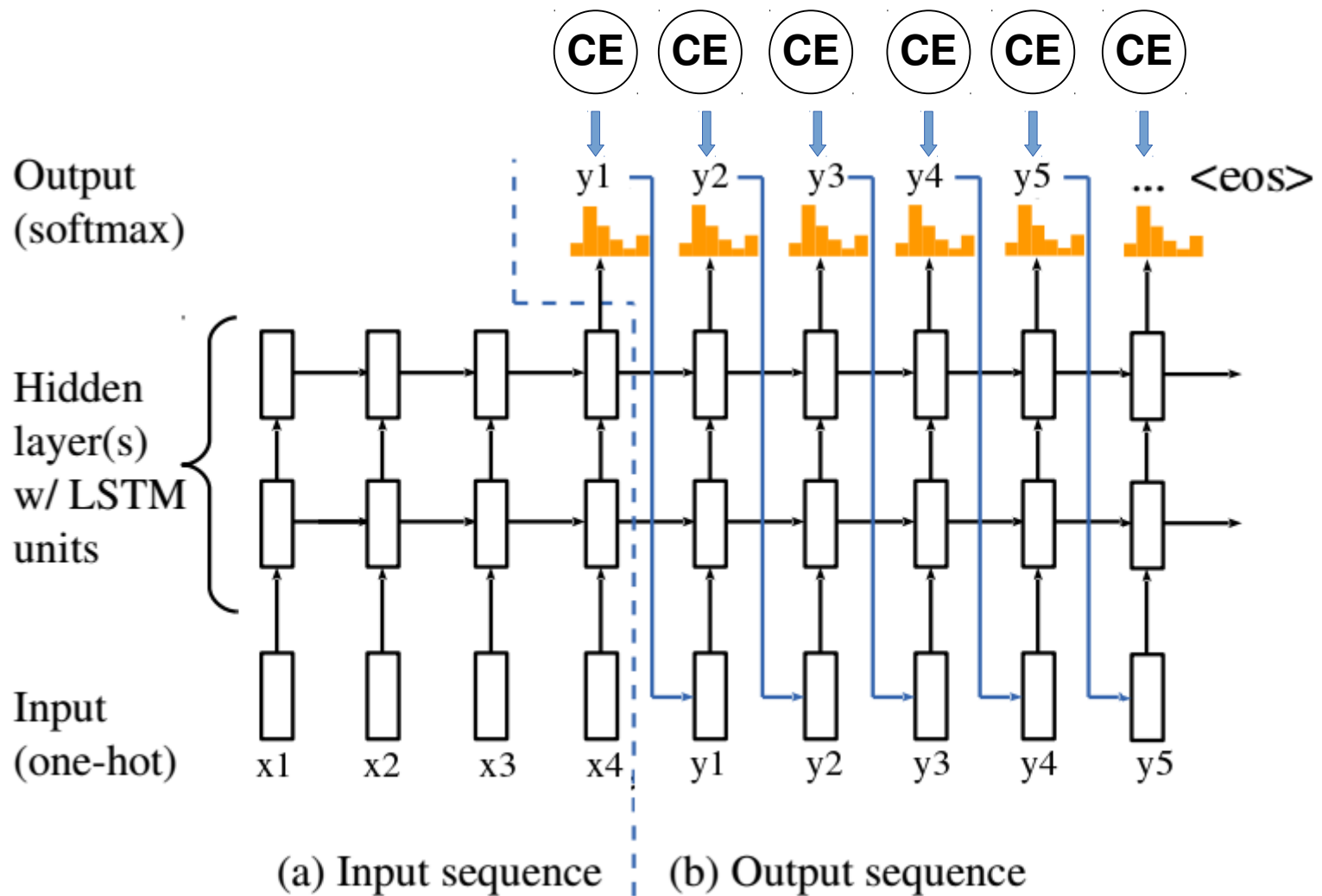
(c)



(d)

# Training

- Step-by-step training



# Indifferentiability

