

# CIPS暑期学校

## 深度学习与机器翻译 (1/2)

张家俊

中国科学院自动化研究所

[www.nlpr.ia.ac.cn/cip/jjzhang.htm](http://www.nlpr.ia.ac.cn/cip/jjzhang.htm)

[jjzhang@nlpr.ia.ac.cn](mailto:jjzhang@nlpr.ia.ac.cn)

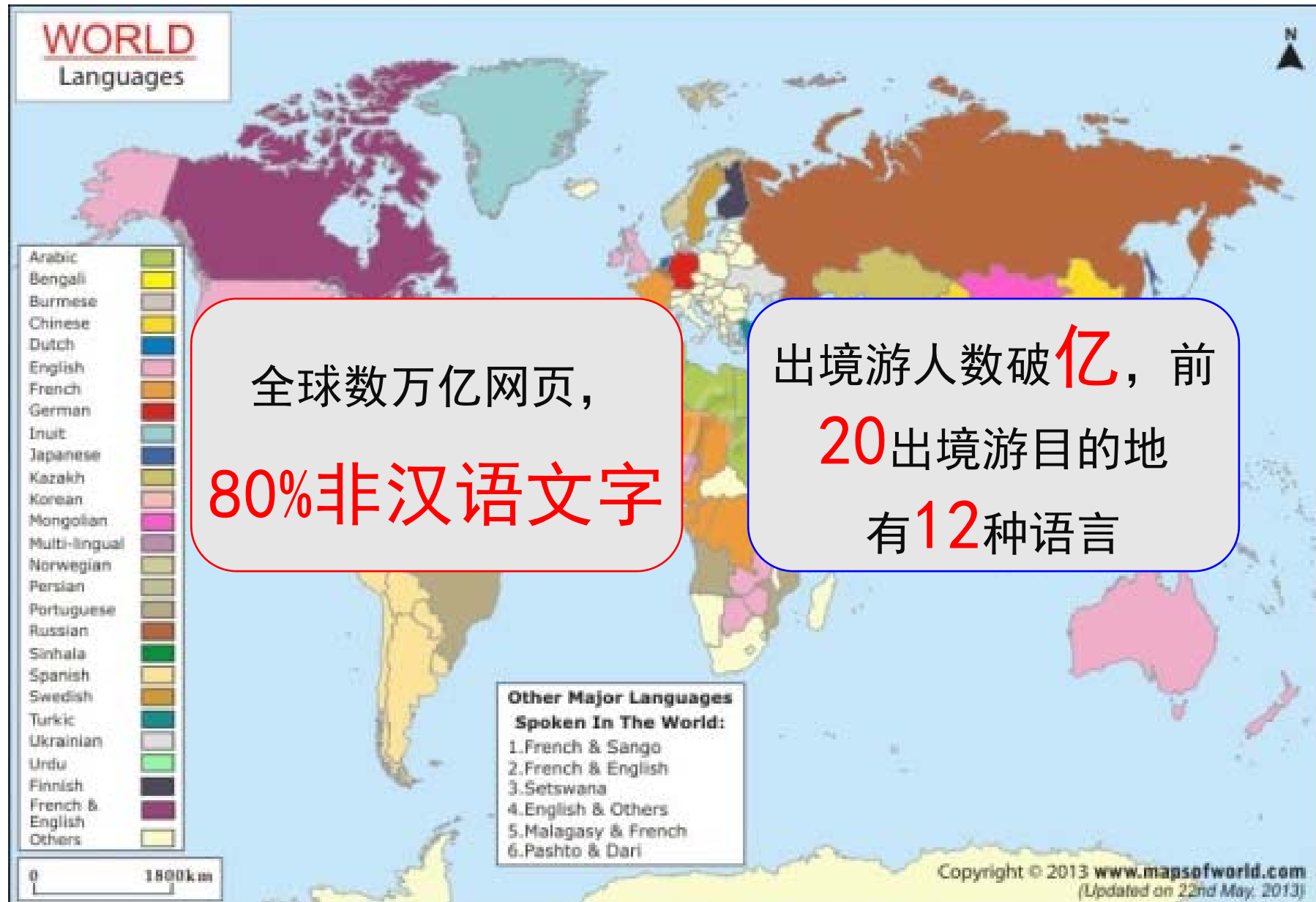


# 机器翻译



**Tower of Babel (巴别塔)**

# 世界语言地图







# 中国语言地图



# 背景：一带一路

64个国家和地区

44亿人口

50多种语言



# 机器翻译

**定义：** 机器翻译是**利用计算机**将一种自然语言（**源语言**）**自动转换**为另一种自然语言（**目标语言**）的技术。







# 机器翻译

**定义：** 机器翻译是**利用计算机**将一种自然语言（**源语言**）**自动转换**为另一种自然语言（**目标语言**）的技术。

The screenshot shows the Google Translate web interface. At the top, the Google logo is on the left, and the word "Translate" is in the center. Below the logo, there are tabs for "Chinese", "English", "Spanish", and "Detect language". The "English" tab is selected. To the right of the tabs, there is a large blue button labeled "文本翻译" (Text Translation). Below the tabs, there is a text input area containing the following Chinese text: "本月13日, 瑞典学院出人意料地将2016年诺贝尔文学奖授予鲍勃·迪伦, 称赞他在'伟大的美国歌曲传统中开创了诗意表达'。鲍勃·迪伦成为诺贝尔奖115年历史上首位获得文学奖的音乐人, 他的官方'推特'和'脸谱网'账号上都提到了获奖一事。不过, 瑞典学院一直联系不到鲍勃·迪伦。美国福克斯新闻网22日发送电邮请求鲍勃·迪伦的发言人拉里·詹金斯回应此事, 但截至23日发稿时仍未得到回复。" Below the text input area, there are icons for "Ä", a microphone, a speaker, and a "拼" (Pinyin) button. To the right of the text input area, there is a text output area containing the following English text: "On the 13th of this month, the Swedish Academy unexpectedly awarded the 2016 Nobel Prize for Literature to Bob Dylan, praising him for creating a poetic expression in the 'great American song tradition.' Bob Dylan became the first Nobel Prize-winning 115-year-old musician to win the prize in his official Twitter and Facebook accounts. However, the Swedish Academy has been linked to Bob Dylan. Fox News Network 22 to send an e-mail request Bob Dylan spokesman Larry Jenkins to respond to the matter, but as of 23 press time has not received a reply." Below the text output area, there are icons for a star, a document, a speaker, and a "Suggest an edit" button.

# 很多人眼中的机器翻译





# 很多人眼中的机器翻译





# 现在的机器翻译

Translate

Chinese English Spanish Detect language

Nine out of the ten Chinese crew members freed by Somali pirates took a flight home on Monday from the Kenyan capital Nairobi, accompanied by officials sent from Beijing.

English Chinese (Simplified) Spanish Translate

在索马里海盗释放的10名中国船员中，有9名在星期一从肯尼亚首都内罗毕飞往家乡，并从北京派出官员陪同。

☆ □ Ä 🔊 ↩

Suggest an edit

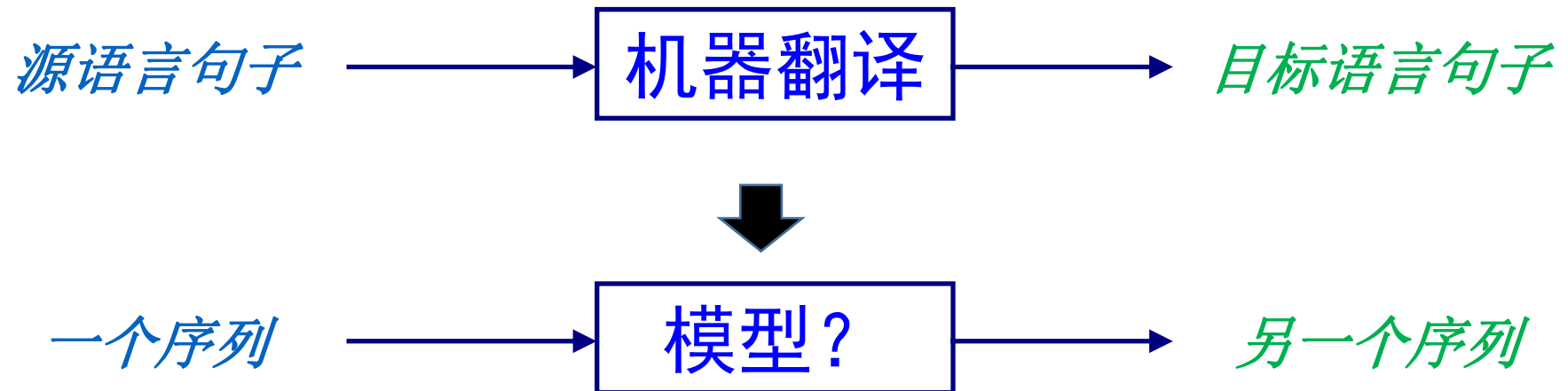
Zài suǒmǎlǐ hǎidào shìfàng de 10 míng zhōngguó chuányuán zhōng, yǒu 9 míng zài xīngqī yī cóng kěnniyǎ shǒudū nèiluóbì fēi wǎng jiāxiāng, bìng cóng běijīng pàichū guānyuán péitóng.



# 机器翻译形式化

今天北京天气不错

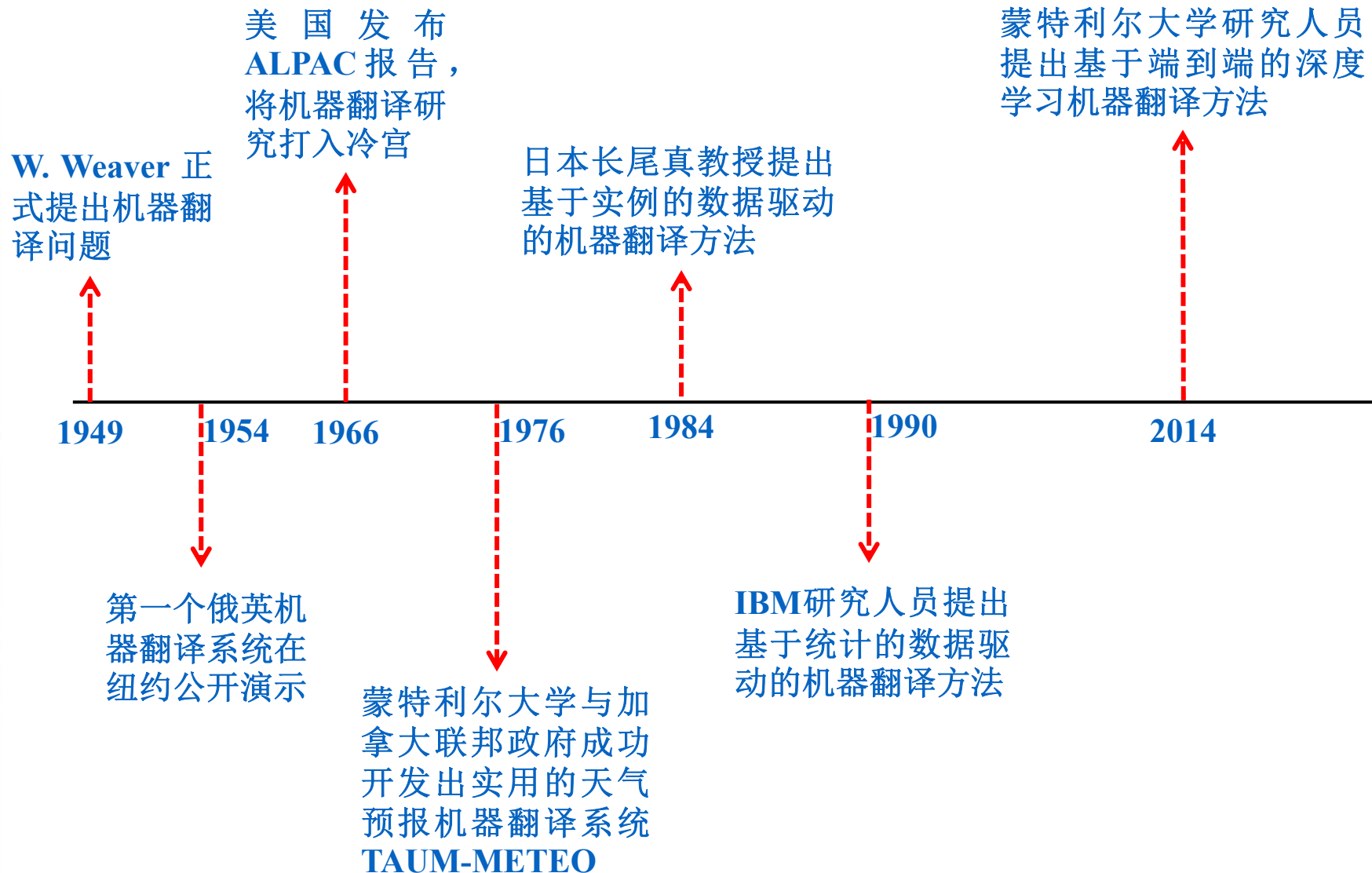
*The weather is fine in Beijing today*







# 机器翻译简史





# 机器翻译的困难

## ◆ 自然语言中普遍存在的歧义和未知现象

- 句法结构歧义/词汇歧义/语用歧义 ...
- 新的词汇、术语、结构、语义 ...

## ◆ 机器翻译不仅仅是字符串的转换

- 不同语言之间文化的差异
- 现有方法无法表示和利用世界知识和常识

## ◆ 机器翻译的解不唯一，而且始终存在的人为的标准

**几乎自然语言处理中的所有问题在机器翻译中都会遇到**



# 机器翻译方法

- ◆ 直接转换法
- ◆ 基于规则的翻译方法
- ◆ 基于中间语言的翻译方法
- ◆ 基于语料库（数据驱动）的翻译方法
  - 基于实例的翻译方法
  - 统计机器翻译
  - 神经网络机器翻译





# 机器翻译方法

- ◆ 直接转换法
- ◆ 基于规则的翻译方法
- ◆ 基于中间语言的翻译方法
- ◆ 基于语料库（数据驱动）的翻译方法
  - 基于实例的翻译方法
  - 统计机器翻译
  - 神经网络机器翻译



# 双语对照数据

人类 共 有 二十三 对 染色体 。

humans have a total of 23 pairs of chromosomes .

澳洲 重新 开放 驻 马尼拉 大使馆

australia reopens embassy in manila

中国 大陆 手机 用户 成长 将 减缓

growth of phone users in mainland china to slow

外交 人员 搭乘 第五 架 飞机 返国

diplomatic staff take the fifth plane home

驻 南韩 美军 三千人 奉命 冻结 调防

us freezes transfer of 3,000 troops in south korea

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-基于词的模型

澳洲<sub>1</sub> 与<sub>2</sub> 北韩<sub>3</sub> 有<sub>4</sub> 邦交<sub>5</sub>

**f<sub>1</sub>**      **f<sub>2</sub>**      **f<sub>3</sub>**      **f<sub>4</sub>**      **f<sub>5</sub>**      **f<sub>6</sub>**      **f<sub>7</sub>**

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

澳洲<sub>1</sub> 与<sub>2</sub> 北韩<sub>3</sub> 有<sub>4</sub> 邦交<sub>5</sub>

**f<sub>1</sub>**      **f<sub>2</sub>**      **f<sub>3</sub>**      **f<sub>4</sub>**      **f<sub>5</sub>**      **f<sub>6</sub>**      **f<sub>7</sub>**

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

澳洲<sub>1</sub> 与<sub>2</sub> 北韩<sub>3</sub> 有<sub>4</sub> 邦交<sub>5</sub>

**Austria<sub>1</sub> has<sub>2</sub> diplomatic<sub>3</sub> relations<sub>4</sub> with<sub>5</sub> North<sub>6</sub> Korea<sub>7</sub>**

$$p(s_j|t_{a_j})$$

IBM Model-2



Garcia and associates .

Garcia y asociados .

Carlos Garcia has three associates .

Carlos Garcia tiene tres asociados .

his associates are not strong .

sus asociados no son fuertes .

Garcia has a company also .

Garcia tambien tiene una empresa .

its clients are angry .

sus clientes estan enfadados .

the associates are also angry .

los asociados tambien estan enfadados .

the clients and the associates are enemies .

los clientes y los asociados son enemigos .

the company has three groups .

la empresa tiene tres grupos .

its groups are in Europe .

sus grupos estan en Europa .

the modern groups sell strong pharmaceuticals .

los grupos modernos venden medicinas fuertes .

the groups do not sell zanzanine .

los grupos no venden zanzanina .

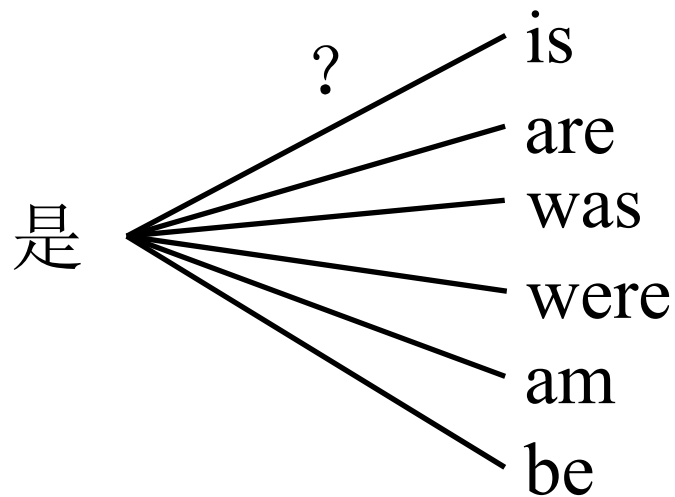
the small groups are not modern .

los grupos pequenos no son modernos .

# 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  $\xrightarrow{?}$  North Korea

邦交  $\xrightarrow{?}$  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



[Koehn, 2003]

短语：连续的词串（非句法意义）

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

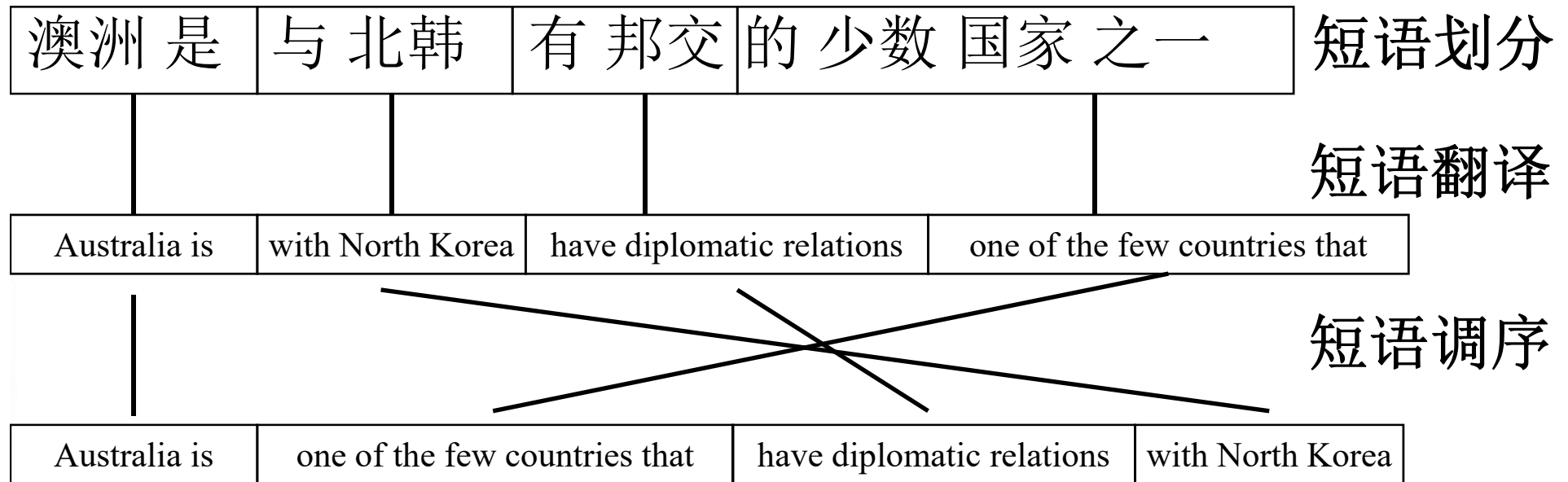


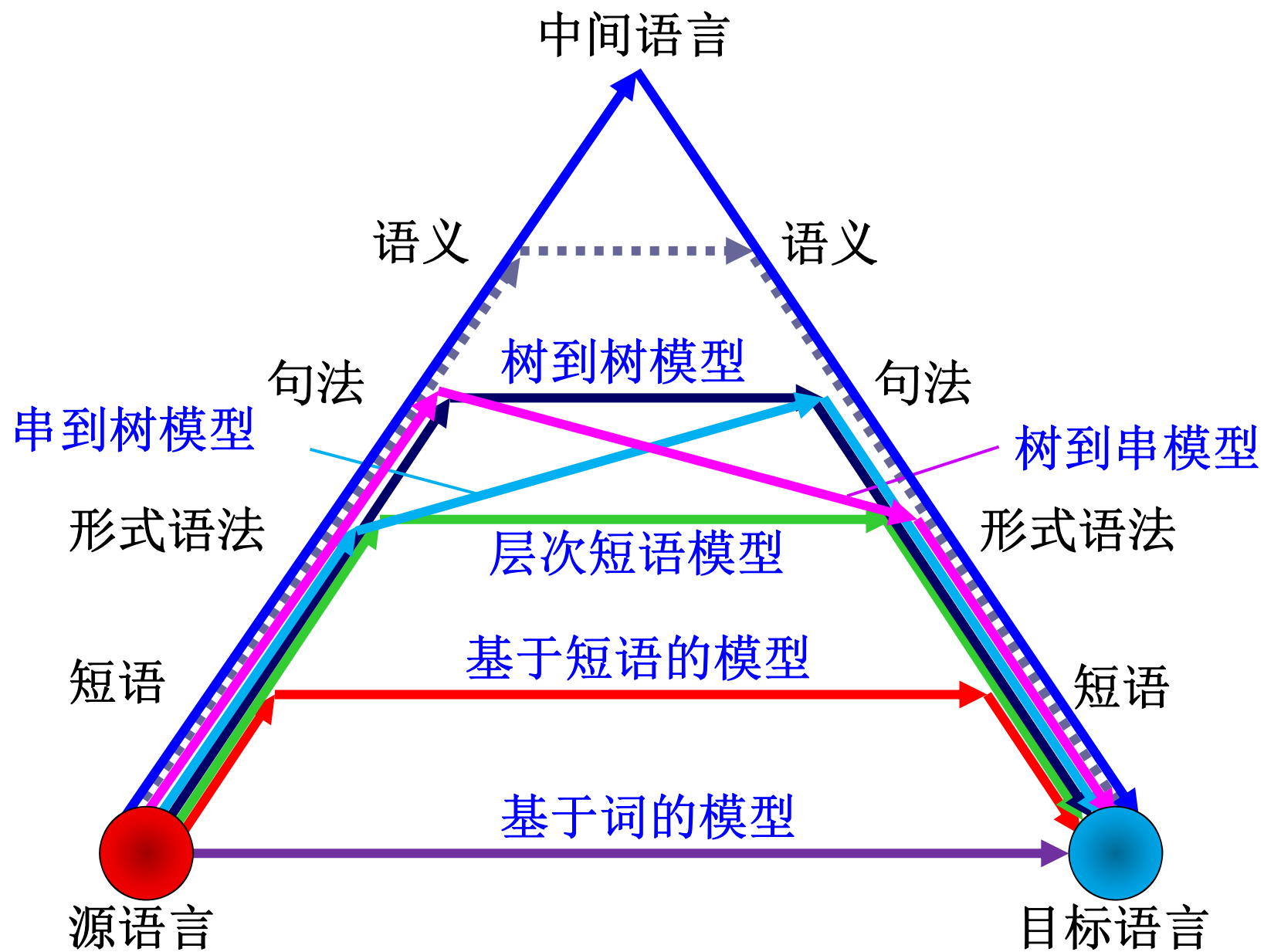
# 基于短语的统计机器翻译

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



# 基于短语的统计机器翻译







# 统计机器翻译

Chinese: 我 在 北京 做了 报告

Phrase Seg:

①可解释性高

做了 报告

Phrase Trans:

人工设定的模块和特征

gave a talk

②模块随便加

③错误易追踪

English:

I

gave

a

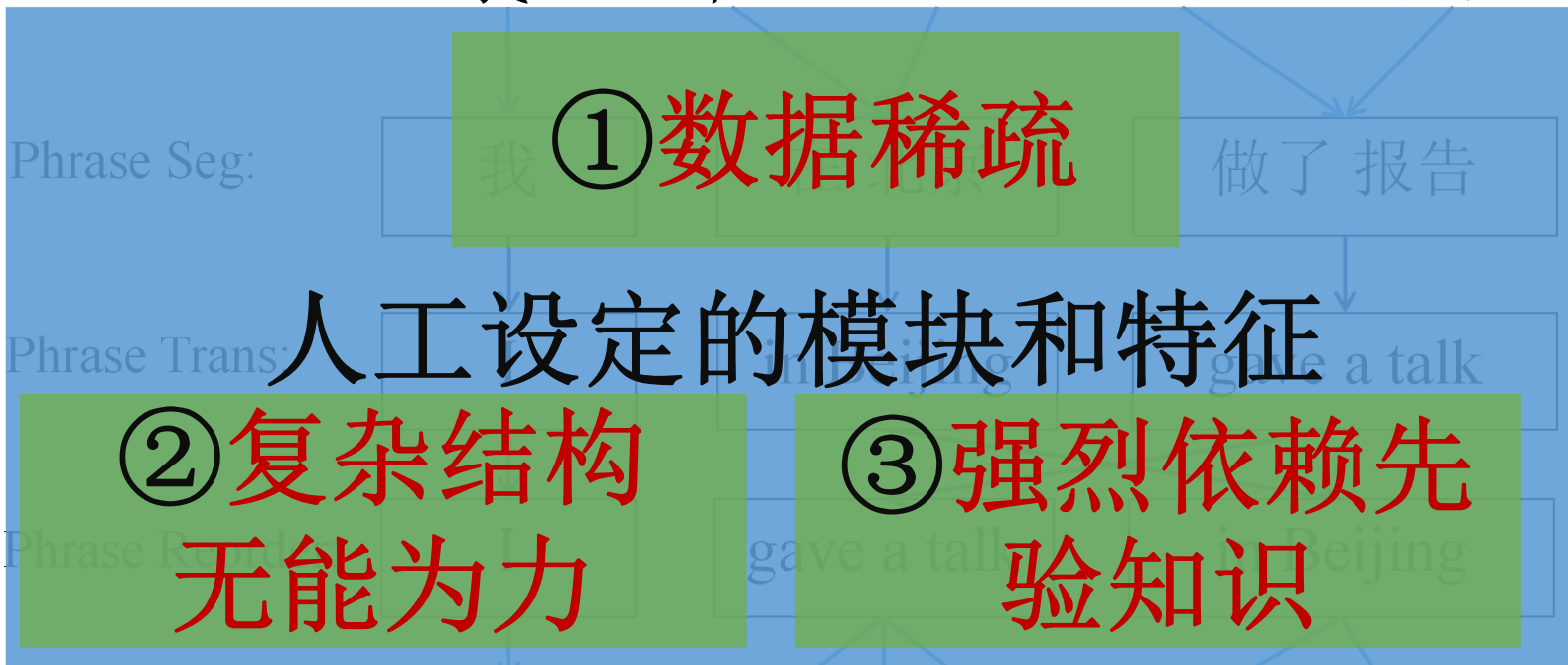
talk

in

Beijing

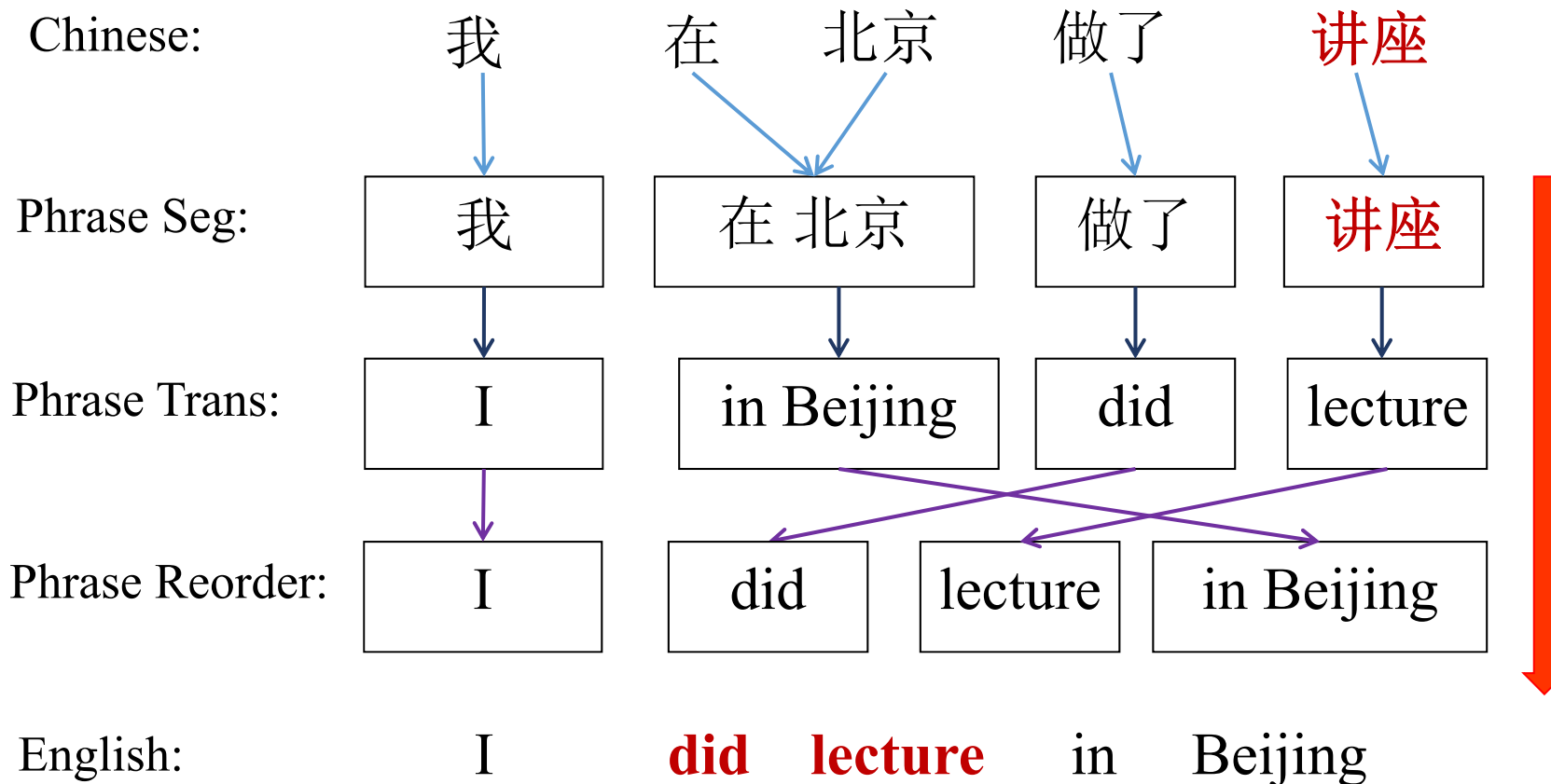
# 统计机器翻译

Chinese: 我 在 北京 做了 报告



English: I gave a talk in Beijing

# 统计机器翻译



①数据稀疏



# 统计机器翻译

Chinese

美国总统布什昨天在白宫与以色列总理沙龙就中东局势 ×  
举行了一个小时的会谈。

English

Yesterday, U.S. President George W. Bush at the White House with Israeli Prime Minister Ariel Sharon on the  
situation in the Middle East held a one-hour talks.

②复杂结构无能为力

# 现实世界 VS. 认知世界

- 现实世界：物体相互独立地存在





# 现实世界 VS. 认知世界

- 认知世界：概念互相联系、语义连续分布



# 统计机器翻译→神经机器翻译

离散符号表示方法  $\Rightarrow$  连续分布式表示方法

讲座  $\otimes$  报告 = 0

讲座

报告

$\begin{bmatrix} 0.48 \\ 0.46 \\ 0.26 \end{bmatrix}$

$\otimes$

$\begin{bmatrix} 0.42 \\ 0.51 \\ 0.21 \end{bmatrix}$

$\approx 1$

分布式语义表示  
是核心和基础

今天  
明天  
昨天

一月  
三月  
五月

低维、稠密的连续实数空间

# 神经机器翻译

Chinese:

我 在 北京 做了 报告

编码网络

仅需要两个神经网络

分布式语义表示

解码网络

English:

I gave a talk in Beijing

# 神经机器翻译

The screenshot displays the Google Translate interface. At the top, the Google logo is visible. Below it, the word "Translate" is written in red. The language selection bar shows "Chinese" as the source language and "English" as the target language. The input text in Chinese is: "美国总统布什昨天在白宫与以色列总理沙龙就中东局势举行了一个小时的会谈。". The output text in English is: "US President George W. Bush held an hour-long meeting with Israeli Prime Minister Ariel Sharon on the situation in the Middle East yesterday at the White House." The English text is segmented by colored bars (green, red, blue, yellow) that correspond to the word alignment with the Chinese text. The interface also includes a "Translate" button, a "Suggest an edit" link, and various utility icons like a star, a document, a speaker, and a share icon.

Google

Translate

Chinese English Spanish Detect language

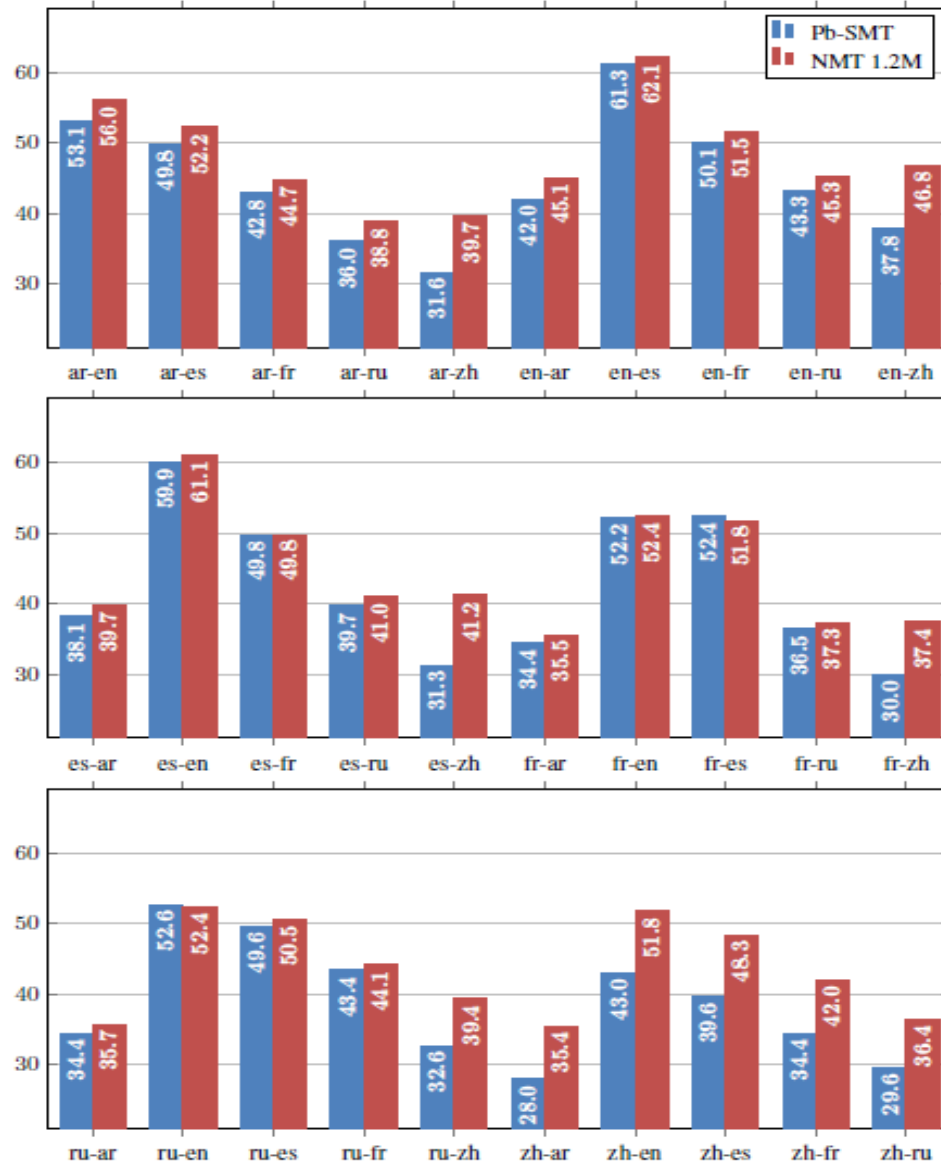
美国总统布什昨天在白宫与以色列总理沙龙就中东局势举行了一个小时的会谈。

English Chinese (Simplified) Spanish Translate

US President George W. Bush held an hour-long meeting with Israeli Prime Minister Ariel Sharon on the situation in the Middle East yesterday at the White House.

☆ 📄 🔊 ➦ Suggest an edit

# 神经机器翻译



**神经机器翻译压倒性胜出！**

[Junczys-Dowmunt et al, 2016]



# 统计机器翻译→神经机器翻译

离散符号表示方法  $\Rightarrow$  连续分布式表示方法

讲座

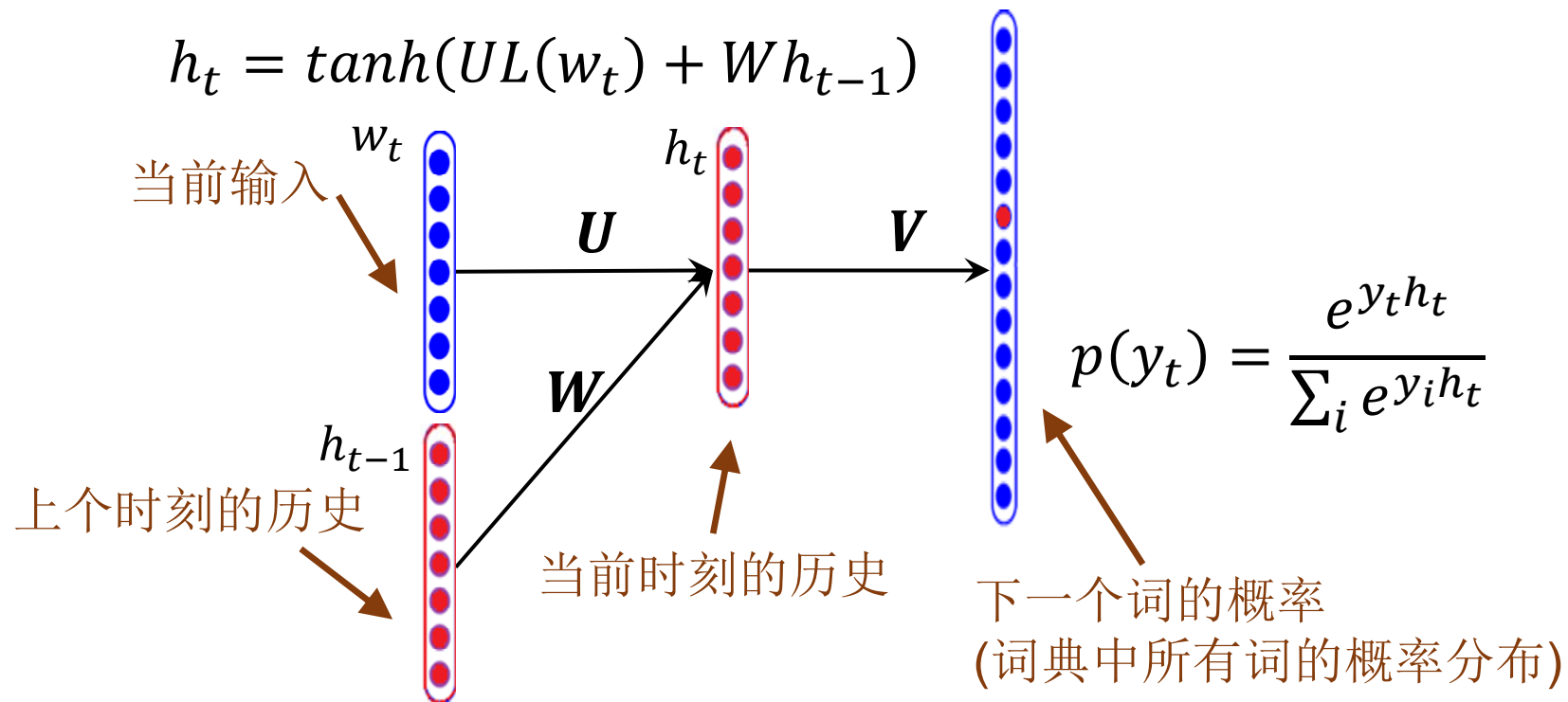
报告

$$\begin{bmatrix} 0.48 \\ 0.46 \\ 0.26 \end{bmatrix} \otimes \begin{bmatrix} 0.42 \\ 0.51 \\ 0.21 \end{bmatrix} \approx 1$$

表示是核心  
运算是关键

# 循环神经网络

- 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  $t + 1$ 时刻所有词的概率分布



# 源语言编码

$$h_t = \tanh(U_S L(w_t) + W_S h_{t-1})$$

$L(w_t):$       $w_t \longrightarrow \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \in R^3$      我  $\longrightarrow \begin{bmatrix} 0.1 \\ 0.9 \\ 0.6 \end{bmatrix}$      随机初始化

<s> 我 在 北 京 做 了 报 告 </s>

0.2	0.1	0.1	0.3	0.2	0.4	0.3
0.3	0.9	0.2	0.8	0.1	0.1	0.1
0.5	0.6	0.4	0.3	0.3	0.2	0.2

# 源语言编码

$$h_t = \tanh(U_S L(w_t) + W_S h_{t-1})$$

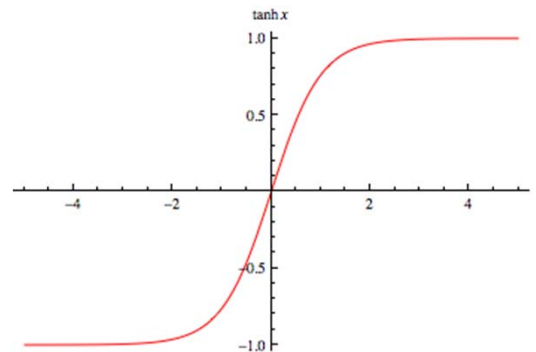
$L(w_t):$        $w_t \longrightarrow \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \in R^3$       我  $\longrightarrow \begin{bmatrix} 0.1 \\ 0.9 \\ 0.6 \end{bmatrix}$       随机初始化

$h_{t-1}:$       上一时刻的历史信息       $h_0 = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.0 \end{bmatrix}$

$$U_S = \begin{bmatrix} 0.1 & 0.2 & 0.1 \\ 0.3 & 0.2 & 0.0 \\ 0.4 & 0.0 & 0.2 \end{bmatrix} \in R^{3 \times 3} \quad W_S = \begin{bmatrix} 0.0 & 0.3 & 0.2 \\ 0.1 & 0.1 & 0.3 \\ 0.0 & 0.4 & 0.1 \end{bmatrix} \in R^{3 \times 3}$$

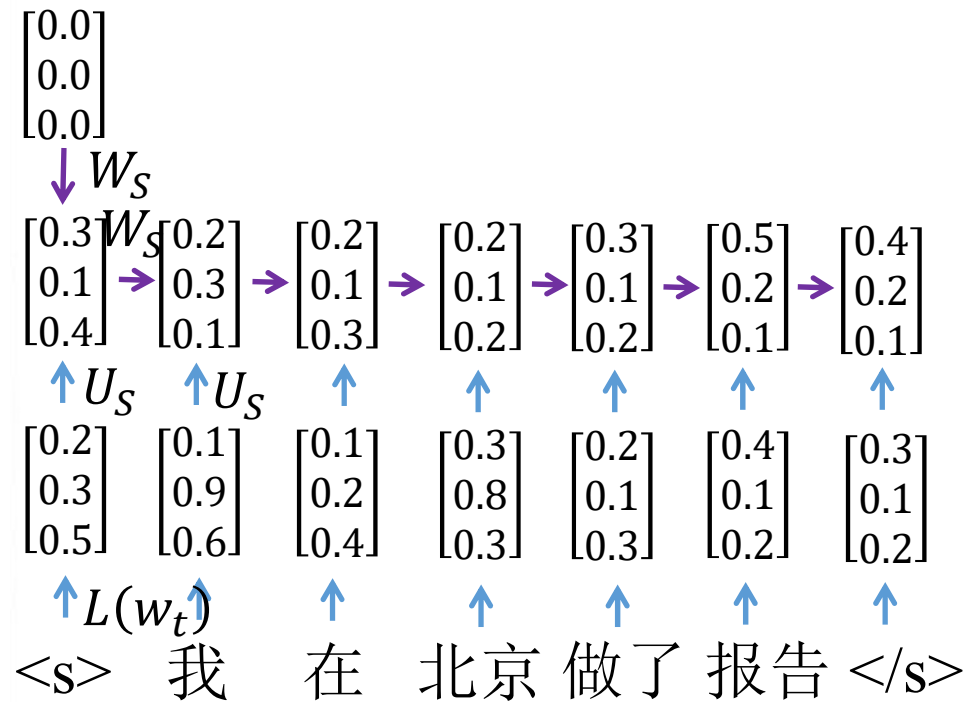
$$z = U_S L(w_t) + W_S h_{t-1} \in R^3$$

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \longrightarrow$$



# 源语言编码

$$h_t = \tanh(U_S L(w_t) + W_S h_{t-1})$$



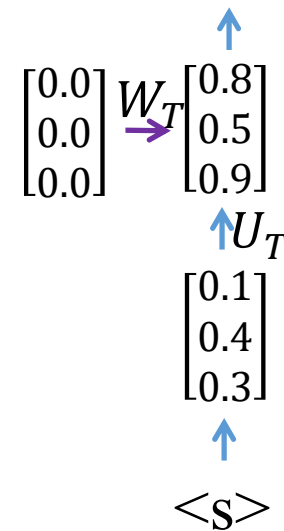




# 目标语言解码

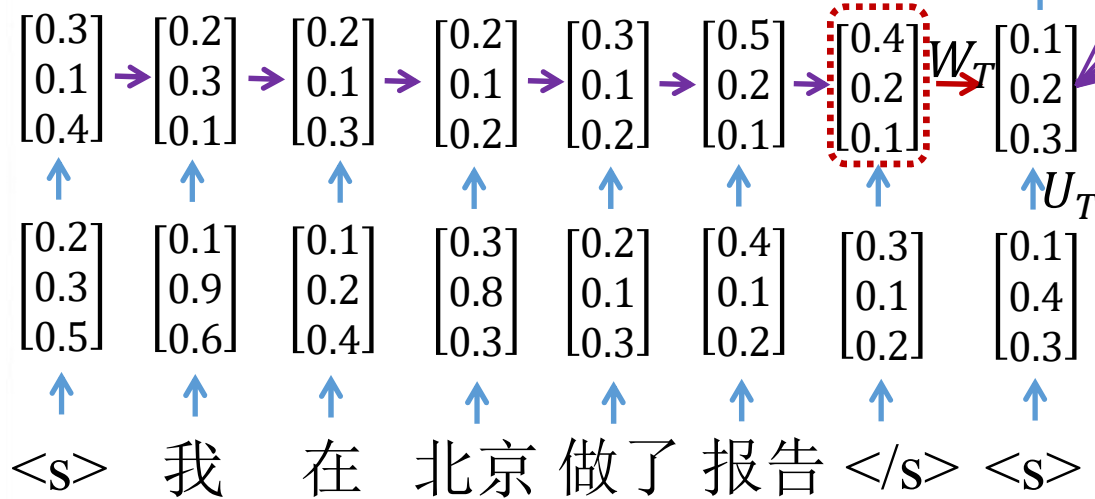
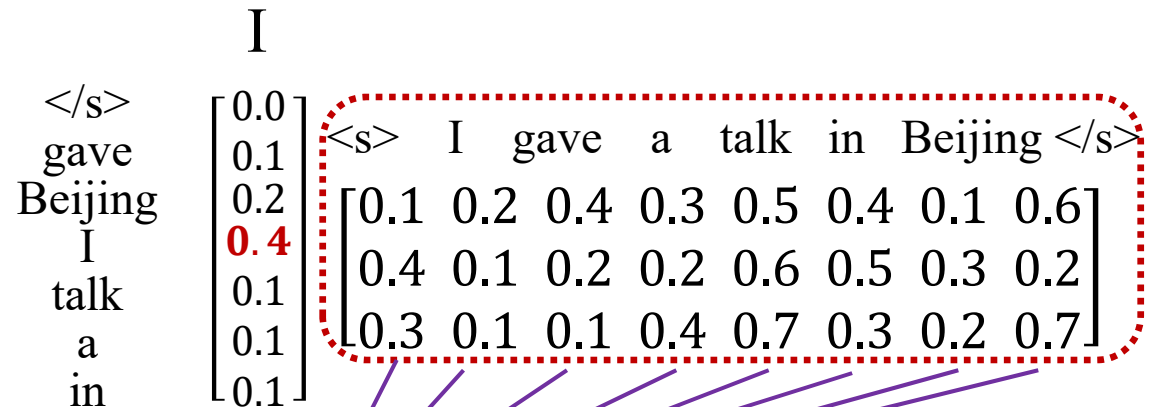
$$h_t = \tanh(U_T L(w_t) + W_T h_{t-1})$$

随机？抽样？



# 目标语言解码

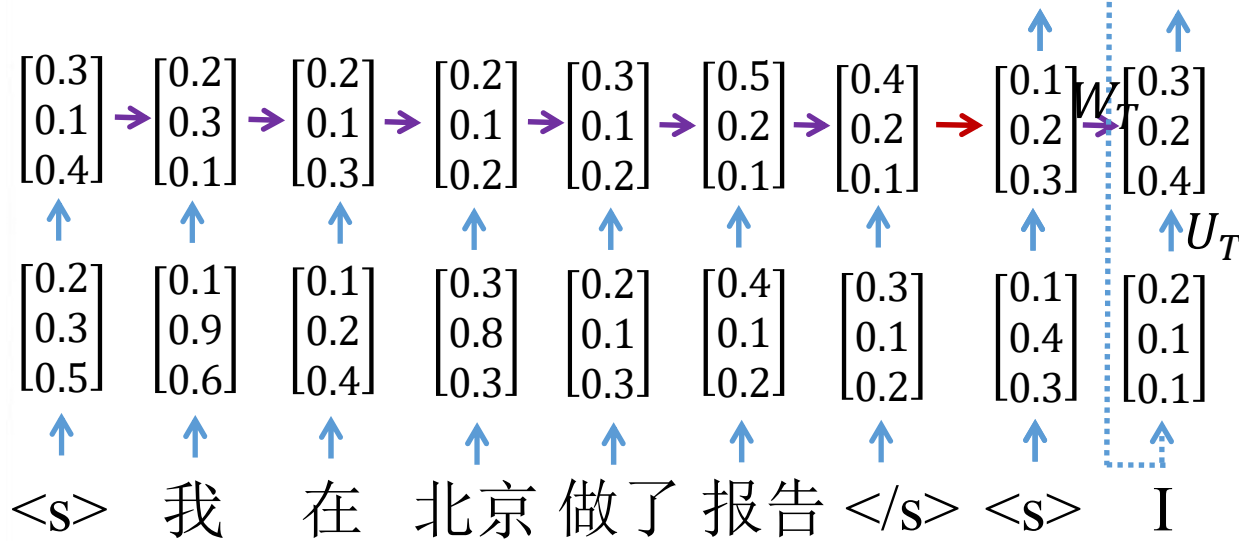
$$p(y_t) = \frac{e^{y_t h_t}}{\sum_i e^{y_i h_t}}$$



# 目标语言解码

$$p(y_t) = \frac{e^{y_t h_t}}{\sum_i e^{y_i h_t}}$$

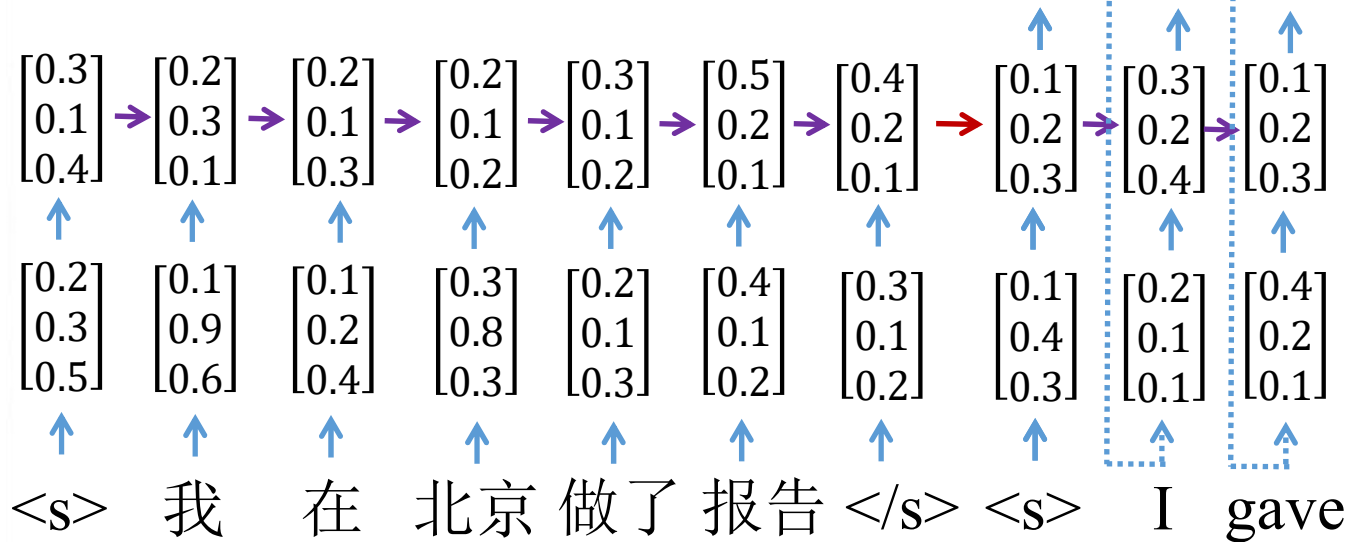
	I	gave
</s>	0.0	0.1
gave	0.1	<b>0.5</b>
Beijing	0.2	0.1
I	<b>0.4</b>	0.1
talk	0.1	0.0
a	0.1	0.1
in	0.1	0.1



# 目标语言解码

$$p(y_t) = \frac{e^{y_t h_t}}{\sum_i e^{y_i h_t}}$$

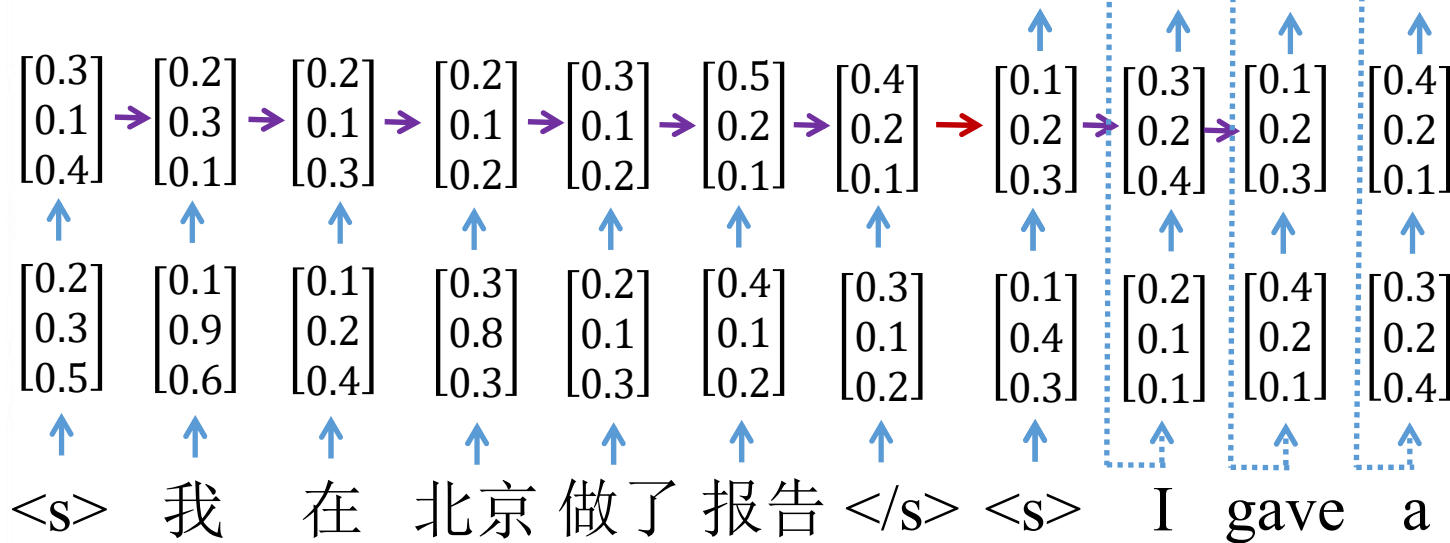
	I	gave	a
</s>	$\begin{bmatrix} 0.0 \\ 0.1 \\ 0.2 \\ \mathbf{0.4} \\ 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$	$\begin{bmatrix} 0.1 \\ \mathbf{0.5} \\ 0.1 \\ 0.1 \\ 0.0 \\ 0.1 \\ 0.1 \end{bmatrix}$	$\begin{bmatrix} 0.1 \\ 0.2 \\ 0.1 \\ 0.1 \\ 0.1 \\ \mathbf{0.3} \\ 0.1 \end{bmatrix}$
gave			
Beijing			
I			
talk			
a			
in			



# 目标语言解码

$$p(y_t) = \frac{e^{y_t h_t}}{\sum_i e^{y_i h_t}}$$

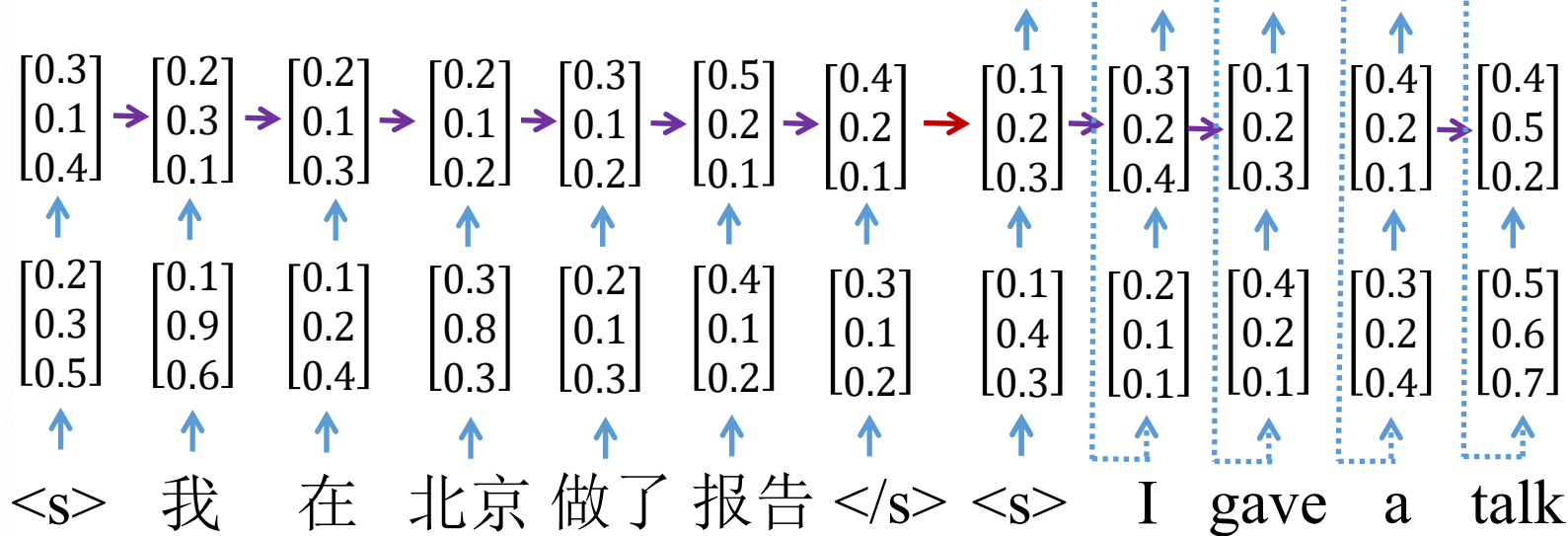
	I	gave	a	talk
</s>	0.0	0.1	0.1	0.1
gave	0.1	<b>0.5</b>	0.2	0.2
Beijing	0.2	0.1	0.1	0.1
I	<b>0.4</b>	0.1	0.1	0.1
talk	0.1	0.0	0.1	<b>0.4</b>
a	0.1	0.1	<b>0.3</b>	0.0
in	0.1	0.1	0.1	0.1



# 目标语言解码

$$p(y_t) = \frac{e^{y_t h_t}}{\sum_i e^{y_i h_t}}$$

	I	gave	a	talk	in
</s>	$\begin{bmatrix} 0.0 \\ 0.1 \\ 0.2 \\ 0.4 \\ 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$	$\begin{bmatrix} 0.1 \\ 0.5 \\ 0.1 \\ 0.1 \\ 0.0 \\ 0.1 \\ 0.1 \end{bmatrix}$	$\begin{bmatrix} 0.1 \\ 0.2 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.3 \\ 0.1 \end{bmatrix}$	$\begin{bmatrix} 0.1 \\ 0.2 \\ 0.1 \\ 0.1 \\ 0.4 \\ 0.0 \\ 0.1 \end{bmatrix}$	$\begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.2 \\ 0.1 \\ 0.1 \\ 0.3 \end{bmatrix}$
gave					
Beijing					
I					
talk					
a					
in					

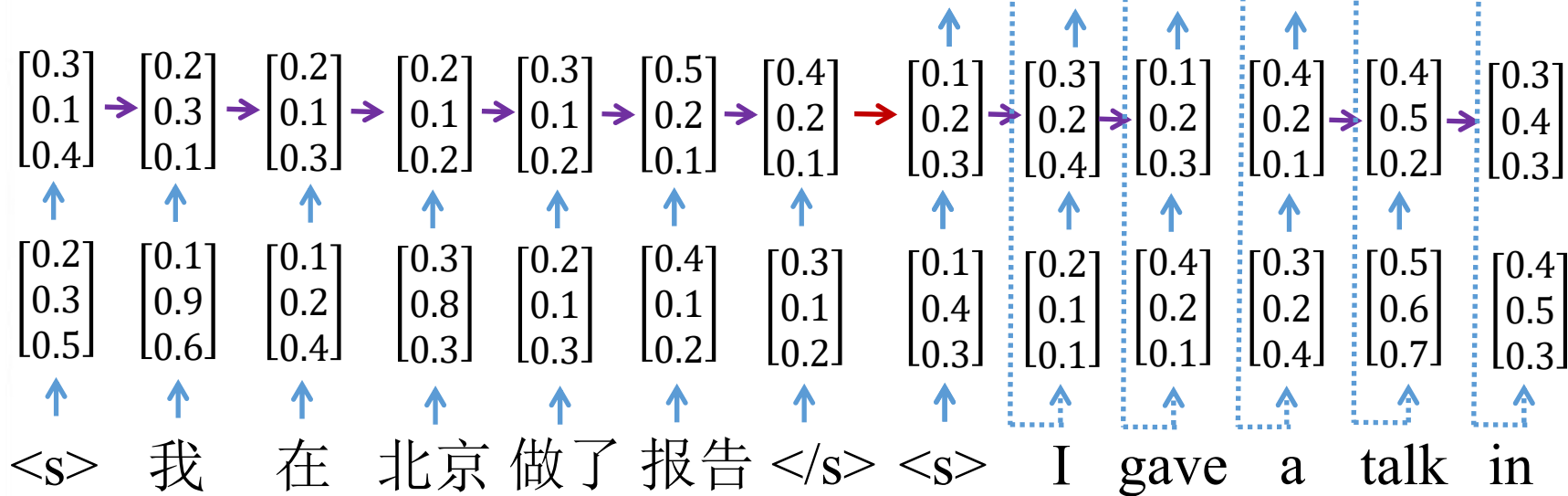




# 目标语言解码

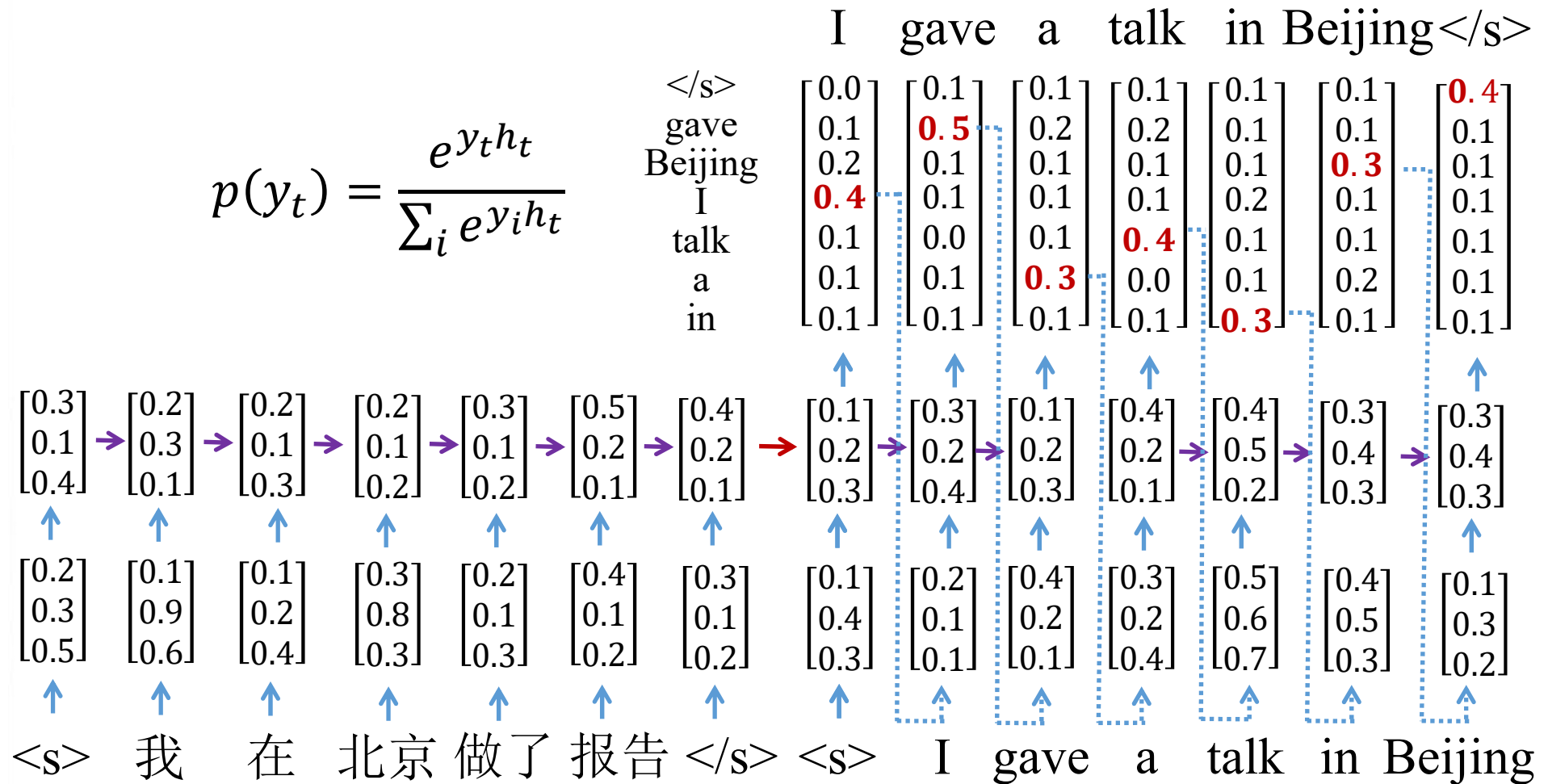
$$p(y_t) = \frac{e^{y_t h_t}}{\sum_i e^{y_i h_t}}$$

	I	gave	a	talk	in	Beijing
</s>	0.0	0.1	0.1	0.1	0.1	0.1
gave	0.1	<b>0.5</b>	0.2	0.2	0.1	0.1
Beijing	0.2	0.1	0.1	0.1	0.1	<b>0.3</b>
I	<b>0.4</b>	0.1	0.1	0.1	0.2	0.1
talk	0.1	0.0	0.1	<b>0.4</b>	0.1	0.1
a	0.1	0.1	<b>0.3</b>	0.0	0.1	0.2
in	0.1	0.1	0.1	0.1	<b>0.3</b>	0.1



# 目标语言解码

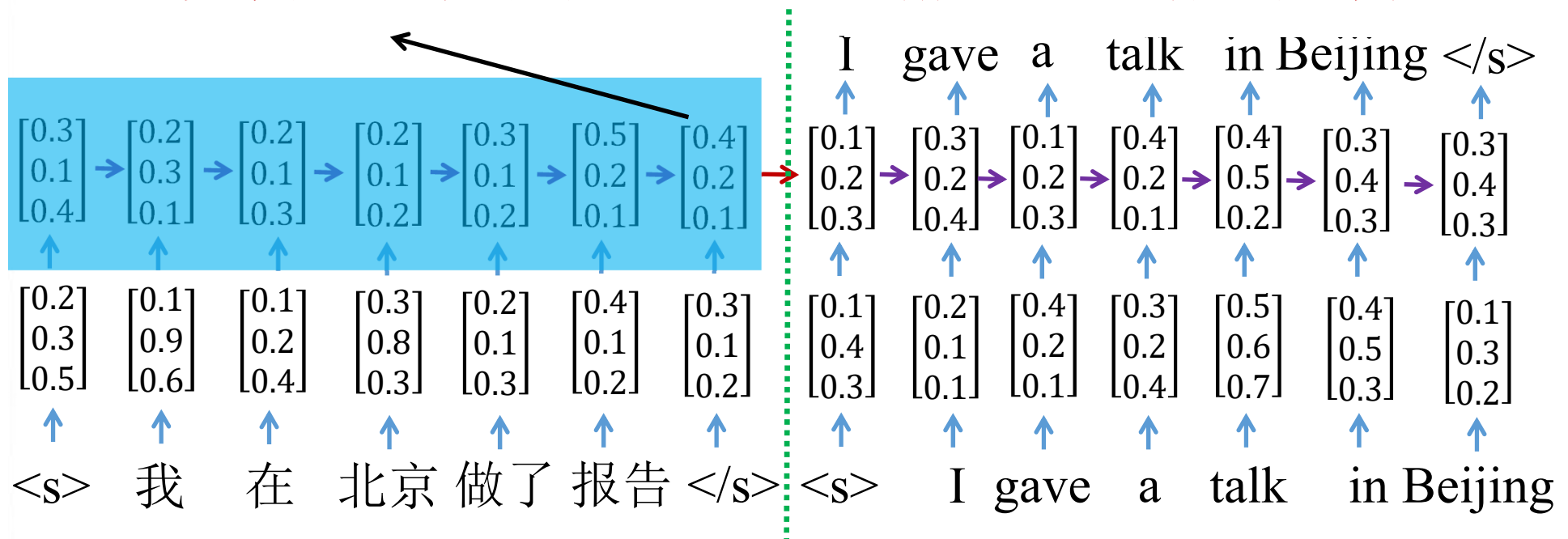
$$p(y_t) = \frac{e^{y_t h_t}}{\sum_i e^{y_i h_t}}$$



# 编码-解码网络

将源语言句子编码成一个  
实数向量语义表示

将源语言句子的语义表示  
解码生成目标语言句子

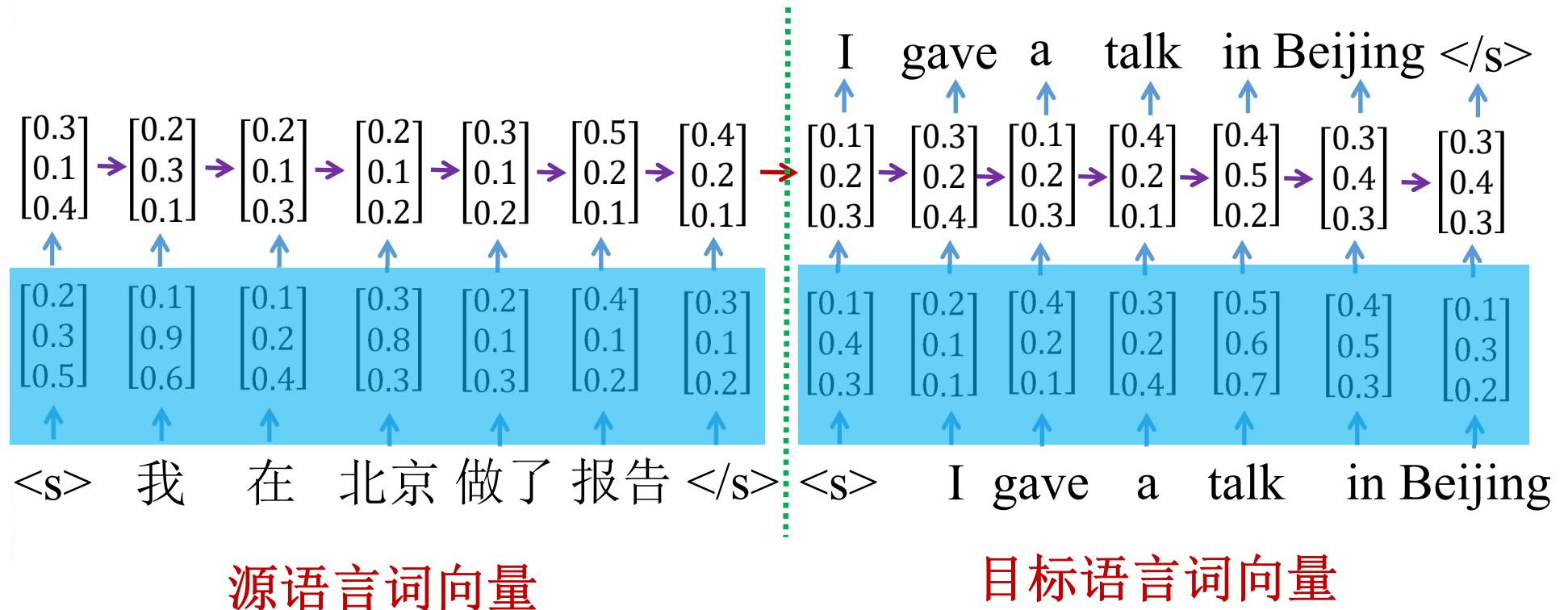


编码器

解码器

# 编码-解码网络

词向量随机初始化，在训练过程中进行优化！





# 测试 vs. 训练

I gave a talk in Beijing</s>

测试解码输入：只有源语言句子

<s> 我在北京做了报告 </s>

0.1  
0.1  
**0.3**  
0.1  
0.1  
0.1  
0.2  
0.1

**0.4**  
0.1  
0.1  
0.1  
0.1  
0.1  
0.1  
0.1

in [0.1] [0.1] [0.1] [0.1] **[0.3]** [0.1]



$\begin{bmatrix} 0.3 \\ 0.1 \\ 0.4 \end{bmatrix} \rightarrow \begin{bmatrix} 0.2 \\ 0.3 \\ 0.1 \end{bmatrix}$



$\begin{bmatrix} 0.2 \\ 0.3 \\ 0.5 \end{bmatrix} \rightarrow \begin{bmatrix} 0.1 \\ 0.9 \\ 0.6 \end{bmatrix}$



模型训练输入：源语言句子和正确译文

<s> 我在北京做了报告 </s>

<s> I gave a talk in Beijing </s>

0.3  
0.4  
0.3

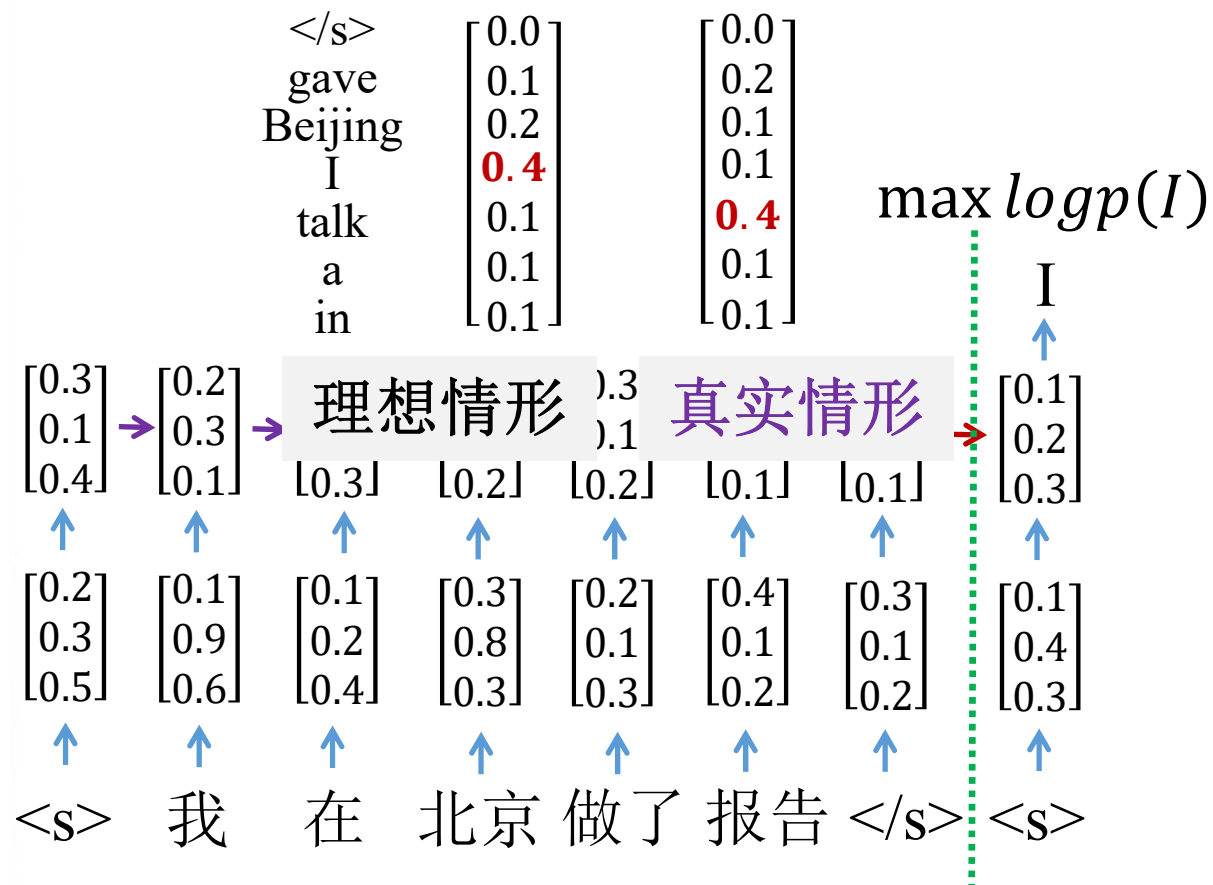
0.3  
0.4  
0.3

0.4  
0.5  
0.3

0.1  
0.3  
0.2

<s> 我在北京做了报告 </s> <s> I gave a talk in Beijing

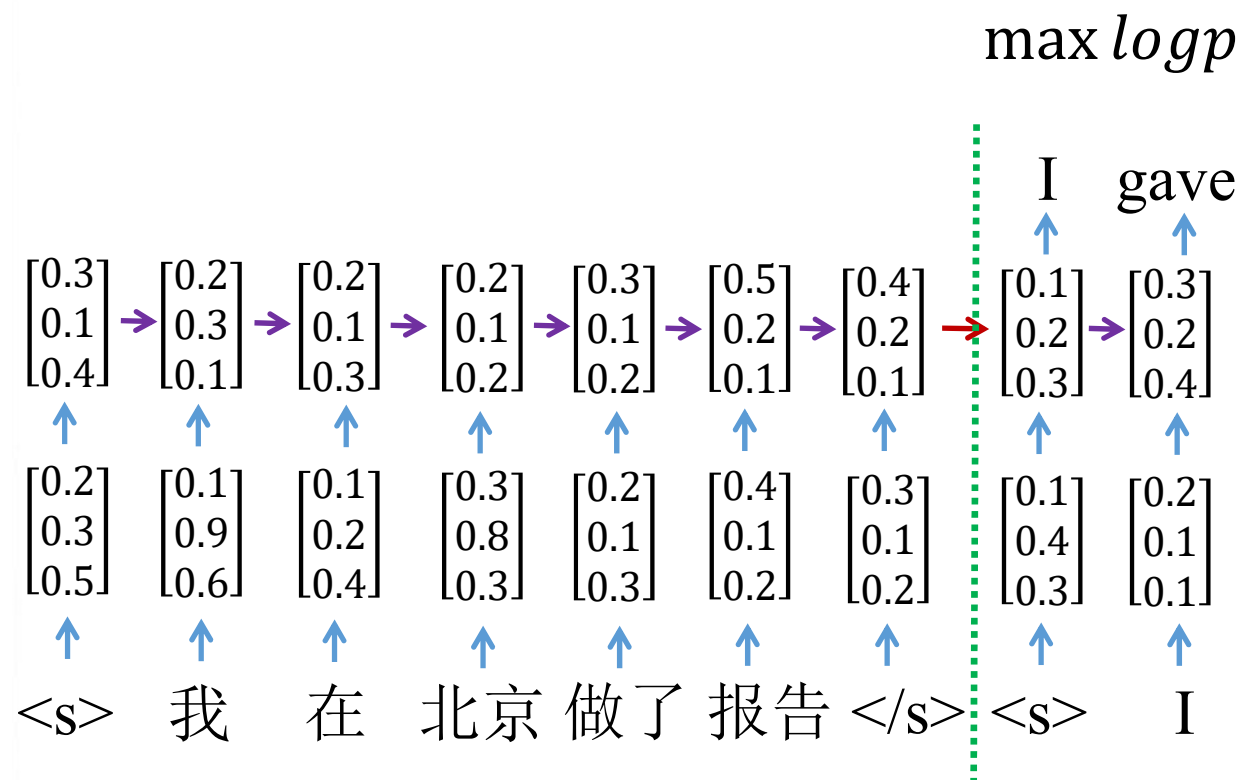
# 神经机器翻译参数优化



最大化  $P(\text{target}|\text{source})$



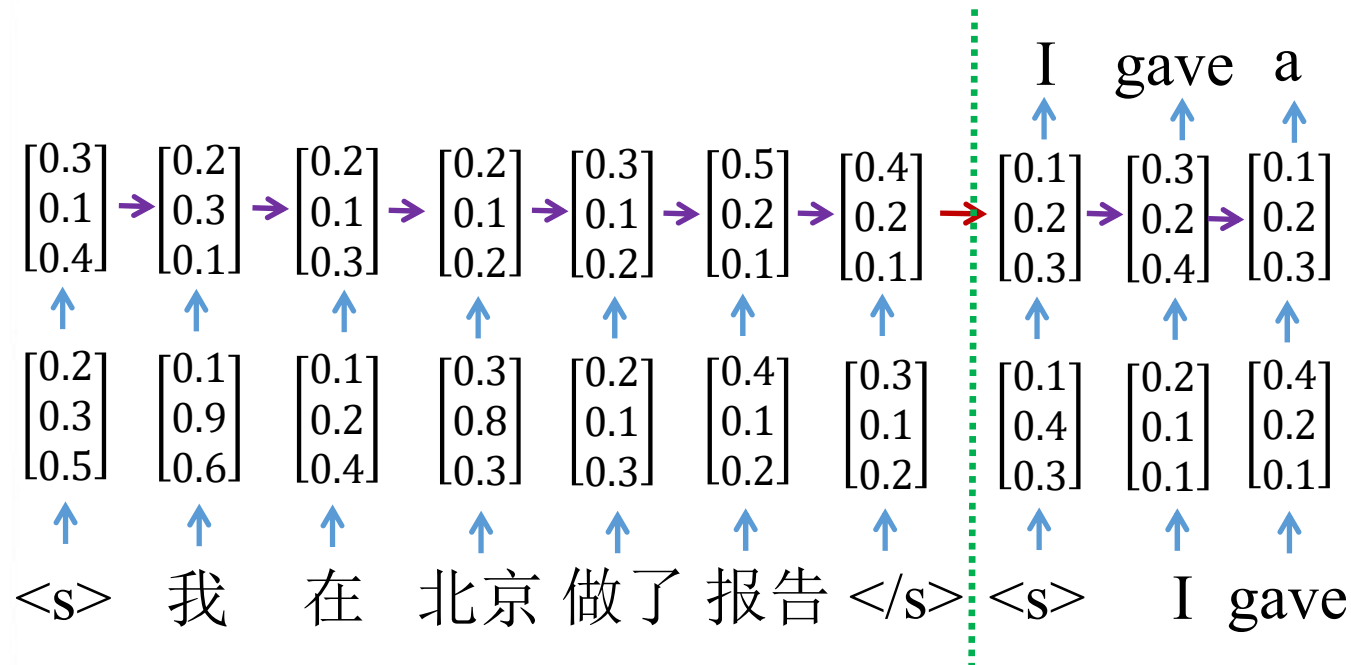
# 神经机器翻译参数优化



最大化  $P(\text{target}|\text{source})$

# 神经机器翻译参数优化

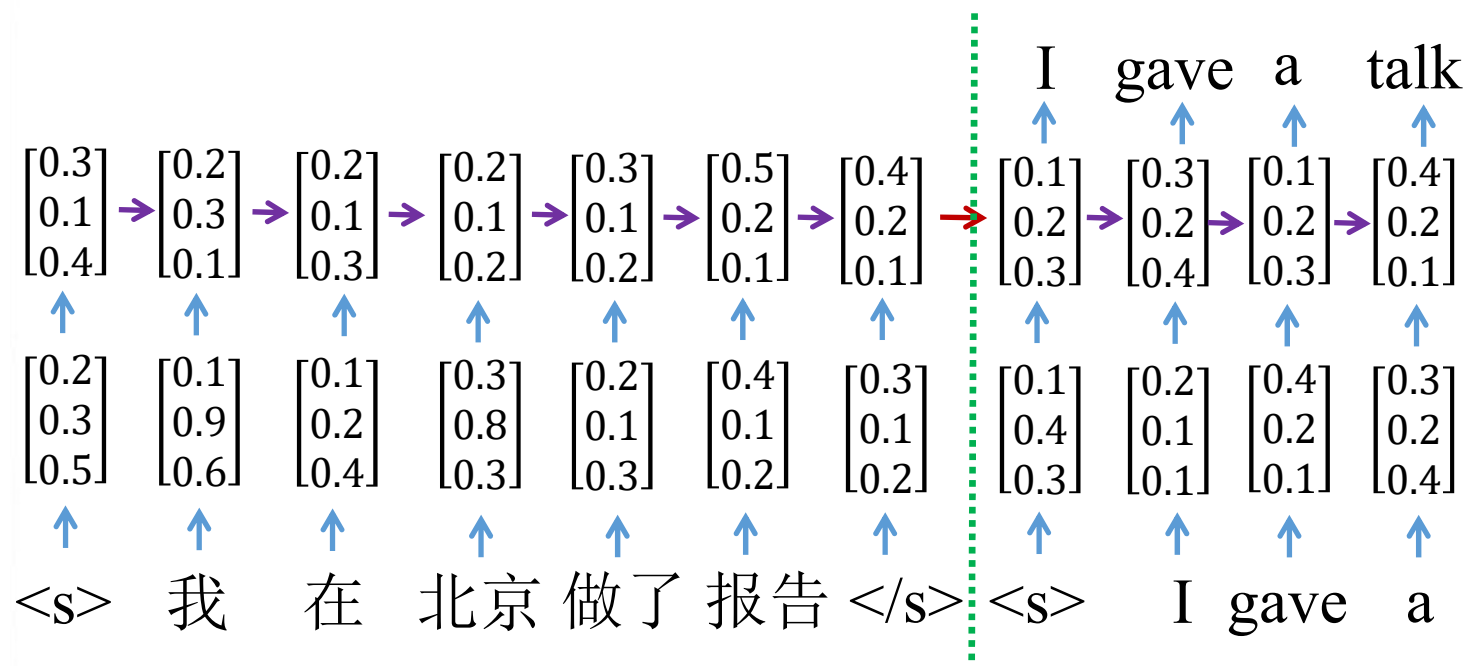
$$\max \log p(a)$$



最大化  $P(\textit{target}|\textit{source})$

# 神经机器翻译参数优化

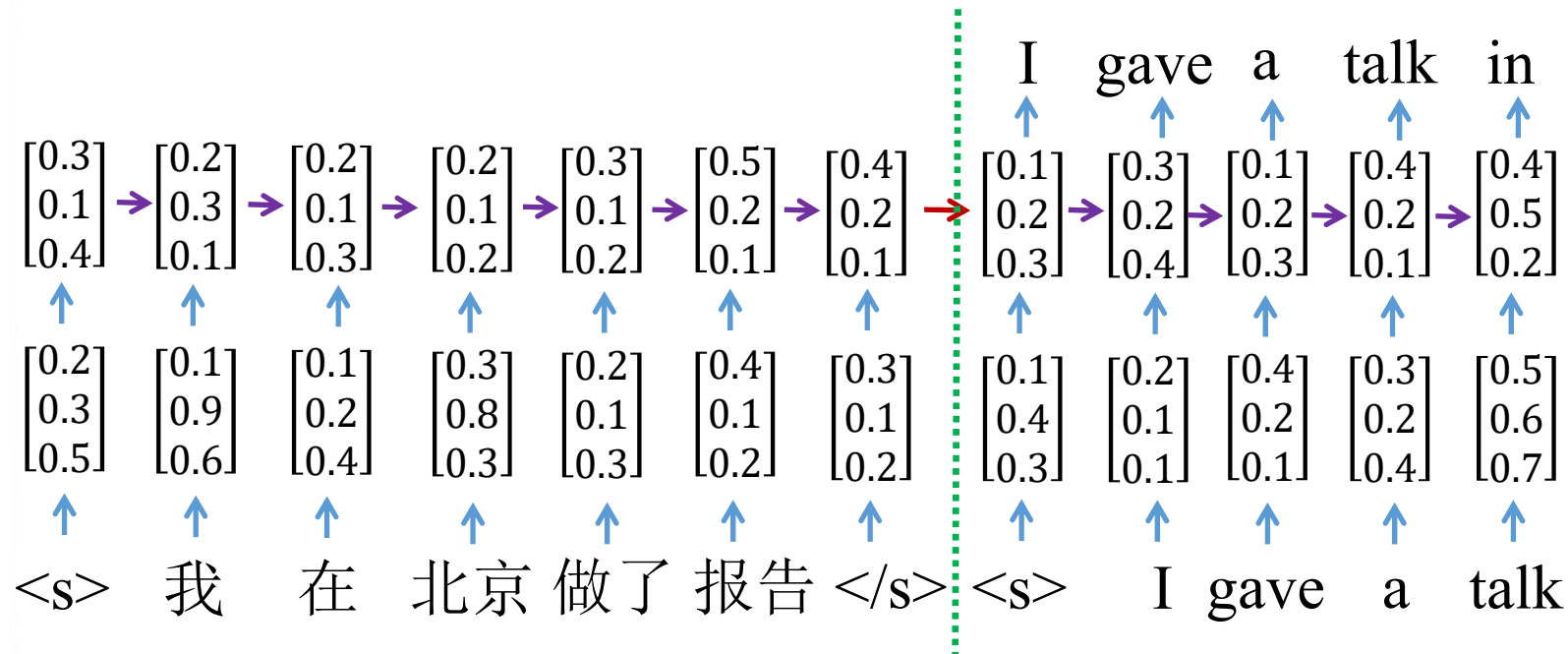
$$\max \log p(\text{talk})$$



最大化  $P(\text{target}|\text{source})$

# 神经机器翻译参数优化

$$\max \log p(in)$$

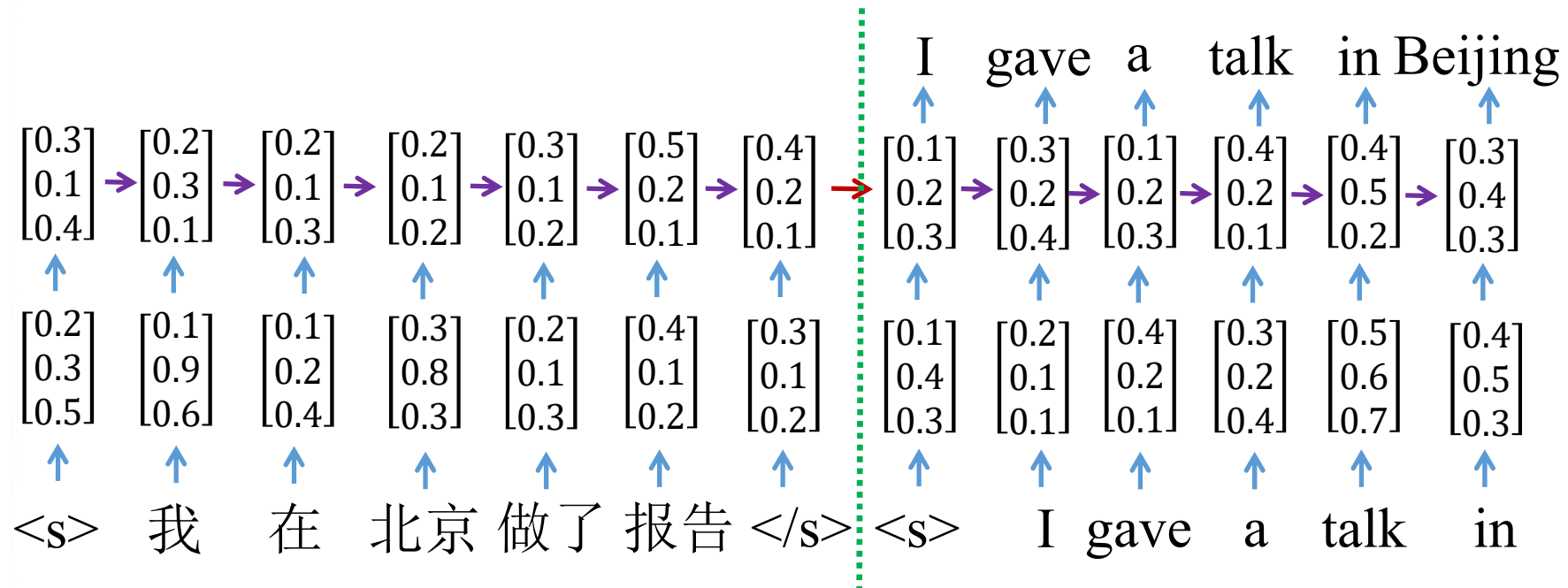


最大化  $P(\text{target}|\text{source})$



# 神经机器翻译参数优化

$$\max \log p(\text{Beijing})$$

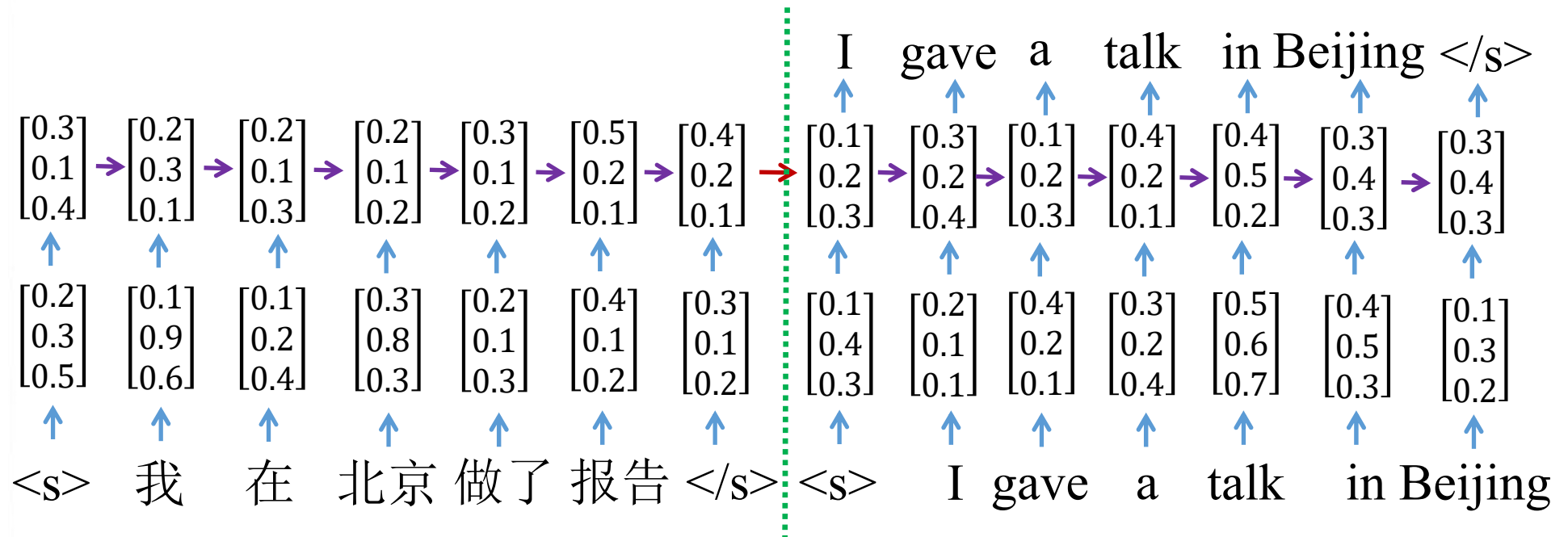


最大化  $P(\text{target}|\text{source})$



# 神经机器翻译参数优化

$$\max \log p(\langle /s \rangle)$$

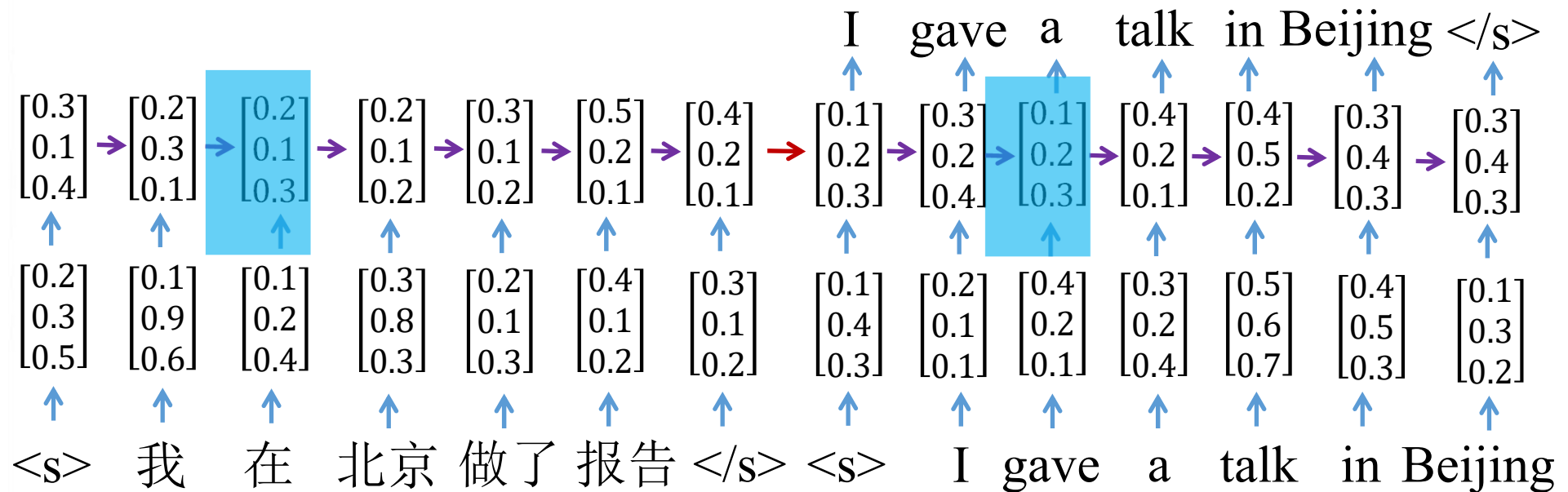


最大化  $P(\textit{target}|\textit{source})$

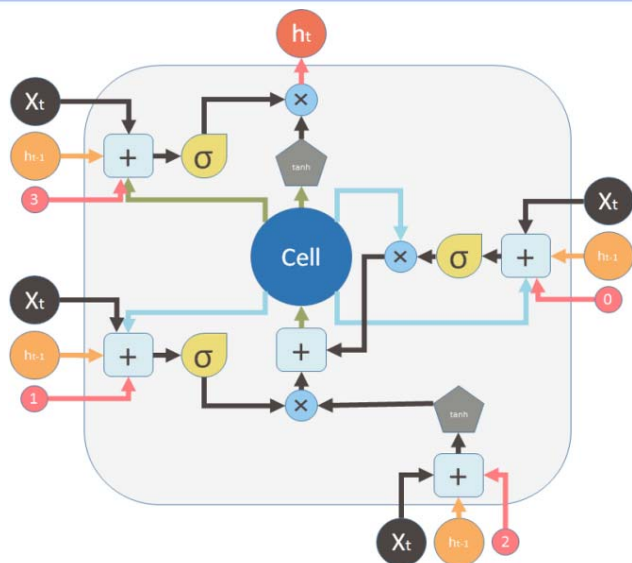


# 神经机器翻译-计算单元

$$h_t = \tanh(U_S L(w_t) + W_S h_{t-1}) \quad h_t = \tanh(U_T L(w_t) + W_T h_{t-1})$$

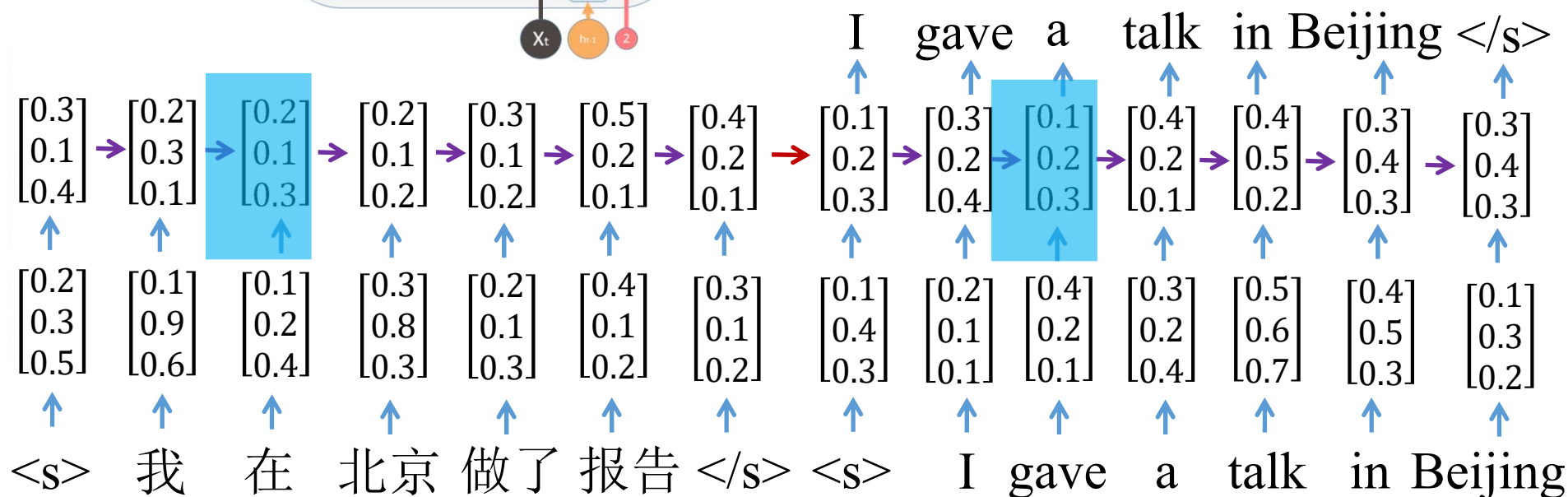


# 神经机器翻译-计算单元



## LSTM计算单元

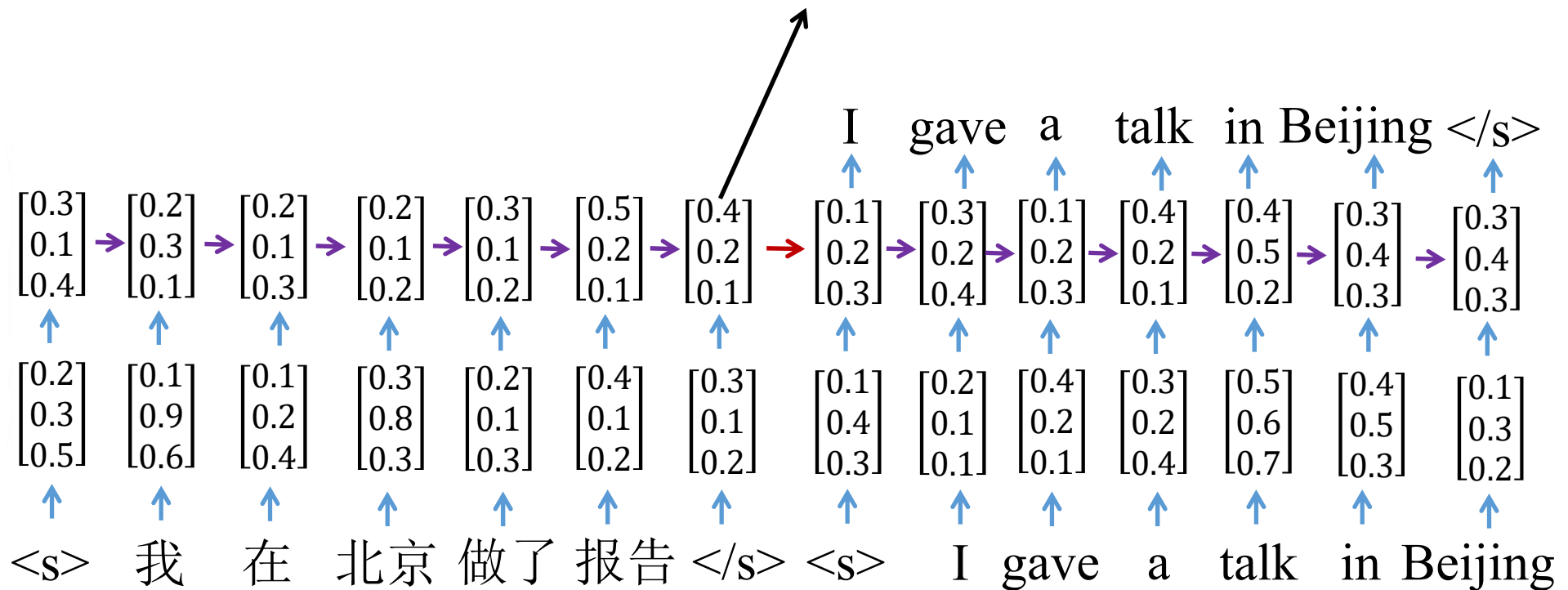
$$h_t = LSTM(w_t, h_{t-1})$$



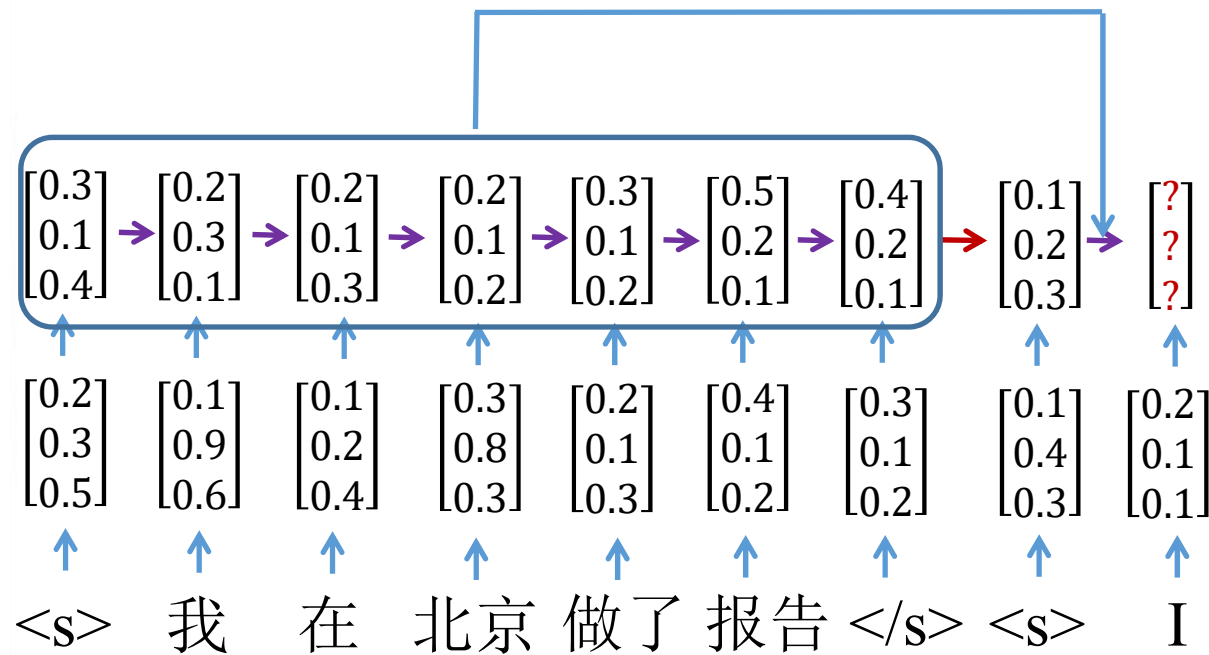


# 神经机器翻译-源端表示问题

一个实数向量无法表示  
源语言句子的完整语义

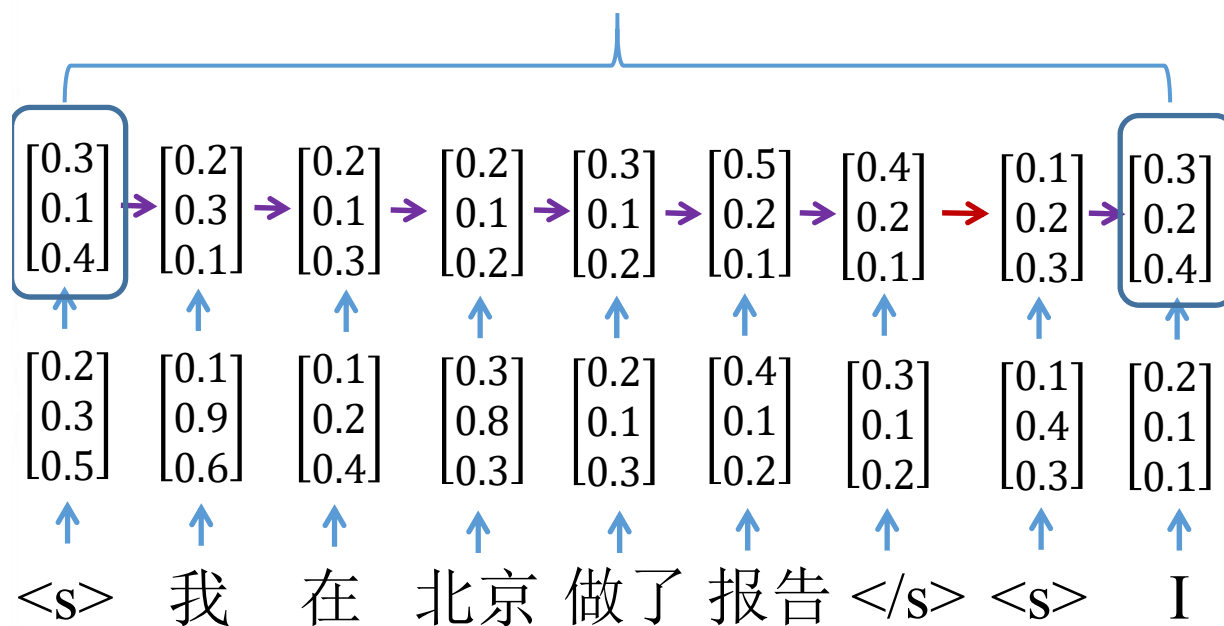


# 神经机器翻译-注意机制



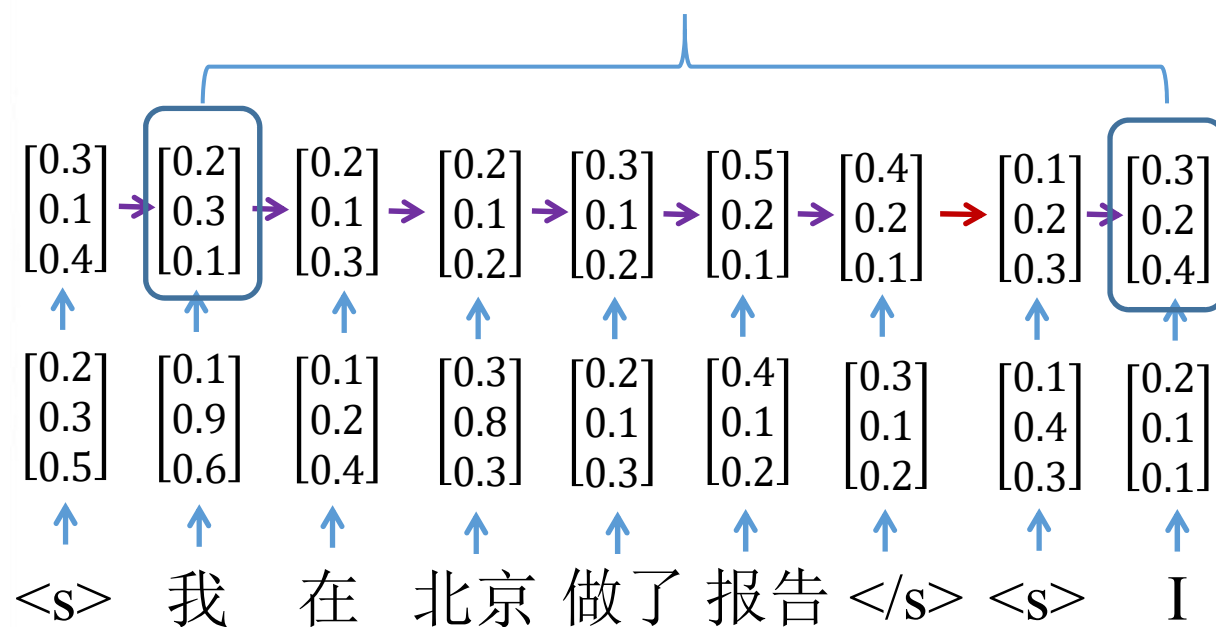
# 神经机器翻译-注意机制

$$score(h_s, h_t) = 1$$



# 神经机器翻译-注意机制

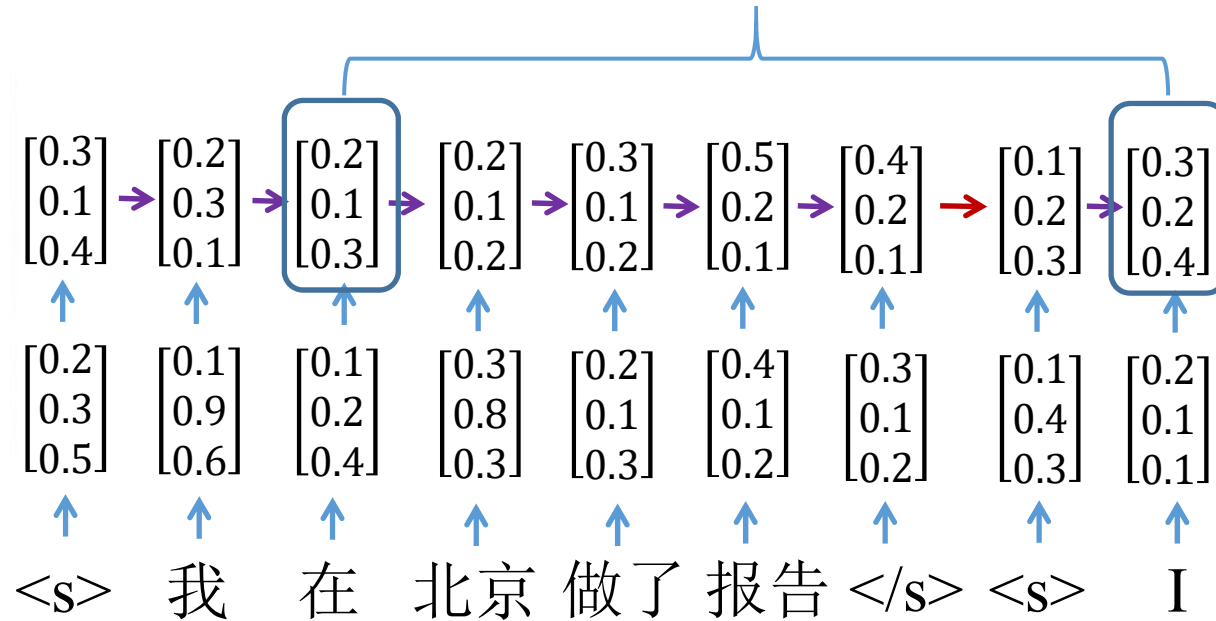
$$score(h_s, h_t) = 1$$





# 神经机器翻译-注意机制

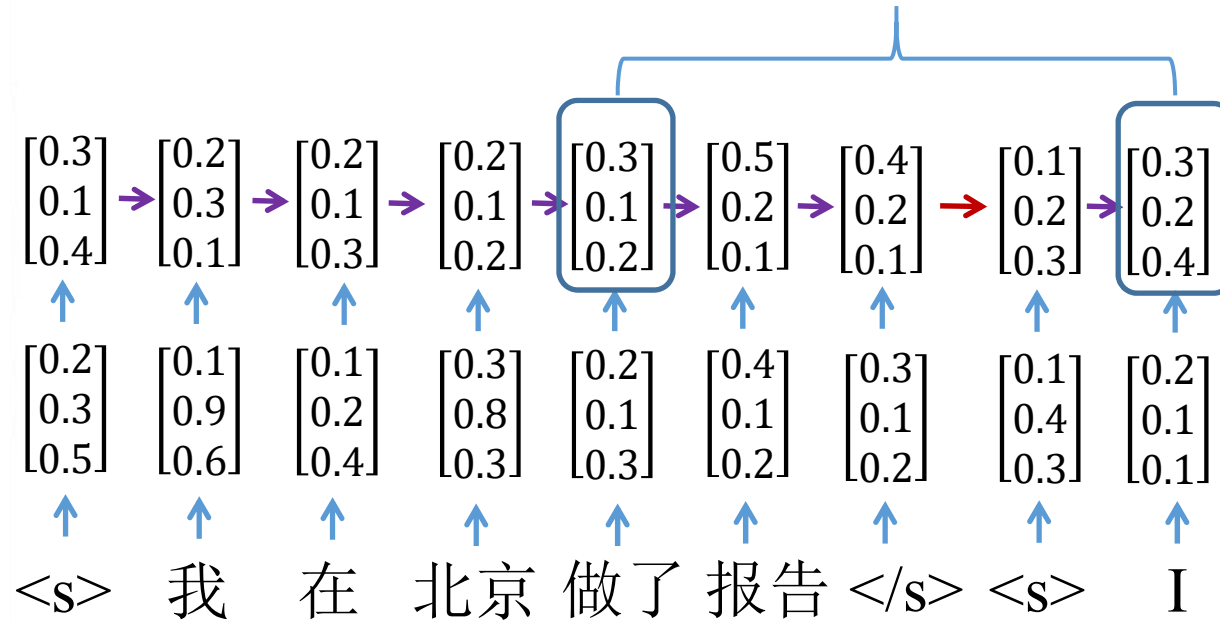
$$\text{score}(h_s, h_t) = 1$$





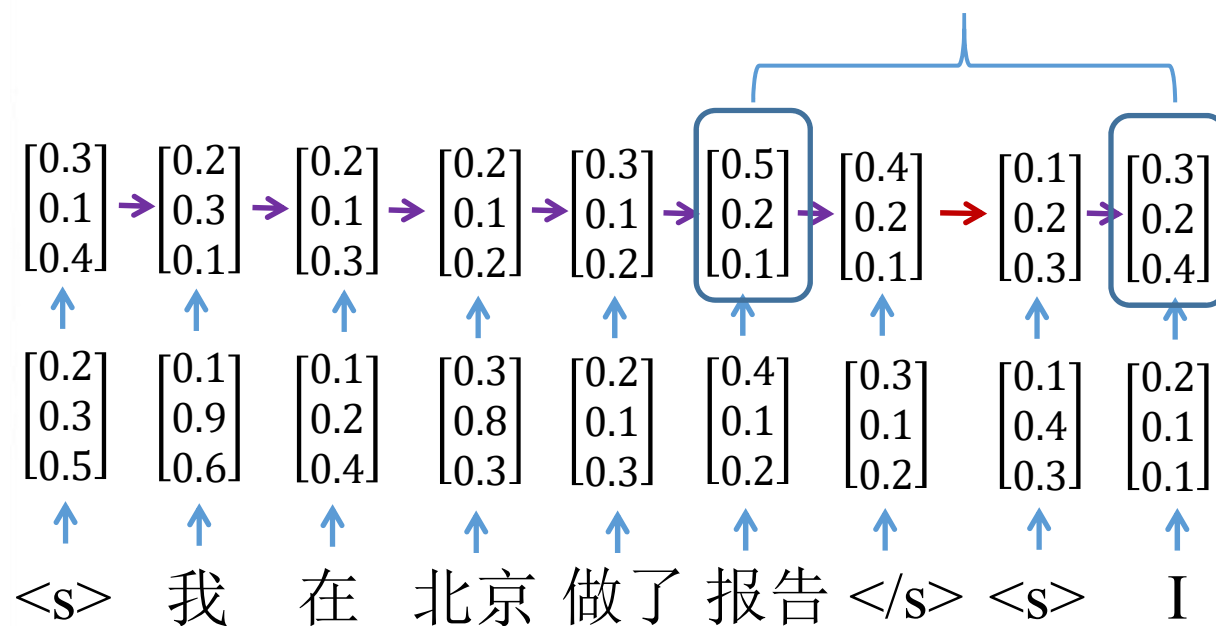
# 神经机器翻译-注意机制

$$\text{score}(h_s, h_t) = 4$$



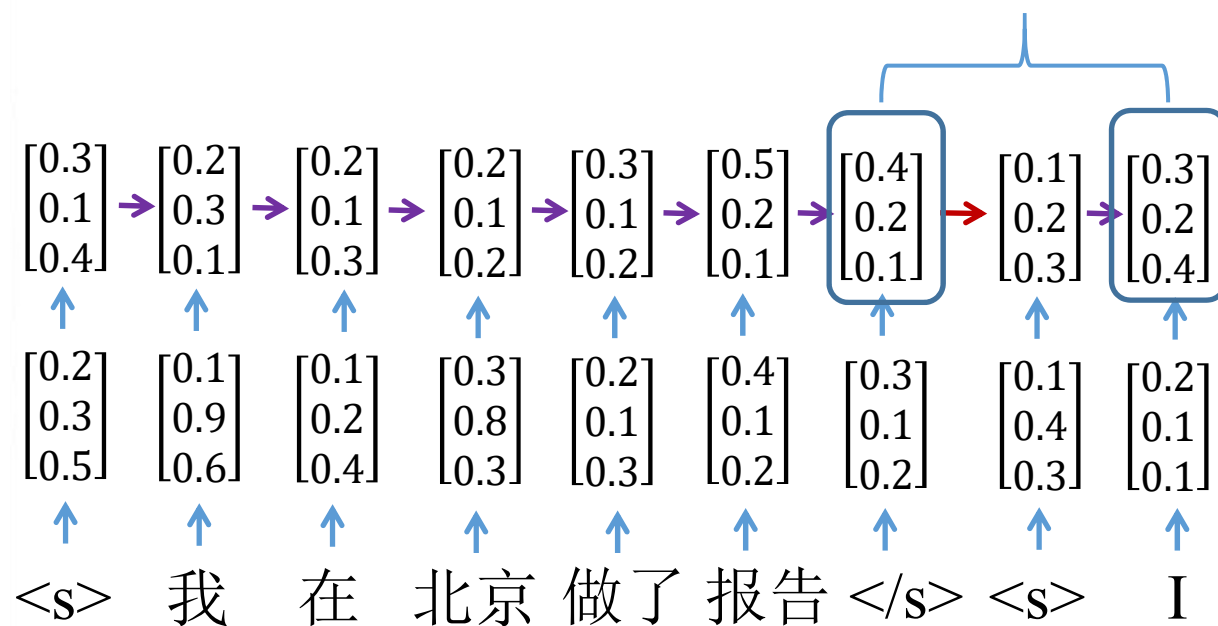
# 神经机器翻译-注意机制

$$score(h_s, h_t) = 2$$



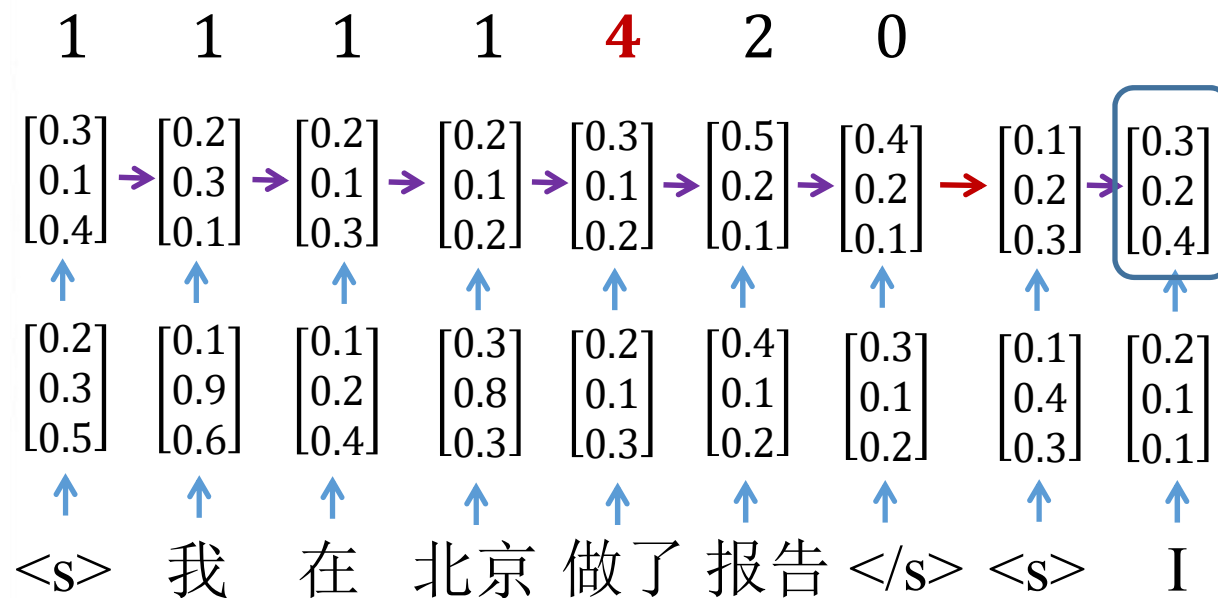
# 神经机器翻译-注意机制

$$score(h_s, h_t) = 0$$

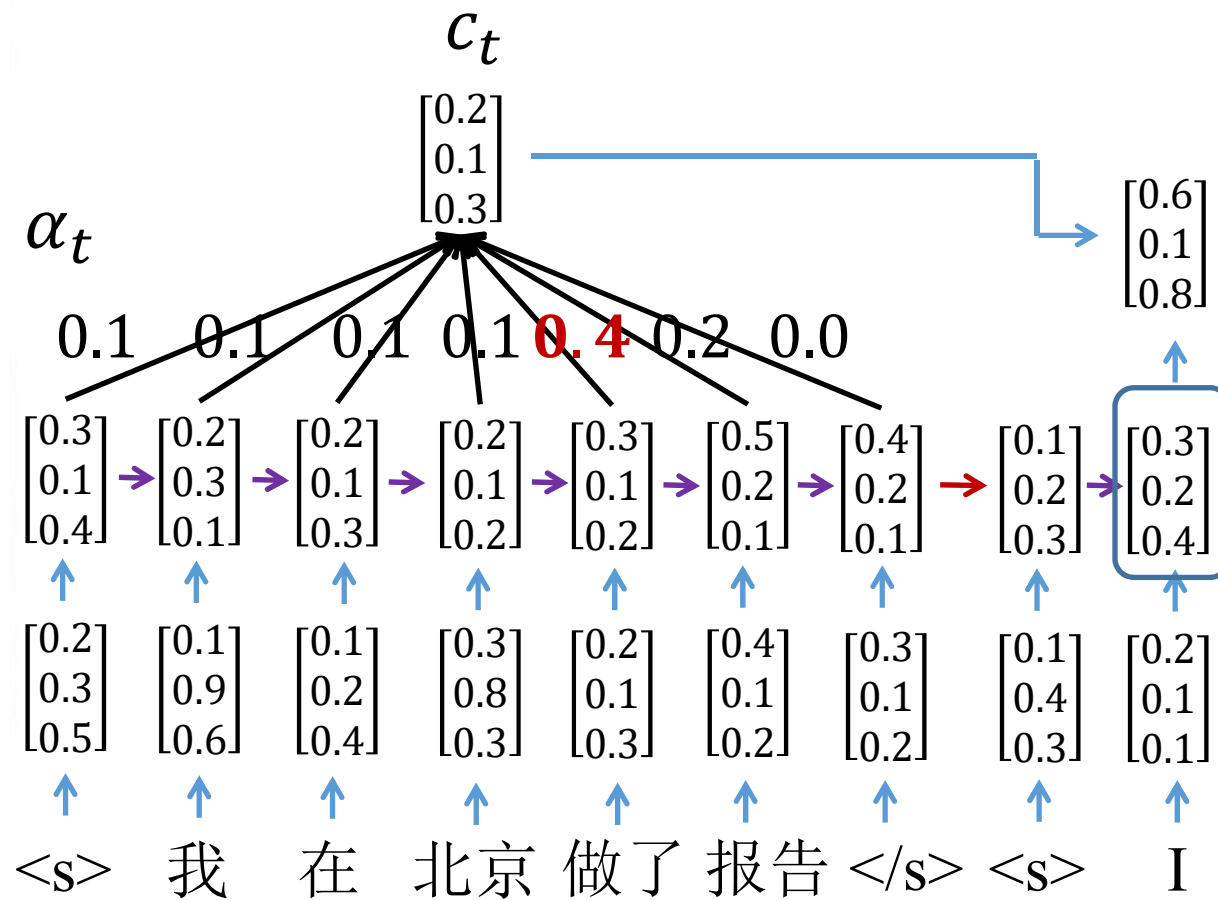


# 神经机器翻译-注意机制

$$\text{score}(h_s, h_t)$$

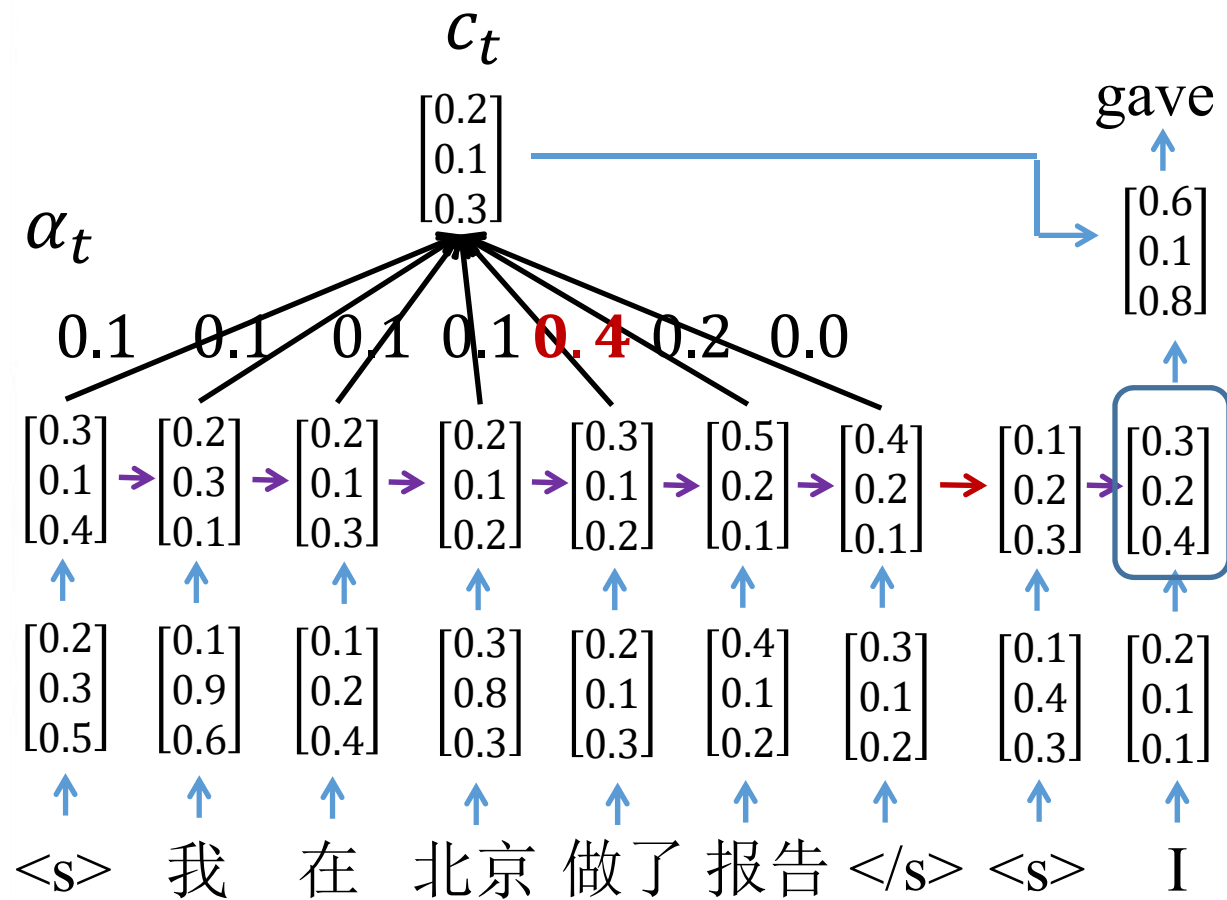
 $\alpha_t$ 


# 神经机器翻译-注意机制

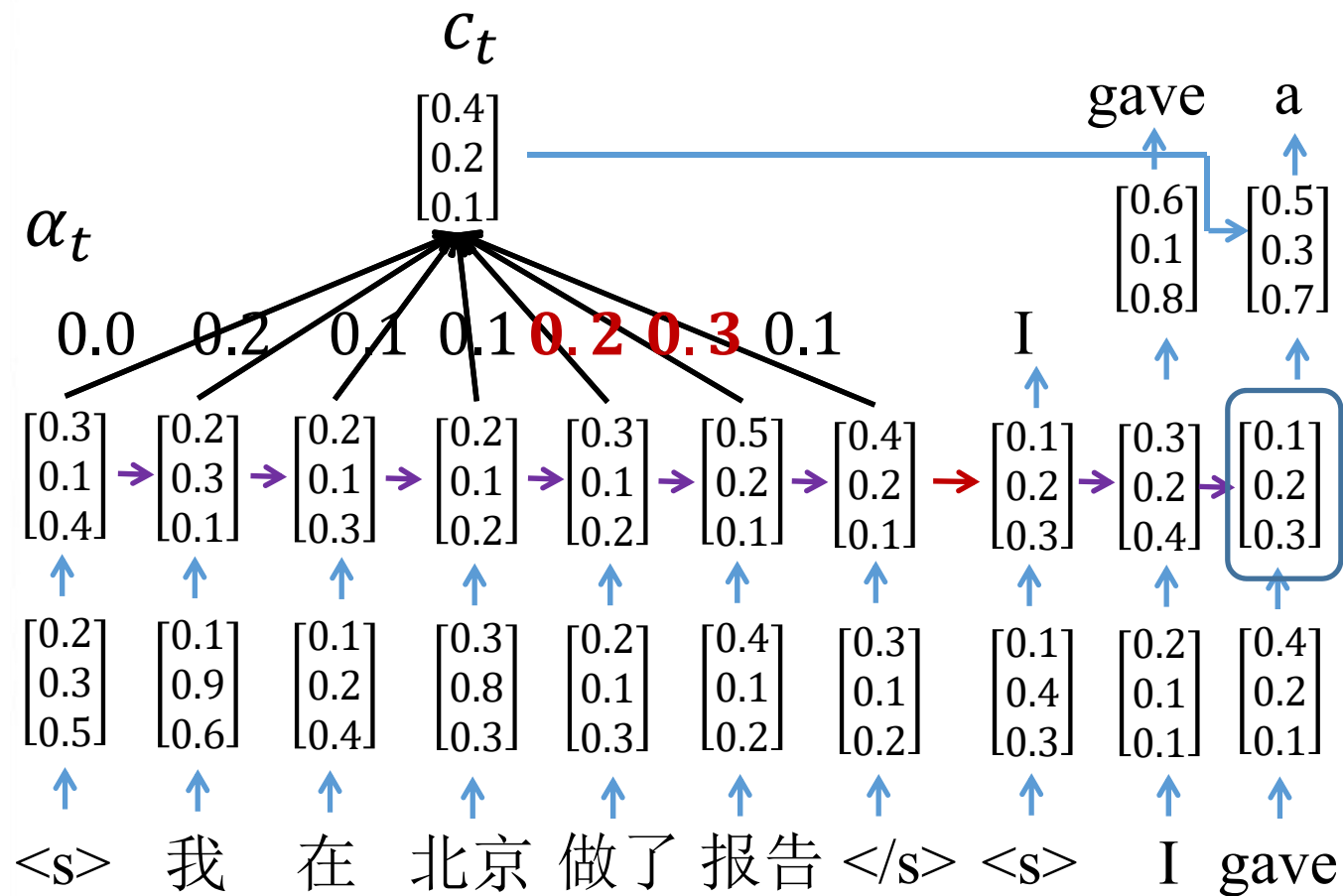




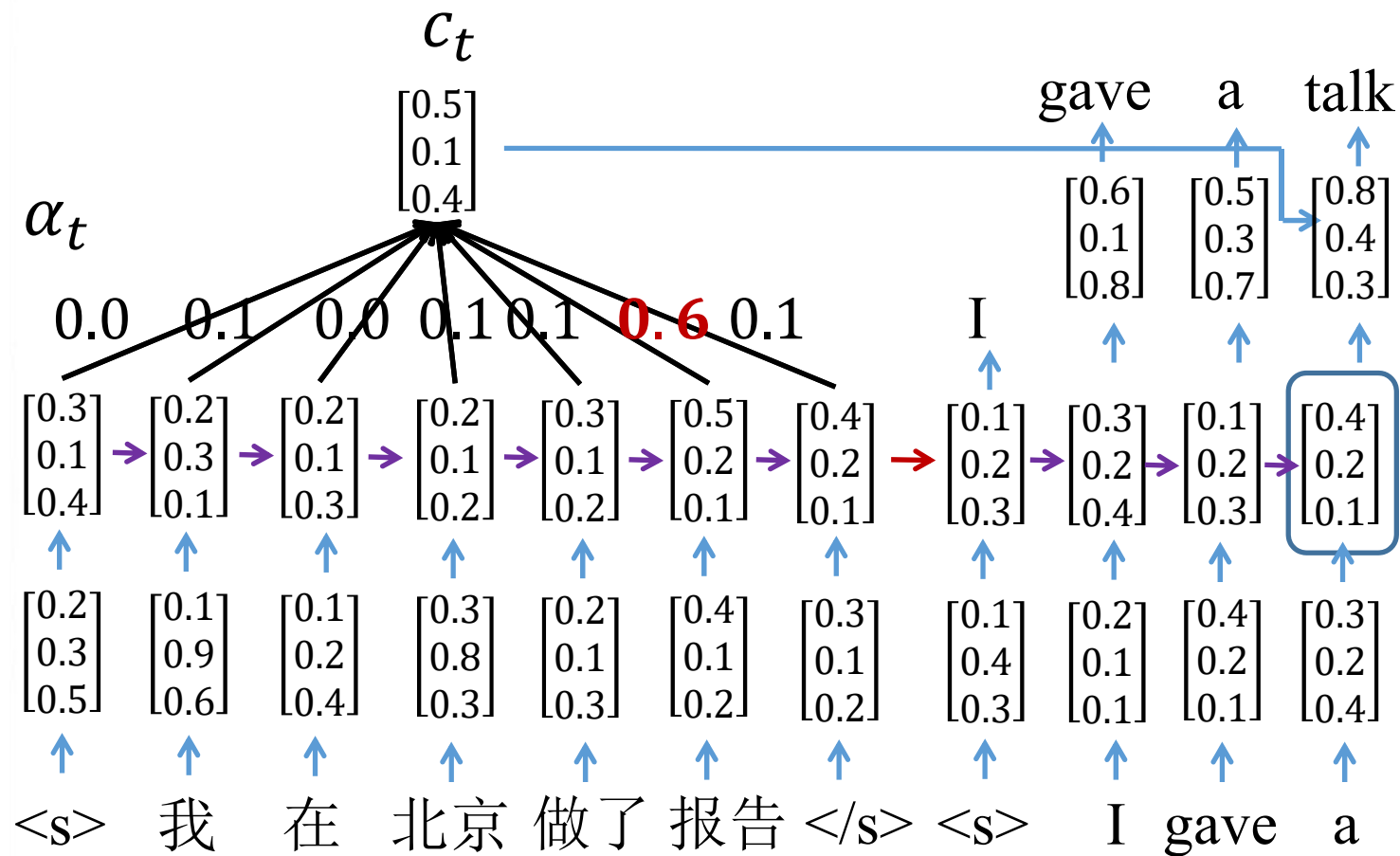
# 神经机器翻译-注意机制



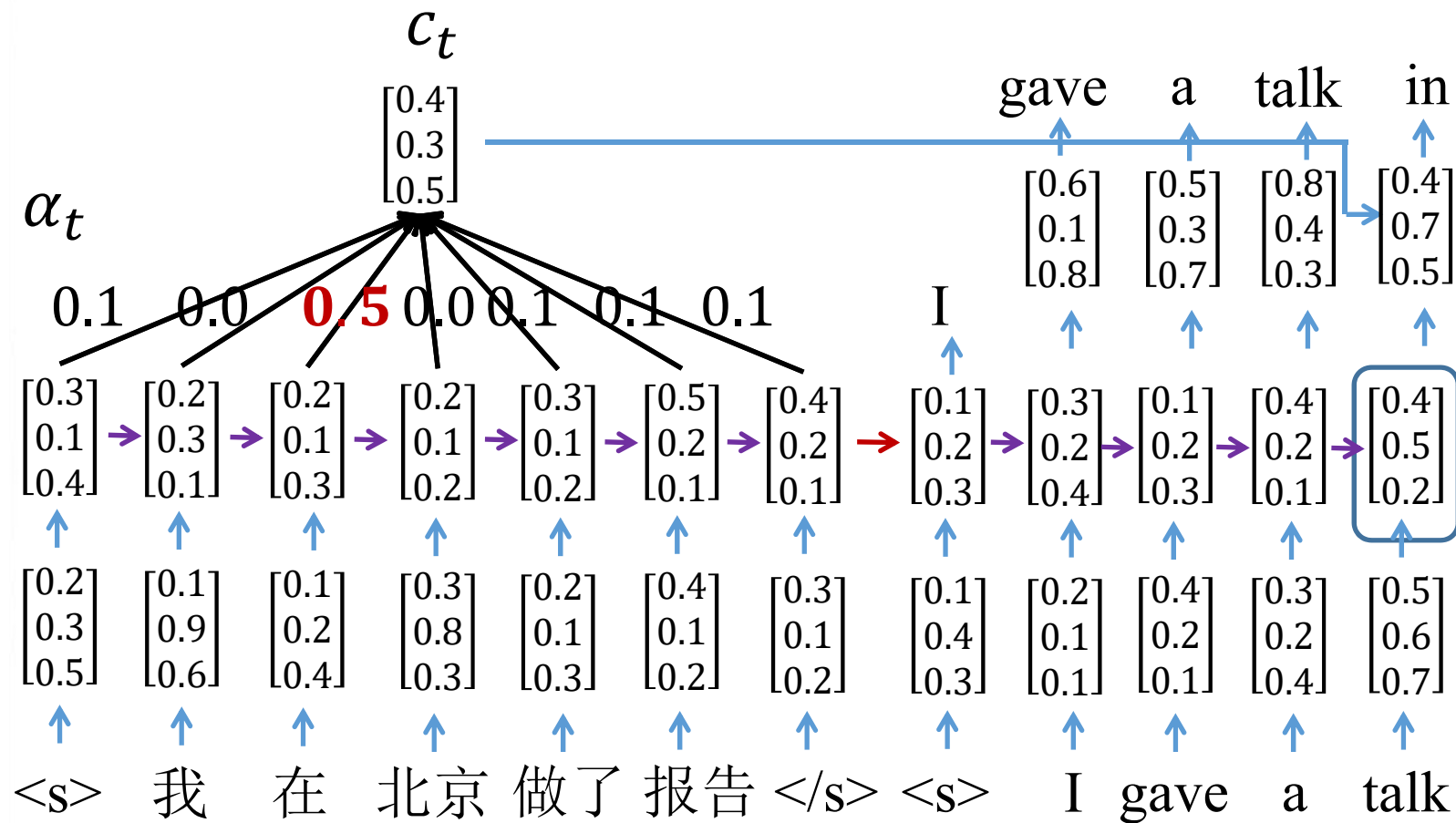
# 神经机器翻译-注意机制



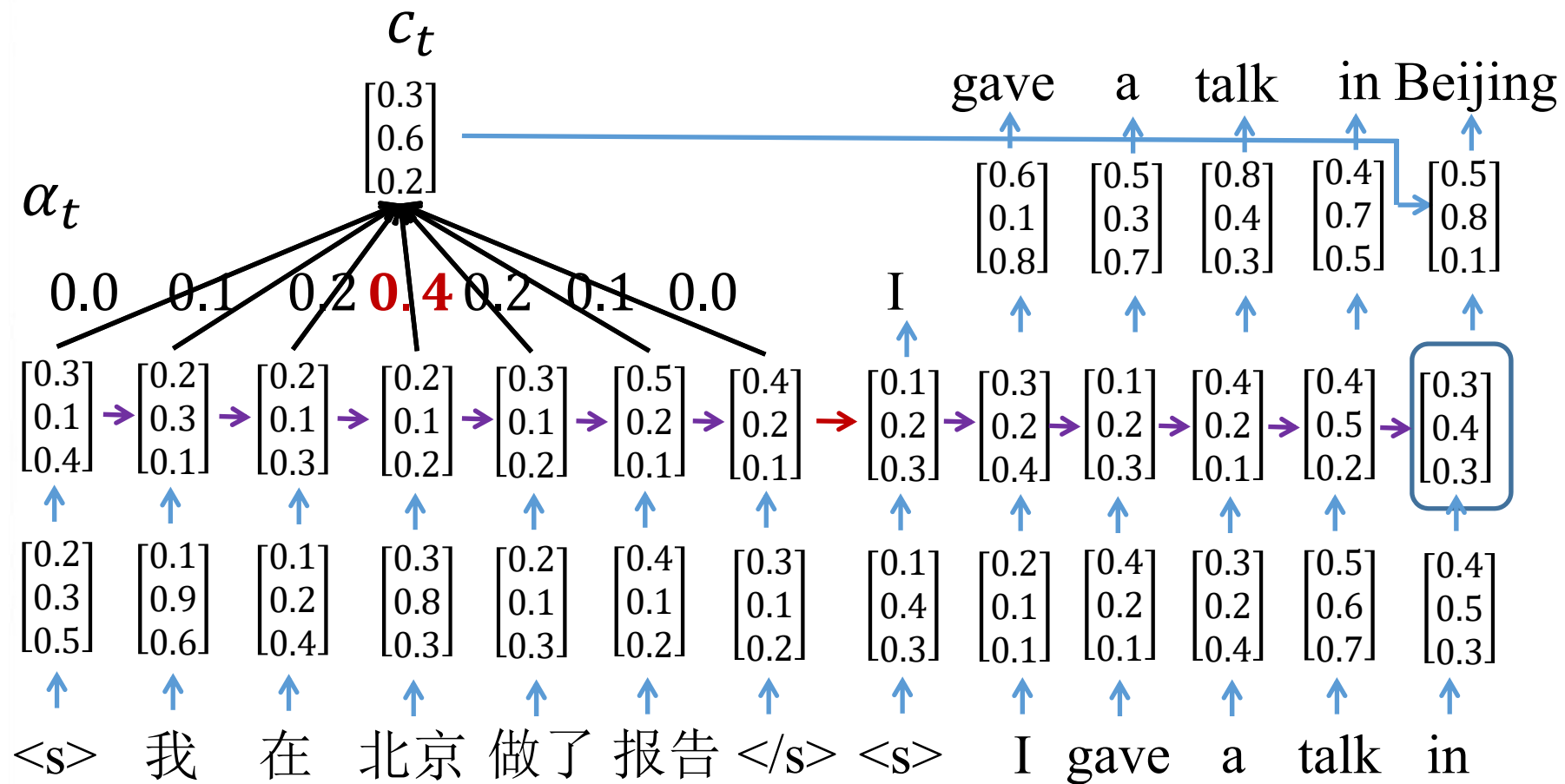
# 神经机器翻译-注意机制



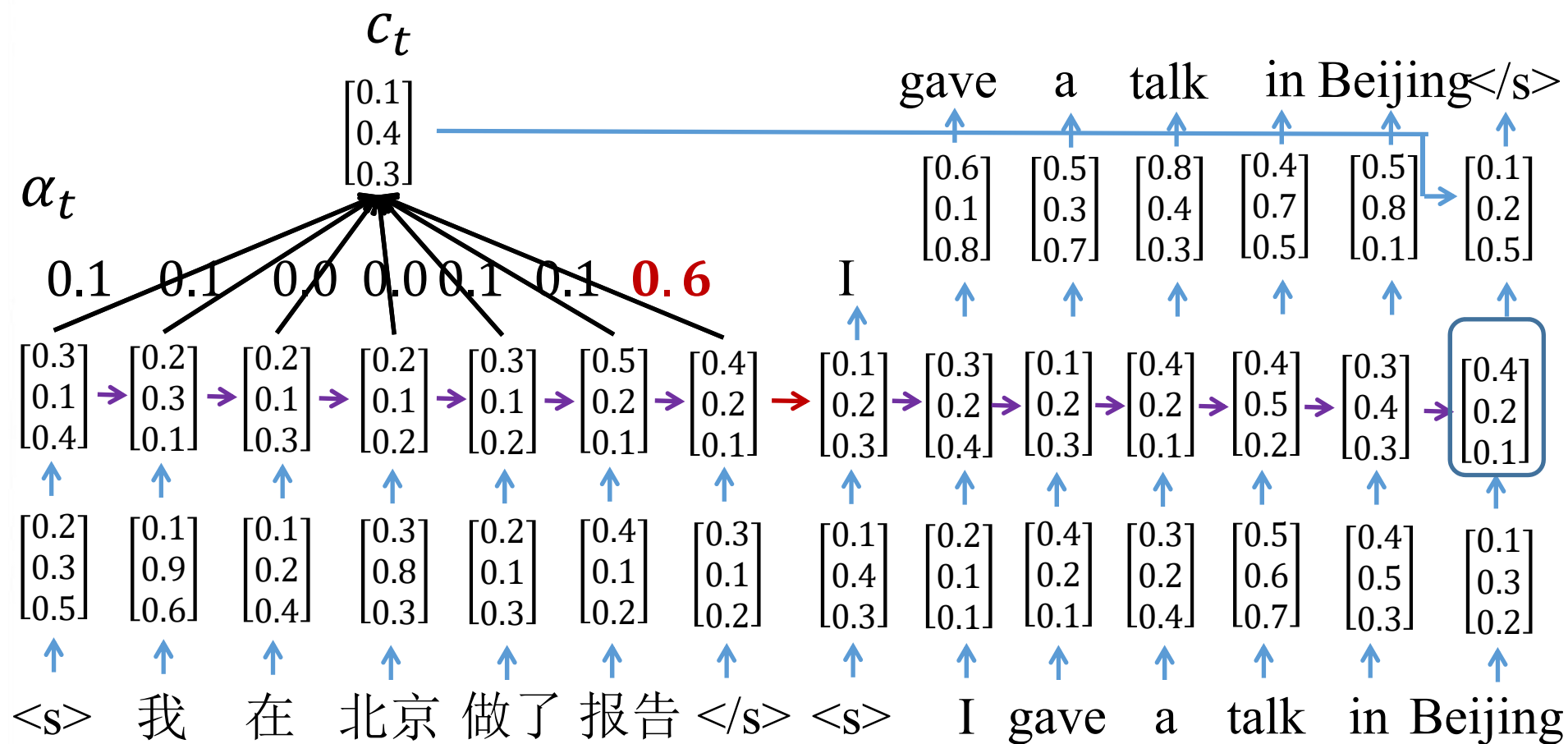
# 神经机器翻译-注意机制



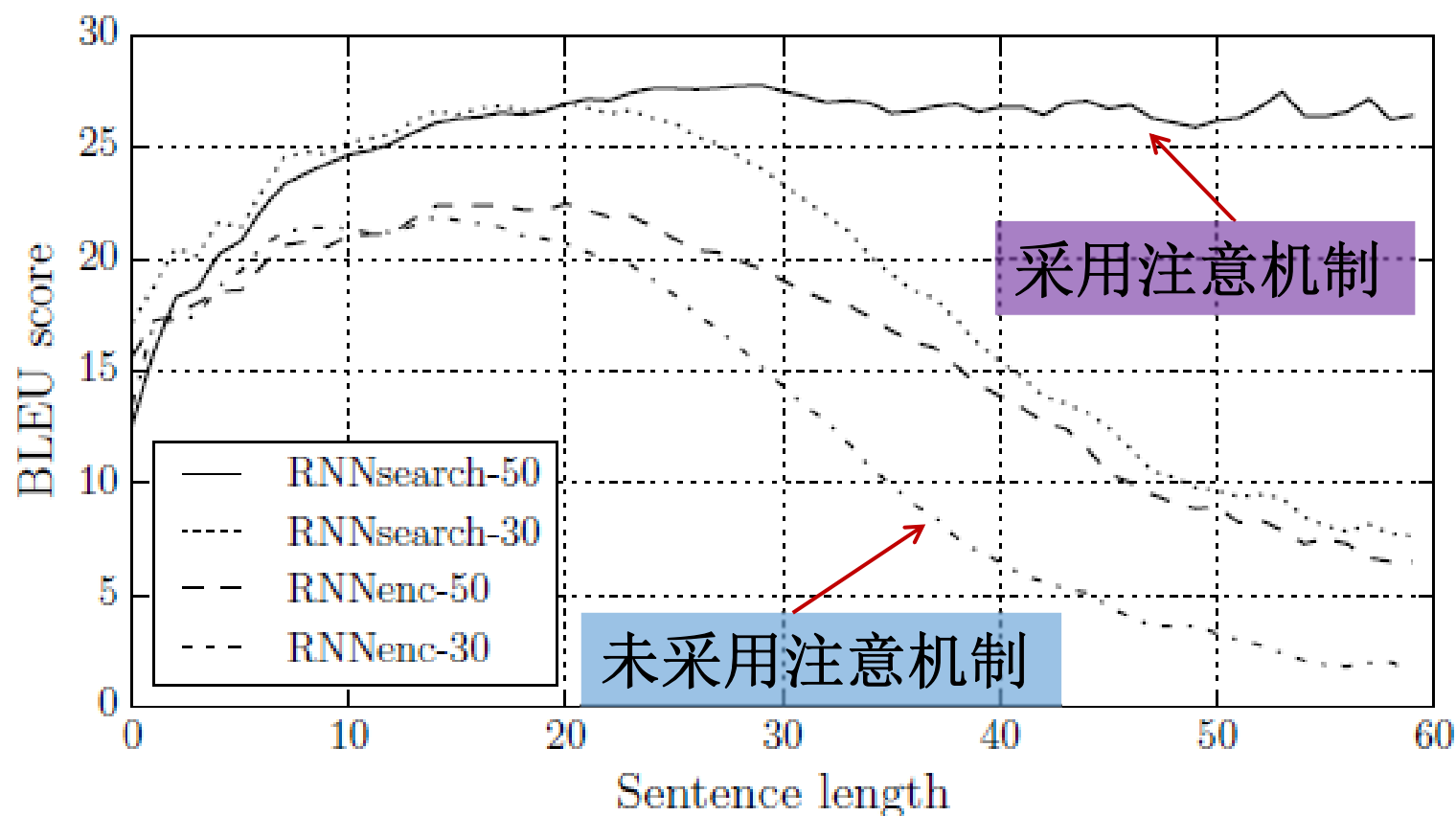
# 神经机器翻译-注意机制



# 神经机器翻译-注意机制



# 神经机器翻译-注意机制



**RNNenc:** 无注意机制, **RNNsearch:** 采用注意机制

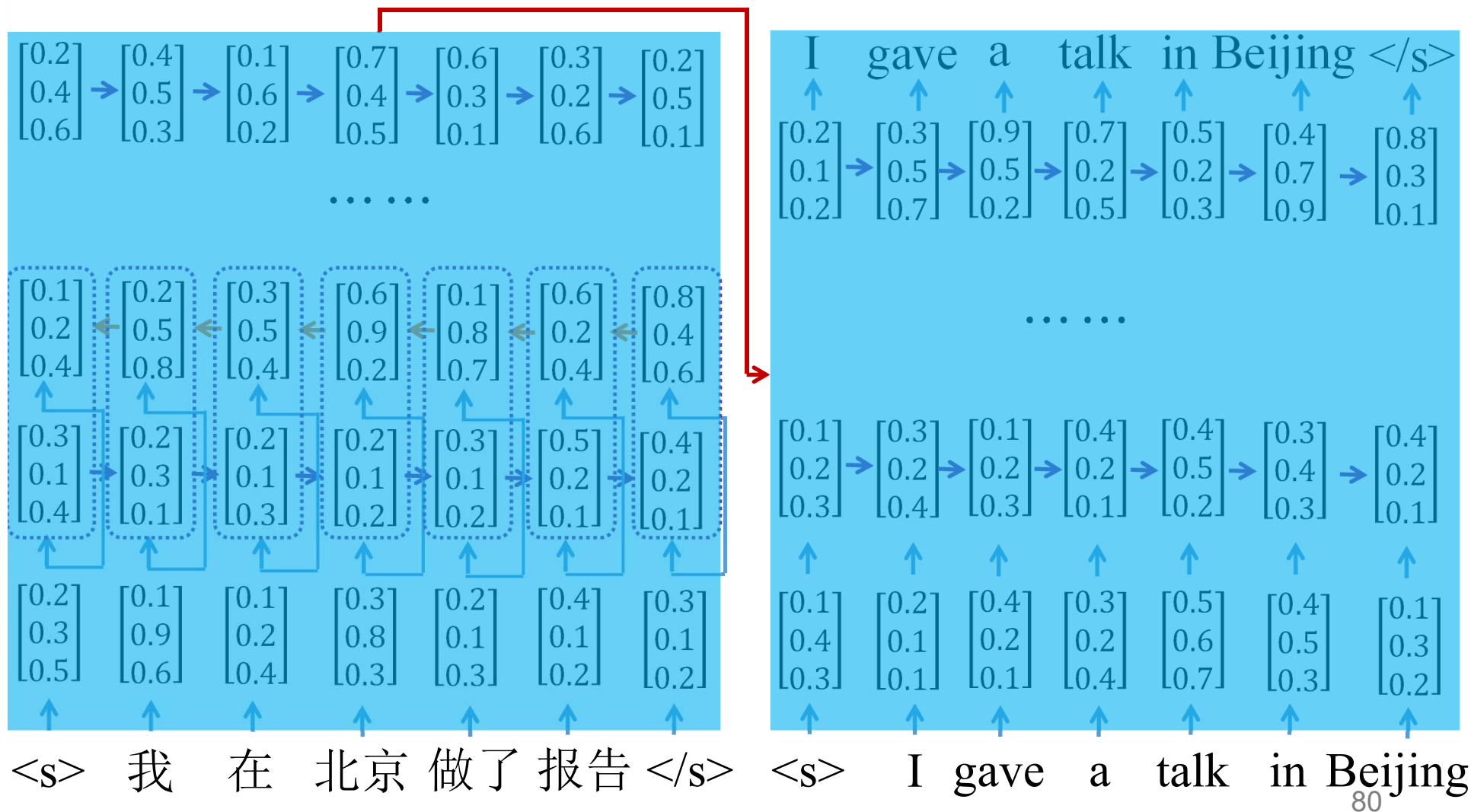


# 翻译实例

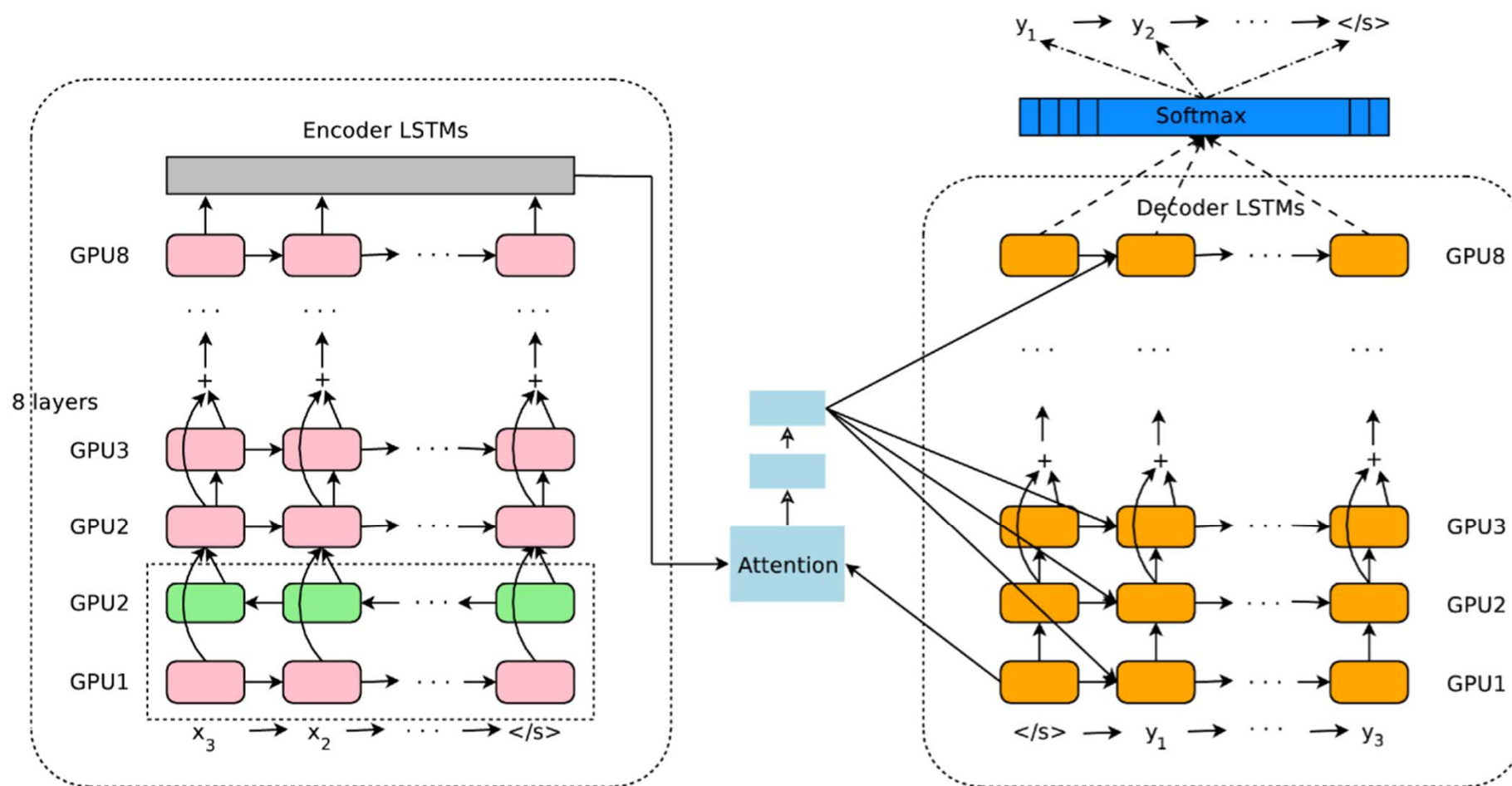
	south	korean	envoy	calls	for	dialogue	between	the	united	states	and	north	korea	.
南韩	Black	Black	White	White	White	White	White	White	White	White	White	White	White	White
特使	White	White	Black	White	White	White	White	White	White	White	White	White	White	White
呼吁	White	White	White	Black	Gray	White	White	White	White	White	White	White	White	Gray
美国	White	White	White	White	Gray	Gray	White	Black	Black	Black	White	White	White	White
与	White	White	White	White	White	Gray	Gray	White	White	White	Gray	White	White	Gray
北韩	White	White	White	White	White	White	White	White	White	White	White	Black	Black	White
对话	White	White	White	White	Gray	Gray	Gray	White	White	White	Gray	White	White	Gray

# 神经机器翻译-隐藏层结构

## 注意机制



# 神经机器翻译-隐藏层结构



**GNMT:** 谷歌神经翻译系统



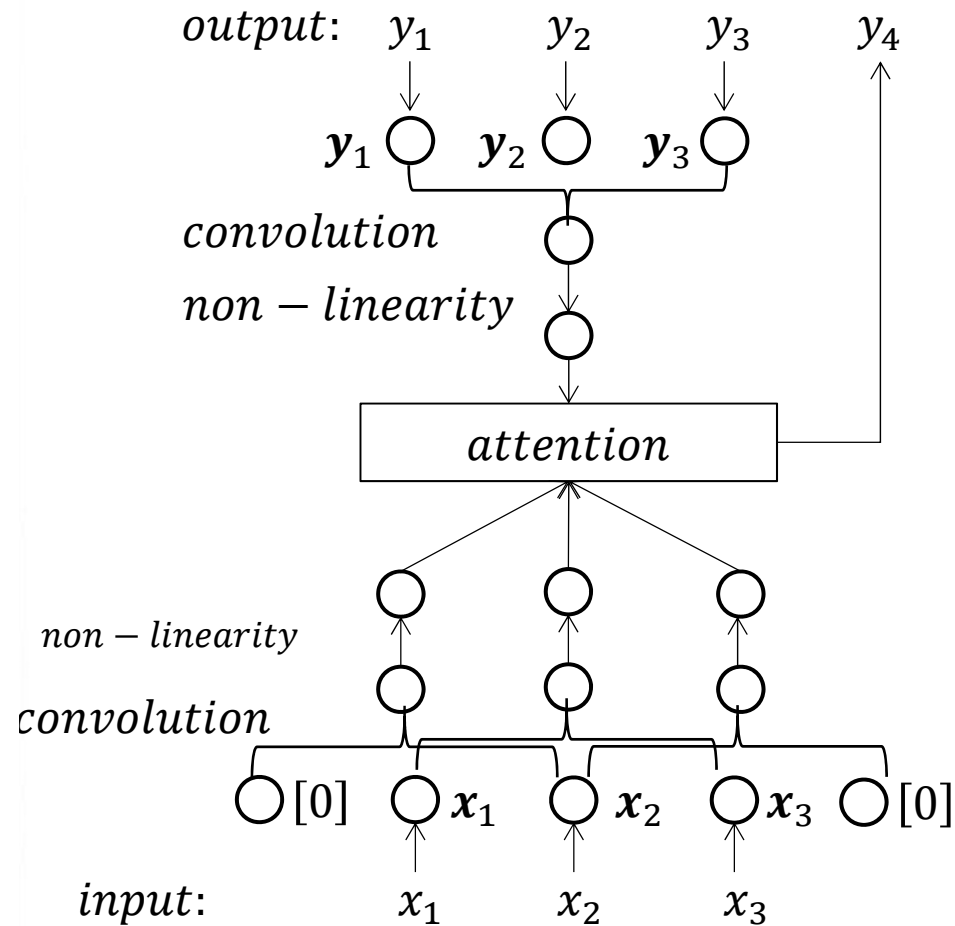
# 神经机器翻译 VS. 统计机器翻译

	PBMT	GNMT	Human	Relative Improvement
English → Spanish	4.885	5.428	5.550	87%
English → French	4.932	5.295	5.496	64%
English → Chinese	4.925	4.594	4.987	58%
Spanish → English	4.872	5.187	5.372	63%
French → English	5.046	5.343	5.404	83%
Chinese → English	3.694	4.263	4.636	60%

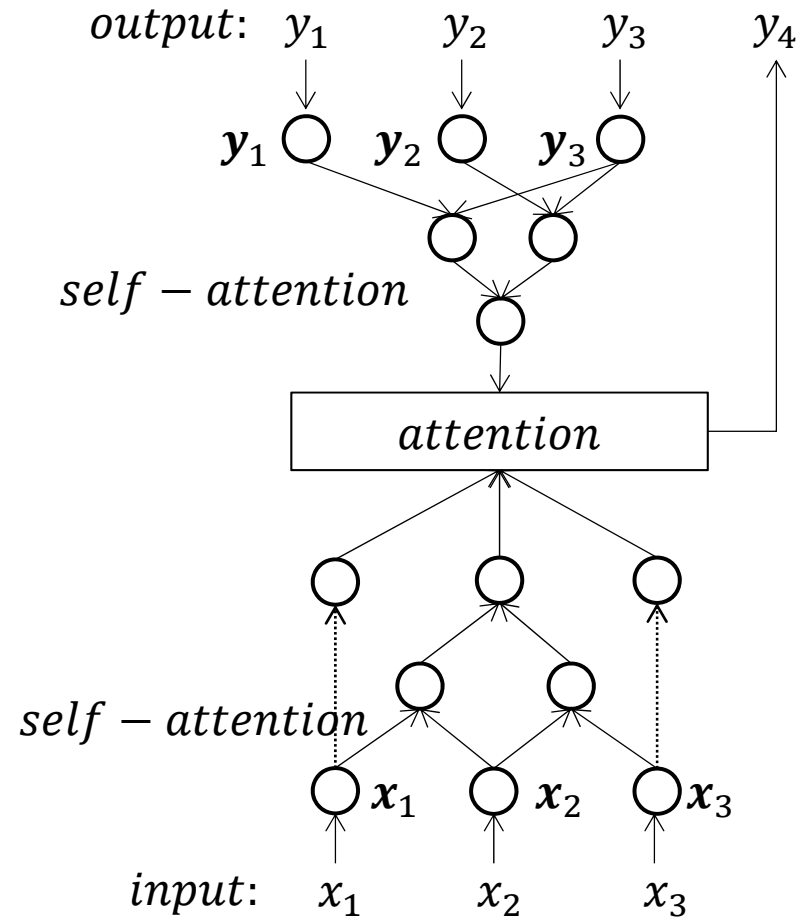
人工评测提升显著！

**GNMT:** 谷歌神经翻译系统

# 其他编码-解码架构



(a) 基于卷积神经网络的翻译模型



(b) 基于纯注意机制的翻译模型



# 两类机器翻译的结果对比

Source: 海珊 也 与 恐怖 组织网 建立 了 联系 。

Reference: Hussein has also established ties with terrorist networks.

PBMT: Hussein also has established relations **and terrorist group** .

HPMT: Hussein **also and** terrorist group established relations .

NMT: Hussein also established relations with **<UNK>** .

模型融合!

- 如何博采众长 ?



# 多粒度融合方法

- 词语级融合 (Arthur et al., 2016; Feng et al., 2017)
  - 以统计机器翻译中的词汇翻译概率为外部知识，影响神经机器翻译的译文选择和输出
- 短语级融合 (Tang et al., 2016; Wang et al., 2017; Zhao et al., 2018)
  - 将统计模型的短语翻译规则作为外部知识，影响和指导神经机器翻译的译文选择
- 句子级融合 (Niehues et al., 2016; Zhou et al., 2017)
  - 以统计机器翻译的译文和神经机器翻译的译文为多源输入，输出兼顾两者优势的译文





# 词汇级融合-NMT利用SMT词汇翻译表

- 动机：NMT对低频词的翻译存在大量错翻

原因：NMT参数的“共享”机制（Feng et al., 2017）

目的：减少参数个数，提高训练效果

副作用：1、高频词有更多地机会去“优化”参数

2、参数对高频词的预测较好，而对低频词则较差

例子：

**Source:** 阿尔卡特 宣称 第四季 销售 成长 近 30%

**Reference:** alcatel says sales in fourth quarter last year grew nearly 30 %

**NMT:** he said sales grew nearly 30 percent in fourth quarter of last year



# 词汇级融合-NMT利用SMT词汇翻译表

- 为了缓解NMT的对低频词的错翻，在解码过程融合SMT的词汇翻译表 (Arthur et al., 2016; Feng et al., 2017)
- 原因：SMT的离散表示和非共享机制，使得SMT对低频单词的预测存在很大优势。例如： $p(\text{alcatel} | \text{阿尔卡特}) = 0.5732$
- Step 1：利用SMT得到一个词汇翻译表
- Step 2：根据词汇翻译表以及输入句子 $X$ ，得到如下矩阵

$$L = \begin{bmatrix} p_l(e=1 | x_1) & \dots & p_l(e=1 | x_{Tx}) \\ \dots & \dots & \dots \\ p_l(e=V_e | x_1) & \dots & p_l(e=V_e | x_{Tx}) \end{bmatrix}$$

列：源语言句子单词的个数  
行：所有的目标语言单词的个数



# 词汇级融合-NMT利用SMT词汇翻译表

- Step 1 : 利用SMT得到一个词汇翻译表
- Step 2 : 根据词汇翻译表以及输入句子X，得到如下矩阵

$$L = \begin{bmatrix} p_l(e=1 | x_1) & \dots & p_l(e=1 | x_{Tx}) \\ \dots & \dots & \dots \\ p_l(e=V_e | x_1) & \dots & p_l(e=V_e | x_{Tx}) \end{bmatrix}$$

列：源语言句子单词的个数  
行：所有的目标语言单词的个数

假设：源语言为“我爱中国”

目标语言一共5个单词：<start>、<end>、I、love、china



$$L = \begin{bmatrix} p_l(start | 我) & p_l(start | 爱) & p_l(e=start | 中国) \\ p_l(end | 我) & p_l(end | 爱) & p_l(e=end | 中国) \\ p_l(I | 我) & p_l(I | 爱) & p_l(e=I | 中国) \\ p_l(love | 我) & p_l(love | 爱) & p_l(e=love | 中国) \\ p_l(china | 我) & p_l(china | 爱) & p_l(e=china | 中国) \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$



# 词汇级融合-NMT利用SMT词汇翻译表

➤ Step 3 : 将矩阵  $L$  与 attention weight 相乘 , 得到如下 SMT 的预测结果

$$P_{SMT} = L * a = \begin{bmatrix} p_l(e=1 | x_1) & \dots & p_l(e=1 | x_{Tx}) \\ \dots & \dots & \dots \\ p_l(e=V_e | x_1) & \dots & p_l(e=V_e | x_{Tx}) \end{bmatrix} \begin{bmatrix} a_{i,1} \\ \dots \\ a_{i,Tx} \end{bmatrix}$$

$$a_{11} = 0.9$$

假设 : 预测第一个目标单词时 , attention weight 为  $a_{12} = 0.05$

$$a_{13} = 0.05$$

$$P_{SMT} = L * a = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0.9 \\ 0.05 \\ 0.05 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0.9 \\ 0.05 \\ 0.05 \end{bmatrix} \begin{matrix} \text{start} \\ \text{end} \\ \text{i} \\ \text{love} \\ \text{china} \end{matrix}$$



# 词汇级融合-NMT利用SMT词汇翻译表

➤ Step 4 : 融合SMT的结果和NMT的结果

$$P_f = \lambda P_{NMT} + (1 - \lambda) P_{SMT}$$

分析：SMT的预测结果只与两个因素有关

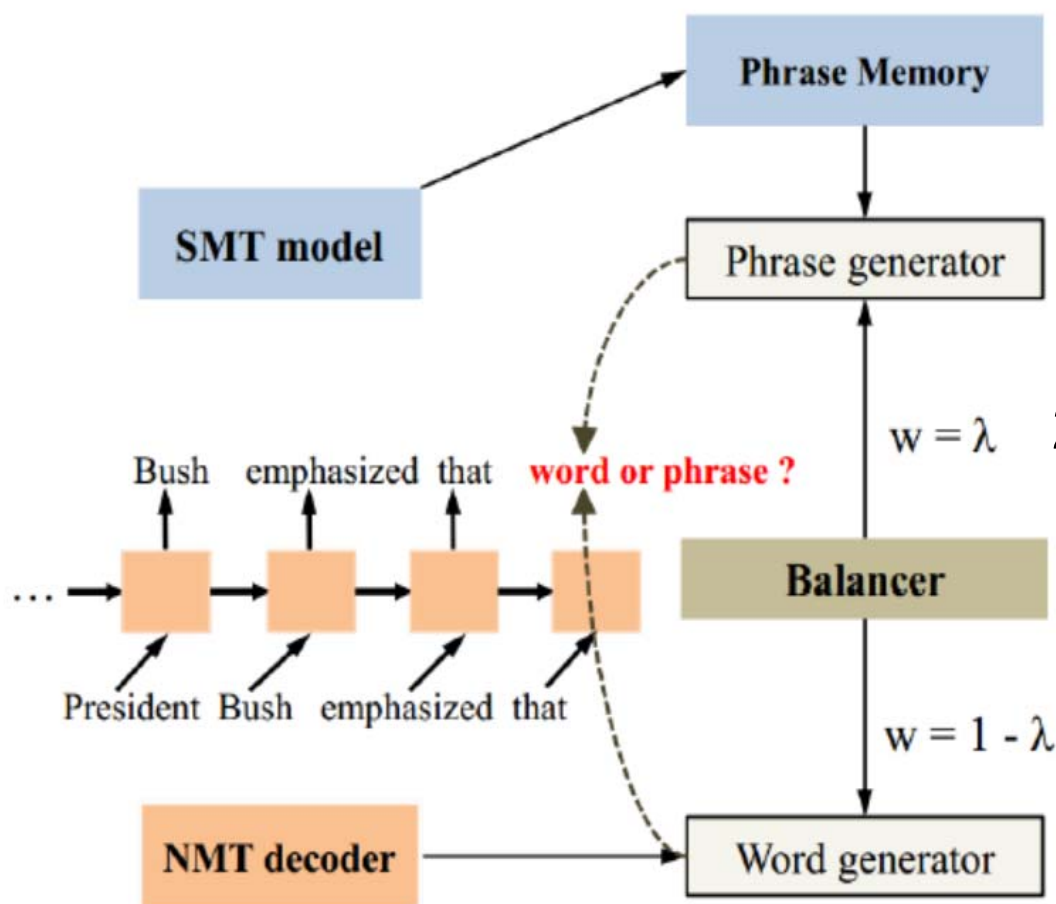
1) attention weight: 翻译系统正在翻译哪个单词

2) 翻译概率 $L$ : 单词翻译概率越大, SMT的预测结果越高



# 短语级融合-NMT利用SMT短语翻译表

方法1. 直接将短语作为翻译单元 (Tang et al., 2016; Wang et al., 2017)



1. 通过权重调节应该是翻译单词但是翻译短语
2. 权重在网络中自动学习



# 短语级融合-NMT利用SMT短语翻译表

---

方法2：短语与部分翻译结果进行匹配（Zhao et al.2018）

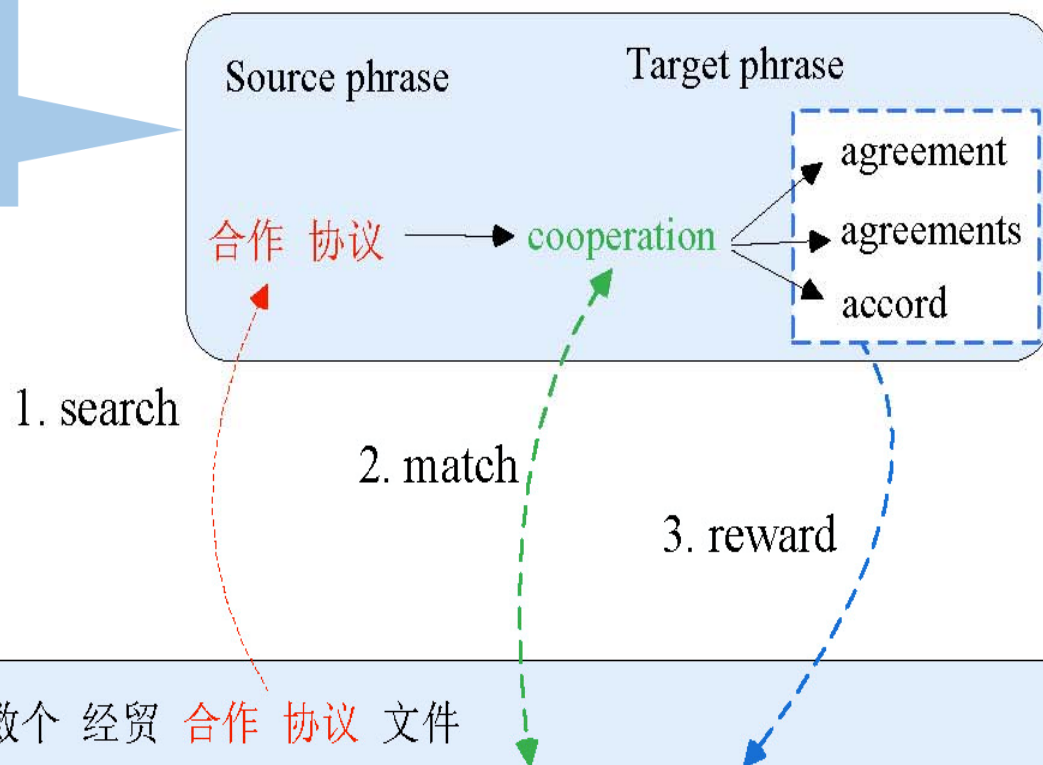
**基本思想：**如果NMT当前译文的后缀能够匹配到某个短语翻译规则的目标语言部分的前缀，那么翻译规则目标语言短语匹配部分的下一个词极有可能就是NMT应该预测的词



# 短语级融合-NMT利用SMT短语翻译表

统计机器翻译使用的翻译知识：  
**短语翻译表**

神经机器翻译解码过程中，当前译文**匹配短语候选译文前缀**时，**鼓励后缀**成为下个时刻的候选译文



Chinese Source: 双方 将 签署 数个 经贸 **合作 协议** 文件

English Reference: The two sides will sign several economic and trade **cooperation agreement** documents

# 短语级融合-NMT利用SMT短语翻译表

$$p(y_i|y_{<i}, C) = p(y_i|y_{<i}, c_i) \\ = \text{softmax} \left( \text{score}(W_{y_i}, h_i) \right)$$



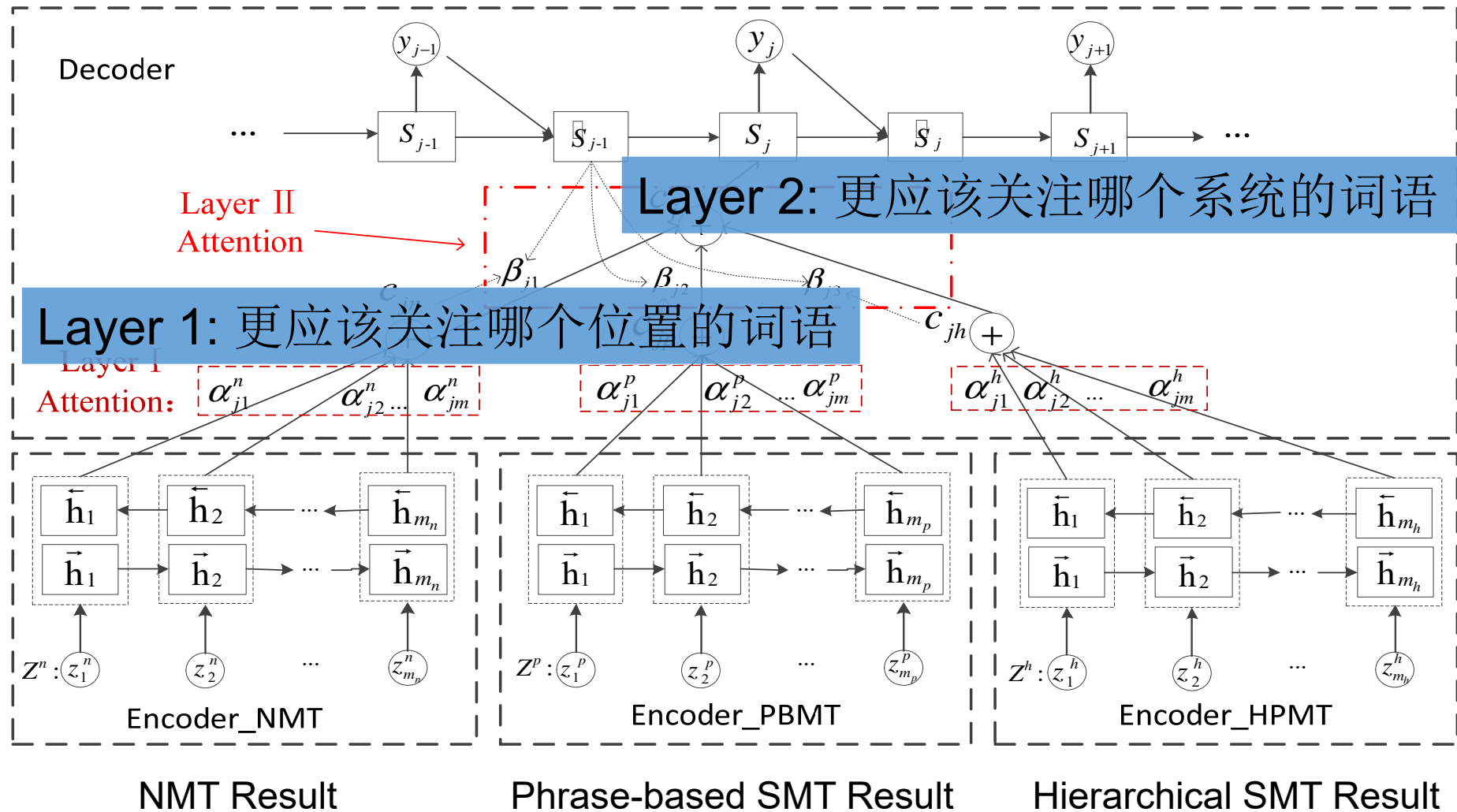
$$p(y_i|y_{<i}, C) = p(y_i|y_{<i}, c_i) \\ = \text{softmax} \left( (1 + \lambda V(R_i)) \text{score}(W_{y_i}, h_i) \right)$$

# 短语级融合-NMT利用SMT短语翻译表

#	Method	CH-EN						EN-JA	
		MT03(dev)	MT04	MT05	MT06	MT08	Ave	dev	test
1	Moses	28.35	30.02	29.10	32.92	23.20	28.72	20.06	22.40
2	Baseline	34.20	36.96	32.60	33.85	25.96	32.71	23.61	25.99
3	Arthur	34.98 <sup>†</sup>	37.96 <sup>†</sup>	33.36 <sup>†</sup>	34.79 <sup>†</sup>	26.53*	33.52	24.33*	26.72 <sup>†</sup>
4	<b>Our method</b>	<b>36.48<sup>†</sup></b>	<b>38.79<sup>†</sup></b>	<b>35.34<sup>†</sup></b>	<b>36.58<sup>†</sup></b>	<b>27.49<sup>†</sup></b>	<b>34.94</b>	<b>25.63<sup>†</sup></b>	<b>27.95<sup>†</sup></b>
5	system(no matching)	34.99 <sup>†</sup>	37.54*	33.32 <sup>†</sup>	34.22*	26.39*	33.29	24.11*	26.47*
6	system(no first)	35.25 <sup>†</sup>	38.07 <sup>†</sup>	34.13 <sup>†</sup>	34.95 <sup>†</sup>	26.67 <sup>†</sup>	33.81	24.37 <sup>†</sup>	26.93 <sup>†</sup>

- 方法2比统计模型平均提升**5个BLEU**值，比神经网络系统平均提升**2个BLEU**值！

# 句子级融合模型 (Zhou et al., 2017)



# 层次注意模型

- **Layer 1: 系统内的注意机制**

第 $j$ 输出与第 $k$ 个系统中第 $i$ 个隐层表示的注意权重

$$c_{jk} = \sum_{i=1}^m \alpha_{ji}^k h_i \rightarrow \text{第 } k \text{ 个系统中第 } i \text{ 个隐层表示}$$

第 $k$ 个系统第 $j$ 个输出对应的上下文

$$\alpha_{ji}^k = \frac{\exp(e_{ji})}{\sum_{l=1}^m \exp(e_{jl})}$$

$$e_{ji} = v_a^T \tanh(W_a \tilde{s}_{j-1} + U_a h_i)$$

# 层次注意模型

- **Layer 2:** 不同系统之间的注意机制

第 $j$ 个输出与第 $k$ 个系统之间的注意权重

$$c_j = \sum_{k=1}^K \beta_{jk} c_{jk}$$

第 $j$ 个输出对应的上下文

第 $k$ 个系统第 $j$ 个输出对应的上下文

$$\beta_{jk} = \frac{\exp(\tilde{s}_{j-1} \cdot c_{jk})}{\sum_{k'} \exp(\tilde{s}_{j-1} \cdot c_{jk'})}$$

# 句子级融合模型的效果

System	MT03	MT04	MT05	MT06	Ave
PBMT	37.47	41.20	36.41	36.03	37.78
HPMT	<b>38.05</b>	<b>41.47</b>	<b>36.86</b>	36.04	<b>38.10</b>
NMT	37.91	38.95	36.02	<b>36.65</b>	37.38
Jane (Freitag et al., 2014)	39.83	42.75	38.63	39.10	40.08
Multi	40.64	44.81	38.80	38.26	40.63
Multi+Source	42.16	45.51	40.28	39.03	41.75
Multi+Ensemble	41.67	45.95	40.37	39.02	41.75
Multi+Source+Ensemble	<b>43.55</b>	<b>47.09</b>	<b>42.02</b>	<b>41.10</b>	<b>43.44</b>

3.4

5.3

- 句子级融合模型比最好的单系统平均提升**5.3**个**BLEU**值，比传统系统融合方法平均提升**3.4**个**BLEU**值！



# 总结

- 机器翻译从统计模型迁移到深度学习模型
- 基于注意机制的神经机器翻译模型具有更优的翻译性能
- 统计方法与神经网络方法融合，对两种范式兼容并包、博采众长





# 参考文献

1. Marcin Junczys-Dowmunt, Tomasz Dwojak and Hieu Hoang, 2016. [Is Neural Machine Translation Ready for Deployment? A Case Study on 30 Translation Directions](https://arxiv.org/pdf/1610.01108.pdf). <https://arxiv.org/pdf/1610.01108.pdf>
2. Nal Kalchbrenner and Phil Blunsom, 2013. [Recurrent Continuous Translation Models](#). *In Proc. of EMNLP 2013*.
3. Dzmitry Bahdanau, Kyunghyun Cho and Yoshua Bengio, 2015. [Neural Translation by Jointly Learning to Align and Translate](http://arxiv.org/pdf/1409.0473.pdf). <http://arxiv.org/pdf/1409.0473.pdf>
4. Ilya Sutskever, Oriol Vinyals and Quoc V. Le, 2014. [Sequence to Sequence Learning with Neural Networks](#). *In Proc. of NIPS 2014*.
5. Marcin Junczys-Dowmunt, Tomasz Dwojak and Hieu Hoang, 2016. [Is Neural Machine Translation Ready for Deployment? A Case Study on 30 Translation Directions](#). [arxiv.org/pdf/1610.01108](https://arxiv.org/pdf/1610.01108.pdf)
6. Philip Arthur, Graham Neubig, and Satoshi Nakamura. [Incorporating discrete translation lexicons into neural machine translation](#). In Proceedings of EMNLP 2016.
7. Yang Feng, Shiyue Zhang, Andi Zhang, Dong Wang, and Andrew Abel. [Memory-augmented neural machine translation](#). In proceedings of EMNLP 2017.



# 参考文献

8. Jiajun Zhang and Chengqing Zong. 2016b. [Bridging Neural Machine Translation and Bilingual Dictionaries](https://arxiv.org/pdf/1610.07272.pdf). <https://arxiv.org/pdf/1610.07272.pdf>
9. Yang Zhao, Yining Wang, Jiajun Zhang and Chengqing Zong, 2018. [Phrase Table as Recommendation Memory for Neural Machine Translation](#). In Proc. of IJCAI 2018.
10. Yaohua Tang, Fandong Meng, Zhengdong Lu, Hang Li, and Philip LH Yu. [Neural machine translation with external phrase memory](#). arXiv preprint arXiv:1606.01792, 2016
11. Xing Wang, Zhaopeng Tu, Deyi Xiong, and Min Zhang. [Translating phrases in neural machine translation](#). In proceedings of EMNLP, 2017.
12. Jan Niehues, Eunah Cho, Thanh-Le Ha and Alex Waibel. [Pre-translation for Neural Machine Translation](#). In Proceedings of COLING 2016.
13. Long Zhou, Wenpeng Hu, Jiajun Zhang and Chengqing Zong, 2017. [Neural System Combination for Machine Translation](#). *In Proc. of ACL 2017*.
14. 宗成庆. 2008. 《统计自然语言处理》，清华大学出版社。

N L P R



谢谢!  
Q&A