



前沿技术讲习班
Advanced Technology Tutorial

深度学习与机器翻译



模式识别国家重点实验室
National Laboratory of Pattern Recognition



前沿技术讲习班
Advanced Technology Tutorial

统计机器翻译中的 深度学习

张家俊

中国科学院自动化研究所

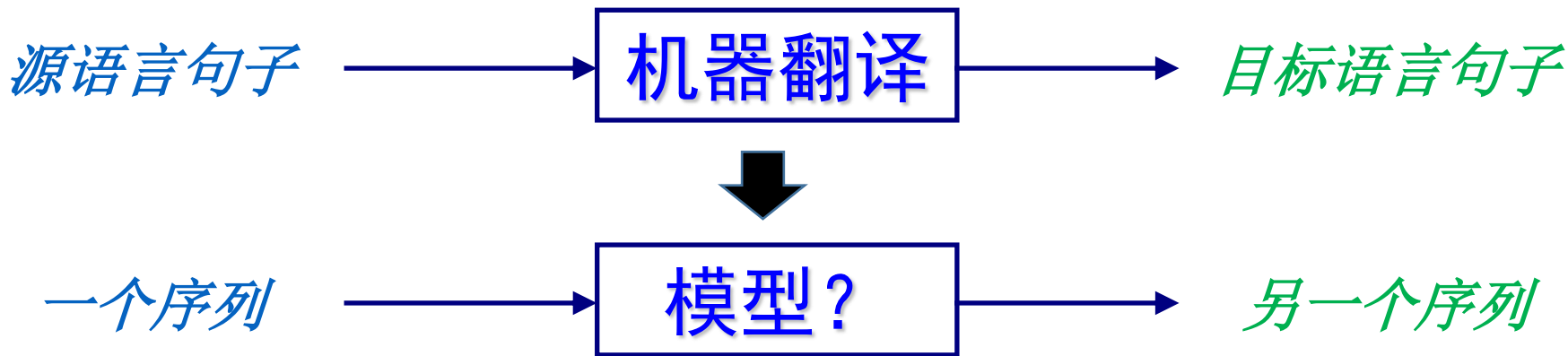
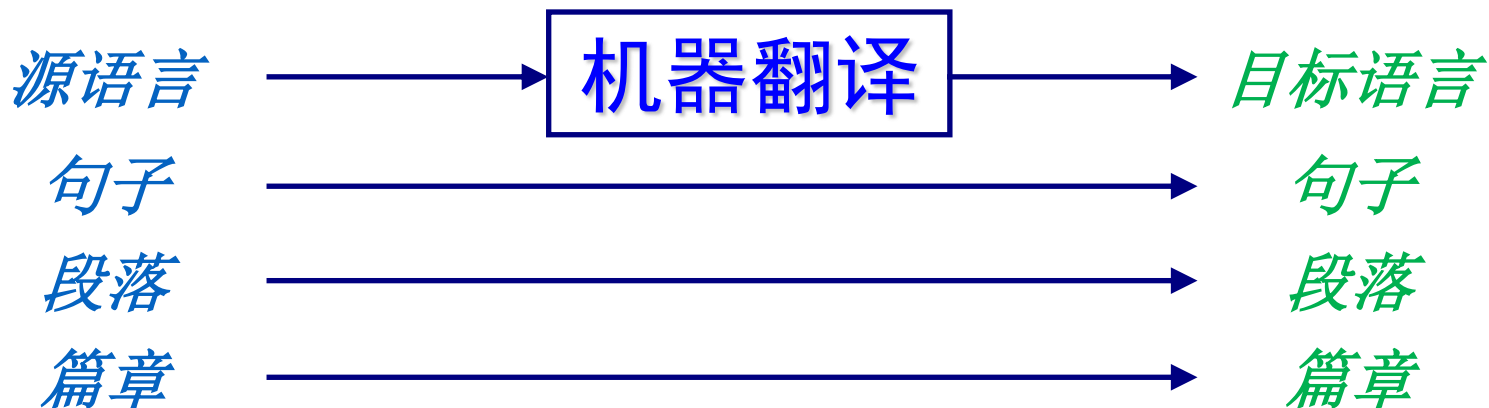
www.nlpr.ia.ac.cn/cip/jjzhang.htm

jjzhang@nlpr.ia.ac.cn

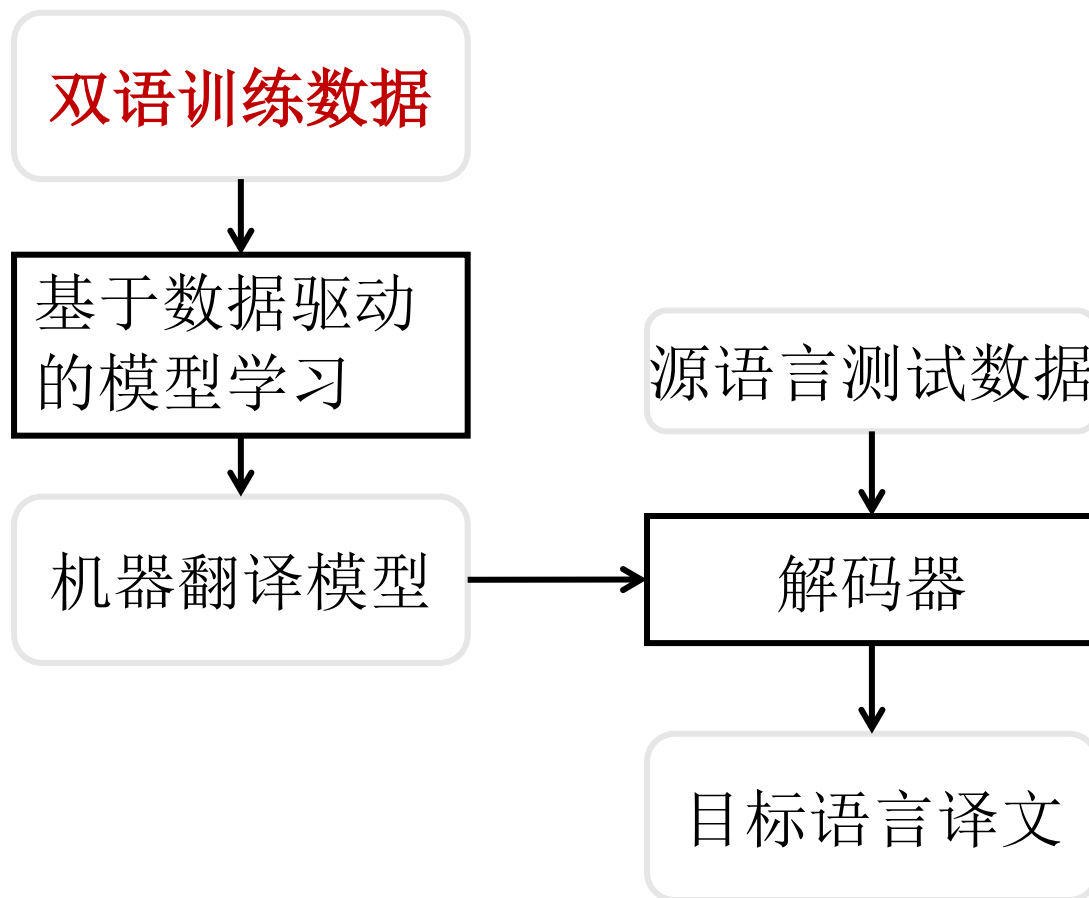
机器翻译

今天烟台天气不错

The weather is fine in Yantai today



统计机器翻译



双语训练数据

人类 共 有 二十三 对 染色体 。

humans have a total of 23 pairs of chromosomes .

中国 大陆 手机 用户 成长 将 减缓

growth of phone users in mainland china to slow

驻 南韩 美军 三千人 奉命 冻结 调防

us freezes transfer of 3,000 troops in south korea

... ..

澳洲 重新 开放 驻 马尼拉 大使馆

australia reopens embassy in manila

外交 人员 搭乘 第五 架 飞机 返国

diplomatic staff take the fifth plane home

姚明 感慨 NBA 的 偶像 来 得 太 快

yao ming feels nba stardom comes too fast

... ..

统计机器翻译-生成模型

源语言句子: $S = s_1^m = s_1 s_2 \cdots s_m$

目标语言句子: $T = t_1^l = t_1 t_2 \cdots t_l$

贝叶斯公式: $P(T|S) = \frac{P(T) \times P(S|T)}{P(S)}$



[Brown et al., 1990, 1993]

$$T' = \operatorname{argmax}_T P(T) \times P(S|T)$$

语言模型

Language model, LM

翻译模型

Translation model, TM

统计机器翻译-生成模型

澳洲₁ 与₂ 北韩₃ 有₄ 邦交₅

f₁ **f₂** **f₃** **f₄** **f₅** **f₆** **f₇**

$\varepsilon \equiv$
 $p(m|T)$

澳洲₁ 与₂ 北韩₃ 有₄ 邦交₅

f₁ **f₂** **f₃** **f₄** **f₅** **f₆** **f₇**

$p(a_j|j, m, l)$

澳洲₁ 与₂ 北韩₃ 有₄ 邦交₅

Austria₁ has₂ diplomatic₃ relations₄ with₅ North₆ Korea₇

$p(s_j|t_{a_j})$

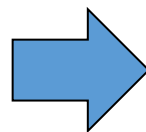
统计机器翻译-判别式模型

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_T \frac{P(T) \times P(S|T)}{P(S)} \\
 &= \operatorname{argmax}_T P(T) \times P(S|T)
 \end{aligned}$$

生成式模型

翻译质量

[Och, 2002]



$$\begin{aligned}
 T' &= \\
 &= \operatorname{argmax}_T P(T) \times P(T|S)
 \end{aligned}$$

?

翻译质量

≈

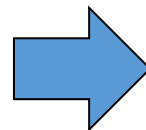


统计机器翻译-判别式模型

$$T' = \operatorname{argmax}_T P(T) \times P(S|T)$$

$$T' = \operatorname{argmax}_T P(T) \times P(T|S)$$

翻译质量



$$T' = \operatorname{argmax}_T P(T) \times P(S|T) \times P(T|S)$$

翻译质量

<

统计机器翻译-判别式模型

$$T' = \operatorname{argmax}_T P(T|S) = \operatorname{argmax}_T \frac{\exp\{\sum_1^M \lambda_m h_m(T, S)\}}{\sum_{T^*} \exp\{\sum_1^M \lambda_m h_m(T^*, S)\}}$$

$$= \operatorname{argmax}_T \left\{ \sum_1^M \lambda_m h_m(T, S) \right\}$$

$$\begin{aligned} h_1(T, S) &= \log P(T) \\ h_2(T, S) &= \log P(S|T) \\ \lambda_1 &= \lambda_2 = 1 \end{aligned}$$



$$T' = \operatorname{argmax}_T P(T|S) = \operatorname{argmax}_T \left\{ \sum_1^M \lambda_m h_m(T, S) \right\}$$

$$= \operatorname{argmax}_T \{ \log P(T) + \log P(S|T) \}$$

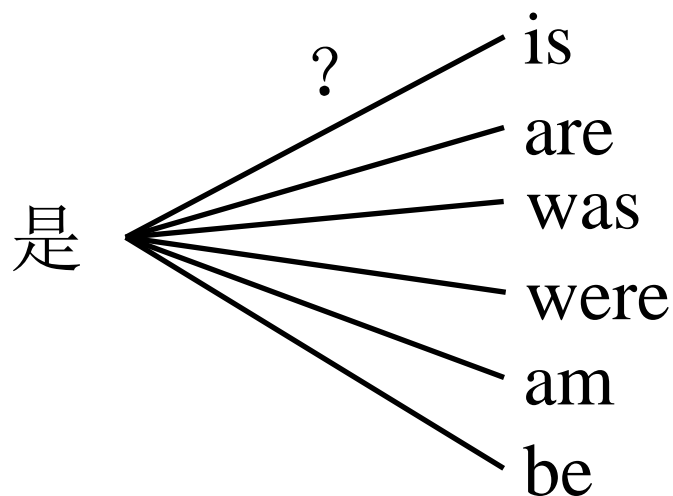
$$= \operatorname{argmax}_T \{ \log (P(T) \times P(S|T)) \}$$

$$= \operatorname{argmax}_T \{ P(T) \times P(S|T) \}$$

基于短语的统计机器翻译

- 基于词的翻译模型的问题：
 - 很难处理词义消歧问题
 - 很难处理一对多、多对一和多对多的翻译问题

澳洲 是 与 北韩 有 邦交 的 少数 国家 之一



基于短语的统计机器翻译

- 基于词的翻译模型的问题：
 - 很难处理词义消歧问题
 - 很难处理一对多、多对一和多对多的翻译问题

澳洲 是 与 北韩 有 邦交 的 少数 国家 之一

北韩 $\overset{?}{\rightarrow}$ North Korea

邦交 $\overset{?}{\rightarrow}$ the diplomatic relations

基于短语的统计机器翻译

➤ 基于短语的统计机器翻译:

澳洲 是 与 北韩 有 邦交 的 少数 国家 之一

(澳洲 是, Australia is)

(与 北韩, with North Korea)

(有 邦交, have the diplomatic relations)

(的 少数 国家 之一, one of the few countries that)

基于短语的统计机器翻译

澳洲 是	与 北韩	有 邦交	的 少数 国家 之一
------	------	------	------------

短语划分

--	--	--	--	--

短语翻译

Australia is	with North Korea	have diplomatic relations	one of the few countries that
--------------	------------------	---------------------------	-------------------------------

--	--	--	--

短语调序

Australia is	one of the few countries that	have diplomatic relations	with North Korea
--------------	-------------------------------	---------------------------	------------------

基于短语的统计机器翻译



[Koehn, 2003]

短语：连续的词串（非句法意义）

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S) \\
 &= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} P(S_1^K | S) \times P(T_1^K | S_1^K, S) \\
 &\quad \times P(T_1^{K'} | T_1^K, S_1^K, S) \times P(T | T_1^{K'}, T_1^K, S_1^K, S)
 \end{aligned}$$

基于短语的统计机器翻译

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S) \\
 &= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} \underbrace{P(S_1^K | S)}_{\substack{\downarrow \\ \text{短语划分模型}}} \underbrace{P(T_1^K | S_1^K, S)}_{\substack{\downarrow \\ \text{短语翻译模型}}} \underbrace{P(T_1^{K'} | T_1^K, S_1^K, S)}_{\substack{\downarrow \\ \text{短语调序模型}}} \underbrace{P(T | T_1^{K'}, T_1^K, S_1^K, S)}_{\substack{\downarrow \\ \text{目标语言模型}}}
 \end{aligned}$$

基于短语的统计机器翻译

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S) \\
 &= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} \underbrace{P(S_1^K | S)}_{\substack{\downarrow \\ \text{短语划分模型}}} \underbrace{P(T_1^K | S_1^K, S)} \underbrace{P(T_1^{K'} | T_1^K, S_1^K, S)} \underbrace{P(T | T_1^{K'}, T_1^K, S_1^K, S)}
 \end{aligned}$$

目标：将一个词序列如何划分为短语序列

方法：一般假设每一种短语划分方式都是等概率的

基于短语的统计机器翻译

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S) \\
 &= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} \underbrace{P(S_1^K | S)}_{\text{短语翻译模型}} \underbrace{P(T_1^K | S_1^K, S)}_{\text{短语调序模型}} \underbrace{P(T_1^{K'} | T_1^K, S_1^K, S)}_{\text{目标语言模型}} \underbrace{P(T | T_1^{K'}, T_1^K, S_1^K, S)}_{\text{目标语言模型}}
 \end{aligned}$$

剩下的三个核心模型：

1. 短语翻译模型： $P(T_1^K | S_1^K, S)$
2. 短语调序模型： $P(T_1^{K'} | T_1^K, S_1^K, S)$
3. 目标语言模型： $P(T | T_1^{K'}, T_1^K, S_1^K, S)$

基于短语的统计机器翻译

短语翻译模型: $P(T_1^K | S_1^K, S)$

1. 如何学习短语翻译规则

2. 如何估计短语翻译概率

双语句对词语对齐

短语翻译规则抽取

(澳洲 是, Australia is)

(与 北韩, with North Korea)

(有 邦交, have the diplomatic relations)

(的 少数 国家 之一, one of the few countries that)

基于短语的统计机器翻译

双语句对词语对齐

澳洲 是 与 北韩 有 邦交 的 少数 国家 之一

Australia is one of the few countries that have diplomatic relations with North Korea

IBM model 1-5

[illegible]

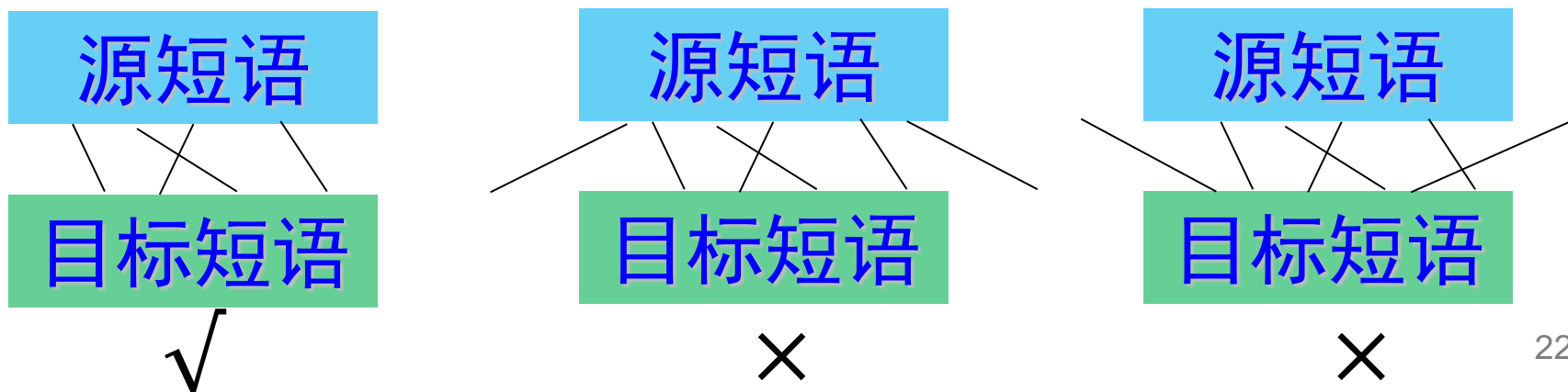
Australia
is
one
of
the
few
countries
that
have
diplomatic
relations
with
North
Korea
.

基于短语的统计机器翻译

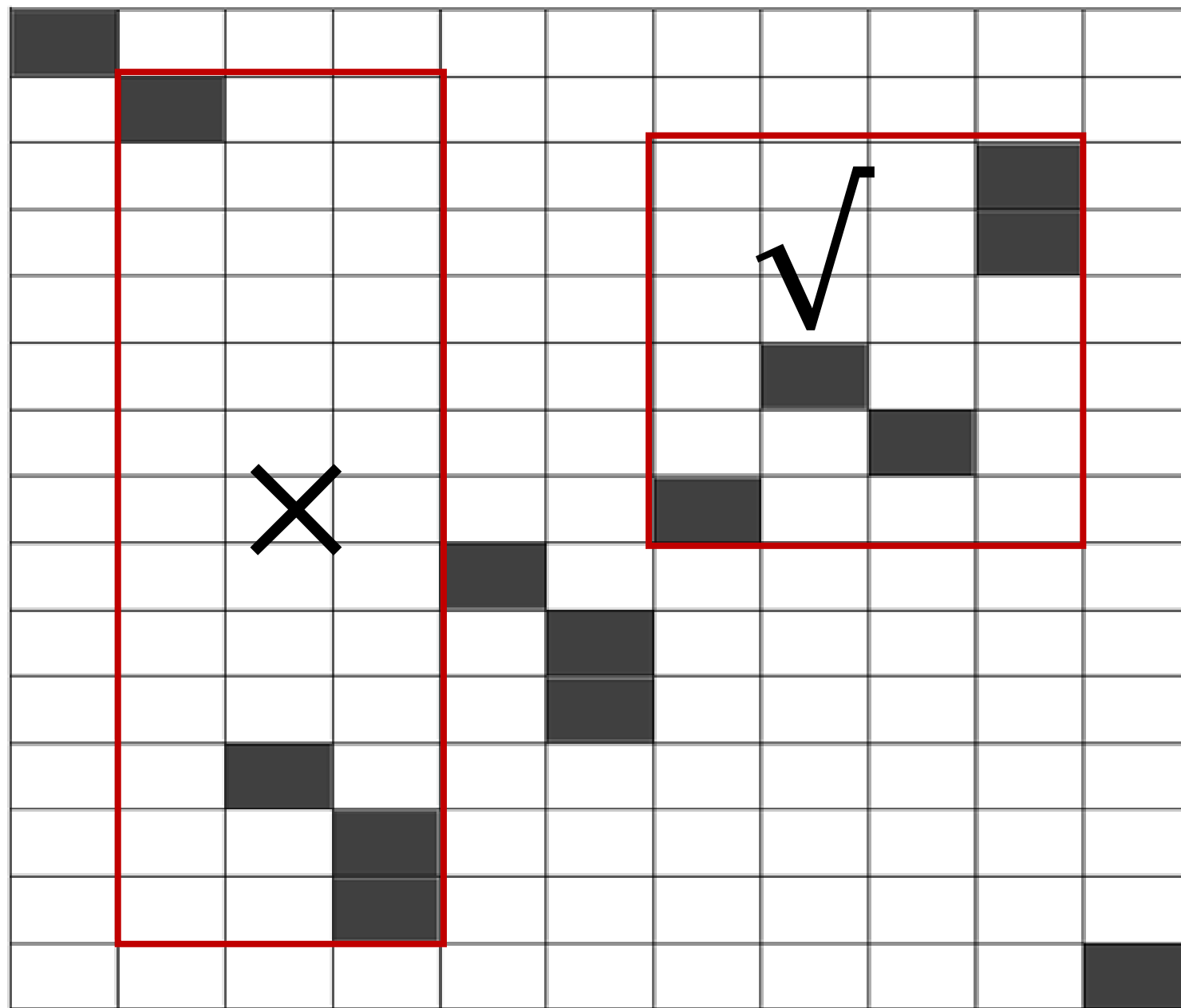
短语翻译规则抽取

算法： 对于源语言句子 S 中的任一短语 S_i^j ，根据词语对齐 A 找到目标语言句子 T 中的对齐片段 $T_{i'}^{j'}$ ，若 S_i^j 与 $T_{i'}^{j'}$ 满足对齐一致性，则 $(S_i^j, T_{i'}^{j'})$ 为一条短语翻译规则。

对齐一致性： S_i^j 中每个词 S_k ，若 $(k, k') \in A$ ，则 $i' \leq k' \leq j'$ ， $T_{i'}^{j'}$ 中每个词 $T_{t'}$ ，若 $(t, t') \in A$ ，则 $i \leq t \leq j$ 。



澳 洲 是 与 北 韩 有 邦 交 的 少 数 国 家 之 一 。



Australia
is
one
of
the
few
countries
that
have
diplomatic
relations
with
North
Korea
.

基于短语的统计机器翻译

短语翻译模型: $P(T_1^K | S_1^K, S)$

1. 如何学习短语翻译规则

2. 如何估计短语翻译概率

双语句对词语对齐

短语翻译规则抽取

(澳洲 是, Australia is)

(与 北韩, with North Korea)

(有 邦交, have the diplomatic relations)

(的 少数 国家 之一, one of the few countries that)

基于短语的统计机器翻译

短语翻译概率估计：4个翻译概率（最大似然）

1. 正向、逆向短语翻译概率 $p(t|s), p(s|t)$
2. 正向、逆向词汇化翻译概率 $p_{lex}(t|s), p_{lex}(s|t)$

(与 北韩, with North Korea)

$$p(t|s) = \frac{1}{2}$$

(与 北韩, and North Korea)

(和 北韩, with North Korea)

$$p(s|t) = \frac{1}{3}$$

(与 朝鲜, with North Korea)

基于短语的统计机器翻译

短语翻译概率估计：4个翻译概率（最大似然）

1. 正向、逆向短语翻译概率 $p(t|s), p(s|t)$
2. 正向、逆向词汇化翻译概率 $p_{lex}(t|s), p_{lex}(s|t)$

(与 北韩, with North Korea)

$$p(with|与) = 0.4$$

$$p(North|北韩) = 0.1$$

$$p(Korea|北韩) = 0.5$$

$$= \prod_{j=1}^{|t|} \frac{1}{|\{i | (j, i) \in A\}|} \sum_{\forall (j, i) \in A} p(t_j | s_i) \quad \begin{aligned} &= \frac{p_{lex}(t|s)}{0.4 \times 0.1 \times 0.5} \\ &= 0.02 \end{aligned}$$

基于短语的统计机器翻译

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S) \\
 &= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} \underbrace{P(S_1^K | S)} \underbrace{P(T_1^K | S_1^K, S)} \underbrace{P(T_1^{K'} | T_1^K, S_1^K, S)} \underbrace{P(T | T_1^{K'}, T_1^K, S_1^K, S)}
 \end{aligned}$$

剩下的三个核心模型：

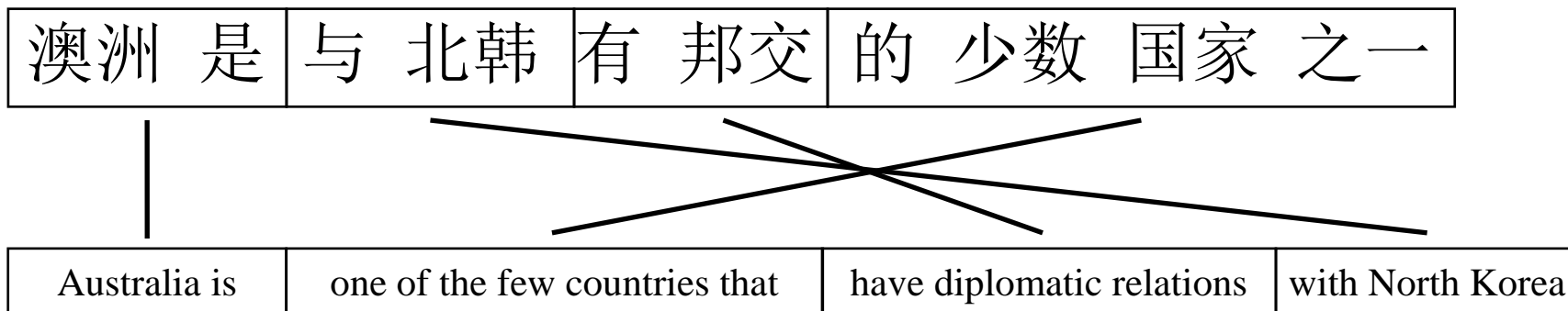
1. 短语翻译模型： $P(T_1^K | S_1^K, S)$
2. 短语调序模型： $P(T_1^{K'} | T_1^K, S_1^K, S)$
3. 目标语言模型： $P(T | T_1^{K'}, T_1^K, S_1^K, S)$

基于短语的统计机器翻译

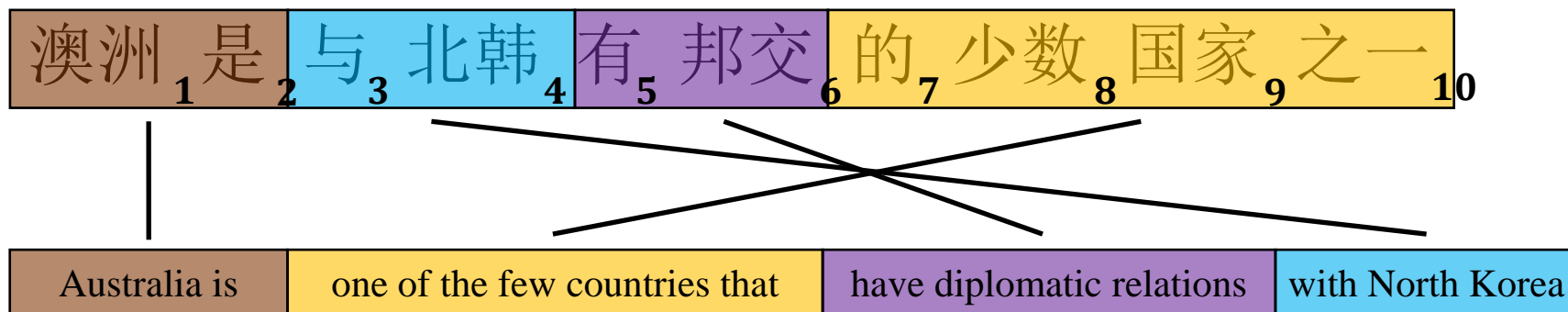
短语调序模型: $(T_1^{K'} | T_1^K, S_1^K, S)$

两种常用方法:

1. 距离跳转模型
2. 分类模型



基于短语的统计机器翻译



Monotone: 顺序拼接

Swap: 交换

Swap: 交换

(与 北韩, with North Korea) (有 邦交, have the diplomatic relations)

↓ 特征提取

f0=与, f1=北韩, f2=有, f3=邦交, f4=with, f5=korea, f6=have, f7=relations

↓ 最大熵模型

Swap

[Xiong et al., 2006]

基于短语的统计机器翻译

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S) \\
 &= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} \underbrace{P(S_1^K | S)}_{\text{短语翻译模型}} \underbrace{P(T_1^K | S_1^K, S)}_{\text{短语调序模型}} \underbrace{P(T_1^{K'} | T_1^K, S_1^K, S)}_{\text{目标语言模型}} \underbrace{P(T | T_1^{K'}, T_1^K, S_1^K, S)}_{\text{目标语言模型}}
 \end{aligned}$$

剩下的三个核心模型：

1. 短语翻译模型： $P(T_1^K | S_1^K, S)$
2. 短语调序模型： $P(T_1^{K'} | T_1^K, S_1^K, S)$
3. 目标语言模型： $P(T | T_1^{K'}, T_1^K, S_1^K, S)$

基于短语的统计机器翻译

$$\begin{aligned}
 & P(w_1 w_2 \cdots w_{t-1} w_n) \\
 = & \prod_{t=1}^n P(w_t | w_{t-1} \cdots w_1) \\
 \approx & \prod_{t=1}^n P(w_t | w_{t-1} \cdots w_{t-n+1}) \\
 & \quad \quad \quad \downarrow \\
 = & \frac{P(w_t | w_{t-1} \cdots w_{t-n+1})}{\text{count}(w_{t-1} \cdots w_{t-n+1} w_t)} \\
 & \quad \quad \quad \text{count}(w_{t-1} \cdots w_{t-n+1})
 \end{aligned}$$

基于短语的统计机器翻译

$$T' = \operatorname{argmax}_T P(T|S)$$


$$= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S)$$

$$= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} P(S_1^K | S) P(T_1^K | S_1^K, S) P(T_1^{K'} | T_1^K, S_1^K, S) P(T | T_1^{K'}, T_1^K, S_1^K, S)$$

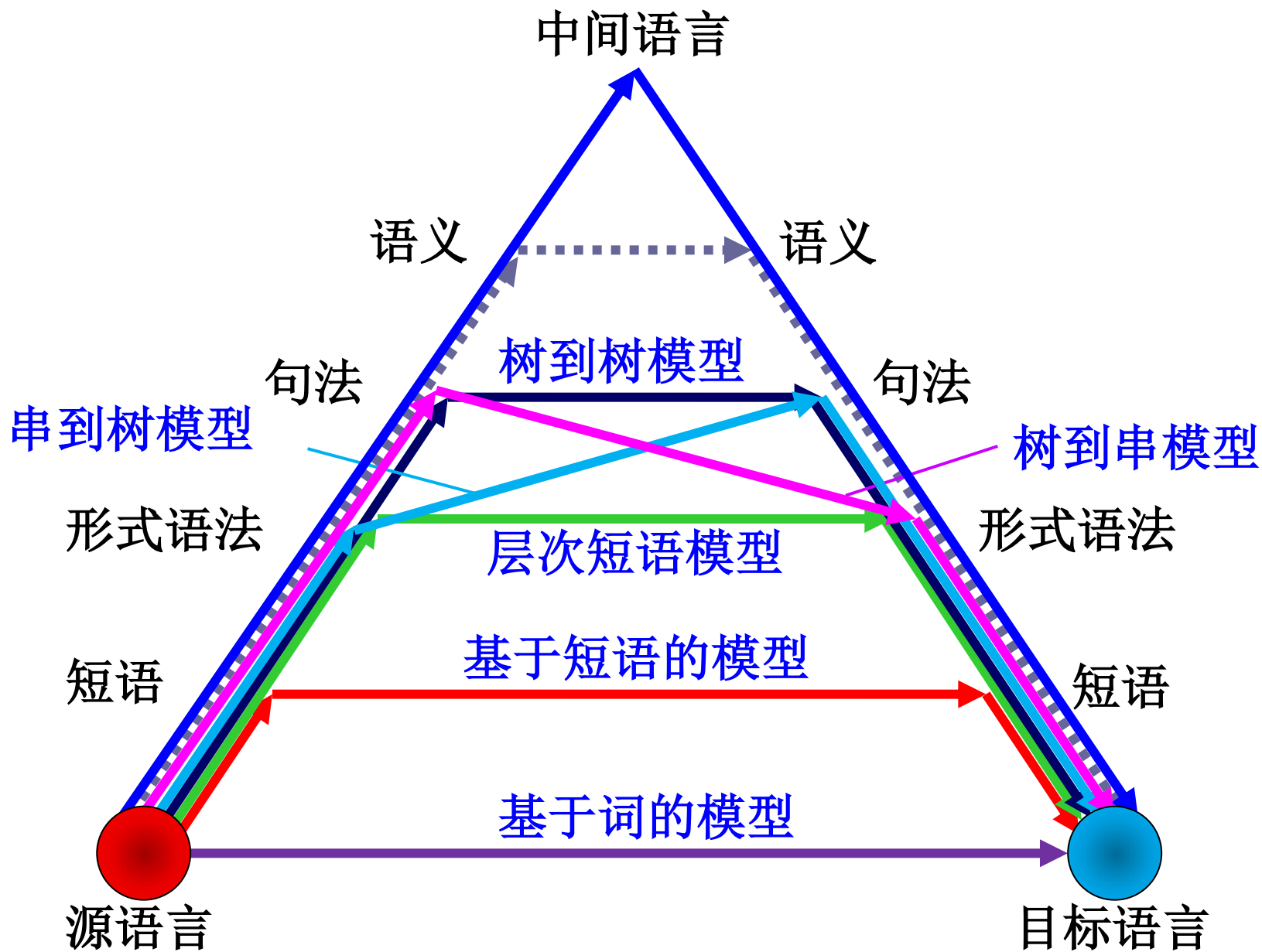


$$T' = \operatorname{argmax}_T P(T|S) = \operatorname{argmax}_T \left\{ \sum_1^M \lambda_m h_m(T, S) \right\}$$

基于短语的统计机器翻译

$$T' = \operatorname{argmax}_T P(T|S) = \operatorname{argmax}_T \left\{ \sum_1^M \lambda_m h_m(T, S) \right\}$$


$$\begin{aligned} h_1(T, S) &= \log p(t|s) & h_5(T, S) &= \log P(T_1^{K'} | T_1^K, S_1^K, S) \\ h_2(T, S) &= \log p(s|t) & h_6(T, S) &= \log P(T | T_1^{K'}, T_1^K, S_1^K, S) \\ h_3(T, S) &= \log p_{lex}(t|s) & h_7(T, S) &= \log \text{len}(T) \\ h_4(T, S) &= \log p_{lex}(s|t) & h_8(T, S) &= \log \text{count}(\text{phrases}) = \log K \end{aligned}$$



统计机器翻译中 深度学习的应用

基于短语的统计机器翻译

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S) \\
 &= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} \underbrace{P(S_1^K | S)}_{\text{短语翻译模型}} \underbrace{P(T_1^K | S_1^K, S)}_{\text{短语调序模型}} \underbrace{P(T_1^{K'} | T_1^K, S_1^K, S)}_{\text{目标语言模型}} \underbrace{P(T | T_1^{K'}, T_1^K, S_1^K, S)}_{\text{目标语言模型}}
 \end{aligned}$$

剩下的三个核心模型：

1. 短语翻译模型： $P(T_1^K | S_1^K, S)$
2. 短语调序模型： $P(T_1^{K'} | T_1^K, S_1^K, S)$
3. 目标语言模型： $P(T | T_1^{K'}, T_1^K, S_1^K, S)$

基于计数的N-元语言模型

该课程很枯燥，大家觉得很无聊。

$$P(w_t | w_{t-1} \cdots w_{t-n+1}) = \frac{P(\text{无聊} | \text{很}) \cdot \text{count}(\text{很 无聊})}{\text{count}(\text{很})}$$

$P(\text{无聊} | \text{很})$
vs.
 $P(\text{枯燥} | \text{很})$

问题①：数据稀疏
N-元组“很 无聊”未出现过，则回退

问题②：忽略语义相似性
“无聊”与“枯燥”虽语义相似，但无法共享信息

基于计数的N-元语言模型

- 典型方法：抽象符号（字符串）

该课程很枯燥，大家觉得很无聊。

w_0 =该 w_1 =课程 w_2 =很 w_3 =枯燥 w_4 =,
 w_5 =大家 w_6 =觉得 w_7 =很 w_8 =无聊 w_9 =。

- 等价表示方法：one-hot表示法

$|V|$

$$\begin{bmatrix} \vdots \end{bmatrix}$$


所有词按照出现的顺序排序



每个词语将对应唯一的下标

枯燥

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

无聊

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

词语表示

- 问题

枯燥

$$\begin{bmatrix} 0 \\ \mathbf{1} \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

无聊

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ \mathbf{1} \\ 0 \end{bmatrix}$$

枯燥 \otimes 无聊

$$\begin{bmatrix} 0 \\ \mathbf{1} \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

\times

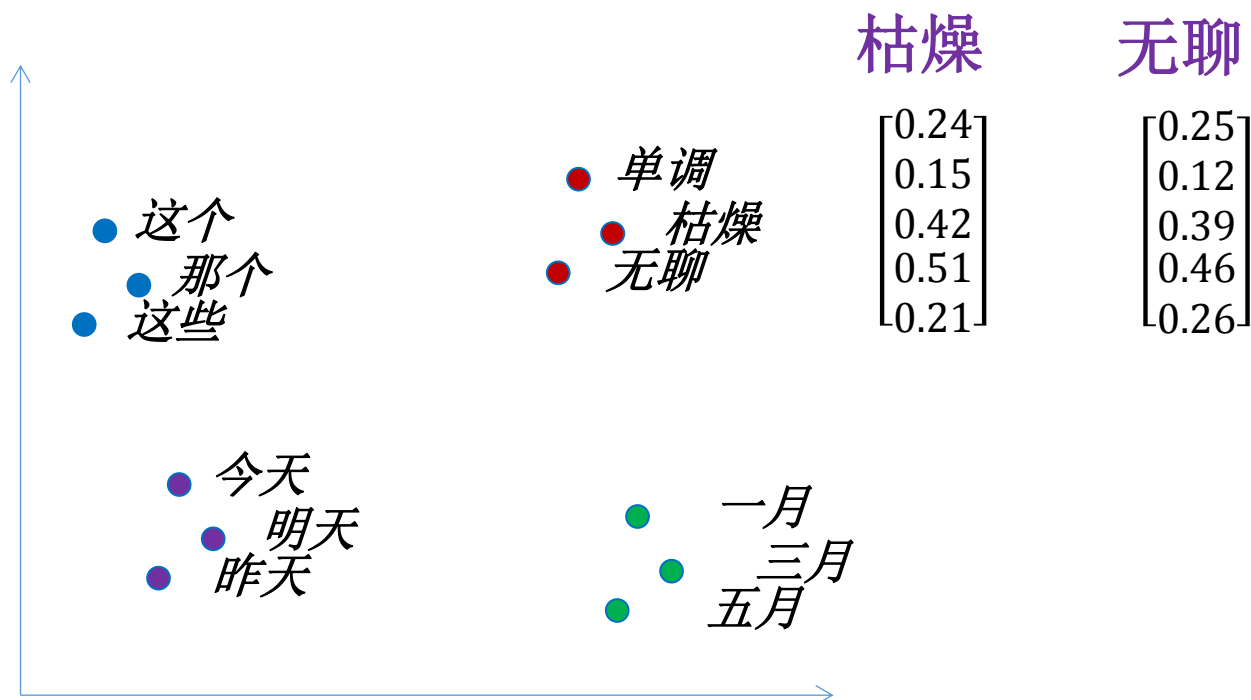
$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ \mathbf{1} \\ 0 \end{bmatrix}$$

$= 0$



任意两个词之间的相似度都为0!

词语表示



低维、稠密的连续实数空间

神经网络语言模型

很

$P(\text{无聊}|\text{很})$

vs.

$P(\text{枯燥}|\text{很})$

$\begin{bmatrix} 0.01 \\ 0.59 \\ 0.18 \\ 0.05 \\ 0.47 \end{bmatrix}$

$$P \left(\begin{bmatrix} 0.25 \\ 0.12 \\ 0.39 \\ 0.46 \\ 0.26 \end{bmatrix} \middle| \begin{bmatrix} 0.01 \\ 0.59 \\ 0.18 \\ 0.05 \\ 0.47 \end{bmatrix} \right)$$

vs.

$$P \left(\begin{bmatrix} 0.24 \\ 0.15 \\ 0.42 \\ 0.51 \\ 0.21 \end{bmatrix} \middle| \begin{bmatrix} 0.01 \\ 0.59 \\ 0.18 \\ 0.05 \\ 0.47 \end{bmatrix} \right)$$

$$f \left(\begin{bmatrix} 0.01 \\ 0.59 \\ 0.18 \\ 0.05 \\ 0.47 \\ 0.25 \\ 0.12 \\ 0.39 \\ 0.46 \\ 0.26 \end{bmatrix} \right)$$

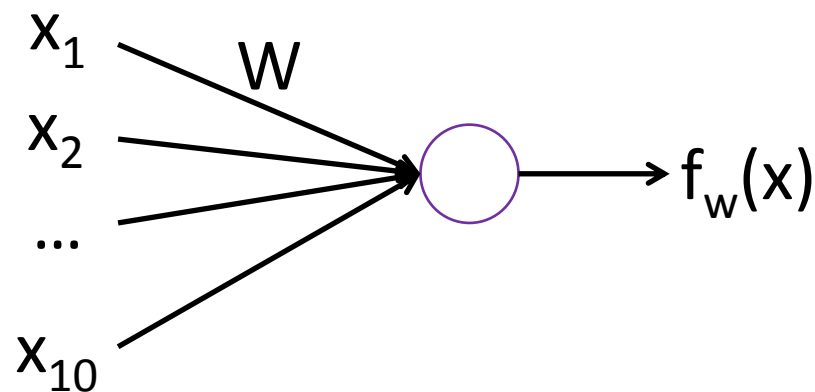
vs.

$$f \left(\begin{bmatrix} 0.01 \\ 0.59 \\ 0.18 \\ 0.05 \\ 0.47 \\ 0.24 \\ 0.15 \\ 0.42 \\ 0.51 \\ 0.21 \end{bmatrix} \right)$$

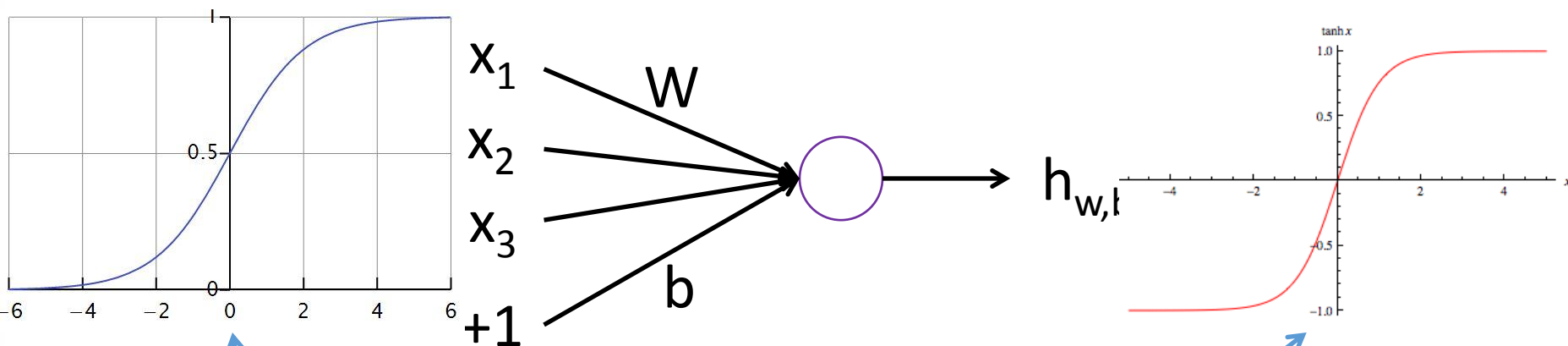
神经网络语言模型

$P(\text{无聊}|\text{很})$

$$f \begin{pmatrix} 0.01 \\ 0.59 \\ 0.18 \\ 0.05 \\ 0.47 \\ 0.25 \\ 0.12 \\ 0.39 \\ 0.46 \\ 0.26 \end{pmatrix} = f \begin{pmatrix} w_1 \times 0.01 \\ w_2 \times 0.59 \\ w_3 \times 0.18 \\ w_4 \times 0.05 \\ w_5 \times 0.47 \\ w_6 \times 0.25 \\ w_7 \times 0.12 \\ w_8 \times 0.39 \\ w_9 \times 0.46 \\ w_{10} \times 0.26 \end{pmatrix} = f(WX)$$



神经网络语言模型



$$h_{W,b}(x) = f(W^T x + b)$$

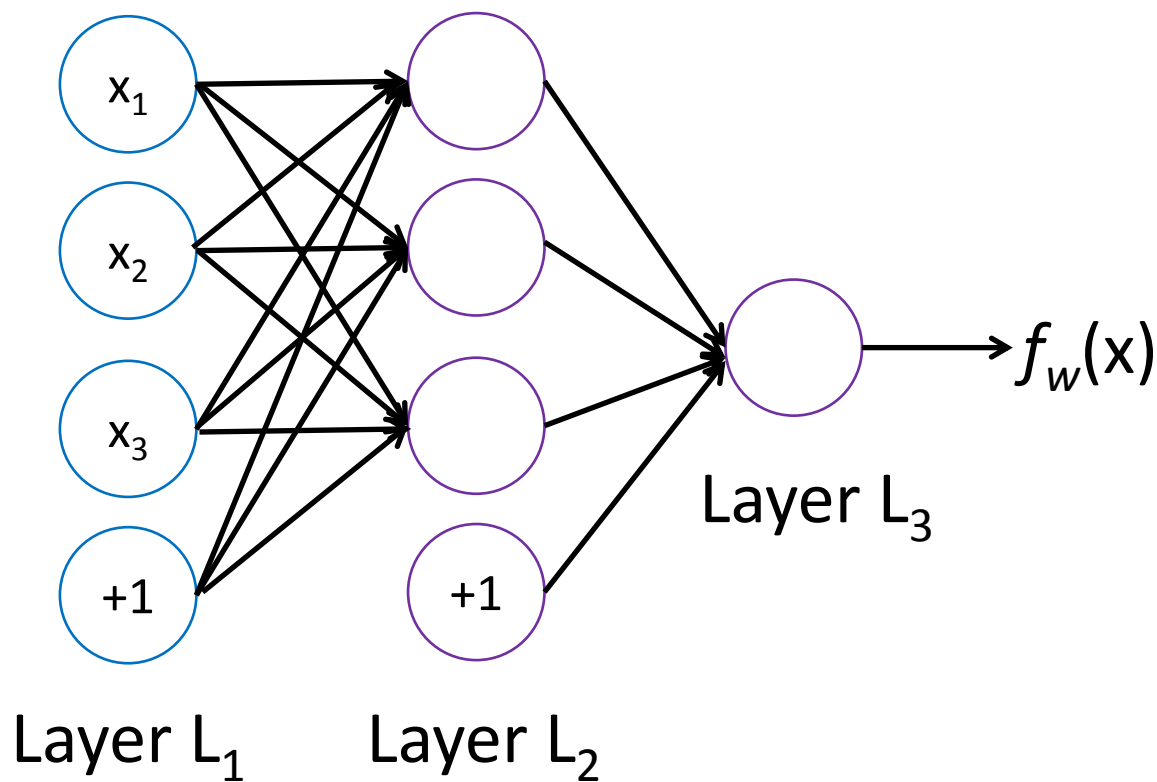
f : 非线性激活函数

$$\begin{cases} f(z) = \frac{1}{1 + \exp(-z)} \\ f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \end{cases}$$

$$f'(z) = f(z)(1 - f(z))$$

$$f'(z) = 1 - f^2(z)$$

神经网络语言模型



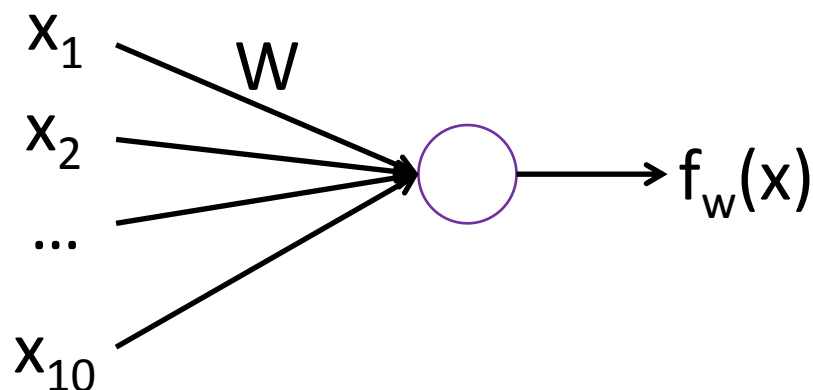
神经网络语言模型

$P(\text{无聊}|\text{很})$

vs.

$P(\text{枯燥}|\text{很})$

$$f \begin{pmatrix} 0.01 \\ 0.59 \\ 0.18 \\ 0.05 \\ 0.47 \\ 0.25 \\ 0.12 \\ 0.39 \\ 0.46 \\ 0.26 \end{pmatrix} = f \begin{pmatrix} w_1 \times 0.01 \\ w_2 \times 0.59 \\ w_3 \times 0.18 \\ w_4 \times 0.05 \\ w_5 \times 0.47 \\ w_6 \times 0.25 \\ w_7 \times 0.12 \\ w_8 \times 0.39 \\ w_9 \times 0.46 \\ w_{10} \times 0.26 \end{pmatrix} = f(WX)$$



问题①：词向量

如何将每个词映射到实数向量空间中的一个点

问题②： f 函数的设计

设计什么样的神经网络结构模拟函数 f

神经网络语言模型-词向量

$$L = \begin{bmatrix} \text{枯燥} & \dots & \text{单调} & \text{无聊} \end{bmatrix} \quad V$$

$$D = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad L \in R^{D \times V}$$

枯燥 ... 单调 无聊

- 通常称为look-up table
 - 我们可以对 L 右乘一个词的one-hot表示 e 得到该词的低维、稠密的实数向量表达： $x = Le$

神经网络语言模型-词向量

$$L = \begin{bmatrix} \text{枯燥} & \dots & \text{单调} & \text{无聊} \end{bmatrix}_D \quad L \in R^{D \times V}$$

The diagram illustrates the word embedding matrix L . It is a matrix of size $D \times V$, where D is the dimensionality of the word vectors and V is the vocabulary size. The matrix is shown as a grid of red dots. The first column is labeled '枯燥' (boring), the second column is labeled '单调' (monotonous), and the third column is labeled '无聊' (boredom). The matrix is enclosed in large blue brackets, with the dimension D indicated on the right. The dimension V is indicated above the matrix. The matrix is labeled L on the left, and the equation $L \in R^{D \times V}$ is shown on the right.

- 词表规模 V 和词向量维度 D 如何确定
 - V 的确定：1, 训练数据中所有词；2, 频率高于某个阈值的所有词；3, 前 V 个频率最高的词
 - D 的确定：超参数，人工设定，一般从几十到几百

神经网络语言模型-词向量

$$L = \begin{bmatrix} \text{枯燥} & \dots & \text{单调} & \text{无聊} \end{bmatrix}_D^V, \quad L \in R^{D \times V}$$

The diagram illustrates the word embedding matrix L . It is a matrix of size $D \times V$, where D is the dimensionality and V is the vocabulary size. The matrix is shown as a grid of red dots. The first column is labeled '枯燥' (boring), the second column is labeled '单调' (monotonous), and the third column is labeled '无聊' (boredom). The matrix is enclosed in large blue brackets, with a purple box highlighting the first column. The matrix is labeled L on the left and $L \in R^{D \times V}$ on the right.

- 如何学习 L
 - 通常先随机初始化，然后通过目标函数优化词的向量表达（e.g. 最大化语言模型似然度）

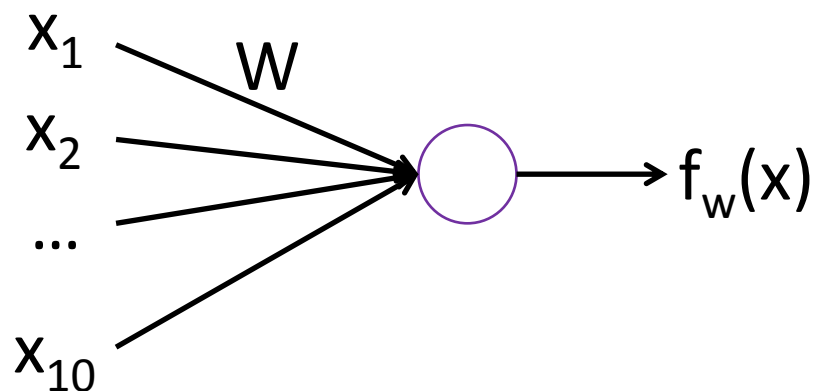
神经网络语言模型

$P(\text{无聊}|\text{很})$

vs.

$P(\text{枯燥}|\text{很})$

$$f \begin{pmatrix} 0.01 \\ 0.59 \\ 0.18 \\ 0.05 \\ 0.47 \\ 0.25 \\ 0.12 \\ 0.39 \\ 0.46 \\ 0.26 \end{pmatrix} = f \begin{pmatrix} w_1 \times 0.01 \\ w_2 \times 0.59 \\ w_3 \times 0.18 \\ w_4 \times 0.05 \\ w_5 \times 0.47 \\ w_6 \times 0.25 \\ w_7 \times 0.12 \\ w_8 \times 0.39 \\ w_9 \times 0.46 \\ w_{10} \times 0.26 \end{pmatrix} = f(WX)$$



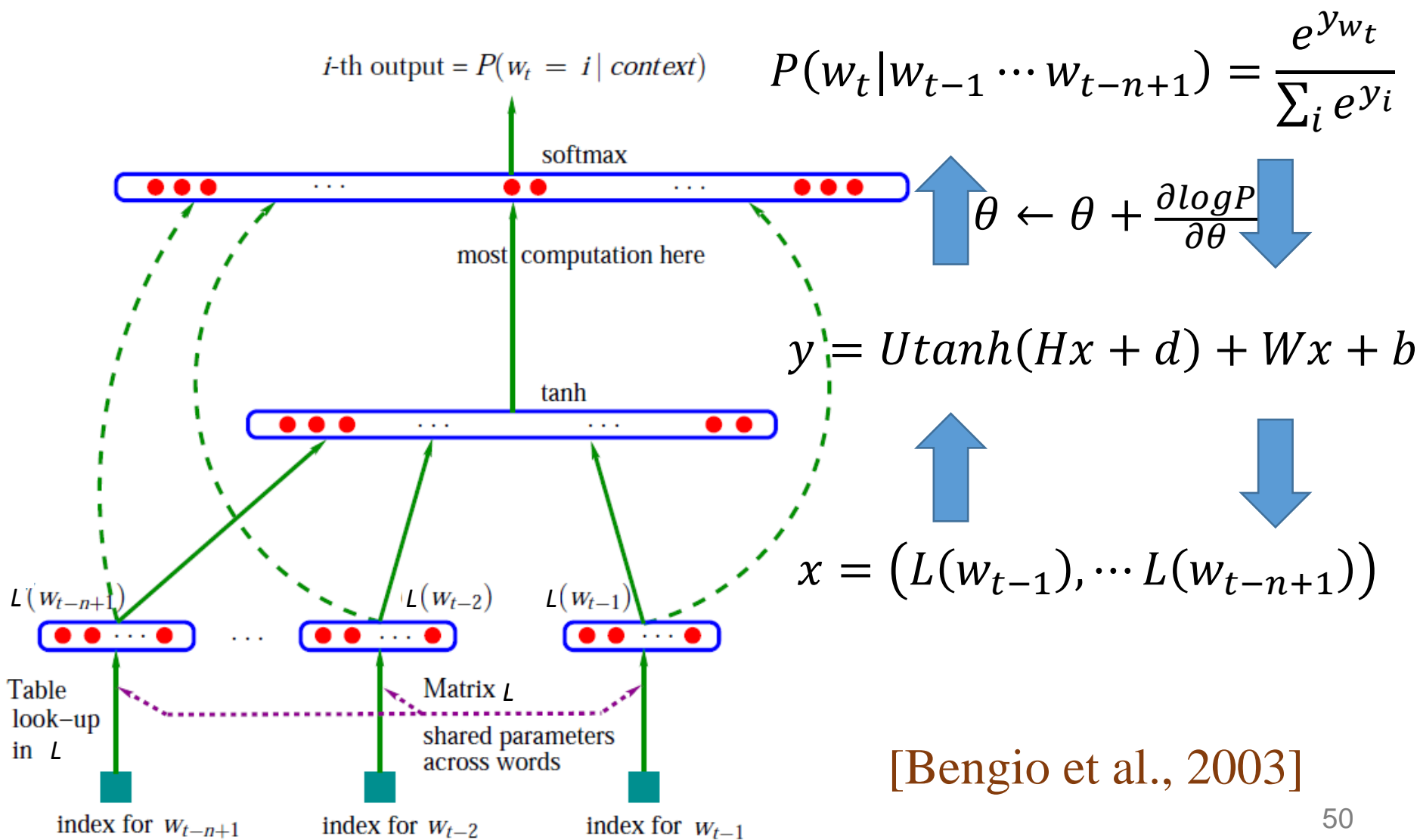
问题①：词向量

如何将每个词映射到实数向量空间中的一个点

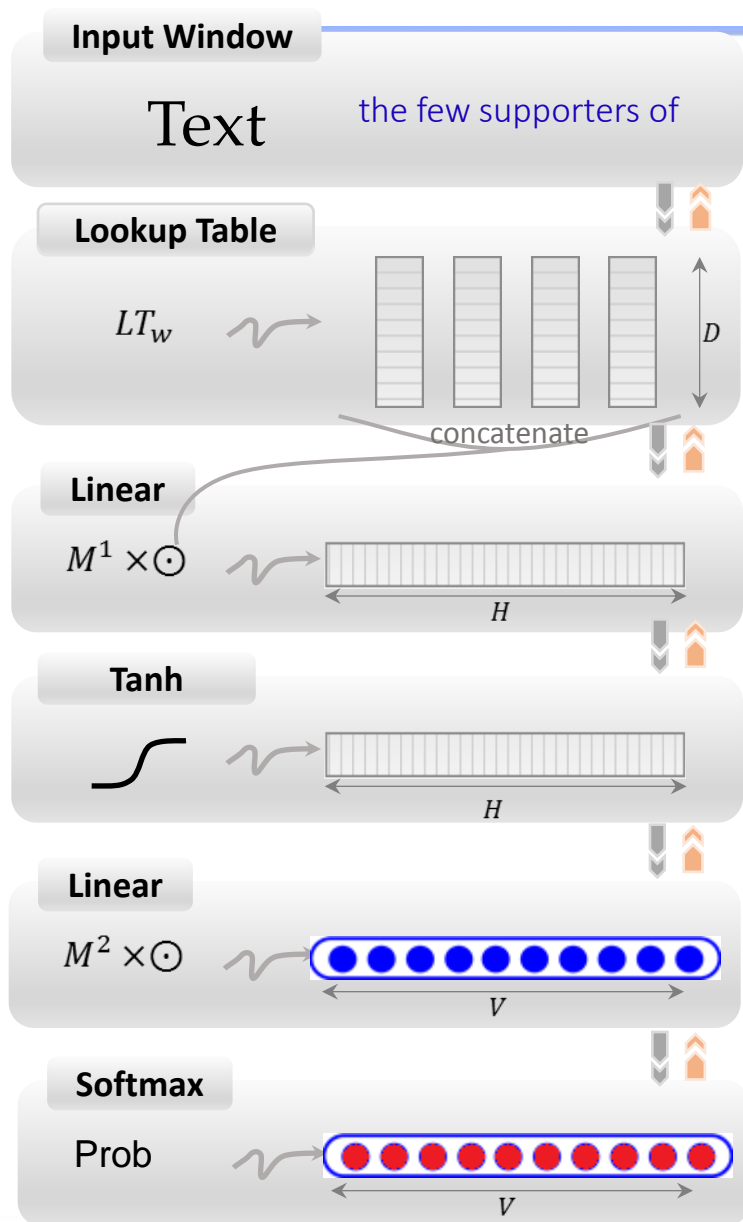
问题②： f 函数的设计

设计什么样的神经网络结构模拟函数 f

前馈神经网络语言模型



前馈神经网络语言模型



$$P(\text{this} | \text{the, few, supporters, of})$$

将每个词通过词向量矩阵 L 映射为低维实数向量

$$\text{of} \rightarrow (0.23, 0.15, 0.08, 0.31, \dots, 0.42)$$

拼接所有词的向量，形成一个向量

隐藏层:

线性映射+非线性变换

⋮

Softmax 输出层:

$$P(\text{this} | \text{the, few, supporters, of})$$

前馈神经网络语言模型

- 问题

仅对小窗口的历史信息建模

例如5-gram语言模型，仅考虑前面4个词的历史信息



能否对所有的历史信息进行建模

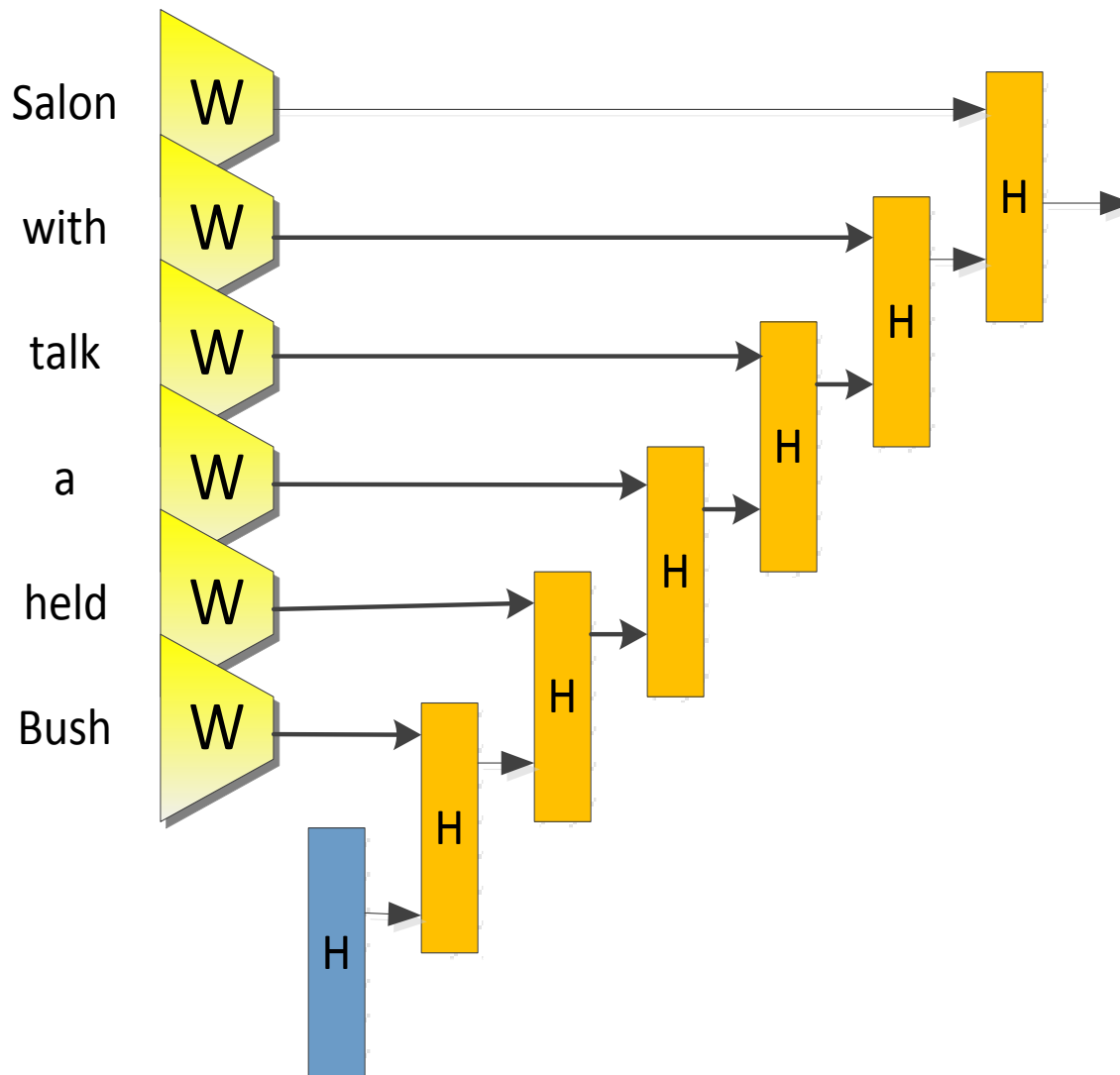
即第t个词的语言模型概率依赖于所有前t-1个词

$$P(w_t | w_{t-1} \cdots w_{t-n+1})$$



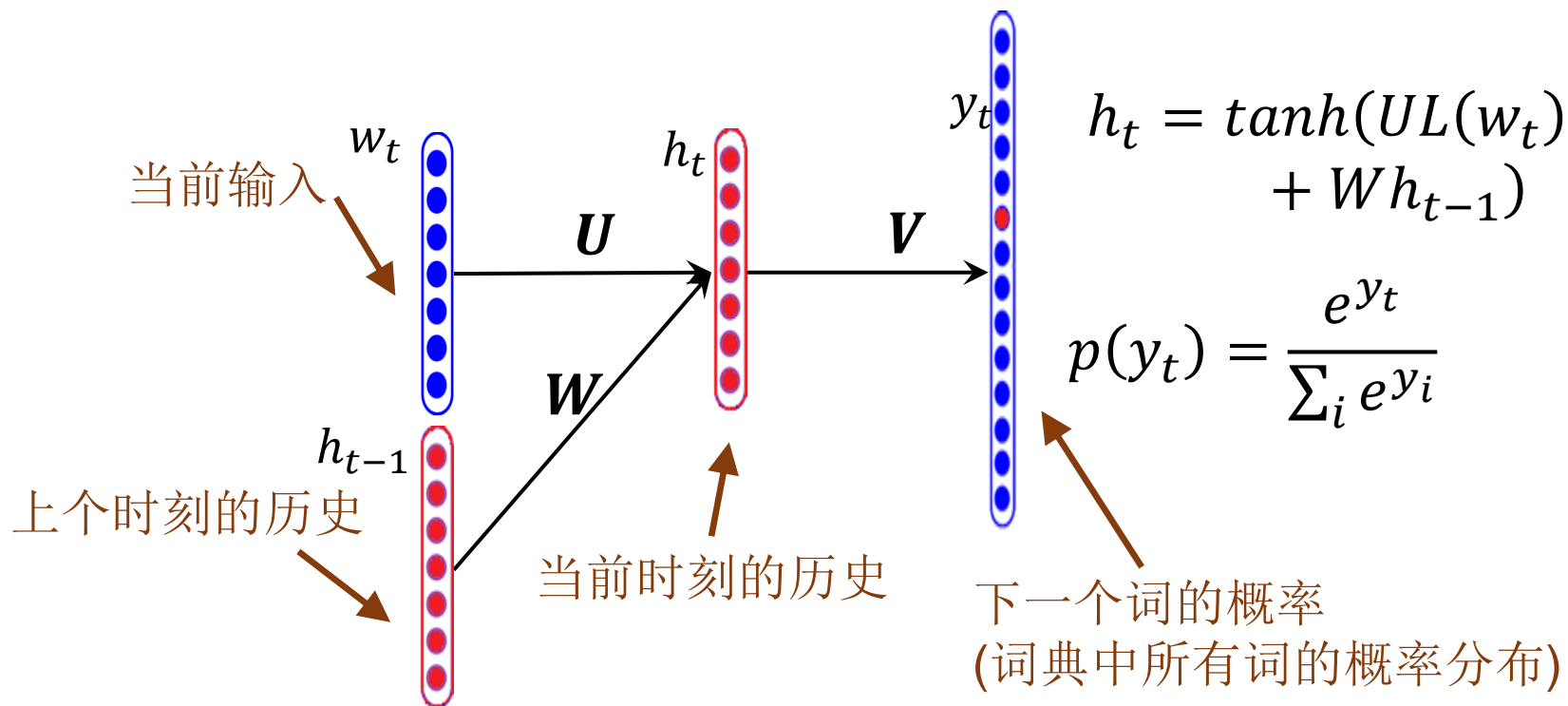
$$P(w_t | w_{t-1} \cdots w_2 w_1)$$

循环神经网络语言模型

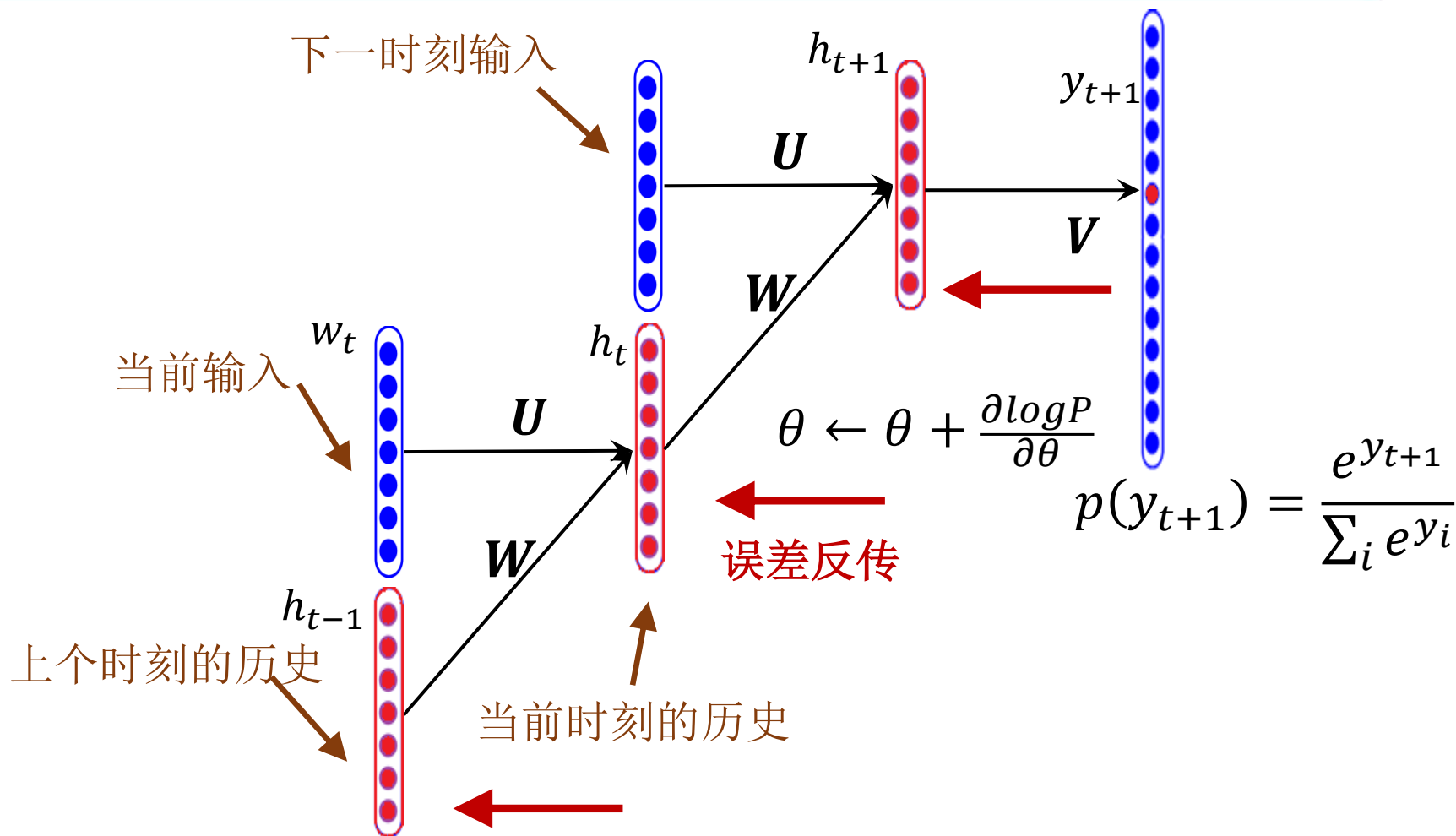


循环神经网络语言模型

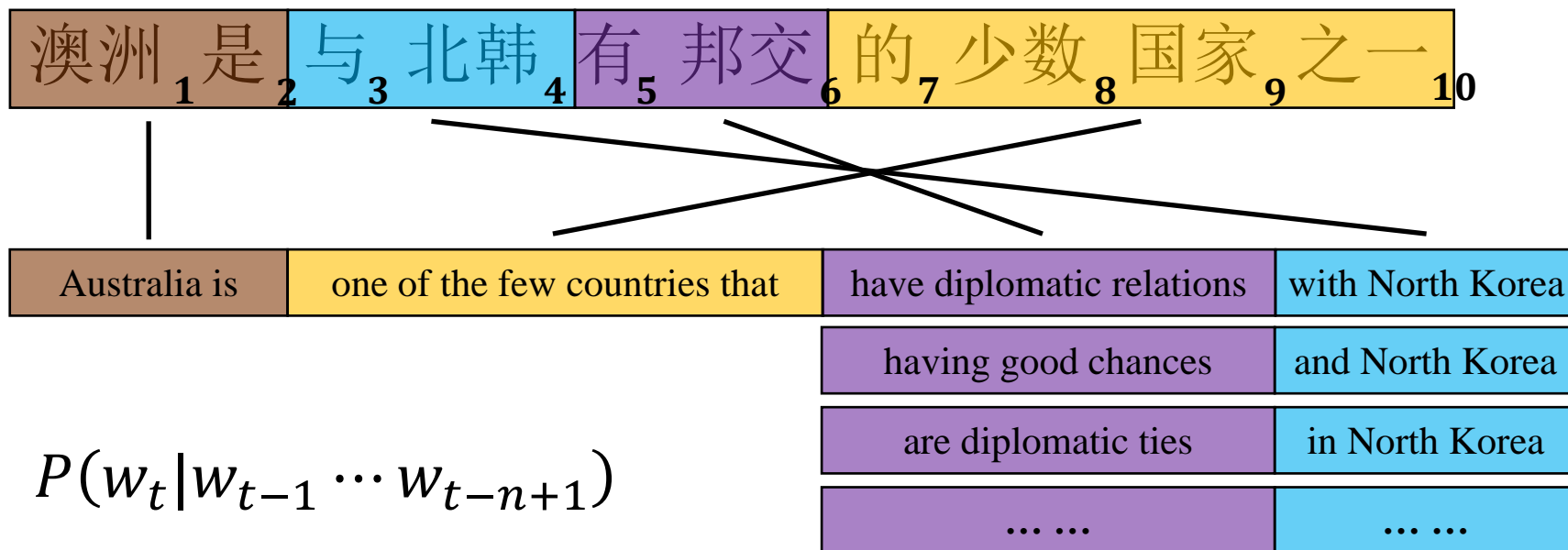
- 输入: $t - 1$ 时刻历史 h_{t-1} 与 t 时刻输入 w_t
- 输出: t 时刻历史 h_t 与 下个时刻 $t + 1$ 输入 y_t 的概率



循环神经网络语言模型



神经网络语言模型-融入解码器



解码中计算语言模型概率：①翻译候选多；②自底往上的解码方式无法知道译文前缀。

➤ 适合前馈神经网络

神经网络语言模型-融入解码器

- 汉语-英语: NIST-2012 受限评测

Setting	dev	2004	2005	2006
baseline	38.2	38.4	37.7	34.3
+NNLM decoding	39.1	39.5	38.8	34.9

[Vaswani et al., EMNLP-2013]

神经网络语言模型-译文重排序

澳洲₁ 是₂ 与₃ 北韩₄ 有₅ 邦交₆ 的₇ 少数₈ 国家₉ 之一₁₀

Australia is	one of the few countries that	have diplomatic relations	with North Korea
Australia are	a few countries	having good chances	and North Korea
Australia was	a country that	are diplomatic ties	in North Korea
...

对多个翻译结果重排序：①翻译候选少；②可以利用译文前缀信息。

➤ 适合循环神经网络

神经网络语言模型-译文重排序

- 法语-英语: WMT-2012

Setting	dev	news2010	news2011	com2011	Ave.
baseline	26.6	27.6	28.3	27.5	27.8
+RNNLM (2m)	27.5	28.1	28.6	28.1	28.3
+RNNLM (50m)	27.7	28.2	29.0	28.1	28.5

[Auli et al., EMNLP-2013]

神经网络语言模型-译文重排序

- 德语-英语: WMT-2012

Setting	dev	news2010	news2011	com2011	Ave.
baseline	21.2	20.7	19.2	20.6	20.0
+RNNLM (2m)	21.8	20.9	19.4	20.9	20.3
+RNNLM (50m)	22.1	21.1	19.7	21.0	20.5

[Auli et al., EMNLP-2013]

神经网络语言模型-译文重排序

- 英语-德语: WMT-2012

Setting	dev	news2010	news2011	com2011	Ave.
baseline	15.2	15.6	14.3	15.7	15.1
+RNNLM (2m)	15.7	15.9	14.6	16.0	15.4
+RNNLM (50m)	15.8	15.9	14.7	16.1	15.5

[Auli et al., EMNLP-2013]

基于短语的统计机器翻译

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S) \\
 &= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} \underbrace{P(S_1^K | S)}_{\text{短语翻译模型}} \underbrace{P(T_1^K | S_1^K, S)}_{\text{短语调序模型}} \underbrace{P(T_1^{K'} | T_1^K, S_1^K, S)}_{\text{目标语言模型}} \underbrace{P(T | T_1^{K'}, T_1^K, S_1^K, S)}_{\text{目标语言模型}}
 \end{aligned}$$

剩下的三个核心模型：

1. 短语翻译模型： $P(T_1^K | S_1^K, S)$
2. 短语调序模型： $P(T_1^{K'} | T_1^K, S_1^K, S)$
3. 目标语言模型： $P(T | T_1^{K'}, T_1^K, S_1^K, S)$

基于最大似然的短语翻译概率估计

短语翻译概率估计：4个翻译概率（最大似然）

1. 正向、逆向短语翻译概率 $p(t|s), p(s|t)$
2. 正向、逆向词汇化翻译概率 $p_{lex}(t|s), p_{lex}(s|t)$

(与 北韩, with North Korea)

100

1

$$p(t|s) = \frac{1}{100}$$

最大似然概率估计无法刻画双语短语之间的语义相似程度！

基于语义向量空间的翻译置信度估计

短语翻译概率估计：4个翻译概率（最大似然）

1. 正向、逆向短语翻译概率 $p(t|s), p(s|t)$
2. 正向、逆向词汇化翻译概率 $p_{lex}(t|s), p_{lex}(s|t)$

(与 北韩, with North Korea)

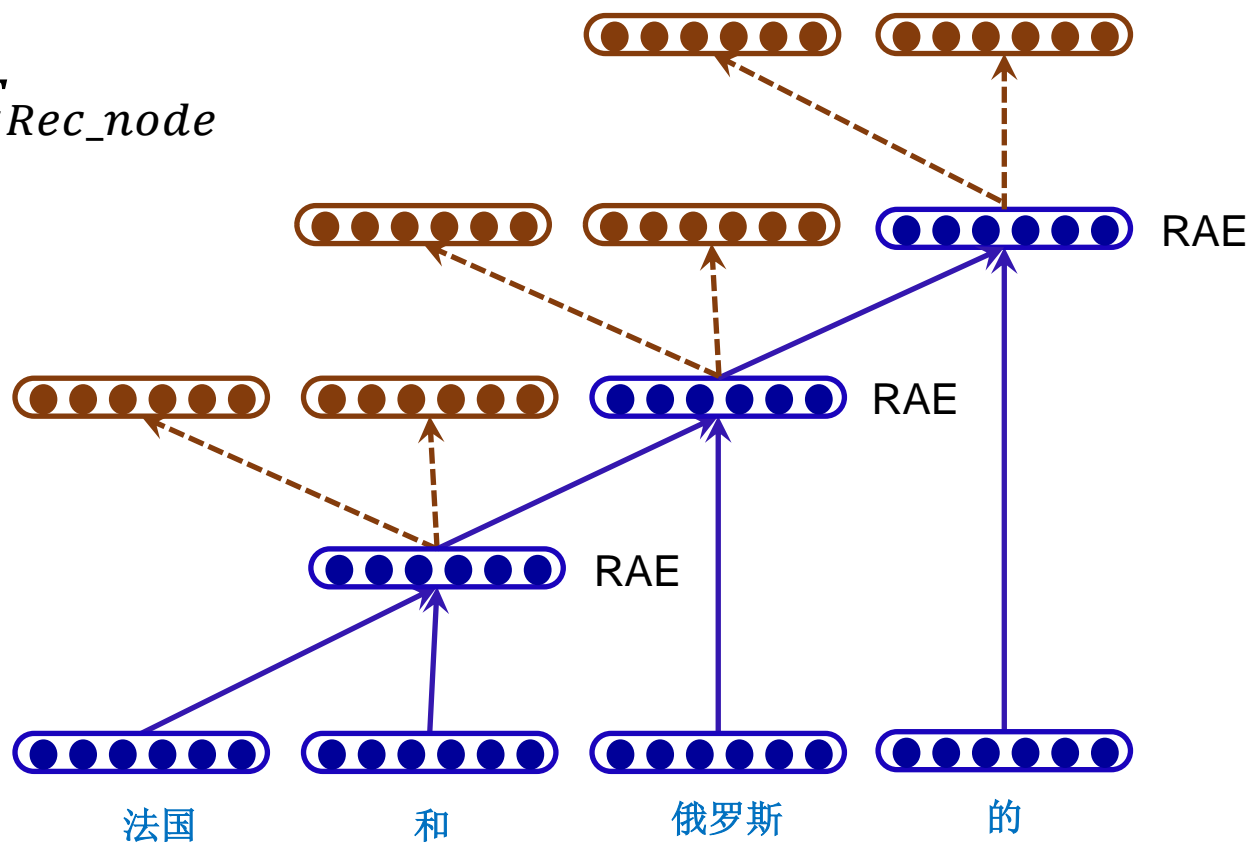


可以在语义向量空间中计算双语短语的语义相似度

基于语义向量空间的短语表示

- 递归自动编码器
- 目标函数：最小化所有节点的重构误差

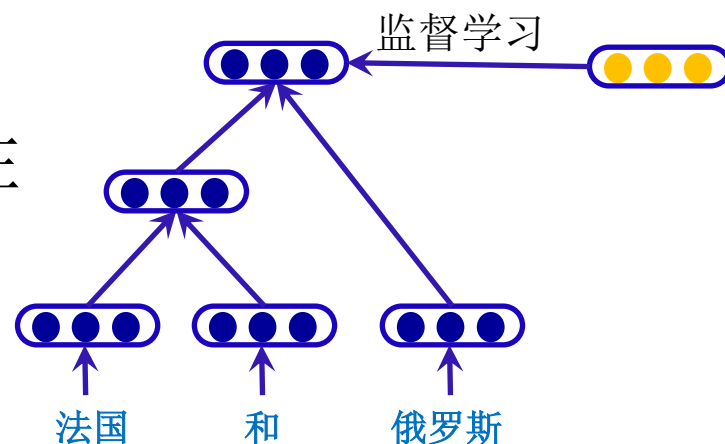
$$E_{total} = \sum_{node} E_{Rec_node}$$



基于语义向量空间的短语表示

- 理想方法：有标注数据

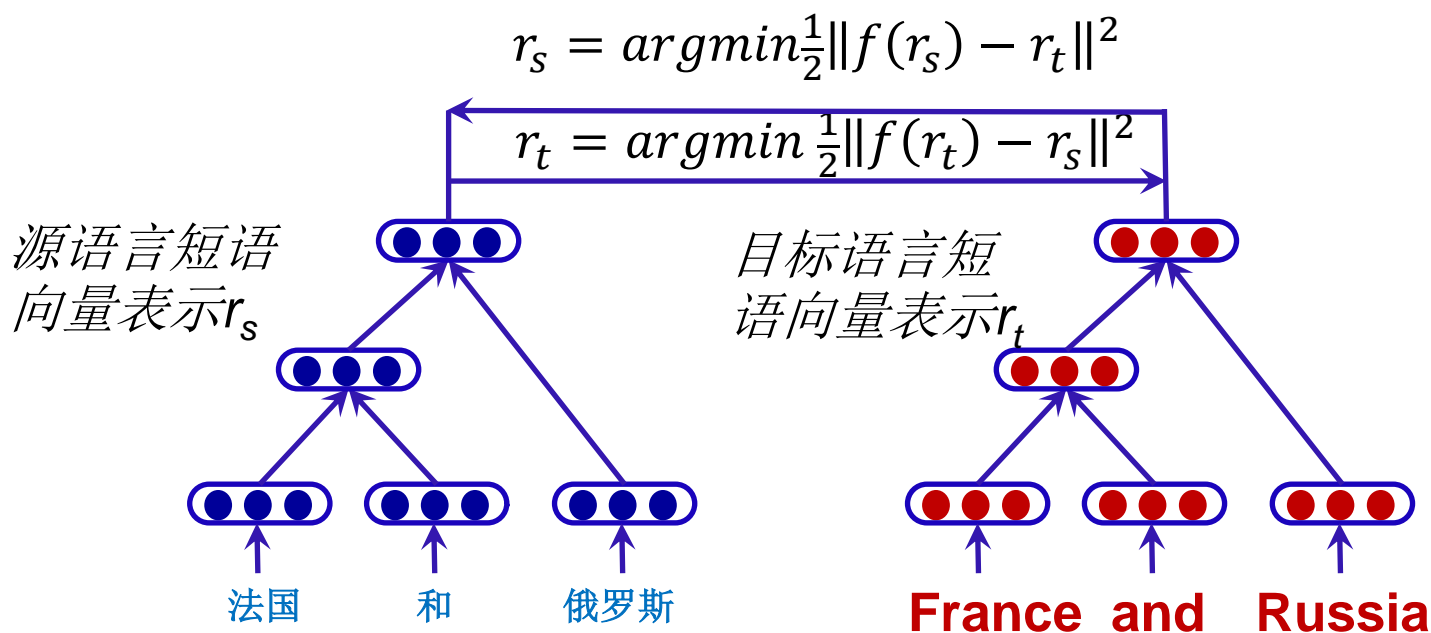
- 但是，现实中不存在正确标注的短语向量



基于语义向量空间的短语表示

- 假设
 - 短语与其翻译具有相同的语义向量表示
- 目标函数
 - 最小化短语翻译对间的语义表示误差
- 模型
 - Pre-training: 无监督递归自编码器学习短语初始表示
 - Fine-tuning: 相互监督学习, 优化短语向量表示

基于语义向量空间的短语表示



基于语义向量空间的短语表示

- 目标函数

正则化项

$$J = E(S, T; \theta) + \frac{1}{2} \|\lambda\|^2$$

重构误差

双语语义误差

$$E(S, T; \theta) = \alpha E_{rec}(S, T; \theta) + (1 - \alpha) E_{regression}(S, T; \theta)$$

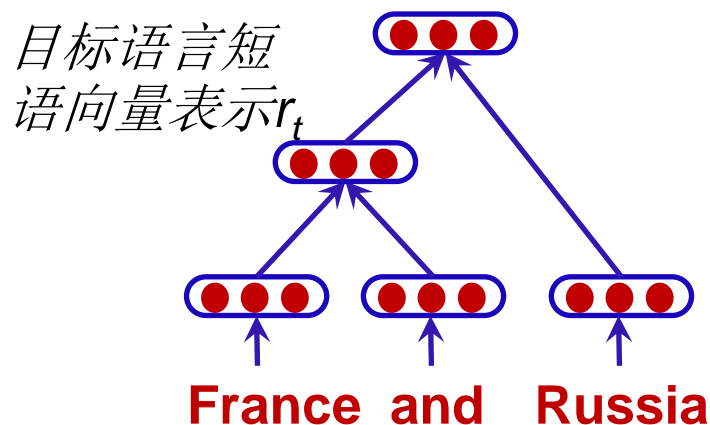
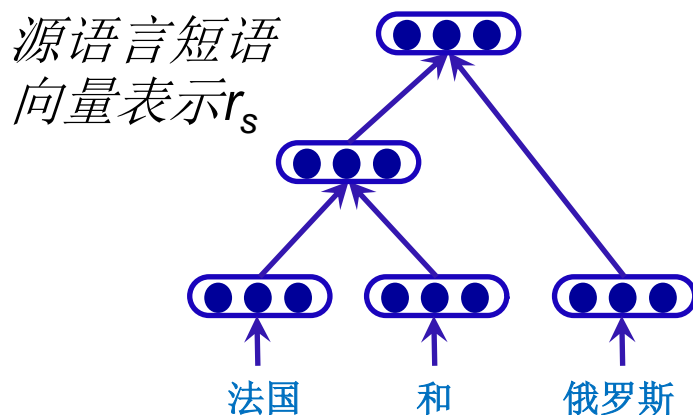
$$E_{rec}(S, T; \theta) = E_{rec}(S; \theta) + E_{rec}(T; \theta)$$

$$E_{regression}(S, T; \theta) = E_{regression}(S|T, \theta) + E_{regression}(T|S, \theta)$$

$$E_{regression}(S|T, \theta) = \sum_{s \in S} \frac{1}{2} \|f(v_{s_root}) - v_{t_root}\|^2$$

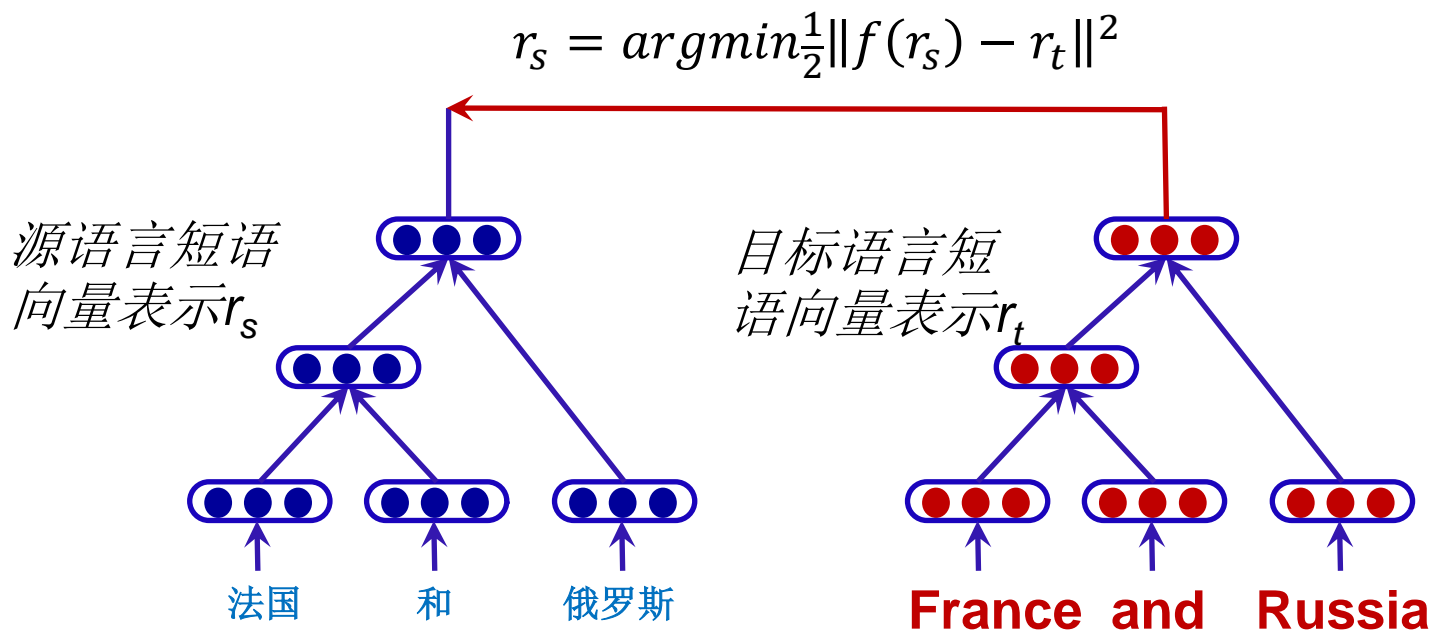
基于语义向量空间的短语表示

- 协同训练：第一步pre-training：利用无监督RAE模型分别学习源语言和目标语言短语的分布式表示



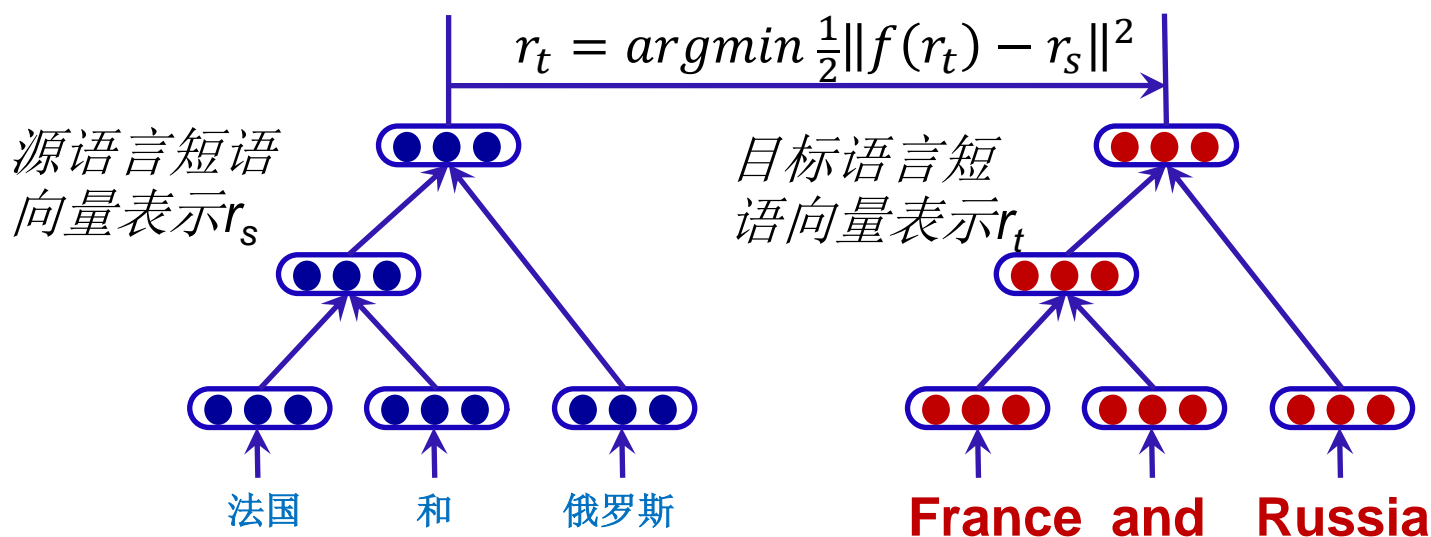
基于语义向量空间的短语表示

- 协同训练：第二步fine-tuning，1) 将目标短语表示视为源语言短语的正确语义表示，有监督地优化源语言短语的分布式表示



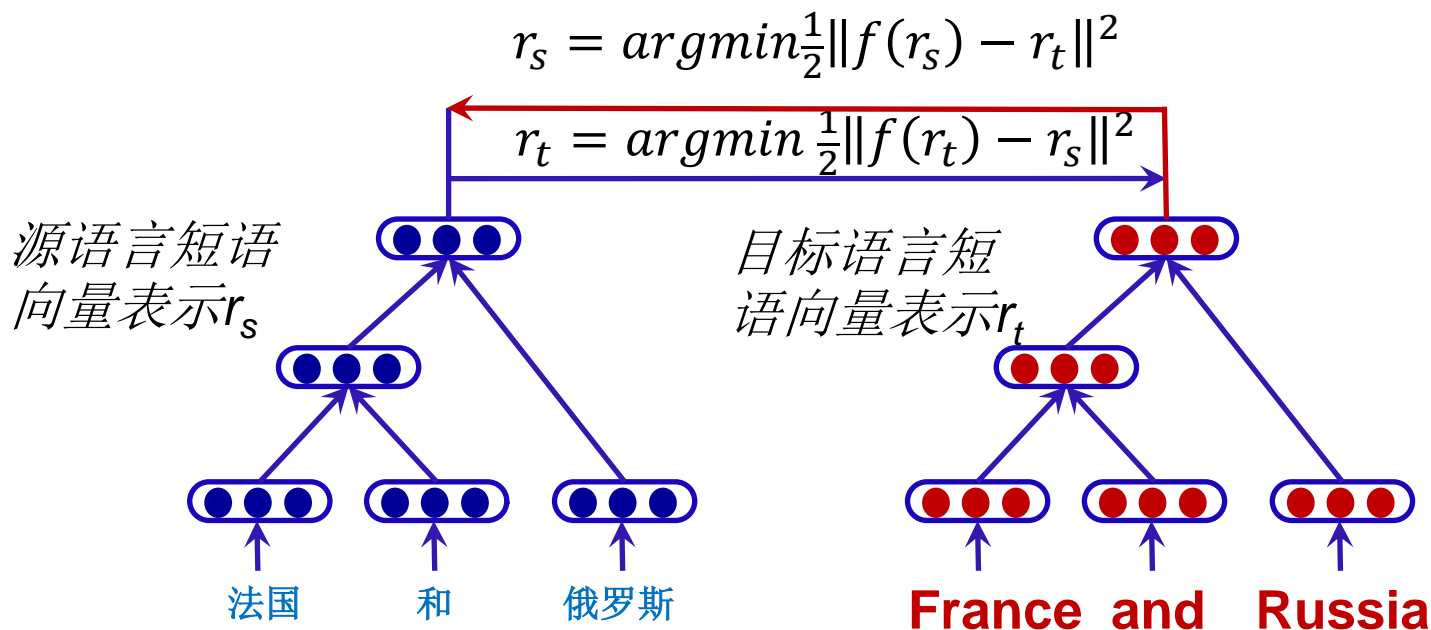
基于语义向量空间的短语表示

- 协同训练：第二步fine-tuning，2) 将源语言短语表示视为目标语言短语的正确语义表示，有监督地优化目标语言短语的分布式表示



基于语义向量空间的短语表示

- 协同训练：第三步收敛检测，若误差小于阈值或者迭代次数超过预设值，则训练结束



基于语义向量空间的短语表示

相似短语: do not agree
 will definitely reject
 will never accept

... ..

相似短语: abstract meaning
 real meaning
 intrinsic logic

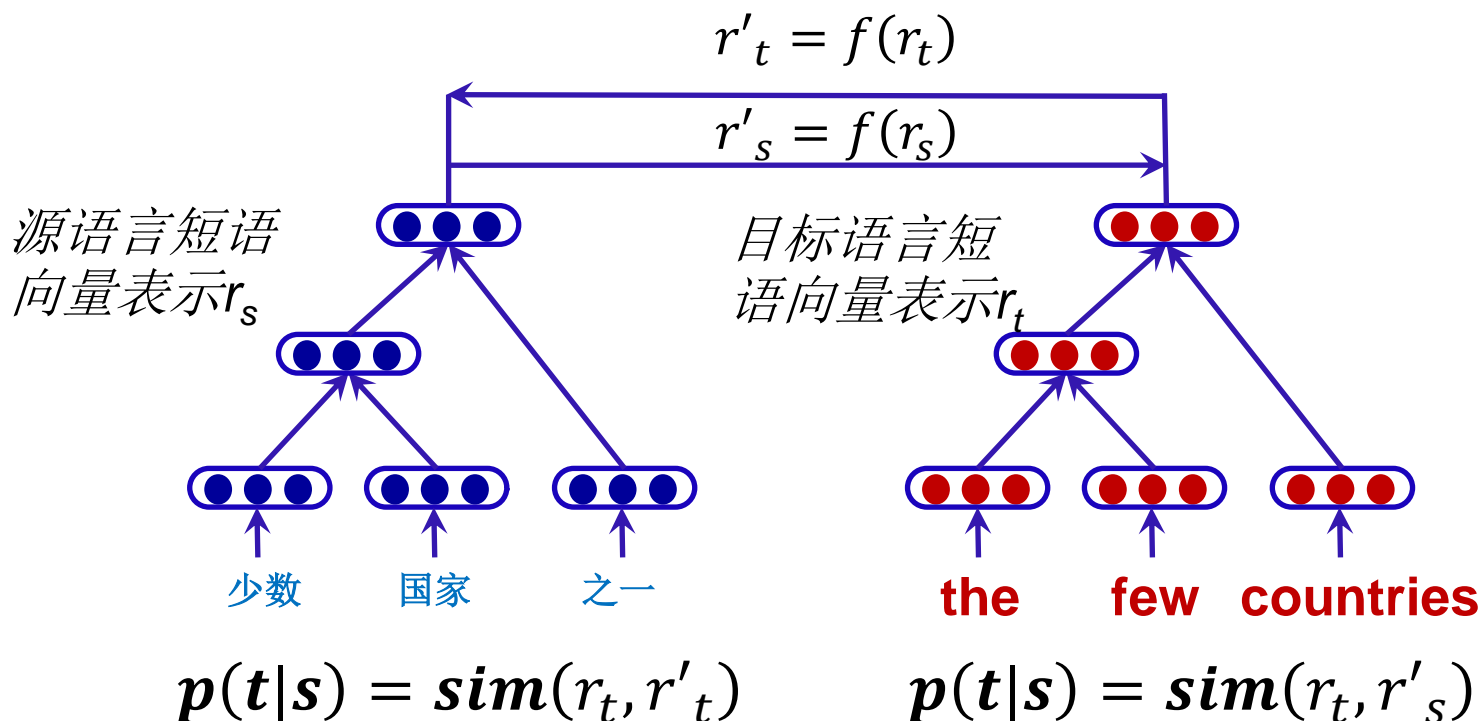
... ..

相似短语: what is your opinion
 what do you think about
 how do you view those

... ..

基于语义向量空间的翻译置信度估计

(少数 国家 之一, the few countries) $p(t|s), p(s|t)$



0.7 \leftarrow $\begin{bmatrix} 0.2 \\ 0.2 \\ 0.5 \\ 0.5 \\ 0.4 \end{bmatrix}$ $\begin{bmatrix} 0.5 \\ 0.5 \\ 0.4 \\ 0.2 \\ 0.5 \end{bmatrix}$

基于语义向量空间的翻译置信度估计

- 汉语-英语：210万训练数据

System	NIST03	NIST04	NIST05	NIST06	NIST08	ALL
MEBTG	35.81	36.91	34.69	33.83	27.17	34.82
+2feats 50-dim	36.43 (0.62↑)	37.64 (0.73 ↑)	35.35 (0.66↑)	35.53 (1.70 ↑)	28.59 (1.42 ↑)	35.84 (1.02↑)
+2feats 100-dim	36.45 (0.64 ↑)	37.44 (0.53↑)	35.58 (0.89↑)	35.42 (1.59↑)	28.57 (1.40↑)	36.03 (1.21 ↑)
+2feats 200-dim	36.34 (0.53↑)	37.35 (0.44↑)	35.78 (1.09 ↑)	34.87 (1.04↑)	27.84 (0.67↑)	35.62 (0.80↑)

[Zhang et al., ACL-2014]

基于短语的统计机器翻译

$$\begin{aligned}
 T' &= \operatorname{argmax}_T P(T|S) \\
 &= \operatorname{argmax}_{T, S_1^K} P(T, S_1^K | S) \\
 &= \operatorname{argmax}_{T, S_1^K, T_1^K, T_1^{K'}} \underbrace{P(S_1^K | S)}_{\text{短语翻译模型}} \underbrace{P(T_1^K | S_1^K, S)}_{\text{短语调序模型}} \underbrace{P(T_1^{K'} | T_1^K, S_1^K, S)}_{\text{目标语言模型}} \underbrace{P(T | T_1^{K'}, T_1^K, S_1^K, S)}_{\text{目标语言模型}}
 \end{aligned}$$

剩下的三个核心模型：

1. 短语翻译模型： $P(T_1^K | S_1^K, S)$
2. 短语调序模型： $P(T_1^{K'} | T_1^K, S_1^K, S)$
3. 目标语言模型： $P(T | T_1^{K'}, T_1^K, S_1^K, S)$

最大熵短语调序模型

(与 北韩, with North Korea) (有 邦交, have the diplomatic relations)



特征提取

f0=与, f1=北韩, f2=有, f3=邦交, f4=with, f5=korea, f6=have, f7=relations



最大熵模型

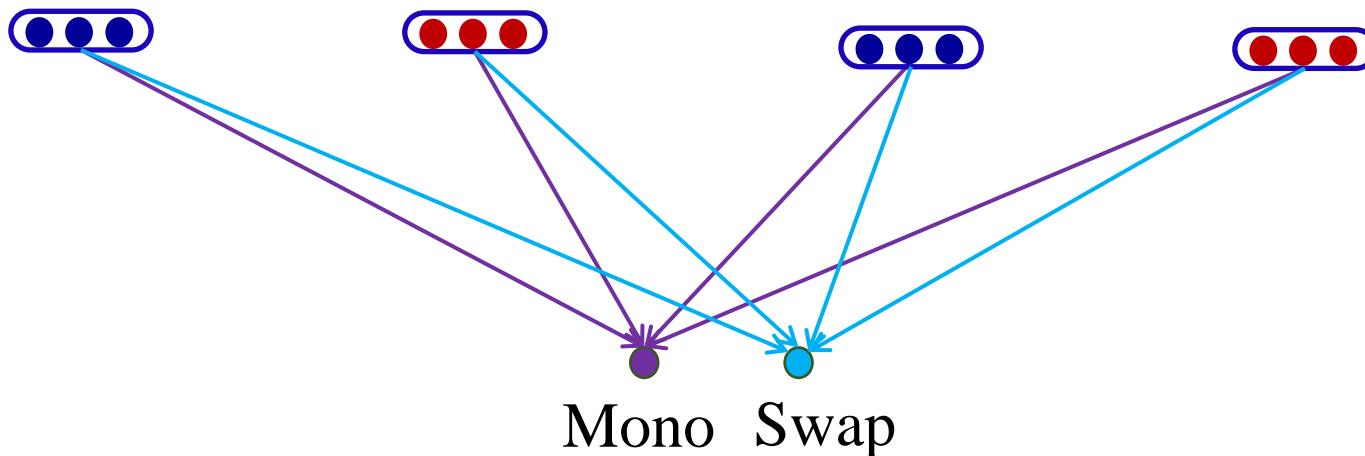
Swap

问题①：信息利用率低，仅利用短语的边界词

问题②：数据稀疏问题严重

基于语义向量空间的短语调序模型

(与 北韩, with North Korea) (有 邦交, have the diplomatic relations)



问题①：短语表示学习

问题②：训练目标函数

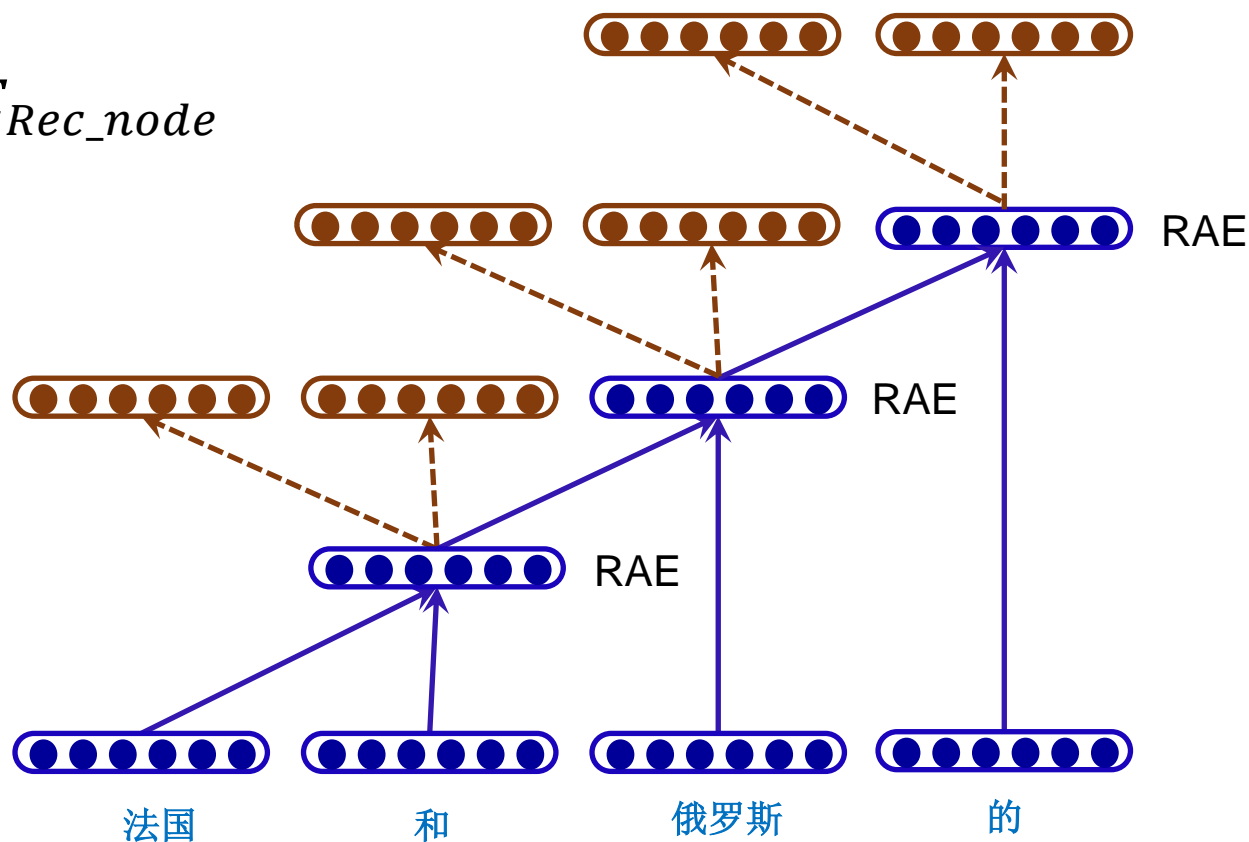
$$P(o|X_{c1}, X_{c2}, X_{e1}, X_{e2}) = \frac{\exp(f(o, X_{c1}, X_{c2}, X_{e1}, X_{e2}))}{\sum_{o'} \exp(f(o', X_{c1}, X_{c2}, X_{e1}, X_{e2}))}$$

$$f(o, X_{c1}, X_{c2}, X_{e1}, X_{e2}) = f(W^o[X_{c1}, X_{c2}, X_{e1}, X_{e2}] + b^o)$$

基于语义向量空间的短语表示

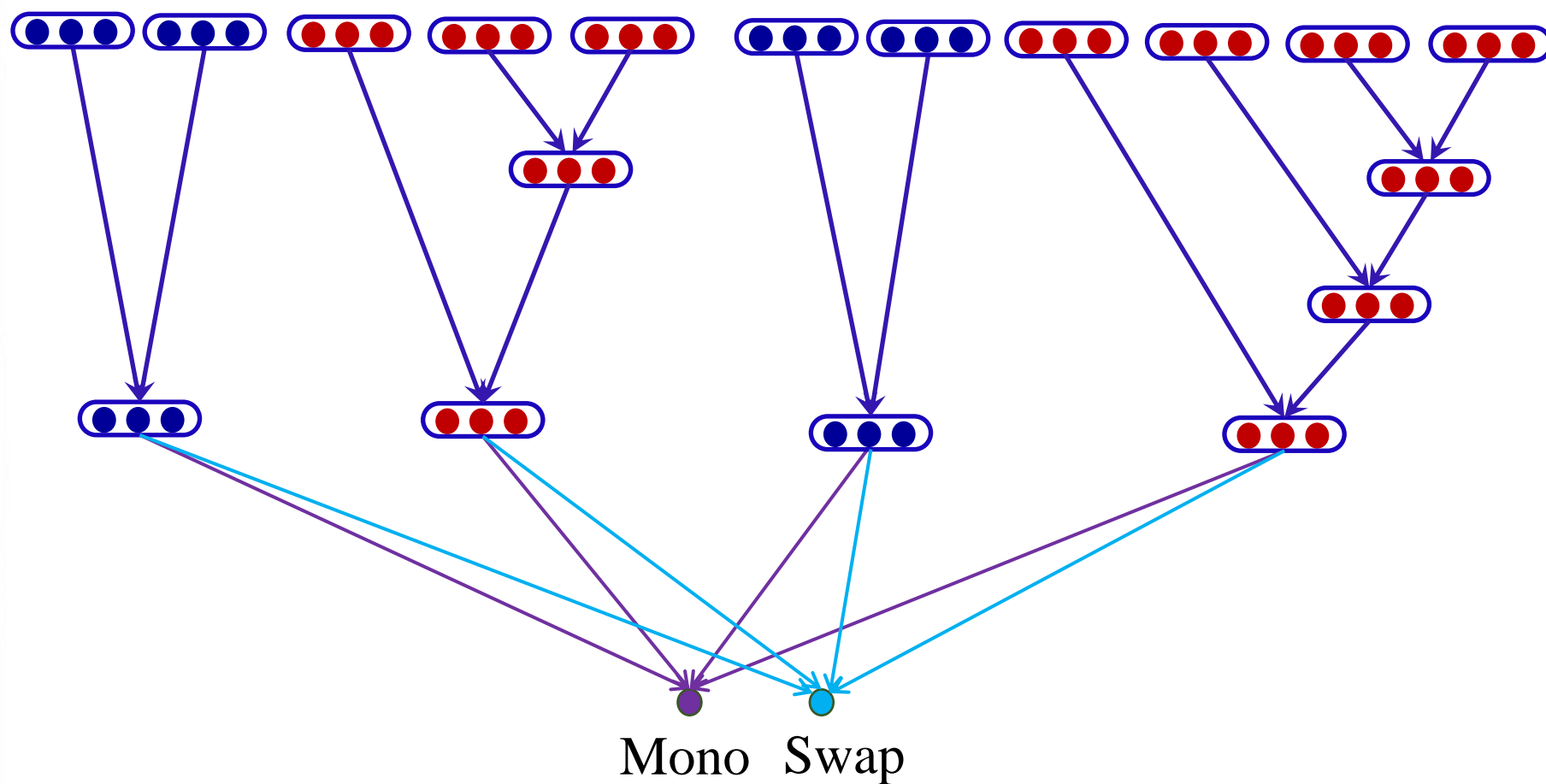
问题①：基于递归自动编码器的短语表示学习

$$E_{total} = \sum_{node} E_{Rec_node}$$



基于语义向量空间的短语调序模型

(与 北韩, with North Korea) (有 邦交, have the diplomatic relations)



基于语义向量空间的短语调序模型

问题②：训练目标函数 $S = \{s = (o, X_{c1}, X_{c2}, X_{e1}, X_{e2})\}$

$$J = \alpha E_{rec}(S; \theta) + (1 - \alpha) E_{reo}(S; \theta) + R(\theta)$$

↓
重构误差

↓
调序误差

↓
正则化项

$$E_{rec} = \sum_{node \in X_{c1}}^{X_{c1}} E_{rec_node}(c_1, c_2; \theta)$$

$$E_{rec_node}(c_1, c_2; \theta) = \frac{1}{2} \|[c_1, c_2] - [c'_1, c'_2]\|^2$$

基于语义向量空间的短语调序模型

问题②：训练目标函数 $S = \{s = (o, X_{c1}, X_{c2}, X_{e1}, X_{e2})\}$

$$J = \alpha E_{rec}(S; \theta) + (1 - \alpha) E_{reo}(S; \theta) + R(\theta)$$

↓
重构误差

↓
调序误差

↓
正则化项

$$E_{reo}(s; \theta) = - \sum_o d(o) \log(P_o(o|X_{c1}, X_{c2}, X_{e1}, X_{e2}))$$

$$d(o) = \{0, 1\}$$

基于语义向量空间的短语调序模型

cluster 1	cluster 2	cluster 3
1.18 accessibility wheelchair candies cough	works for verify on tunnels from transparency in opinion at	alternative duties one-day conference armed groups chinese language works eating habits

cluster 4	cluster 5
these people who the reasons why the story of how the system which the trend towards	of the three on the fundamental over the entire through its own with the best

短语表示的效果

基于语义向量空间的短语调序模型

- 汉语-英语：123万双语训练数据

Setting	NIST06	NIST08
maxent	30.40	23.75
neural	31.61	24.82

[Li et al., EMNLP-2013; COLING-2014]

联合翻译模型与语言模型

原文：

我	就	取	钱	给	了	她们
---	---	---	---	---	---	----

短语划分：

我 就	取 钱 给 了	她们
-----	---------	----

翻译：

I will	draw the money to	them
--------	--------------------------	------

短语调序：

I will	draw the money to	them
--------	--------------------------	------



$$P(the) \approx P(the | \text{取 钱 给 了}, draw\ the\ money\ to)$$

选择某个译文单词的概率仅取决于局部短语翻译规则，以及已生成的译文！

联合翻译模型与语言模型

S: 我 ³就 ⁴取 ⁵钱 ⁶给 ⁷了 她们
i will get money to perf. them

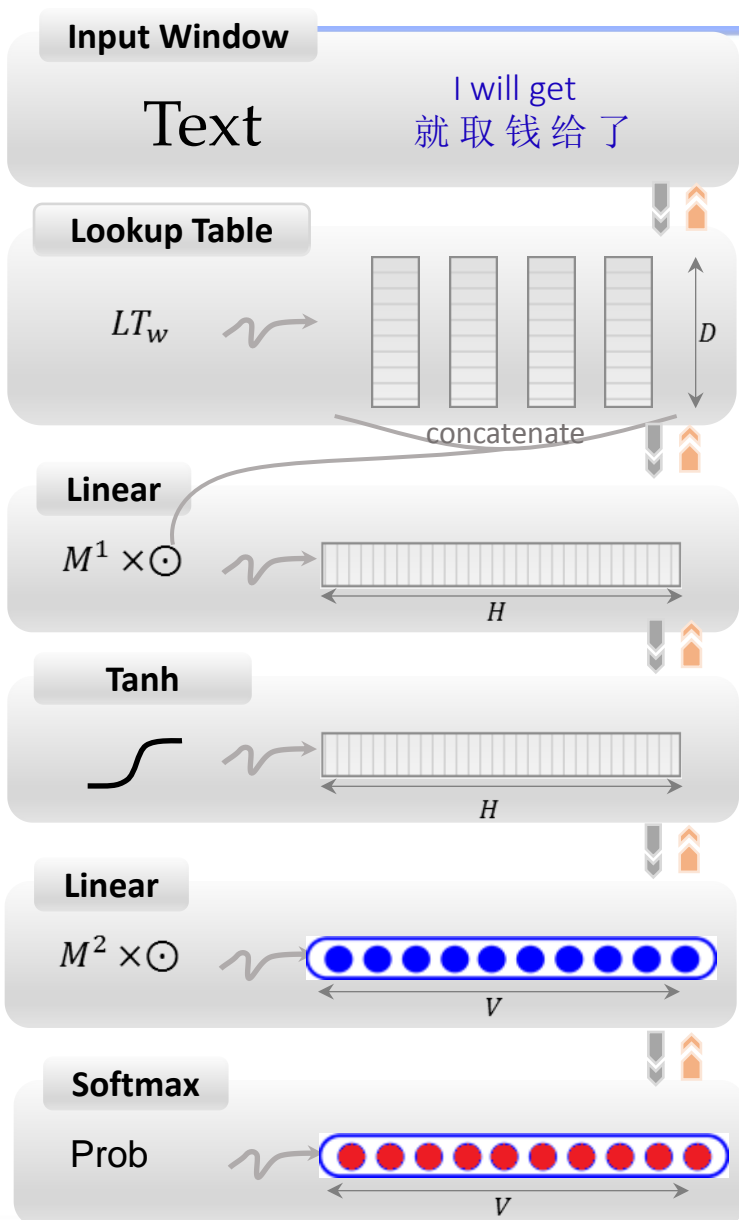
T: ²i ¹will ⁰get the money to them

$P(\text{the} \mid \text{get, will, i, 就, 取, 钱, 给, 了})$

$\begin{bmatrix} 0.1 \\ 0.2 \\ 0.1 \\ 0.5 \\ 0.4 \end{bmatrix}$	$\begin{bmatrix} 0.1 \\ 0.5 \\ 0.2 \\ 0.1 \\ 0.2 \end{bmatrix}$	$\begin{bmatrix} 0.2 \\ 0.1 \\ 0.3 \\ 0.1 \\ 0.4 \end{bmatrix}$	$\begin{bmatrix} 0.1 \\ 0.9 \\ 0.6 \\ 0.5 \\ 0.7 \end{bmatrix}$	$\begin{bmatrix} 0.1 \\ 0.2 \\ 0.4 \\ 0.2 \\ 0.1 \end{bmatrix}$	$\begin{bmatrix} 0.3 \\ 0.8 \\ 0.3 \\ 0.2 \\ 0.5 \end{bmatrix}$	$\begin{bmatrix} 0.4 \\ 0.5 \\ 0.2 \\ 0.3 \\ 0.1 \end{bmatrix}$	$\begin{bmatrix} 0.2 \\ 0.1 \\ 0.3 \\ 0.4 \\ 0.6 \end{bmatrix}$	$\begin{bmatrix} 0.4 \\ 0.1 \\ 0.2 \\ 0.5 \\ 0.7 \end{bmatrix}$
---	---	---	---	---	---	---	---	---

$$\begin{aligned}
 P(e_i) &\approx P(e_i \mid e_1 \cdots e_{i-1}, f) \\
 &\approx P(e_i \mid e_{i-3} \cdots e_{i-1}, f_{j-c} \cdots f_j \cdots f_{j+c})
 \end{aligned}$$

联合翻译模型与语言模型



- 上下文
 - 目标语言 4-gram
 - 源语言中心词左右5个词
- 词向量 (192 维)
- 两个隐藏层 (512维)
- 输出层 softmax

$$P(e_i | e_{i-3} \cdots e_{i-1}, f_{j-c} \cdots f_j \cdots f_{j+c})$$

联合翻译模型与语言模型

	Ar-En	Ch-En
	BLEU	BLEU
OpenMT12 - 1st Place	49.5	32.6
OpenMT12 - 2nd Place	47.5	32.2
OpenMT12 - 3rd Place	47.4	30.8
...
OpenMT12 - 9th Place	44.0	27.0
OpenMT12 - 10th Place	41.2	25.7
Baseline (w/o RNNLM)	48.9	33.0
Baseline (w/ RNNLM)	49.8	33.4
+ S2T/L2R NNJM (Dec)	51.2	34.2
+ S2T NNLTm (Dec)	52.0	34.2
+ T2S NNLTm (Resc)	51.9	34.2
+ S2T/R2L NNJM (Resc)	52.2	34.3
+ T2S/L2R NNJM (Resc)	52.3	34.5
+ T2S/R2L NNJM (Resc)	52.8	34.7
“Simple Hier.” Baseline	43.4	30.1
+ S2T/L2R NNJM (Dec)	47.2	31.5
+ S2T NNLTm (Dec)	48.5	31.8
+ Other NNJMs (Resc)	49.7	32.2

- 汉语-英语：
NIST-2012 受限评测

[Devlin et al., ACL-2014]

句子表示学习在统计机器翻译中的应用

这样 不仅 防止 不 安全 的 分子 ， 也 是 防止 那些 非法 移民 进入 日本



this not only prevents dangerous **people**

三个问题

Q1: 我们为什么要学习整个句子的语义表示?

这样 不仅 防止 不 安全 的 分子 , 也 是 防止 那些 非法 移民 进入 日本

**Q2: 如何学习句子的
语义表示?**



this not only prevents dangerous people

**Q3: 如何将句子的
语义表示融入统计
翻译模型?**

三个问题

Q1: 我们为什么要学习整个句子的语义表示?

这样 不仅 防止 不 安全 的 分子 , 也 是 防止 那些 非法 移民 进入 日本

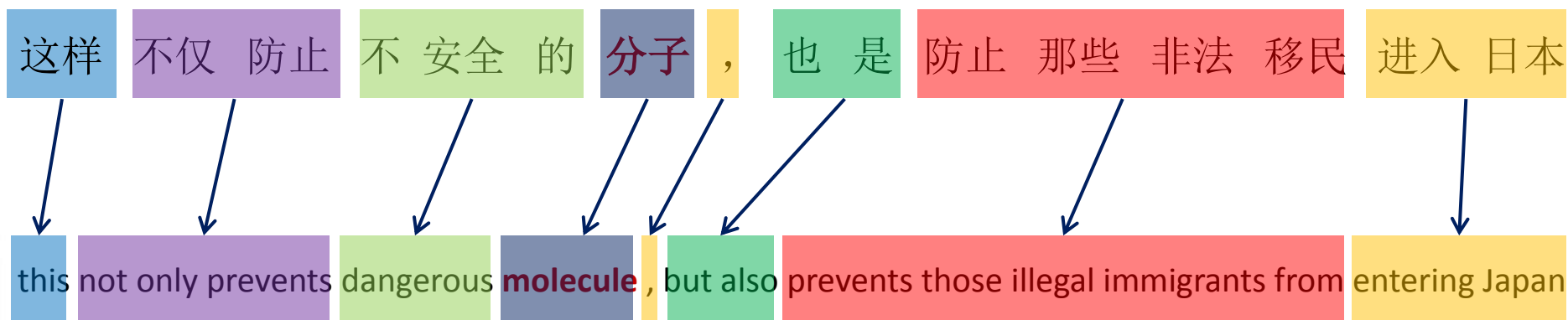
Q2: 如何学习句子的语义表示?



this not only prevents dangerous **people**

Q3: 如何将句子的语义表示融入统计翻译模型?

句子表示的必要性



句子表示的必要性

这样 不仅 防止 不 安全 的 分子 , 也 是 防止 那些 非法 移民 进入 日本



this not only prevents dangerous **molecule** , but also prevents those illegal immigrants from entering Japan

三个问题

Q1: 我们为什么要学习整个句子的语义表示?

这样 不仅 防止 不 安全 的 分子 , 也 是 防止 那些 非法 移民 进入 日本

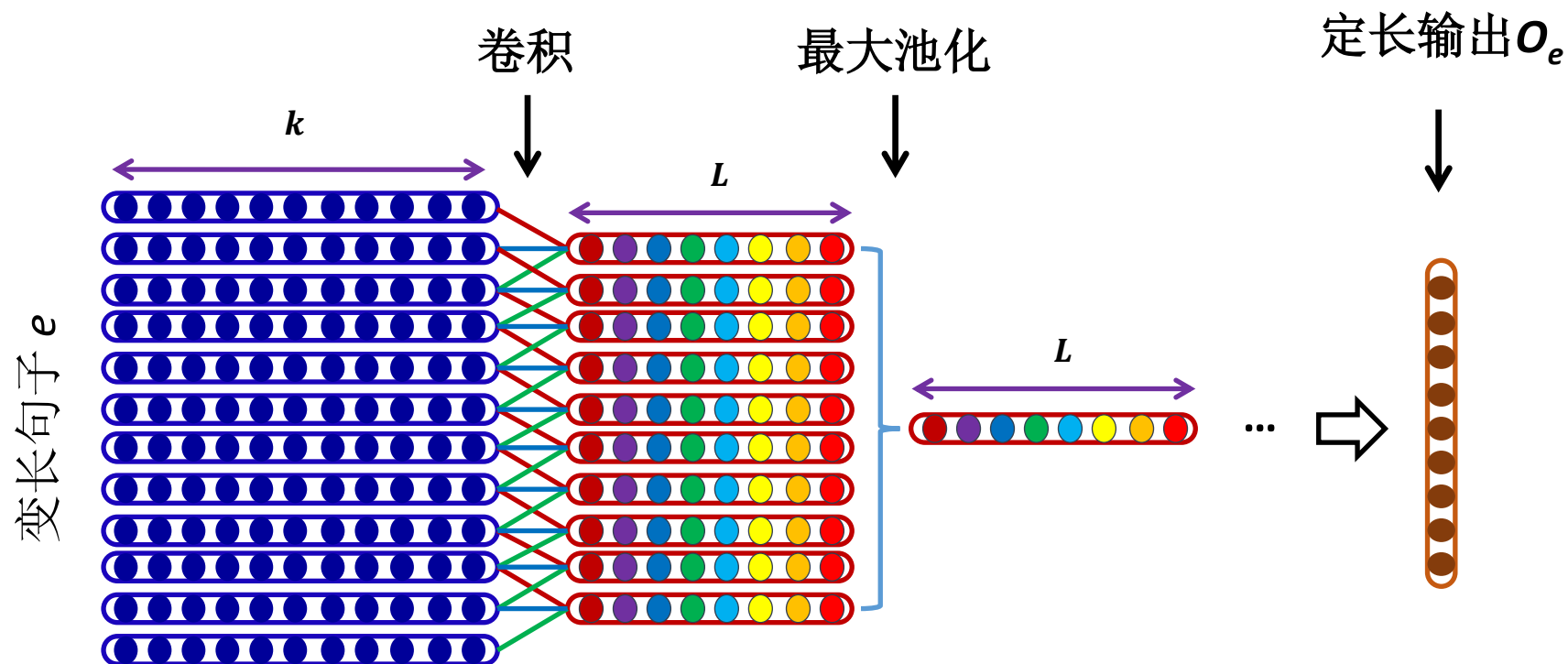
Q2: 如何学习句子的语义表示?



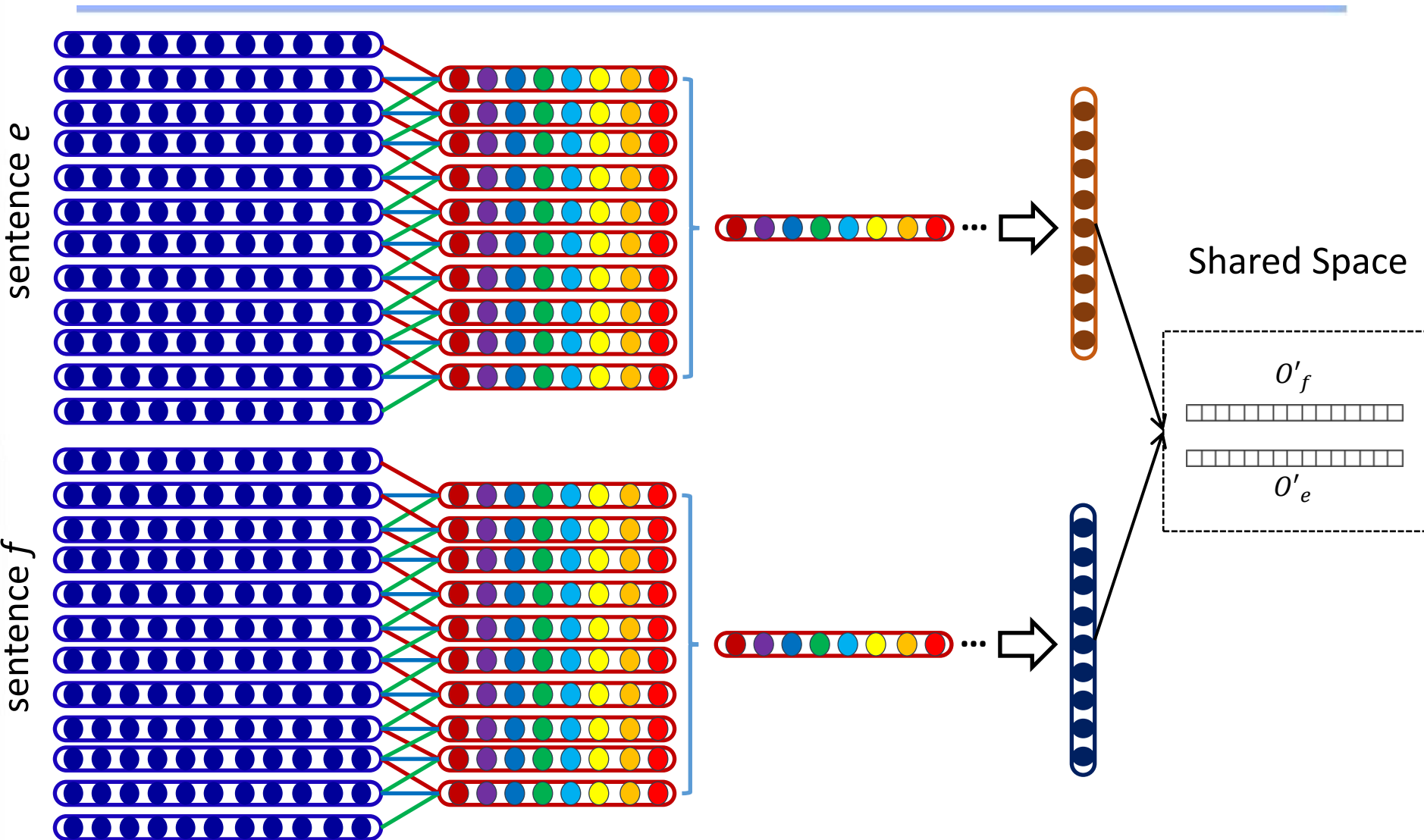
this not only prevents dangerous people

Q3: 如何将句子的语义表示融入统计翻译模型?

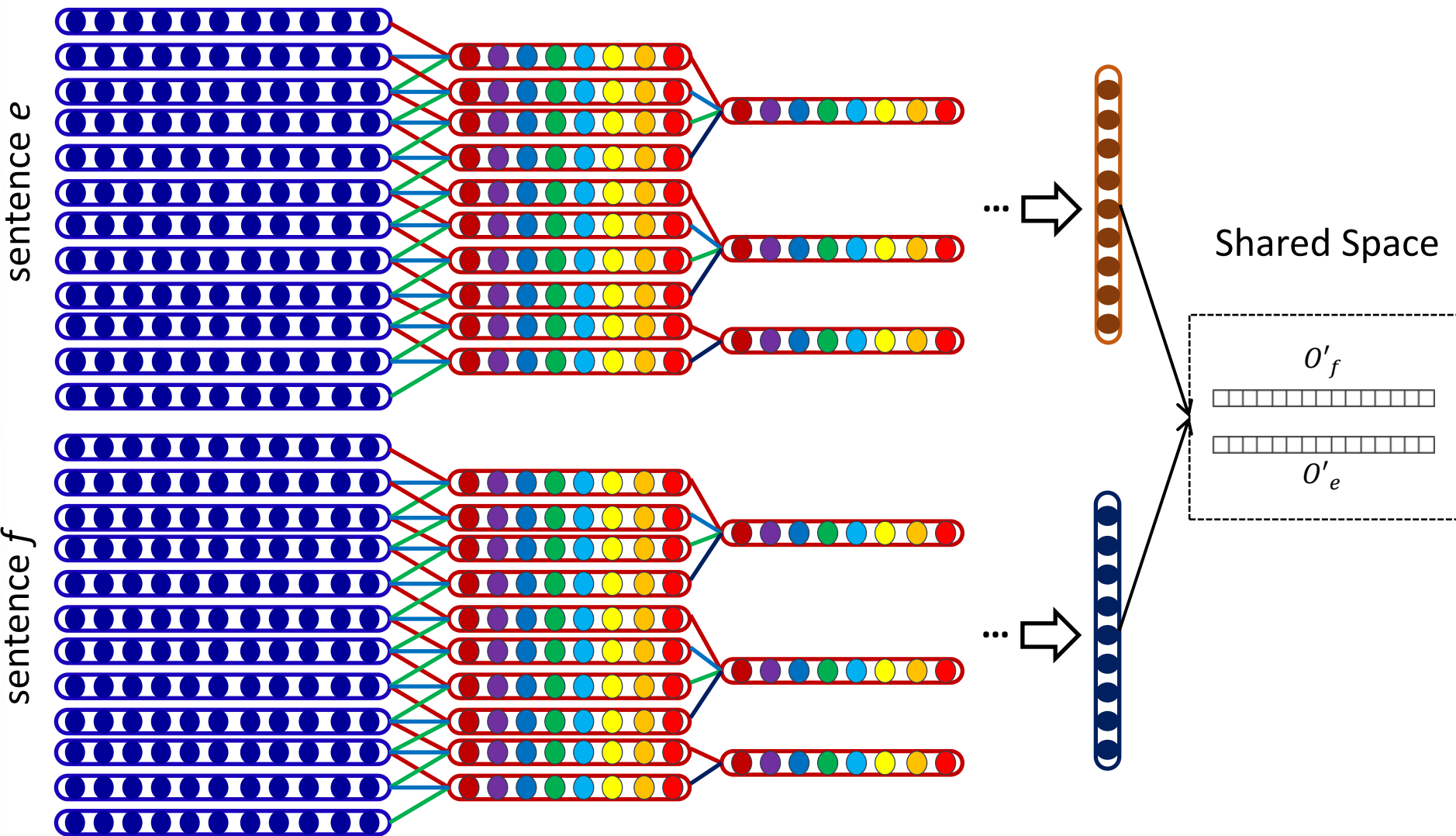
句子表示学习-卷积神经网络



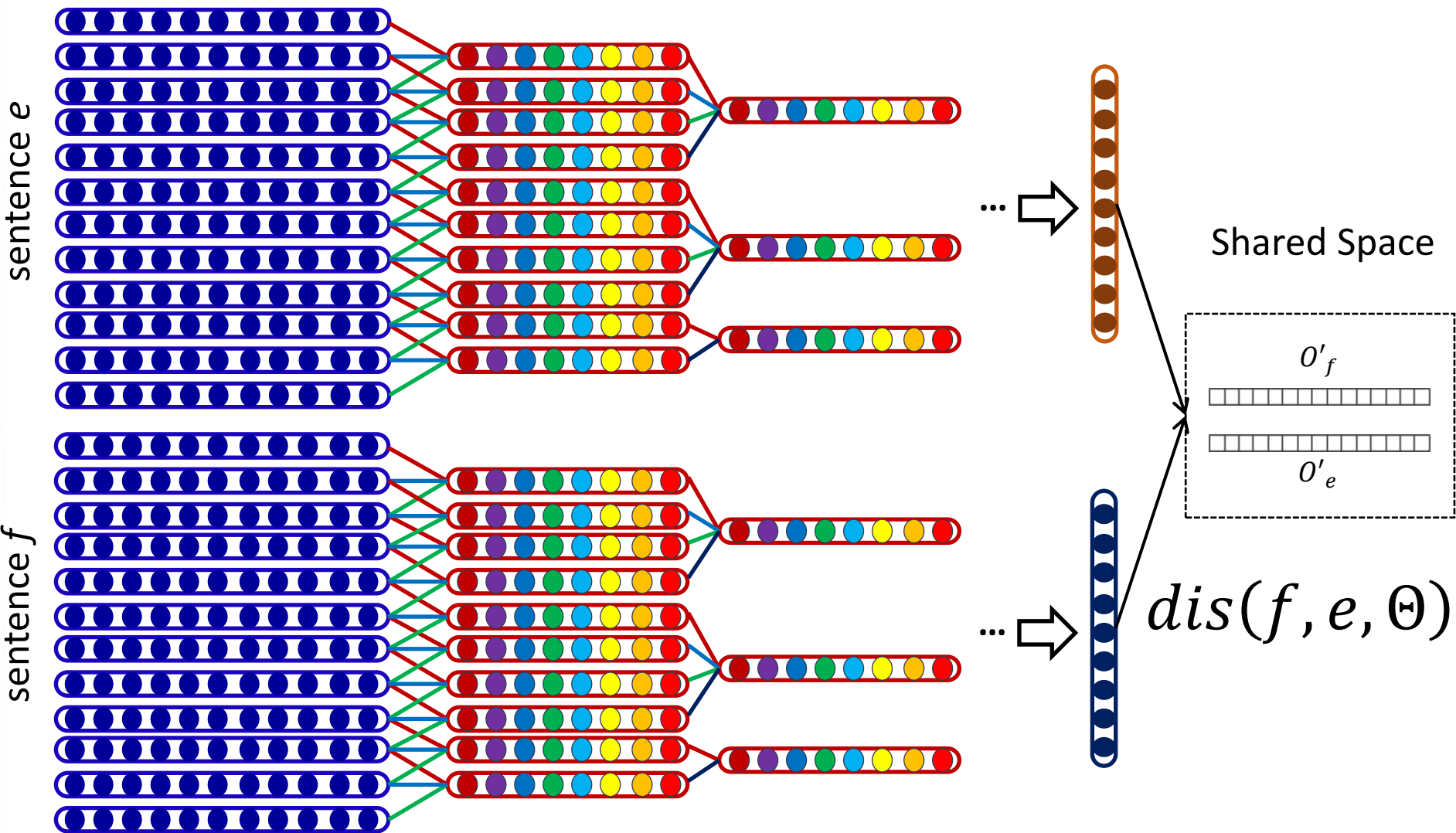
双语约束的卷积神经网络



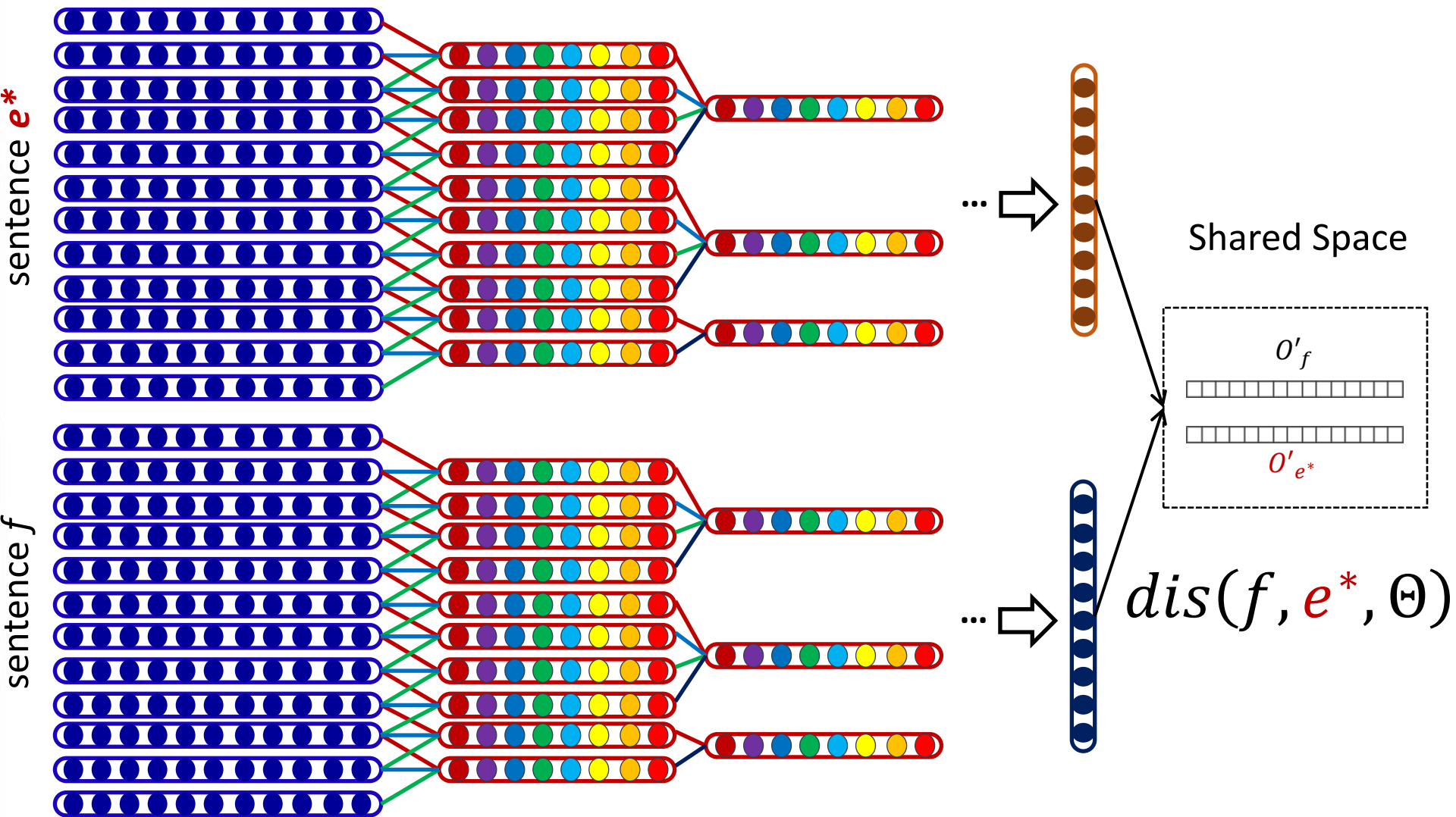
基于语块的卷积神经网络



最大间隔训练



最大间隔训练



最大间隔训练

$$j(f, e, e^*, \Theta) \\ = \max(0, 1 + \text{dis}(f, e, \Theta) - \text{dis}(f, e^*, \Theta))$$

三个问题

Q1: 我们为什么要学习整个句子的语义表示?

这样 不仅 防止 不 安全 的 分子 , 也 是 防止 那些 非法 移民 进入 日本

**Q2: 如何学习句子的
语义表示?**



this not only prevents dangerous people

**Q3: 如何将句子的
语义表示融入统计
翻译模型?**

基线系统

S: 我 ³就 ⁴取 ⁵钱 ⁶给 ⁷了 她们
i will get money to perf. them

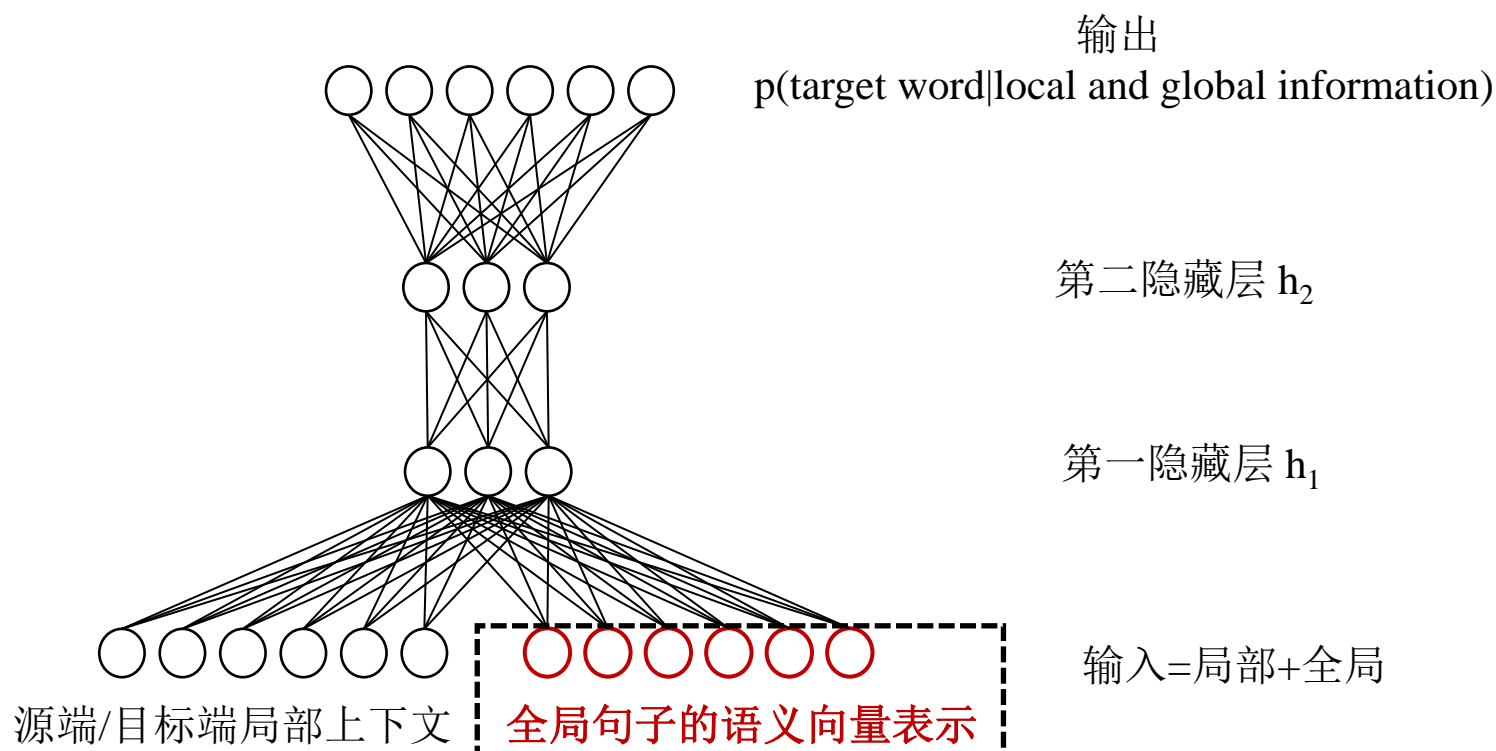
T: ²i ¹will ⁰get the money to them
P(the | get, will, i, 就, 取, 钱, 给, 了)

$$\begin{aligned}
 P(e_i) &\approx P(e_i | e_1 \cdots e_{i-1}, f) \\
 &\approx P(e_i | e_{i-3} \cdots e_{i-1}, f_{j-c} \cdots f_j \cdots f_{j+c})
 \end{aligned}$$

卷积神经网络-统计机器翻译

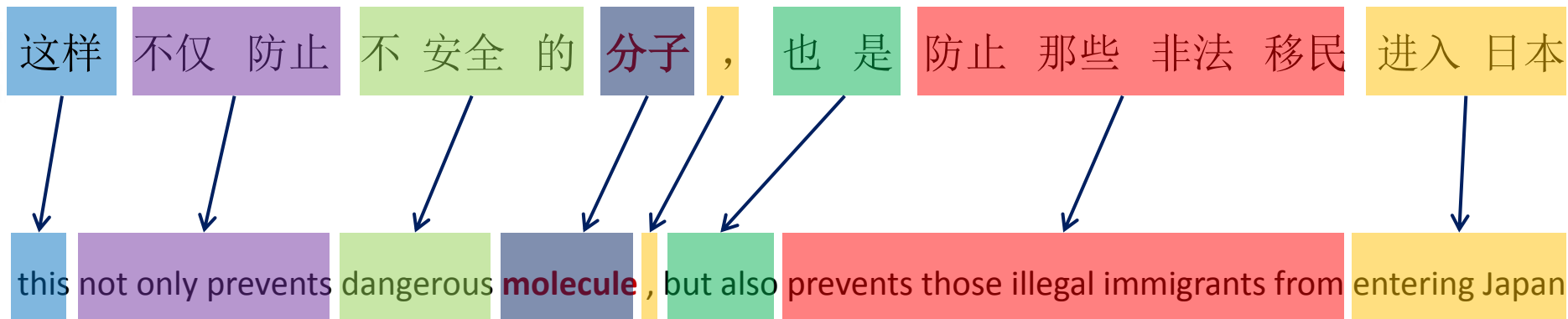
$$P(e_i) \approx P(e_i | e_{i-3} \cdots e_{i-1}, f_{j-c} \cdots f_j \cdots f_{j+c})$$

$$\approx P(e_i | e_{i-3} \cdots e_{i-1}, f_{j-c} \cdots f_j \cdots f_{j+c}, \mathbf{f})$$



融入统计机器翻译

系统: 层次短语翻译系统

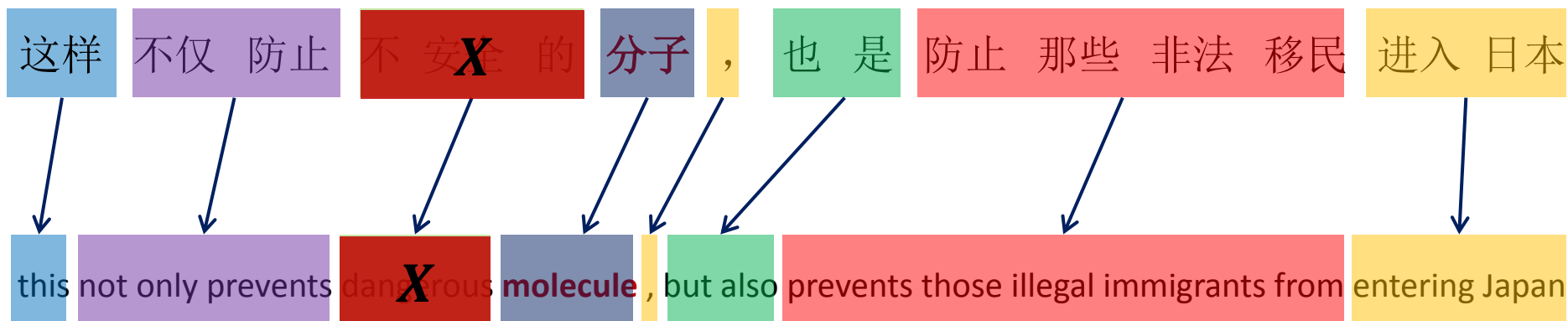


$X = \langle \text{不 安全 的}, \text{dangerous} \rangle$

$X = \langle X_0 \text{ 分子}, X_0 \text{ molecule} \rangle$

融入统计机器翻译

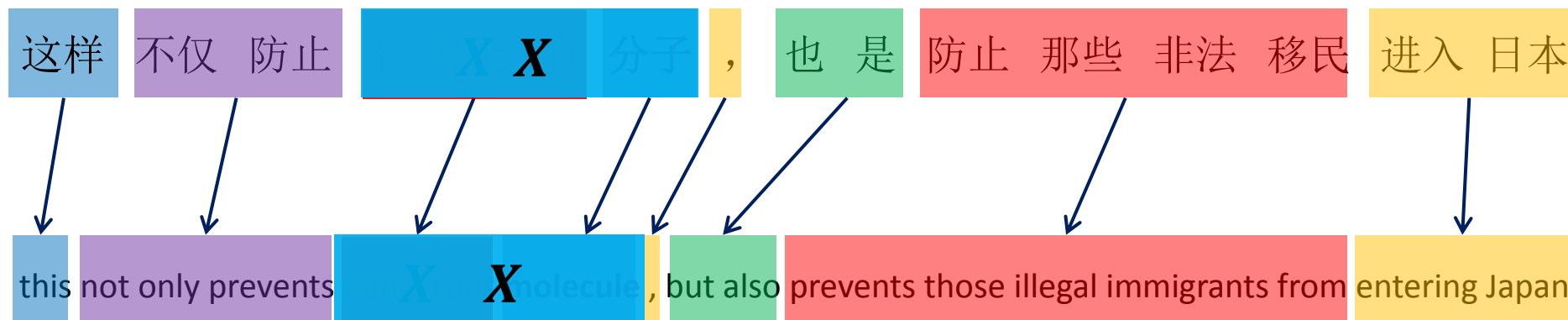
系统: 层次短语翻译系统



$X = \langle \text{不安全的}, \text{dangerous} \rangle$

融入统计机器翻译

系统: 层次短语翻译系统



$X = \langle \text{不 安全 的, dangerous} \rangle$
 $X = \langle X_0 \text{ 分子, } X_0 \text{ molecule} \rangle$

$X = \langle \text{不 安全 的 分子, dangerous molecule} \rangle$

Log-linear Model:

$$P(e_i | e_{i-3} \cdots e_{i-1}, f_{j-c} \cdots f_j \cdots f_{j+c})$$

正向反向概率

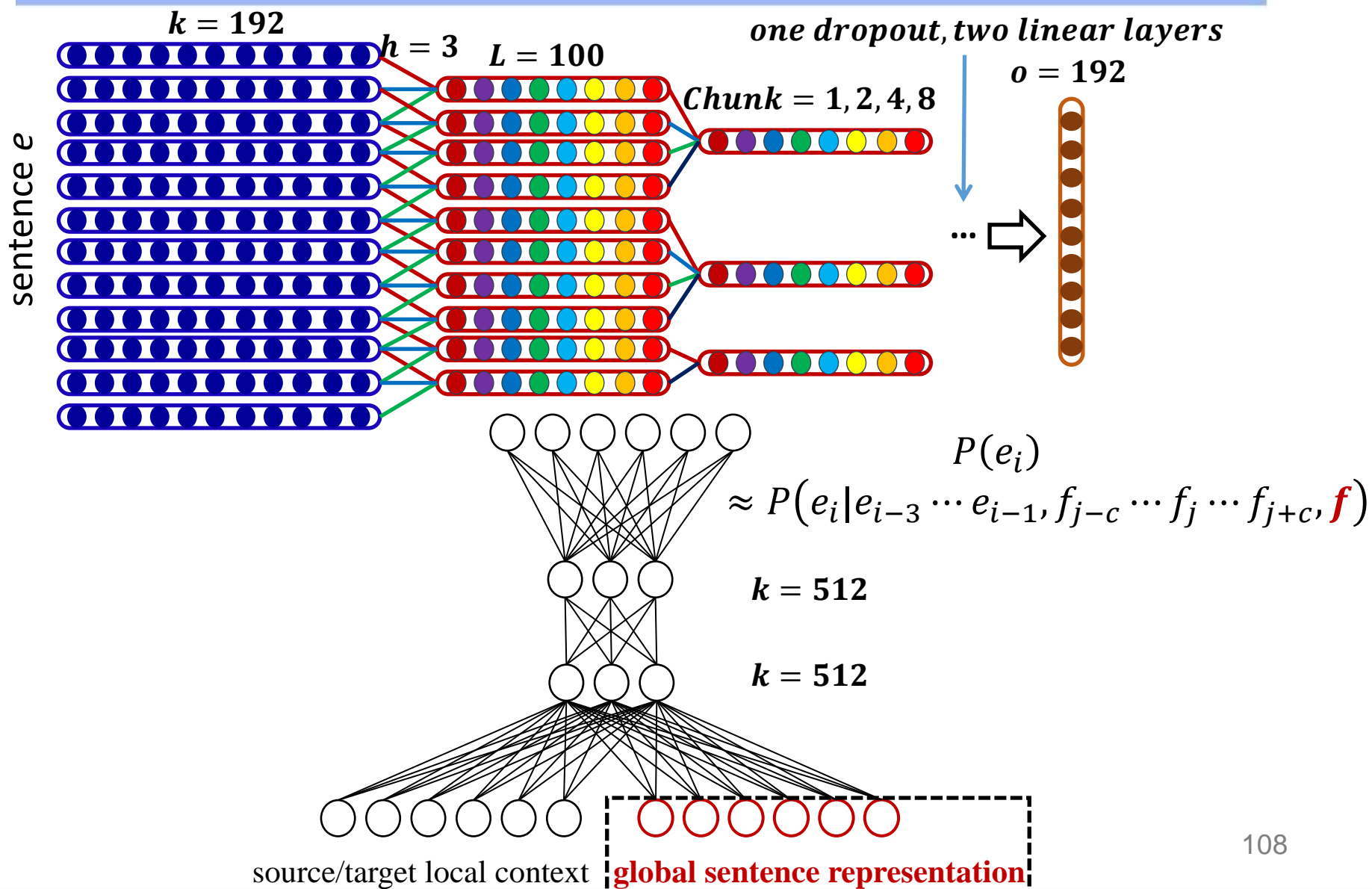
目前语言模型特征

神经网络翻译概率

正向反向词汇化概率

规则数目和译文长度特征

网络参数



句子表示在统计机器翻译中的效果

- 汉语-英语：210万双语训练数据

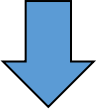
System	MT03	MT05	MT06	MT08
HPB	35.98	34.66	35.25	27.80
+NNJM	36.93	35.55 ⁺	35.77	28.64 ⁺
+AVE_SENT	37.16	35.88 ⁺	36.07 ⁺	29.19 ⁺
+BCCNN-1	37.32	36.06 ⁺	36.42 ⁺	29.35 ⁺
+BCCNN-2	37.75	36.24 ⁺	36.65 ⁺ *	29.97 ⁺ *
+BCCNN-4	37.98	36.22 ⁺	36.78⁺*	30.02⁺*
+BCCNN-8	37.64	36.29⁺*	36.49 ⁺	29.98 ⁺ *

0.95↑

1.38↑

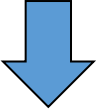
[Zhang et al., IJCAI-2015]

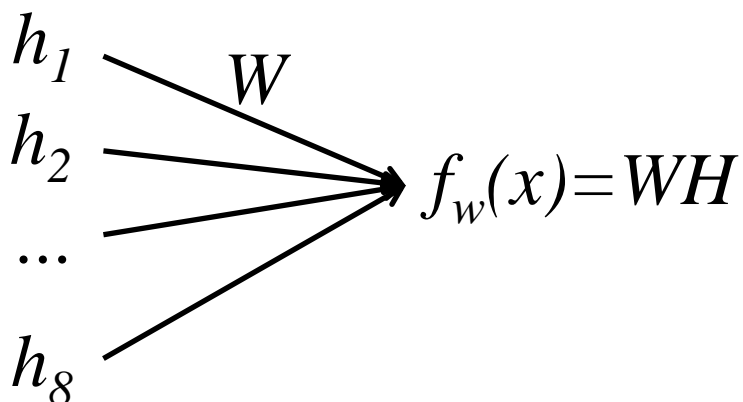
统计机器翻译

$$T' = \operatorname{argmax}_T P(T|S) = \operatorname{argmax}_T \left\{ \sum_1^M \lambda_m h_m(T, S) \right\}$$


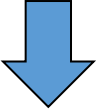
$$\begin{aligned} h_1(T, S) &= \log p(t|s) & h_5(T, S) &= \log P(T_1^{K'} | T_1^K, S_1^K, S) \\ h_2(T, S) &= \log p(s|t) & h_6(T, S) &= \log P(T | T_1^{K'}, T_1^K, S_1^K, S) \\ h_3(T, S) &= \log p_{lex}(t|s) & h_7(T, S) &= \log \text{len}(T) \\ h_4(T, S) &= \log p_{lex}(s|t) & h_8(T, S) &= \log \text{count}(\text{phrases}) = \log K \end{aligned}$$

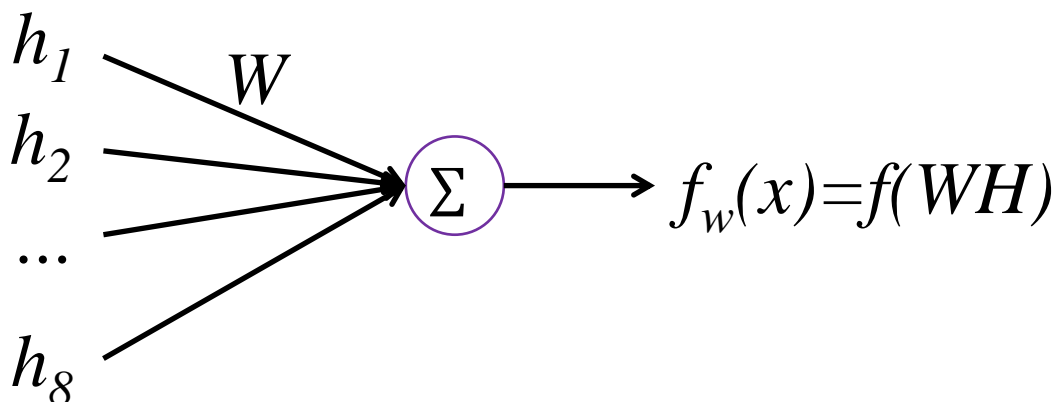
统计机器翻译

$$T' = \operatorname{argmax}_T P(T|S) = \operatorname{argmax}_T \left\{ \sum_1^M \lambda_m h_m(T, S) \right\}$$




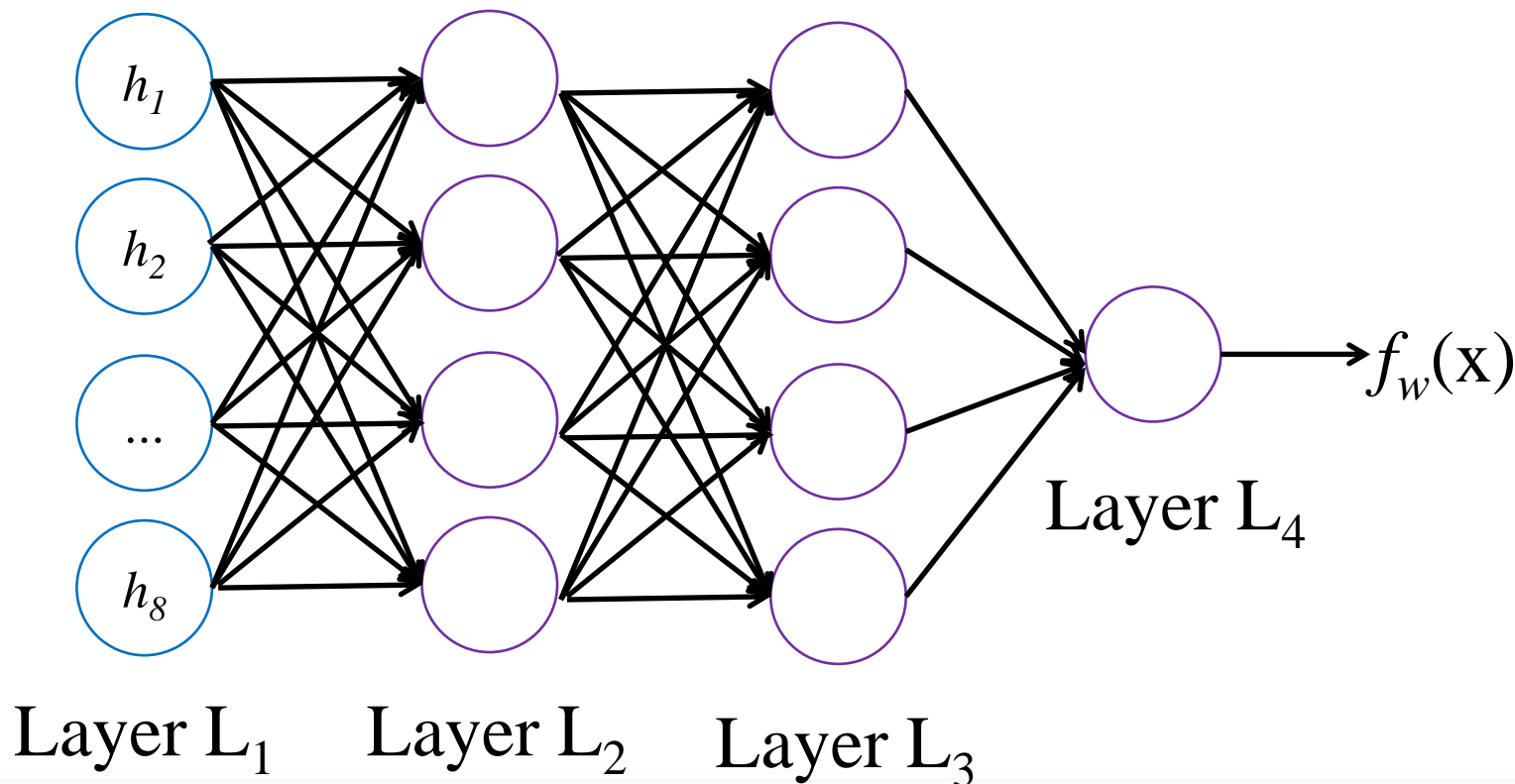
非线性统计机器翻译

$$T' = \operatorname{argmax}_T P(T|S) = \operatorname{argmax}_T \left\{ \sum_1^M \lambda_m h_m(T, S) \right\}$$




非线性统计机器翻译

$$T' = \operatorname{argmax}_T P(T|S) = \operatorname{argmax}_T \left\{ \sum_1^M \lambda_m h_m(T, S) \right\}$$



非线性统计机器翻译

- 汉语-英语：820万双语训练数据

Criteria	MT03(train)	MT02(dev)	MT04	MT05
BR _c	35.02	36.63	34.96	34.15
BR	38.66	40.04	38.73	37.50
BW	39.55	39.36	38.72	37.81
PW	38.61	38.85	38.73	37.98

[Huang et al., ACL-2015]

参考文献

1. Yoshua Bengio, Rejean Ducharme, Pascal Vincent and Christian Jauvin. [A neural probabilistic language model](#). *Journal of Machine Learning Research*, 3:1137–1155..
2. David Chiang. [Hierarchical phrase-based translation](#). *Computational Linguistics*, 2007.
3. Jacob Devlin, Rabih Zbib, Zhongqiang Huang, Thomas Lamar, Richard Schwartz, and John Makhoul. [Fast and robust neural network joint models for statistical machine translation](#). *In Proc. of ACL 2014*.
4. Jianfeng Gao, Xiaodong He, Yih Wen-tao and Li Deng. [Learning continuous phrase representations for translation modeling](#). *In Proc. of ACL 2014*.
5. Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. [Moses: Open source toolkit for statistical machine translation](#). *In Proc. of ACL 2007*.
6. Peng Li, Yang Liu, and Maosong Sun. [Recursive autoencoders for itg-based translation](#). *In Proc. of EMNLP 2013*.
7. Peng Li, Yang Liu, Maosong Sun, Tatsuya Izuha and Dakun Zhang. [A neural reordering model for phrase-based translation](#). *In Proc. of COLING 2014*.
8. Fandong Meng, Zhengdong Lu, Mingxuan Wang, Hang Li, Wenbing Jiang and Qun Liu. [Encoding source language with convolutional neural network for machine translation](#). *In Proc. of ACL 2015*.
9. Franz Josef Och and Hermann Ney. [Discriminative training and maximum entropy models for statistical machine translation](#). *In Proc. of ACL 2002*.
10. Shujian Huang, Huadong Chen, Xinyu Dai and Jiajun Chen. [Non-linear Learning for Statistical Machine Translation](#). *In Proc. of ACL 2015*.

参考文献

11. Jinsong Su, Deyi Xiong, Biao Zhang, Yang Liu, Junfeng Yao and Min Zhang. [Bilingual Correspondence Recursive Autoencoder for Statistical Machine Translation](#). *In Proc. of EMNLP 2015*.
12. Akihiro Tamura, Taro Watanabe and Eiichiro Sumita. [Recurrent neural networks for word alignment model](#). *In Proc. of ACL 2014*.
13. Ashish Vaswani, Yingdong Zhao, Victoria Fossum and David Chiang. [Decoding with large-scale neural language models improves translation](#). *In Proc. of EMNLP 2013*.
14. Deyi Xiong, Qun Liu and Shouxun Lin. [Maximum entropy based phrase reordering model for statistical machine translation](#). *In Proc. of ACL 2006*.
15. Nan Yang, Shujie Liu, Mu Li, Ming Zhou and Nenghai Yu. [Word alignment modeling with context dependent deep neural network](#). *In Proc. of ACL 2013*.
16. Jiajun Zhang, Shujie Liu, Mu Li, Ming Zhou and Chengqing Zong. [Bilingually-constrained phrase embeddings for machine translation](#). *In Proc. of ACL 2014*.
17. Jiajun Zhang, Dakun Zhang and Jie Hao. [Local Translation Prediction with Global Sentence Representation](#). *In Proc. of IJCAI 2015*.
18. Jiajun Zhang and Chengqing Zong. [Deep Neural Networks in Machine Translation: an Overview](#). *IEEE Intelligent Systems 2015*.
19. Will Y Zou, Richard Socher, Daniel Cer and Christopher D Manning. [Bilingual word embeddings for phrase-based machine translation](#). *In Proc. of EMNLP 2013*.

开源工具

- 1, Google Word2Vec, <http://code.google.com/p/word2vec/>
- 2, NPLM, 前馈神经网络, 便于在统计机器翻译中应用前馈神经网络语言模型, <http://nlg.isi.edu/software/nplm/>
- 3, RWTHLM, 前馈和循环神经网络, 便于在统计机器翻译中应用循环神经网络语言模型, <http://www-i6.informatik.rwth-aachen.de/web/Software/index.html>

... ..

N L P R



谢谢!
Q&A