

中国中文信息学会《前沿技术讲习班》第二十五期 (CIPS ATT25)

机器同传

何中军

2021-10-8

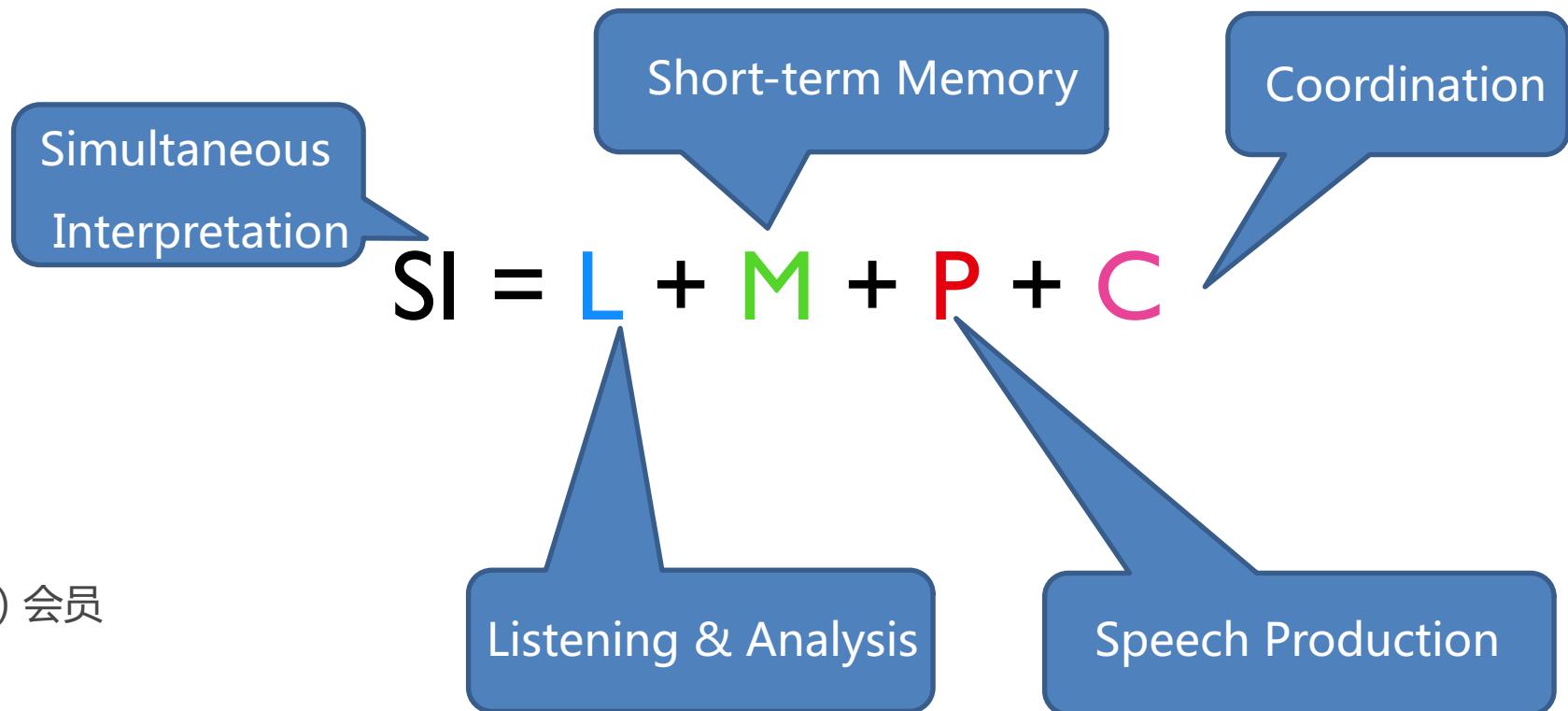
提纲

- 同声传译
- 主要挑战
- 主流方法
 - 级联模型
 - 端到端模型
- 数据、鲁棒性模型及评价方法
- 产品及应用
- 总结及展望

同声传译

- 极具挑战性的翻译方式

吉尔精力分配模式 (Gile's Effort Model)



欧洲翻译研究学会前主席
国际会议口译员协会 (AIIC) 会员

GILE, D. "Basic Theoretical Components in Interpreter and Translator Training." DOLLERUP, C. and LODDEGAARD, A. (eds). *Teaching Translation and Interpreting: Training, Talent and Experience*. Amsterdam: John Benjamins, 1992.

同声传译

准确度及完整性：连续工作30分钟后，每过5分钟，准确率下降10%

工作方式：一般2人一组，轮流工作和休息

AIIC会员：全球约3000人



同传常用技巧

skipping, approximation, filtering, comprehension omission, and substitution

同传常用技巧

skipping, approximation, filtering, comprehension omission, and substitution

然后这个粉色的方框是算法自动预测出来的结果。

Translation: While the pink **boxes** were by an algorithm.

SI: And the pink **one** is **through** algorithms.

同传常用技巧

skipping, approximation, filtering, comprehension omission, and substitution

我们在 2015 年 就 已 经 取 得 了 世 界 第 一 的 水 平 ， 错 误 率 只 有 0.23%。

Translation: We ranked No.1 in the world **in 2015**, with an error rate of only **0.23%**.

SI: We rank the first in the world. And the error rate is **really low**

同传常用技巧

顺句驱动、适当重复

We are most pleased that we share identical views on a wide range of issues.

Translation : 我们在很多问题上观点一致，非常高兴

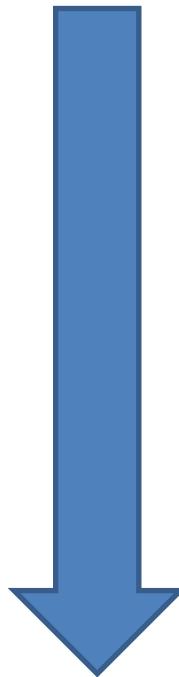
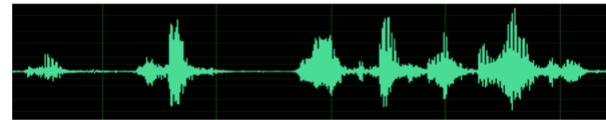
SI : 我们非常高兴，我们的观点一致，在很多问题上 观点一致

提纲

- 同声传译
- 主要挑战
- 主流方法
 - 级联模型
 - 端到端模型
- 数据、鲁棒性模型及评价方法
- 产品及应用
- 总结及展望

机器同传

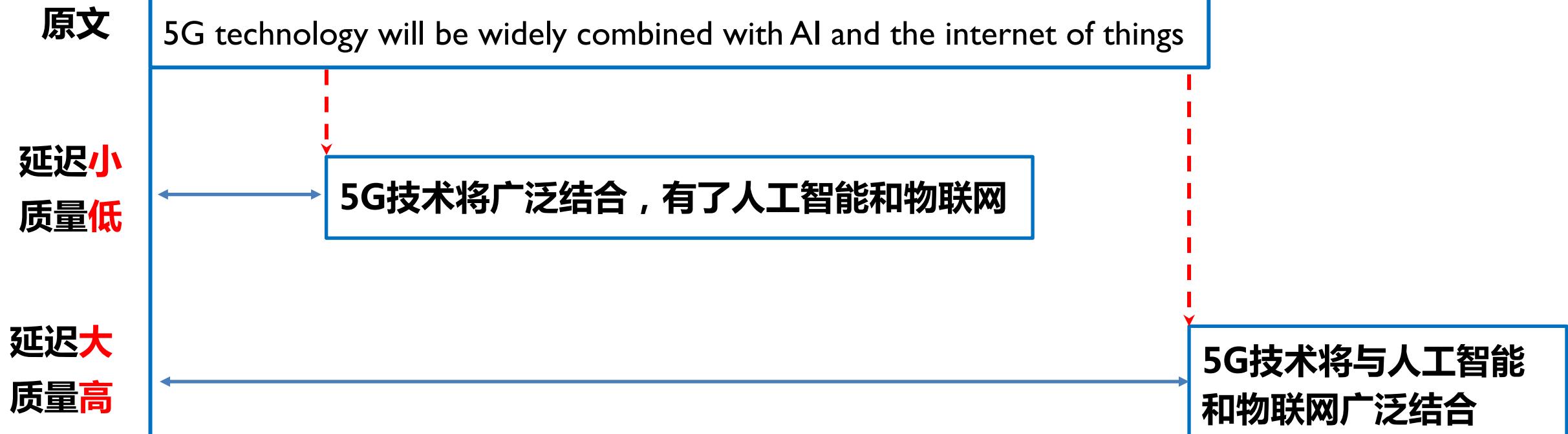
源语言



目标语言



挑战一：翻译质量与时延难以兼顾



挑战二：无断句标点

源语言



|那么大家知道这个庄稼最怕的是出现病虫害一旦就是一个病虫害出现防治不及时的话会造成大量的这个减产|

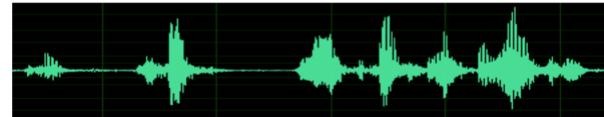
从哪里开始？到哪里结束？

目标语言



挑战二：语音识别错误

源语言



|那么大家知道这个**重庄家**最怕的是出现病虫害一旦就是一个病虫害出现防治不及时的话会造成大量的这个减产|

错误传递和放大：差之毫厘，谬以千里

目标语言



挑战三：口语化成分（冗余、停顿等）

源语言



那么大家知道这个重庄家最怕的是出现病虫害一旦就是一个病虫害出现防治不及时的话会造成大量的这个减产

影响质量，增加时延

目标语言



挑战四：数据稀缺

MT

Source Text

Bilingual

Target Text

我们齐心搭建世界联通新桥梁

Together, we have built new bridges for global connectivity.

ASR

Audio Signal

Monolingual

Transcription



那么我们今天呢就希望，从一个二十年的AI工作者来说，如何从专业的角度去解读一下 ...

SimulTrans

Source Audio

Bilingual

Target Text / Audio



So today, as one who has been working on AI for twenty years, I wish I could give you a professional interpretation ...

挑战五：无统一评价标准

然后这个粉色的方框是算法自动预测出来的结果。

Translation: While the pink **boxes** were **predicted automatically** by an algorithm.

SI: And the pink **one** is **through** algorithms.

评价文本翻译的指标（关注翻译的准确、语法等），

不适合评价同传（关注信息传递的质量和效率）

提纲

- 同声传译
- 主要挑战
- 主流方法
 - 级联模型
 - 端到端模型
- 数据、鲁棒性模型及评价方法
- 产品及应用
- 总结及展望

主流方法

级联模型

语音识别(ASR)



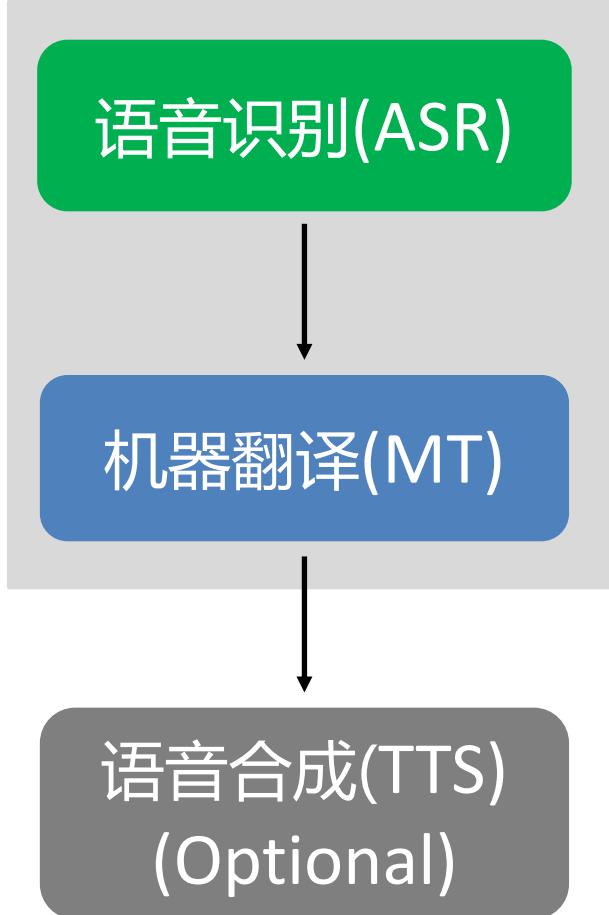
机器翻译(MT)



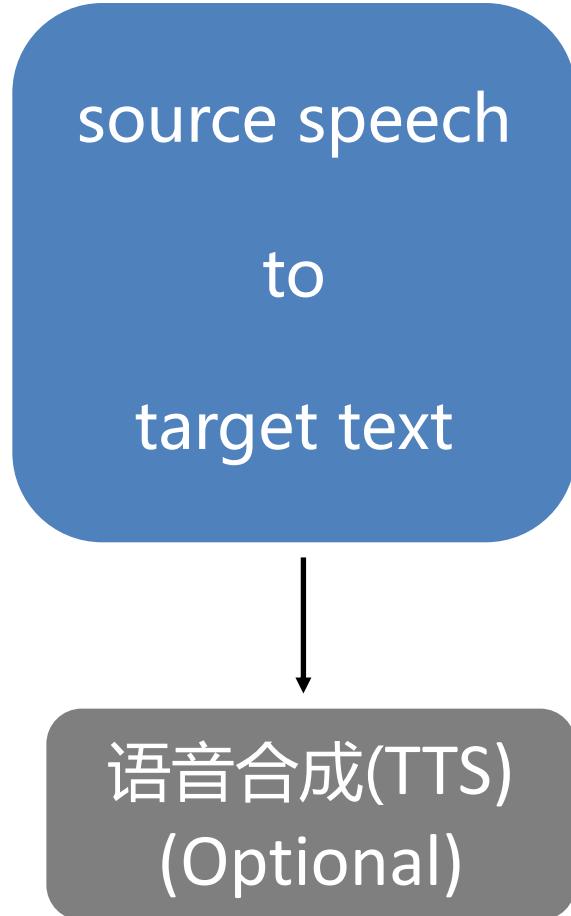
语音合成(TTS)
(Optional)

主流方法

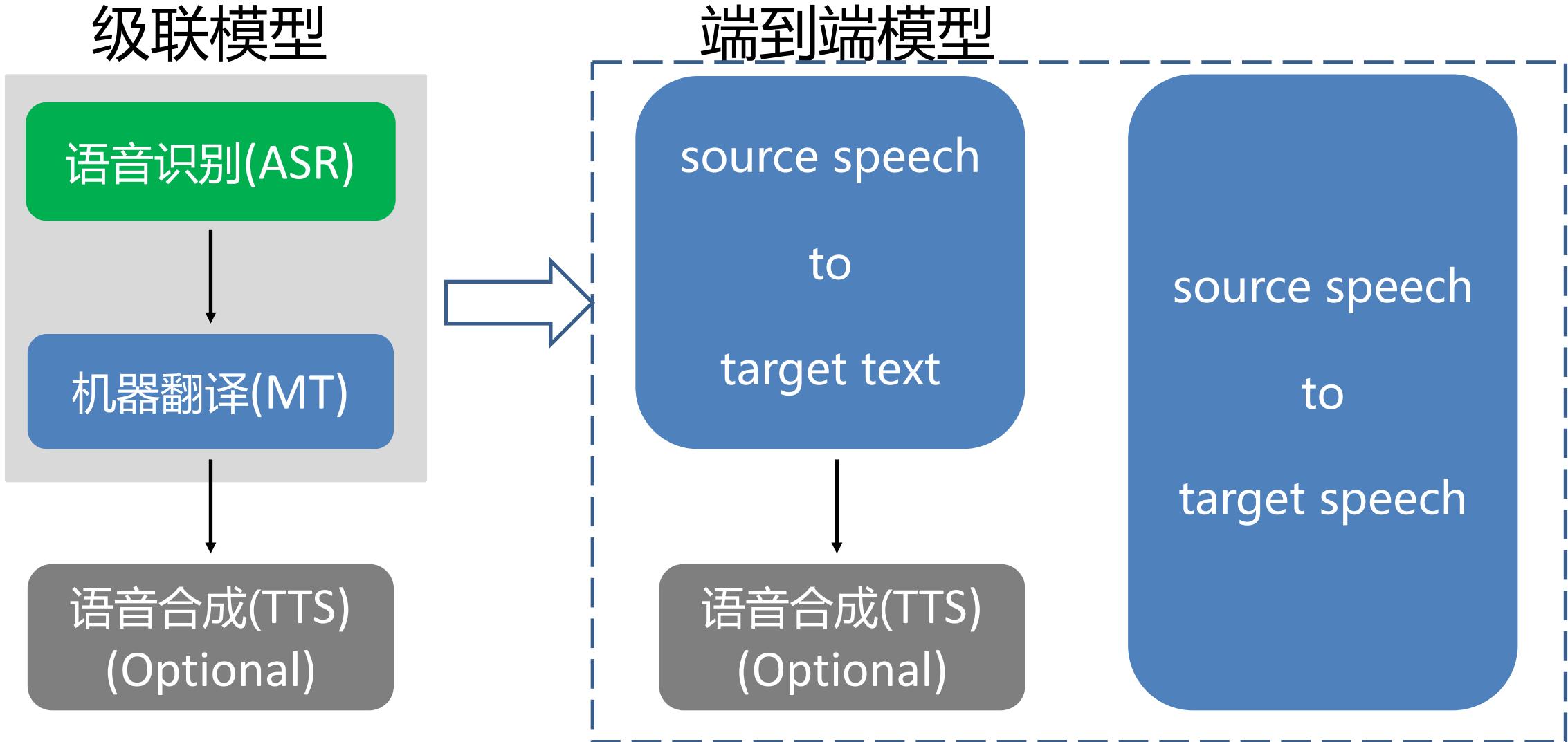
级联模型



端到端模型



主流方法



级联模型

ASR



MT

那么大家知道这个种庄稼最怕的是出现病虫害一旦就是一个病虫害出现防治不及时的话会造成大量的这个减产

ASR: 无断句标点



MT: 以句子为单位进行翻译

As we all know, The biggest fear of planting crops is the emergence of pests and diseases. Once there is a disease and pest, if the prevention and control is not timely, it will cause a lot of production reduction.

级联模型

ASR

那么大家知道这个种庄稼最怕的是出现病虫害一旦就是一个病虫害出现防治不及时的话会造成大量的这个减产



同传策略

读入 (READ) : 读入流文本

翻译 (WRITE) : 翻译已经读入的文本



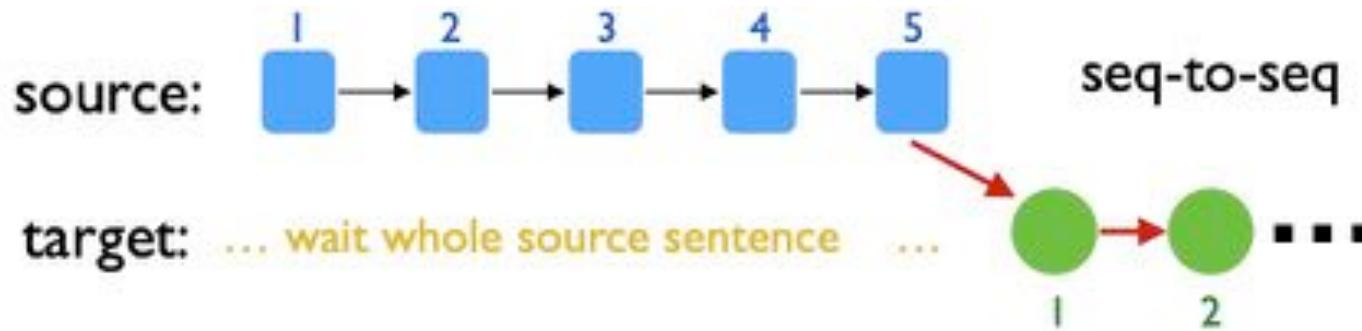
MT

As we all know, The biggest fear of planting crops is the emergence of pests and diseases. Once there is a disease and pest, if the prevention and control is not timely, it will cause a lot of production reduction.

级联模型 – 同传策略

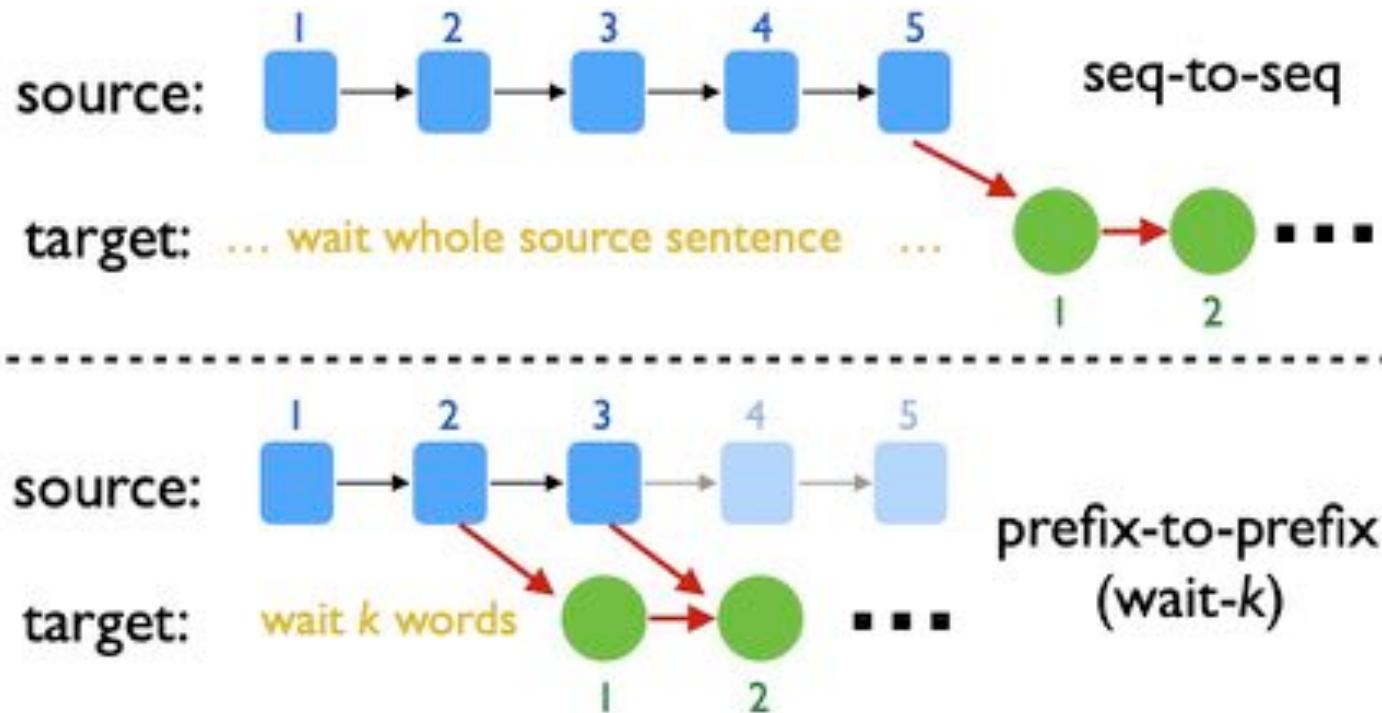
- 固定策略 (Fixed Policy) : 读入长度 (翻译单元) 固定
 - Static read-write (Dalvi et al., 2018)
 - STACL (Ma et al., 2018), etc.
- 自适应策略 (Adaptive Policy) : 动态调整翻译单元长度
 - Rule-based (Cho et al., 2016)
 - RL-based (Gu et al., 2017)
 - Supervised policy (Zheng et al., 2019)
 - Meaningful unit (Zhang et al., 2020)
 - MILk (Arivazhagan et al., 2019)
 - Multihead monotonic attention (Ma et al., 2020), etc.

wait-k策略



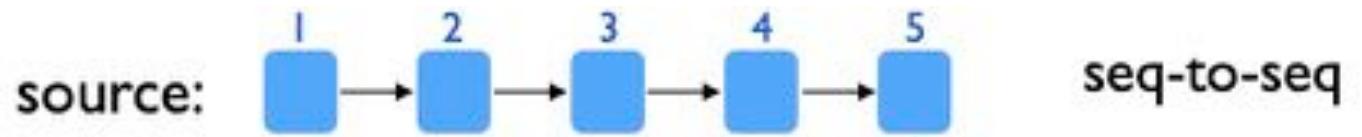
wait- k 策略

先读入 k 个词，然后开始翻译

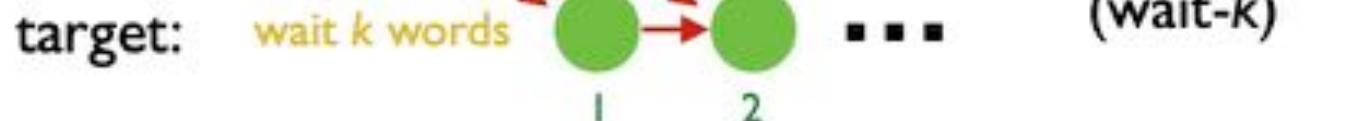
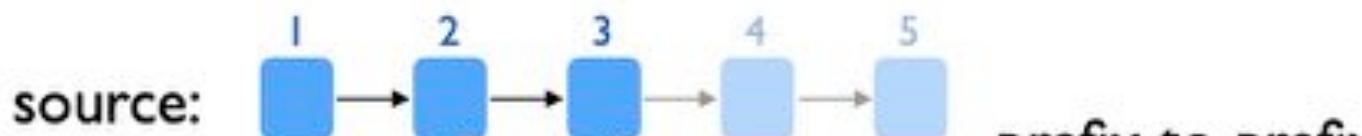


wait- k 策略

先读入 k 个词，然后开始翻译



$$p(y_t | \{x_1, \dots, x_n\}, \{y_1, \dots, y_{t-1}\})$$



$$p(y_t | \{x_1, \dots, x_{t+k-1}\}, \{y_1, \dots, y_{t-1}\})$$

$g(t)$: number of source words used to predict y_t

$$p(y_t | \{x_1, \dots, x_{g(t)}\}, \{y_1, \dots, y_{t-1}\})$$

wait-k策略

两国 领导人 将 在 十月 举行 会谈

wait-2

The leaders of the two countries will hold talks in October

wait-k策略

两国 领导人 将 在 十月 举行 会谈

wait-2

The leaders of the two countries will hold talks in October

wait-k策略

两国 领导人 将 在 十月 举行 会谈

wait-2

The leaders of the two countries will hold talks in October

The	
两国	R
领导人	R W
将	
在	
十月	
举行	
会谈	

$$g(1) = 2$$

wait-k策略

两国 领领导人 将 在 十月 举行 会谈

wait-2

The leaders of the two countries will hold talks in October

The leaders	
两国	R
领导人	R W
将	R
在	W
十月	
举行	
会谈	

$$g(1) = 2 \quad g(2) = 3$$

wait-k策略

两国 领领导人 将 在 十月 举行 会谈

wait-2

The leaders of the two countries will hold talks in October

The leaders of	
两国	R
领导人	R W
将	R W
在	R W
十月	
举行	
会谈	

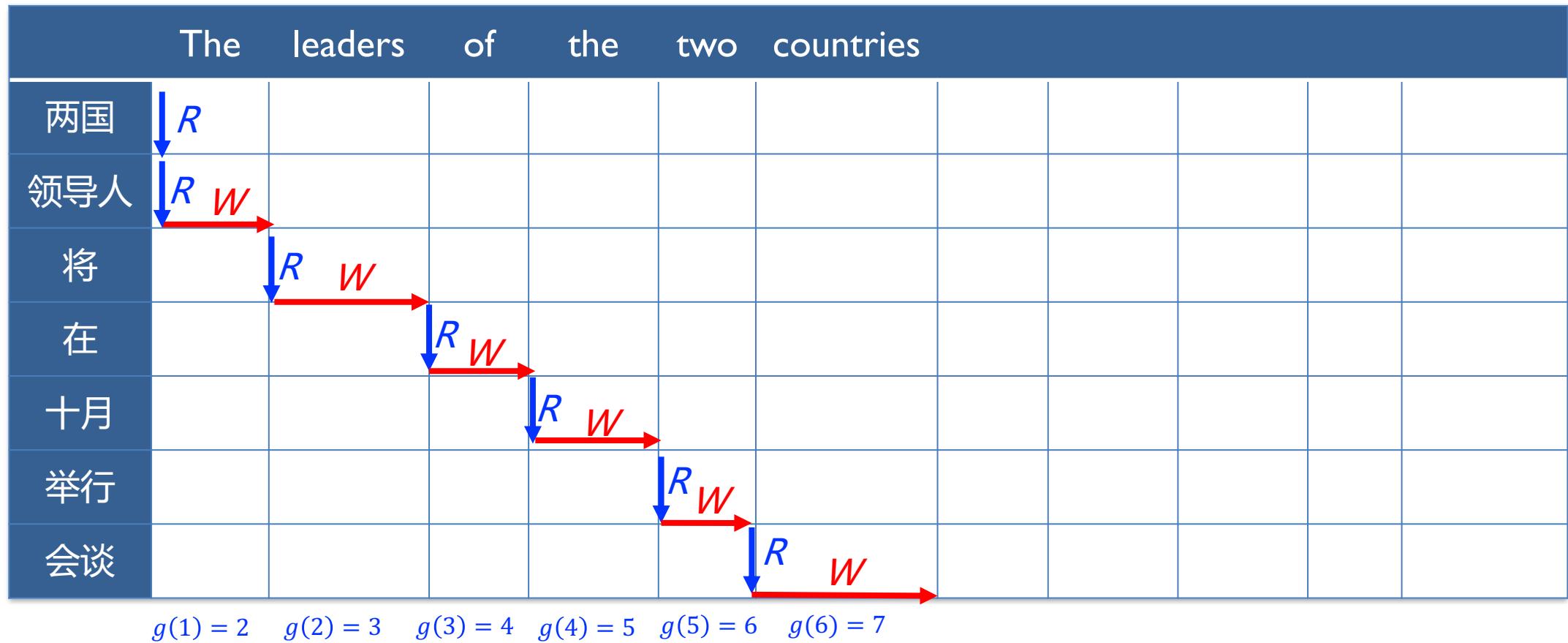
$$g(1) = 2 \quad g(2) = 3 \quad g(3) = 4$$

wait-k策略

两国 领导人 将 在 十月 举行 会谈

wait-2

The leaders of the two countries will hold talks in October

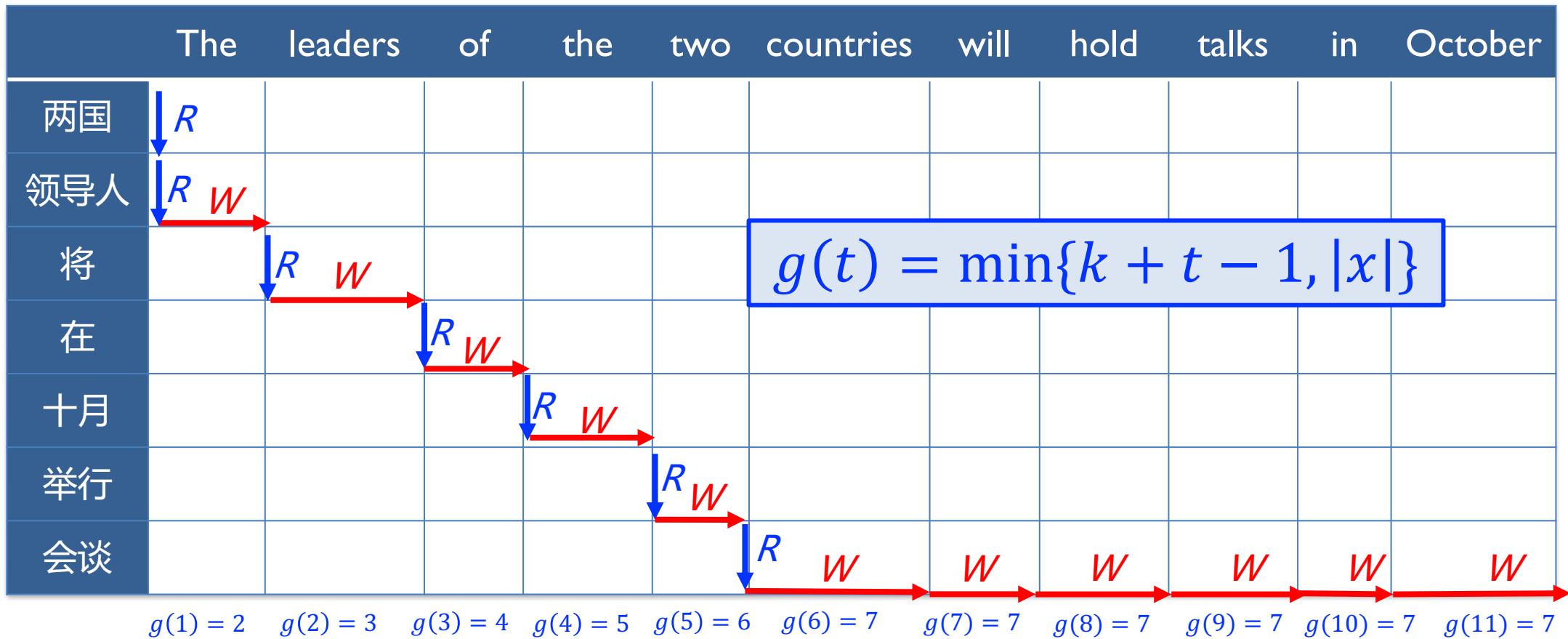


wait-k策略

两国 领领导人 将 在 十月 举行 会谈

wait-2

The leaders of the two countries will hold talks in October



wait-k存在的问题

目标句子尾部信息量大，增加听众负担

两国 领导人 将 在 十月 举行 会谈

wait-2

The leaders of the two countries will hold talks in October

wait-k存在的问题

信息不充分导致译文预测错误

大家 对 他 取得 的 成绩 表示 怀疑

wait-2

Everyone is **satisfied** with his achievement

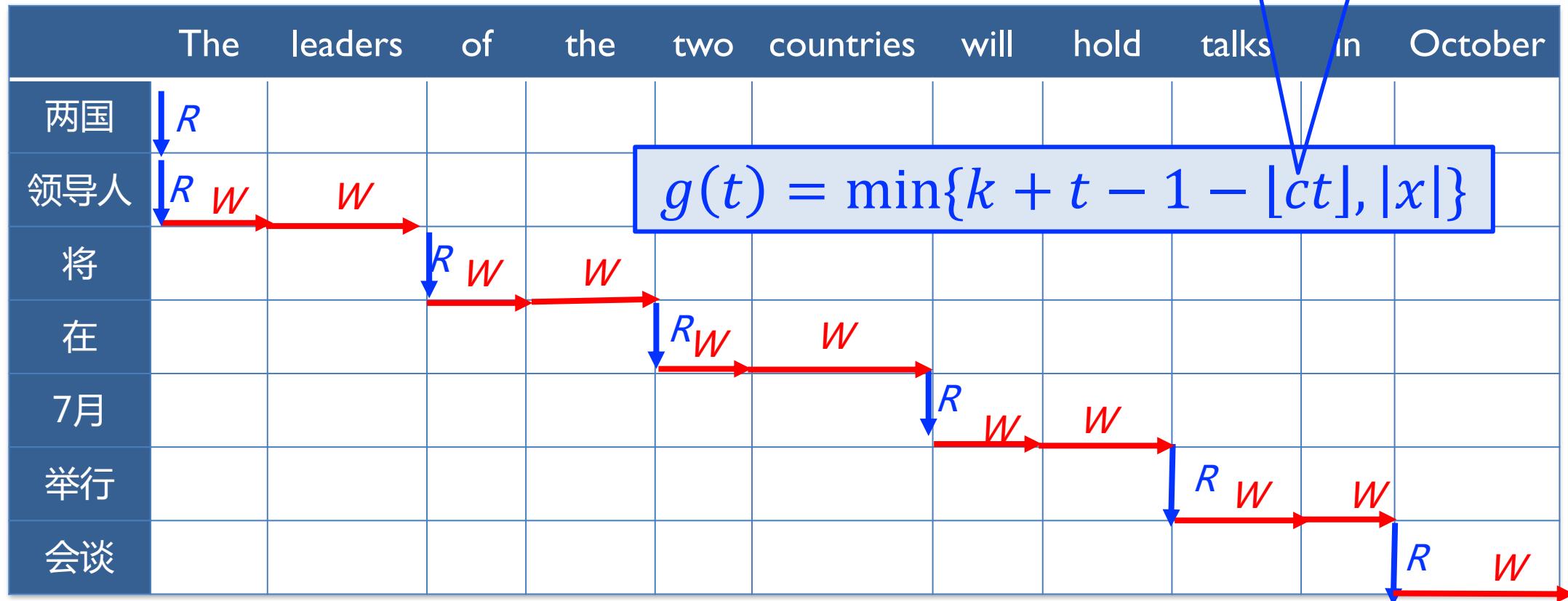
wait-5

Everyone is **doubt** with his achievement

wait-k catchup

-根据目标语言和源语言长度比，调整翻译速度

$$c = \frac{|y^*|}{|x|} - 1$$



wait-k策略 —— Demo

Welcome to Simultaneous Translation(text+audio)

PaddleNLP

Chinese input: I

Rec

Jieba+BPE:

Simultaneous Translation (wait 1):

Simultaneous Translation (wait 3):

Simultaneous Translation (wait 5):

Full Sentence Translation (wait -1):

CLEAR

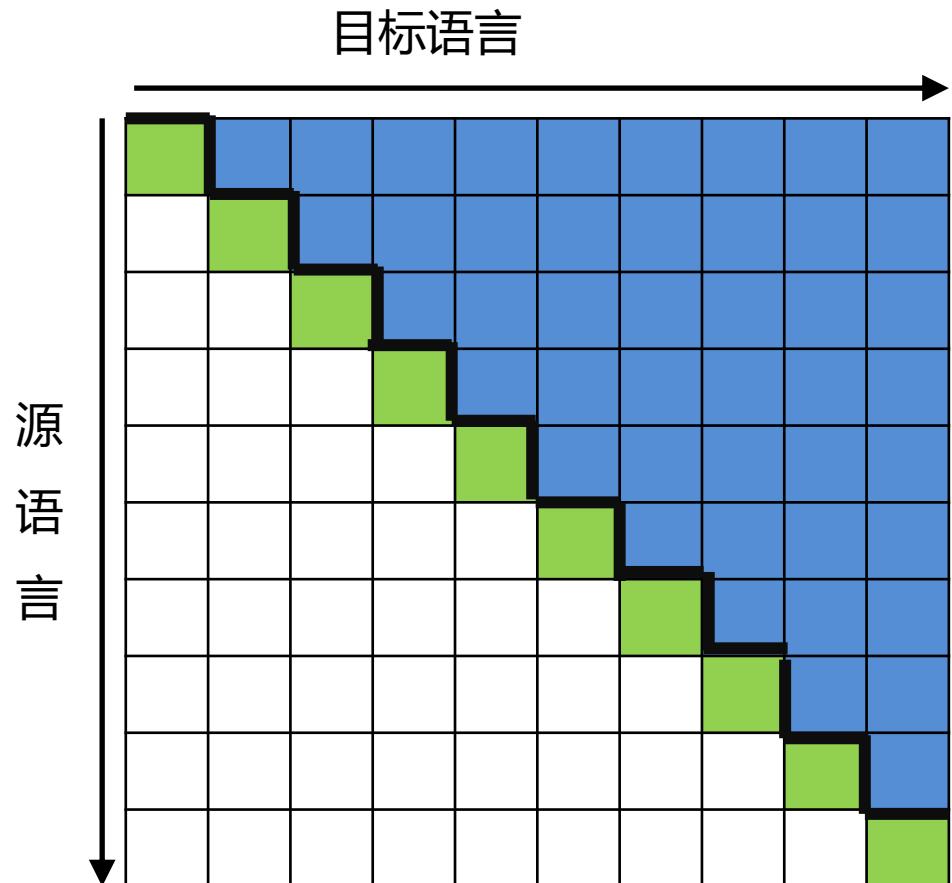
使用说明: 1. 在 Chinese input 输入中文, 按【回车键】开始实时翻译, 遇到【。!?】结束整句, 按【CLEAR】清空所有的输入和输出;
2. 按【Rec】开始录音并开始实时翻译, 遇到【。!?】结束整句, 按【CLEAR】清空所有的输入和输出。

https://github.com/PaddlePaddle/PaddleNLP/blob/develop/examples/simultaneous_translation/stacl/



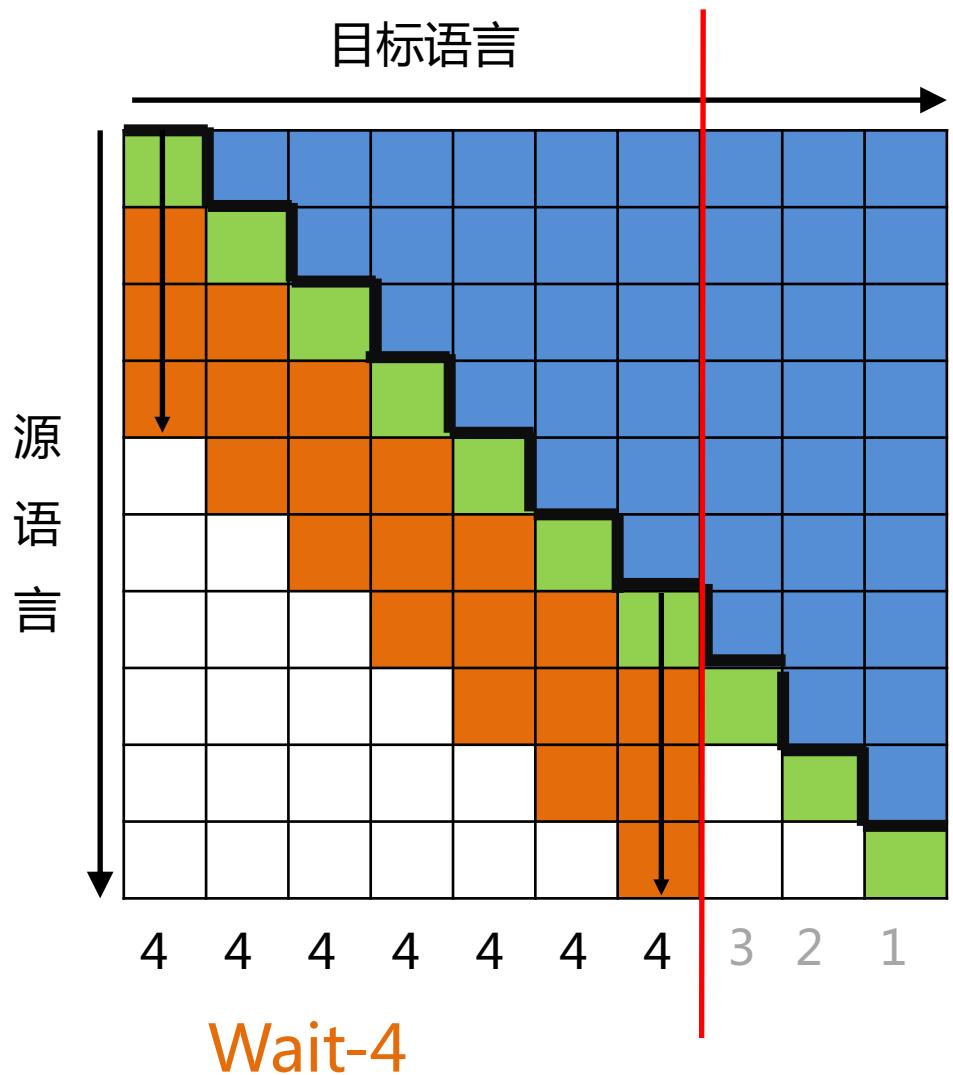
同传时延的评价

时延评价方法 – Average Lagging



理想情况 (oracle) : $AL = 0$
先写1个词，然后再读1个词，依次进行

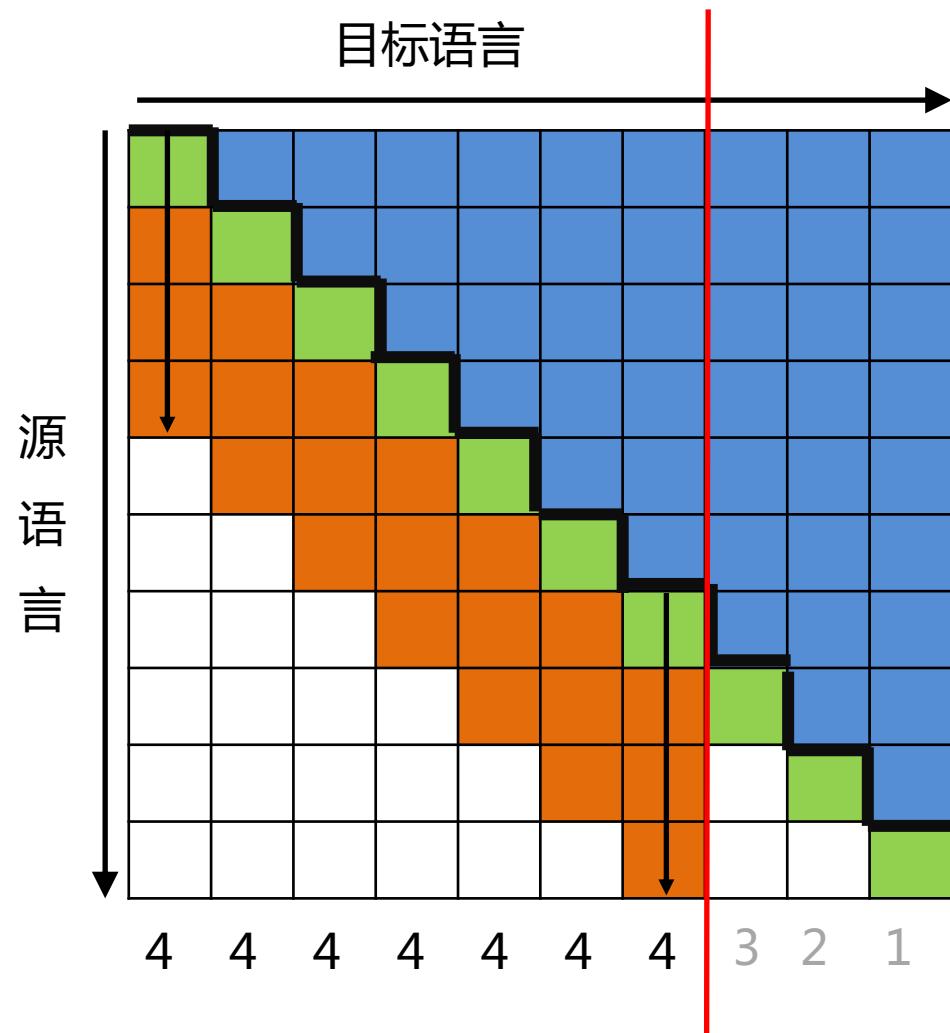
时延评价方法 – Average Lagging



$$AL = 34/10 = 3.4 ?$$

需要在源语言句子读完时，停止计算

时延评价方法 – Average Lagging

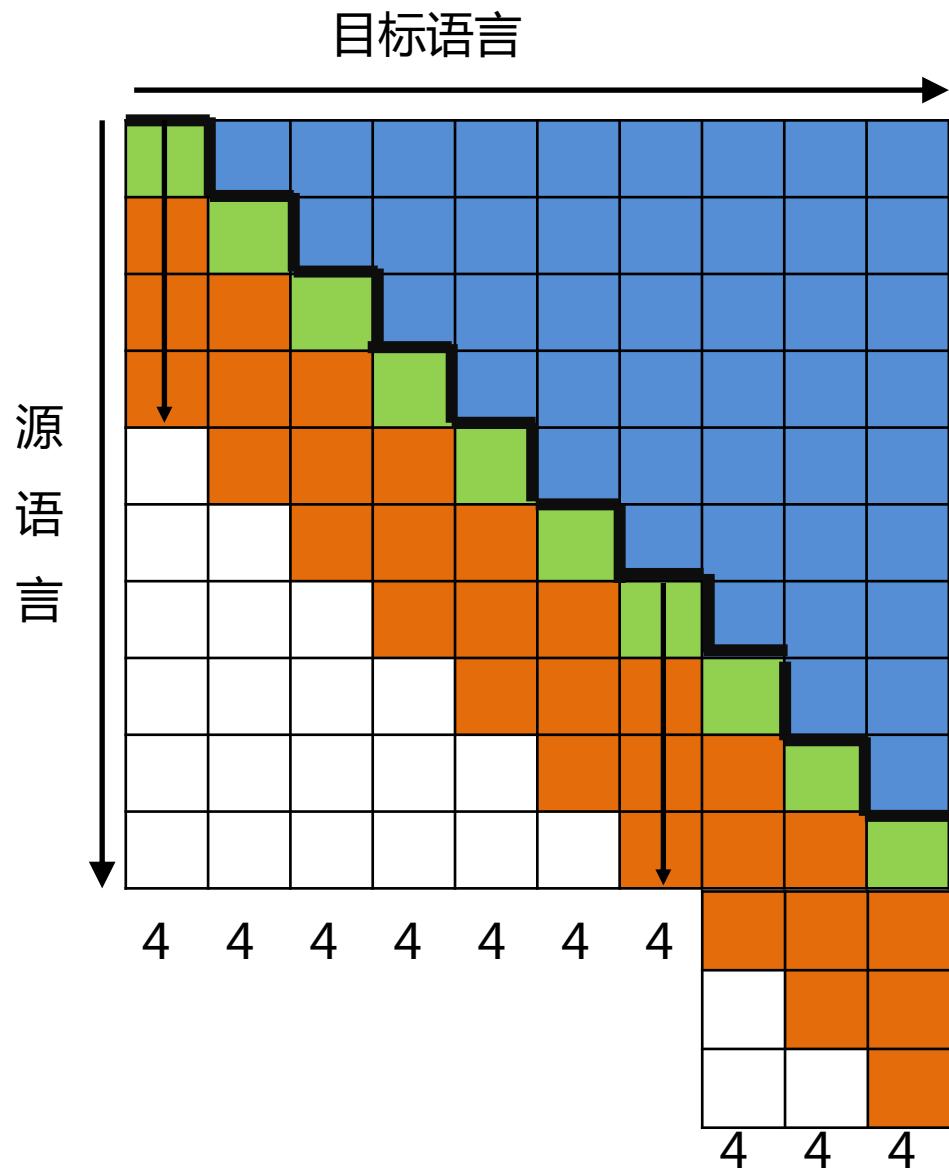


$$AL = \frac{1}{\tau} \sum_{i=1}^{\tau} g_i - \frac{i-1}{\gamma}$$

$$\tau = \operatorname{argmin}_i g_i = |\mathbf{x}|$$

$$\gamma = \frac{|y|}{|\mathbf{x}|}$$

时延评价方法 – Differentiable Average Lagging (DAL)

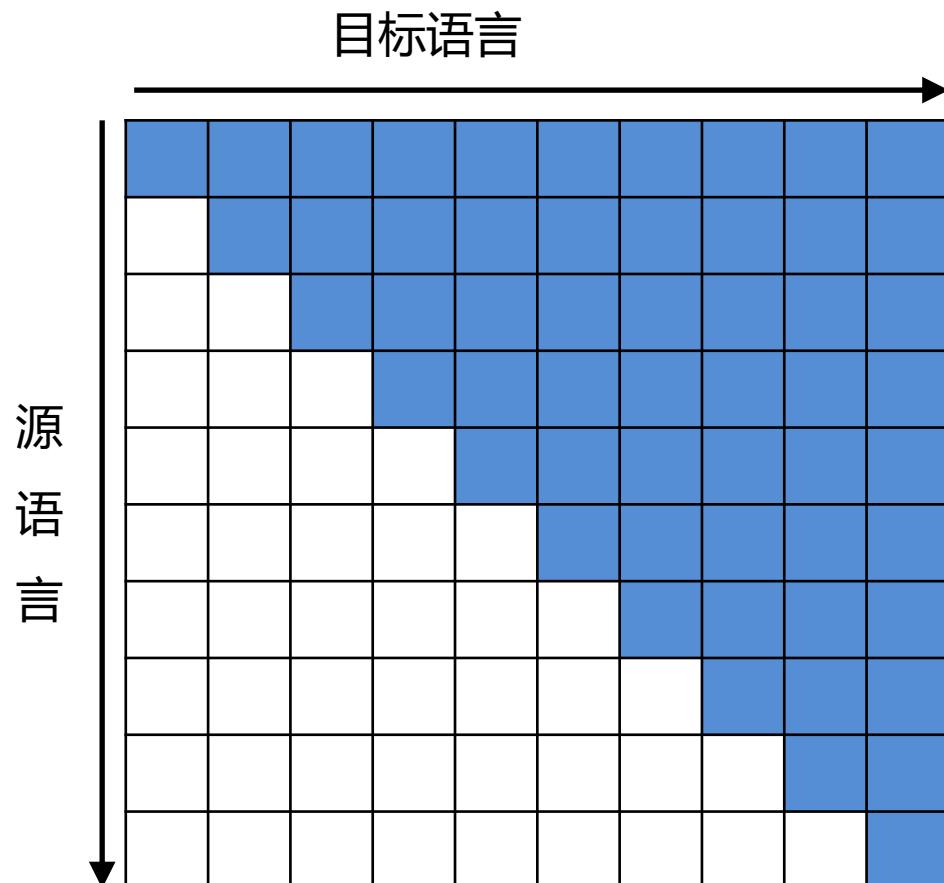


$$DAL = \frac{1}{|Y|} \sum_{t=1}^{|Y|} g'_i - \frac{i-1}{\gamma}$$

$$g'_i = \max(g_i, g'_{i-1} + \frac{1}{\gamma})$$

可集成到损失函数，用于模型训练

时延评价方法 – Average Proportion (AP)

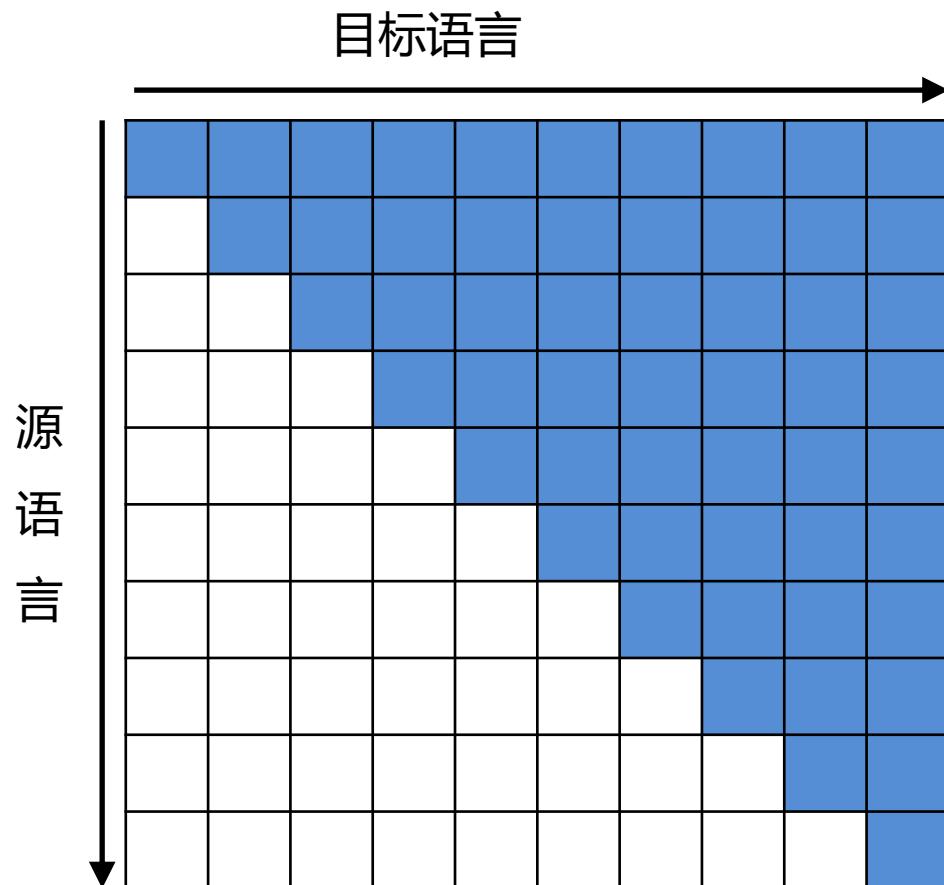


$$AP = \frac{1 + 2 + \dots + 10}{10 \times 10} = 0.55$$

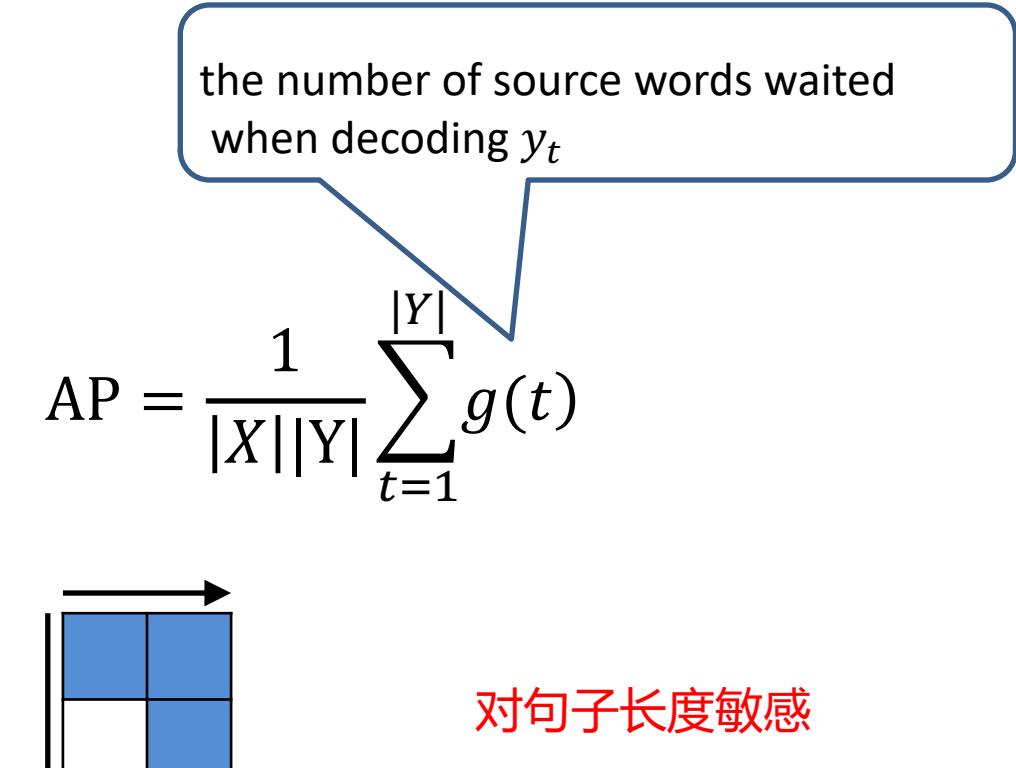
the number of source words waited
when decoding y_t

$$AP = \frac{1}{|X||Y|} \sum_{t=1}^{|Y|} g(t)$$

时延评价方法 – Average Proportion (AP)



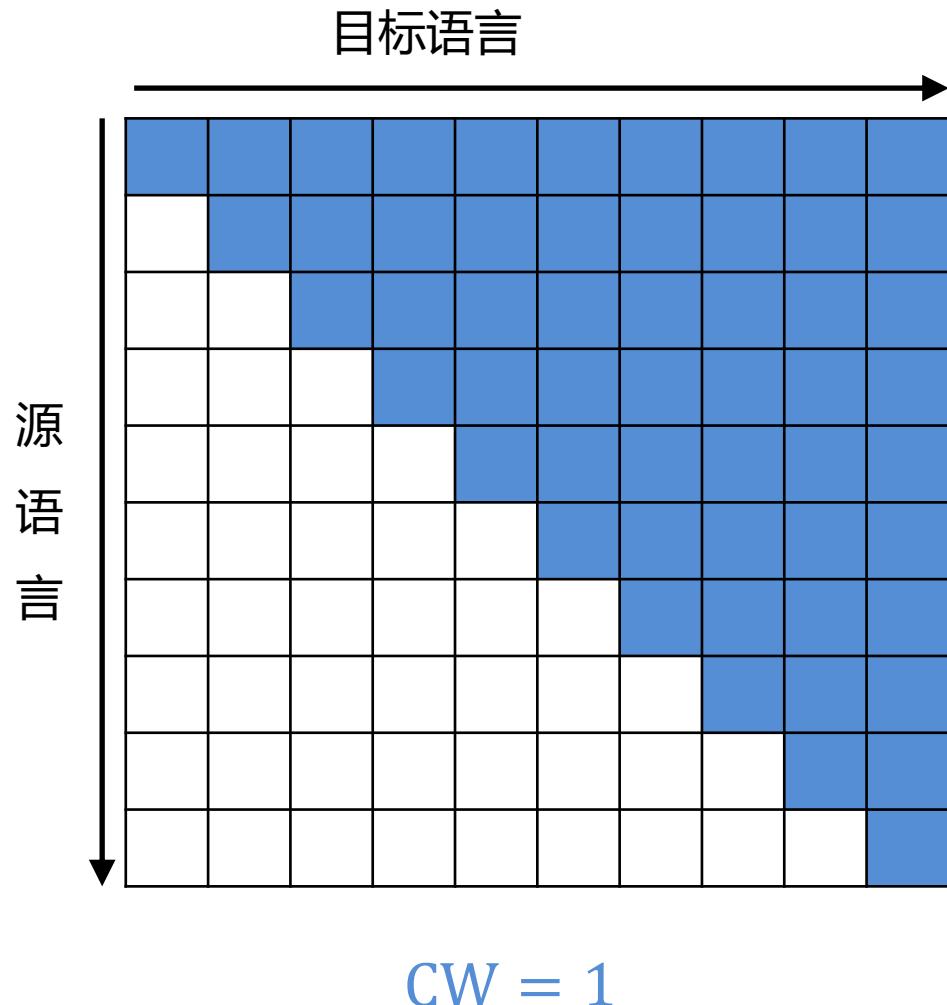
$$AP = \frac{1 + 2 + \dots + 10}{10 \times 10} = 0.55$$



对句子长度敏感

$$AP = \frac{1 + 2}{2 \times 2} = 0.75$$

时延评价方法 – Consecutive Wait (CW)



$$c_t = \begin{cases} c_{t-1} + 1 & a_t = \text{READ} \\ 0 & a_t = \text{WRITE} \end{cases}$$

$$CW = \frac{\sum_{t=1}^{|Y|} c_t}{\sum_{t=1}^{|Y|} 1_{c_t > 0}}$$

实际衡量的是源语言句子平均切分长度

级联模型 – 同传策略

- 固定策略 (Fixed Policy) : 读入长度 (翻译单元) 固定
 - Static read-write (Dalvi et al., 2018)
 - STACL (Ma et al., 2018), etc.
- 自适应策略 (Adaptive Policy) : 动态调整翻译单元长度
 - Rule-based (Cho et al., 2016)
 - RL-based (Gu et al., 2017)
 - Supervised policy (Zheng et al., 2019)
 - Meaningful unit (Zhang et al., 2020)
 - MILk (Arivazhagan et al., 2019)
 - Multihead monotonic attention (Ma et al., 2020), etc.

预先定义规则

- Wait-If-Worse (score comparison): 读入更多原文，预测概率下降

$$\Lambda(C, \mathcal{C} \cup \mathcal{C}') : \log p(\hat{y} | \hat{y}_{<t}, C) > \log p(\hat{y} | \hat{y}_{<t}, \mathcal{C} \cup \mathcal{C}')$$

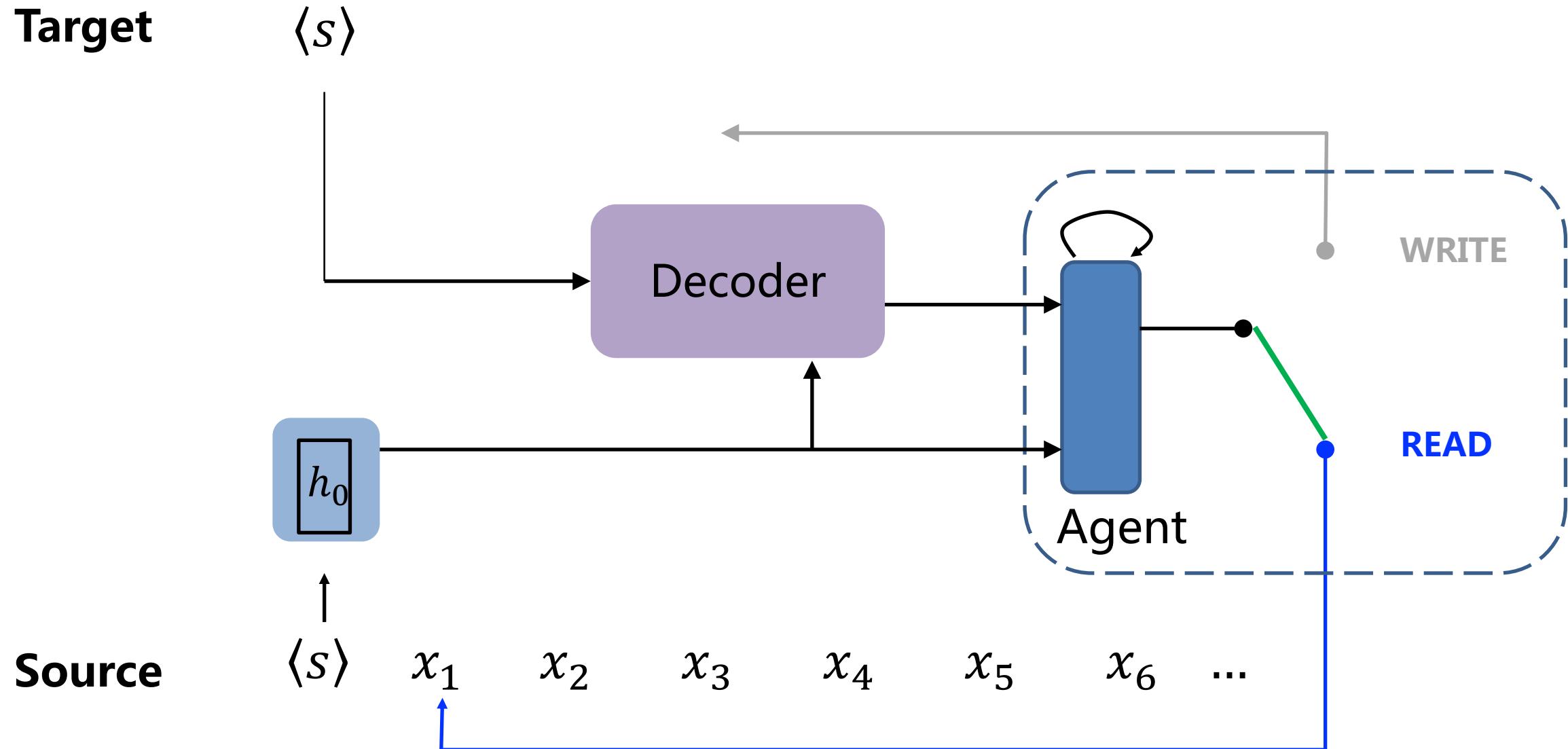
$$\hat{y} = \operatorname{argmax}_y \log p(y | \hat{y}_{<t}, C)$$

- Wait-If-Diff (Hypothesis comparison): 读入更多原文，预测单词不一致

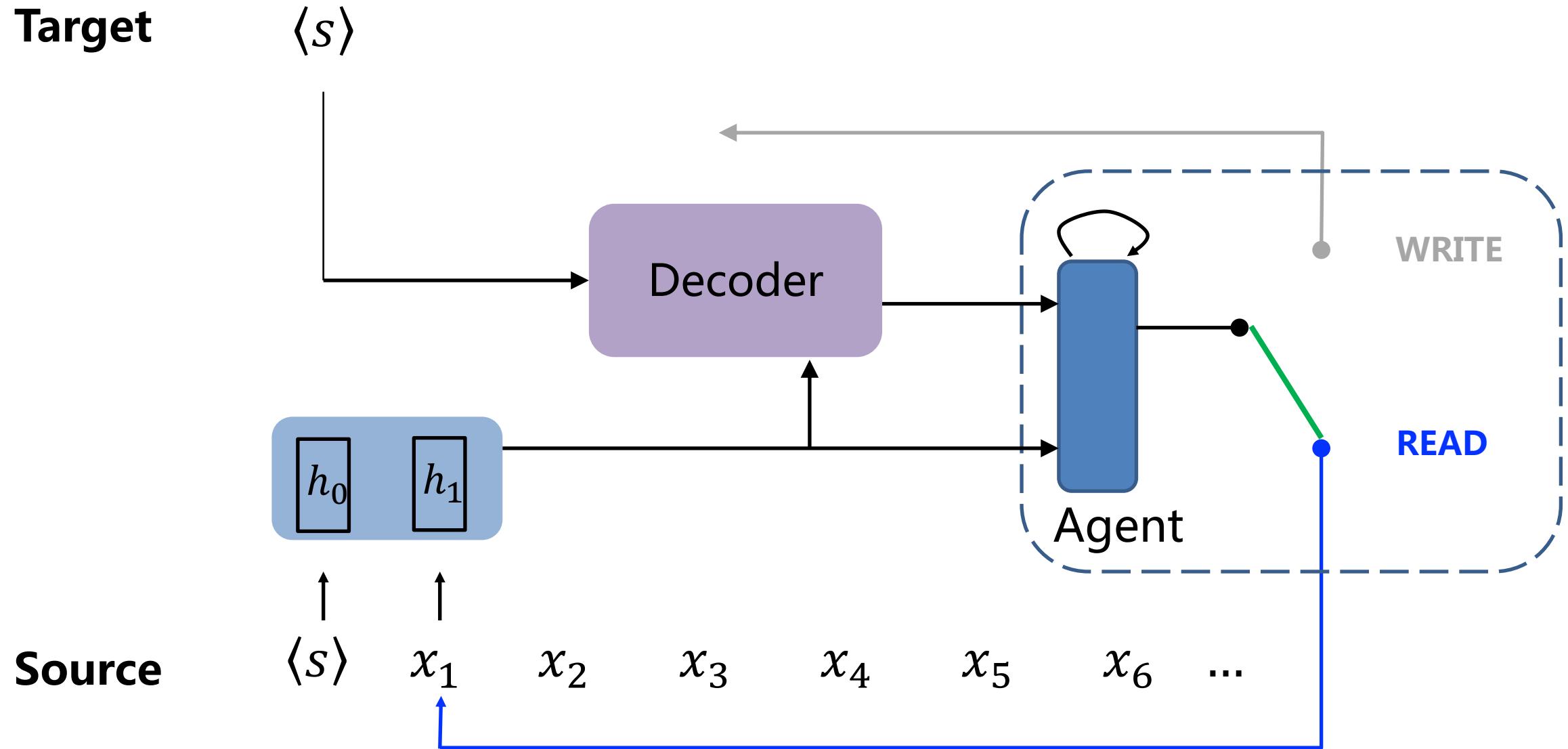
$$\Lambda(C, \mathcal{C} \cup \mathcal{C}') : \hat{y} \neq \hat{y}'$$

$$\hat{y}' = \operatorname{argmax}_y \log p(y | \hat{y}_{<t}, \mathcal{C} \cup \mathcal{C}')$$

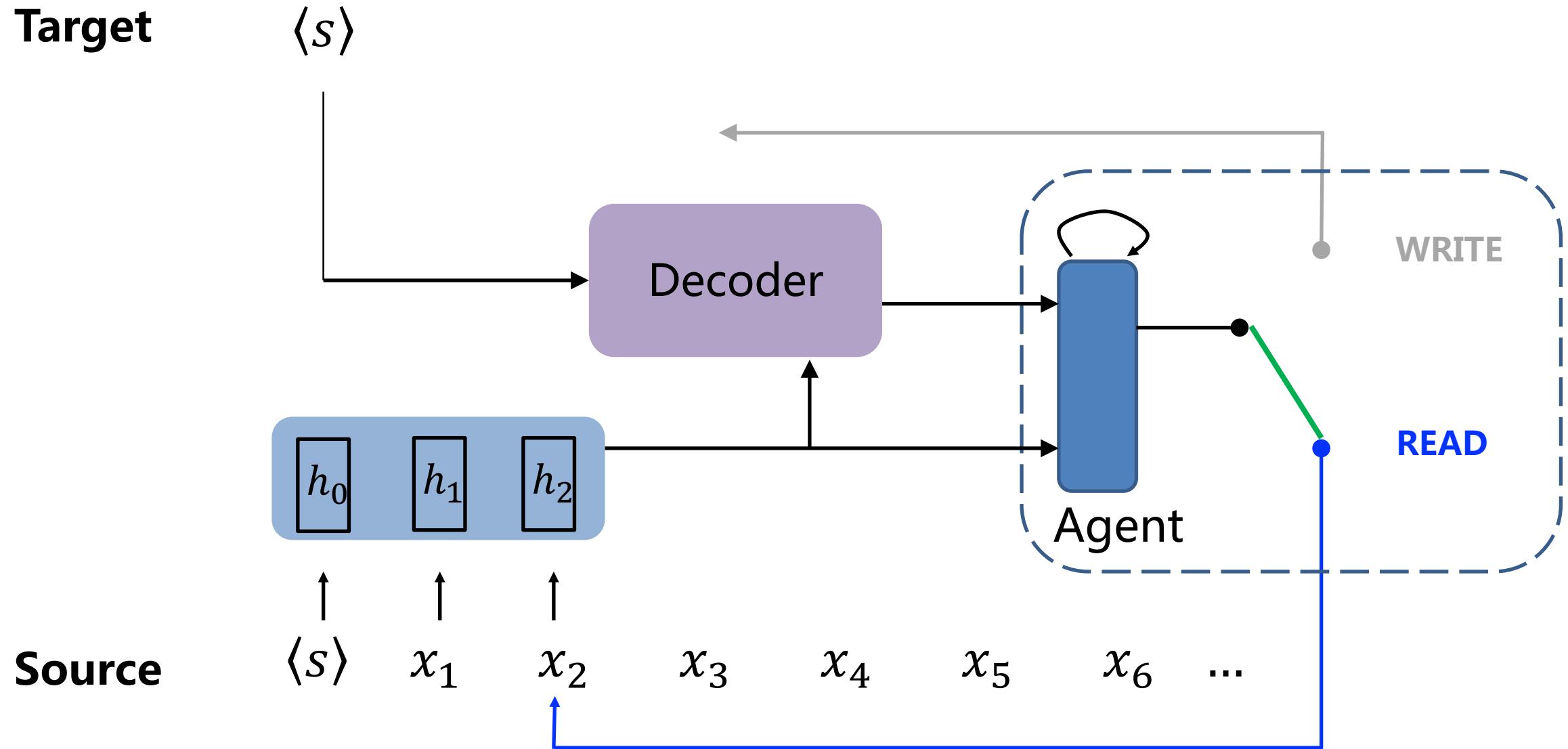
基于强化学习的策略：READ/WRITE



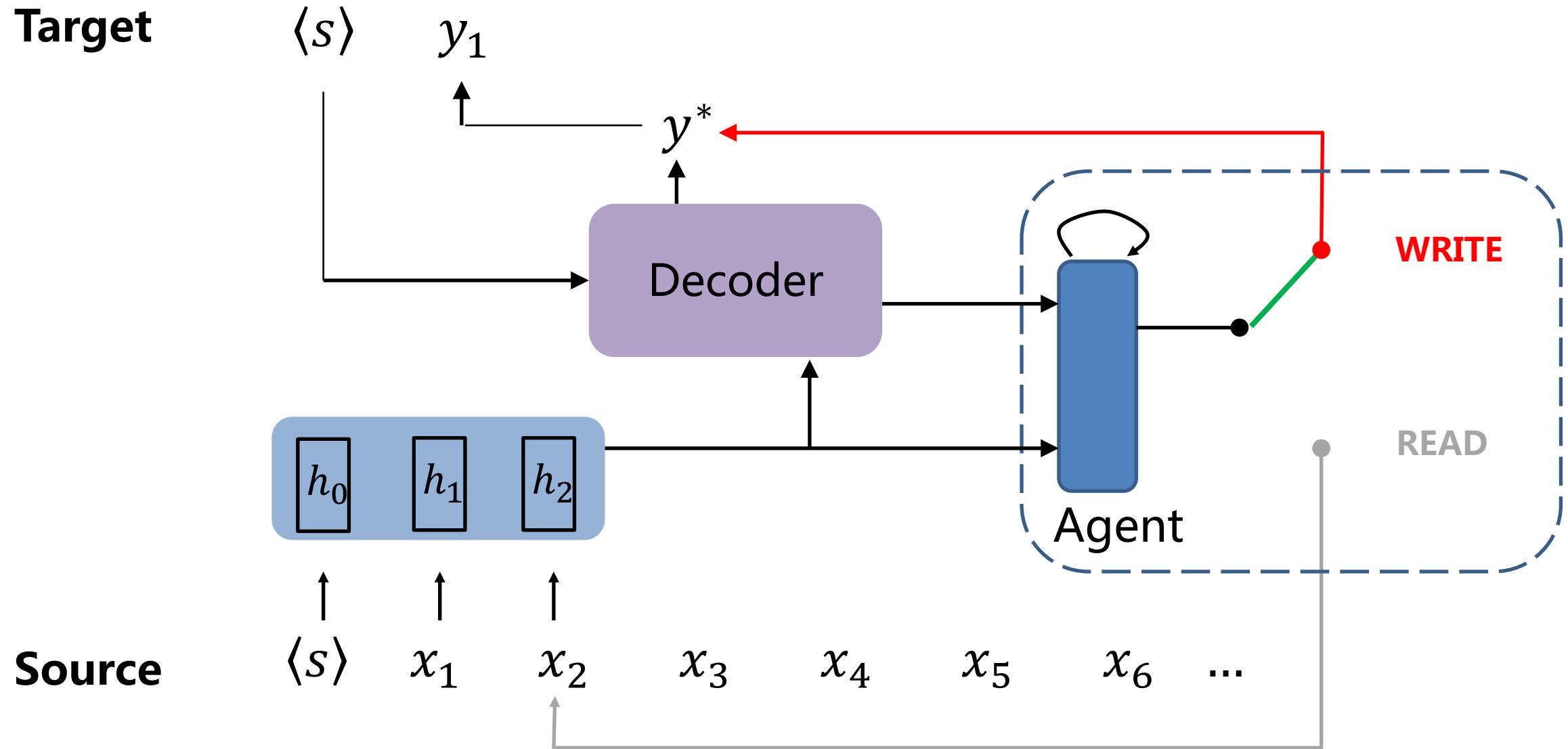
基于强化学习的策略：READ/WRITE



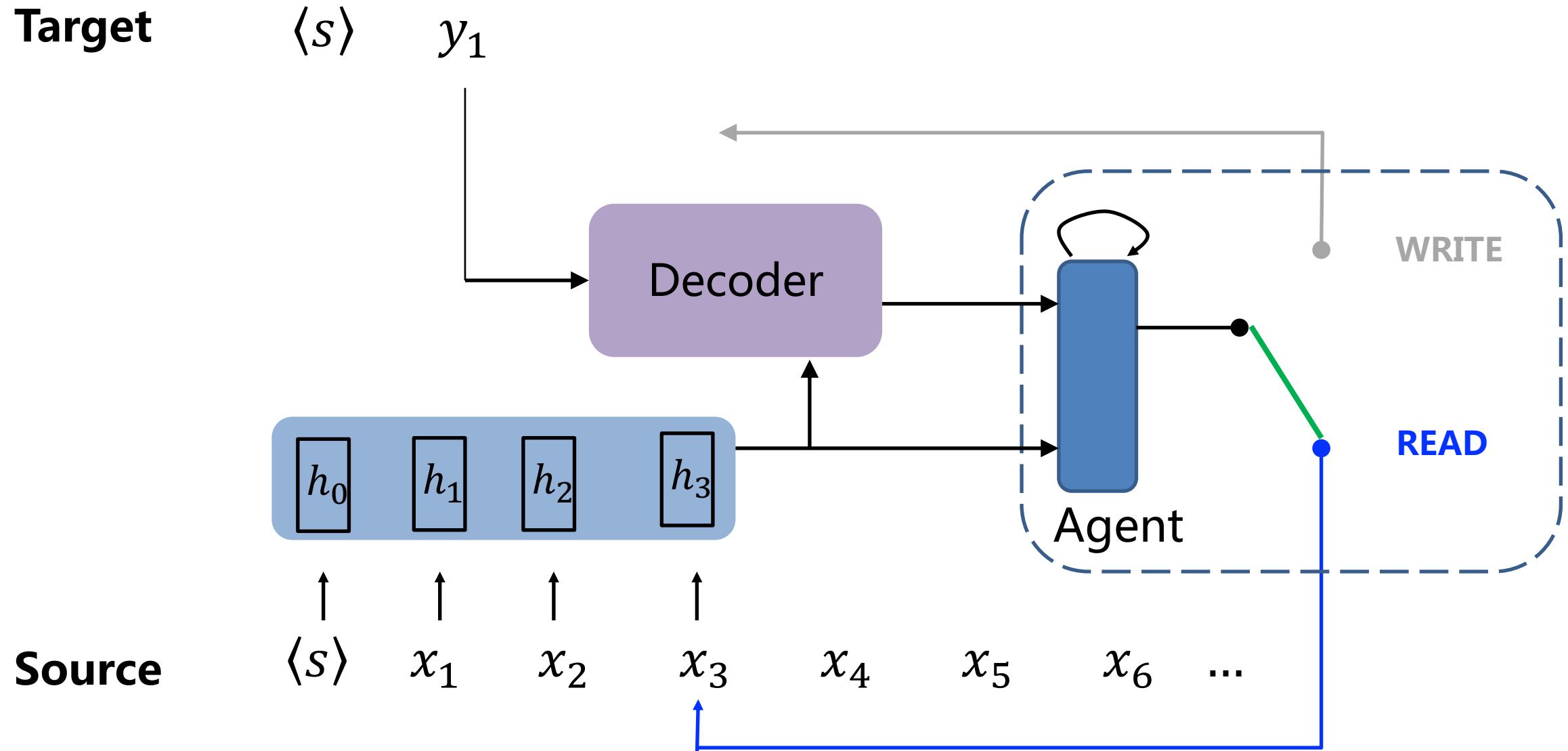
基于强化学习的策略：READ/WRITE



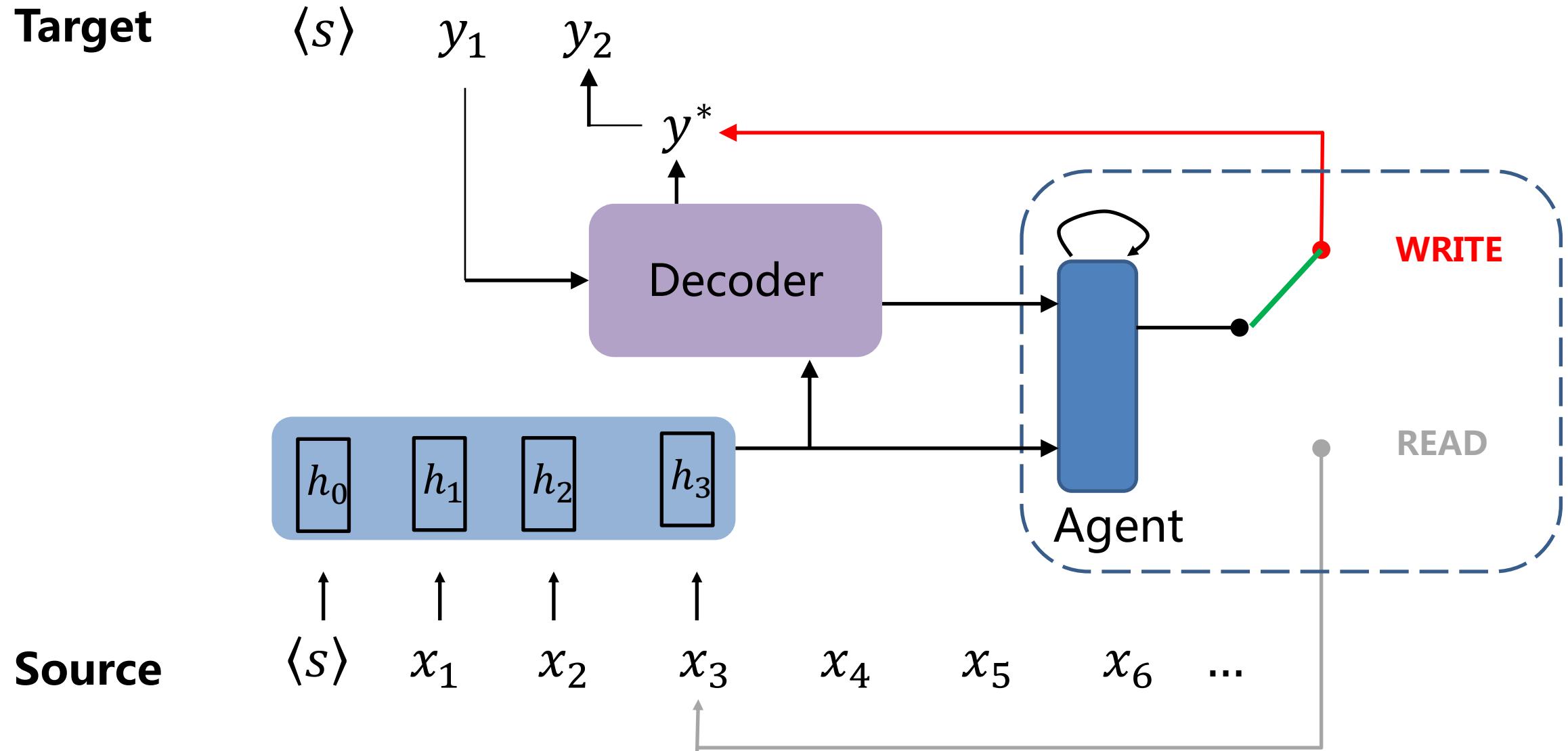
基于强化学习的策略：READ/WRITE



基于强化学习的策略：READ/WRITE



基于强化学习的策略：READ/WRITE



奖励函数 (Reward Function)

翻译质量

$$r_t^Q = \begin{cases} \Delta BLEU^0(Y, Y^*, t) & t < T \\ BLEU(Y, Y^*) & t = T \end{cases}$$

$$r_t = r_t^Q + r_t^D$$

时间延迟

$$r_t^D = \underbrace{\alpha \cdot [\text{sgn}(c_t - c^*) + 1]}_{CW} + \underbrace{\beta \cdot [d_t - d^*]_+}_{AP}$$

级联模型 – 同传策略

- 固定策略 (Fixed Policy) : 读入长度 (翻译单元) 固定
 - Static read-write (Dalvi et al., 2018)
 - STACL (Ma et al., 2018), etc.
- 自适应策略 (Adaptive Policy) : 动态调整翻译单元长度
 - Rule-based (Cho et al., 2016)
 - RL-based (Gu et al., 2017)
 - Supervised policy (Zheng et al., 2019)
 - Meaningful unit (Zhang et al., 2020)
 - MILk (Arivazhagan et al., 2019)
 - Multihead monotonic attention (Ma et al., 2020), etc.

Supervised Learning

目标：预测动作序列

问题：没有标注数据，也没有标注标准

会议 将 在 十月 举行

Action

R W W R W W R R R W W W

Translation

The meeting will be held in October

动作序列标注

id_s

会议

将 在 十月 举行

翻译模型

$rank(\text{The} | x_{\leq id_s}) > r ? \text{READ} : \text{WRITE}$

R

The meeting will be held in October



动作序列标注

$i d_s$

会议

将 在 十月 举行

翻译模型

$rank(\text{The} | x_{\leq id_s}) > r ? \text{READ} : \text{WRITE}$

R W

The meeting will be held in October



动作序列标注

$i d_s$

会议

将 在 十月 举行

翻译模型

$rank(\text{meeting} | x_{\leq i d_s}) > r ? \text{READ: WRITE}$

R W W

The meeting will be held in October



动作序列标注

$i d_s$

会议

将 在 十月 举行

翻译模型

$rank(\text{will} | x_{\leq id_s}) > r ? \text{READ} : \text{WRITE}$

R W W

The meeting will be held in October



动作序列标注

id_s

会议

将

在

十月

举行

翻译模型

$rank(\text{will} | x_{\leq id_s}) > r ? \text{READ} : \text{WRITE}$

R W W R

The meeting will be held in October



动作序列标注

id_s

会议

将

在

十月

举行

翻译模型

$rank(\text{will} | x_{\leq id_s}) > r ? \text{READ} : \text{WRITE}$

R W W R W

The meeting will be held in October



动作序列标注

$i d_s$

会议 将 在 十月 举行

翻译模型

$rank(October|x_{\leq id_s}) > r ? READ : WRITE$

R W W R W W R R R W W W

The meeting will be held in October



训练与解码

$$p(Action) = \operatorname{argmax} \sum p(a_t | o_{\leq t}, a_{<t})$$

$$p(a_t | o_{\leq t}, a_{<t}) > \rho? \text{READ: WRITE}$$

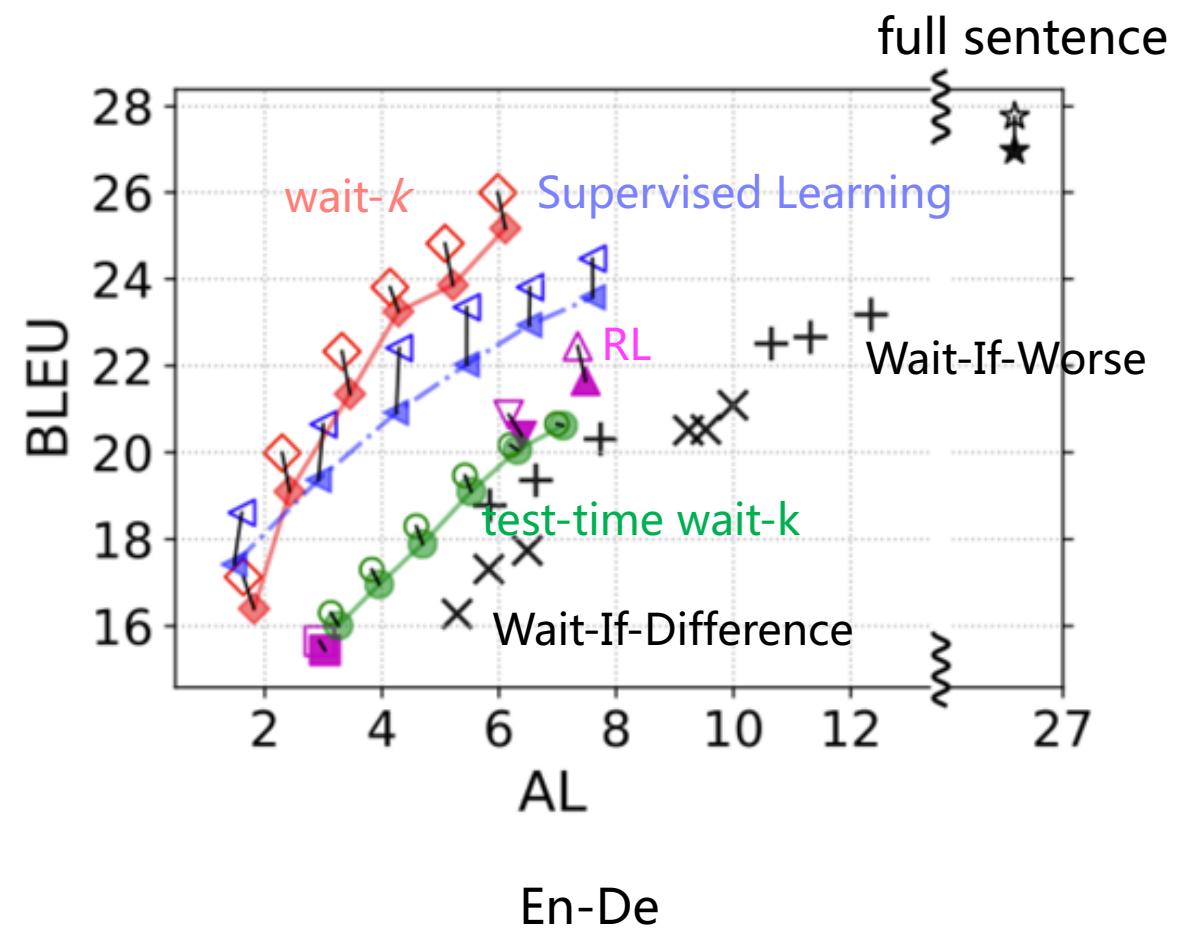
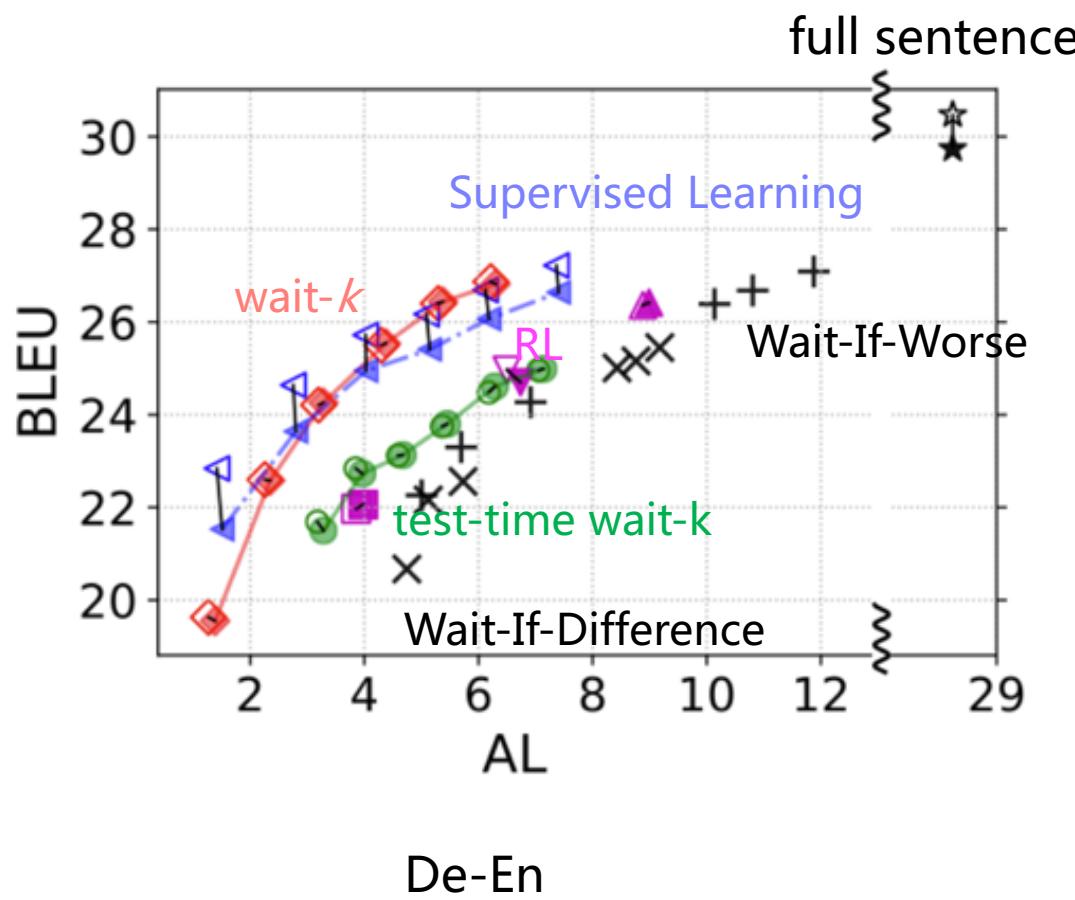
trade off between quality and latency

ρ 越大，倾向于WRITE操作，时延低

ρ 越小，倾向于READ操作，时延高

实验 (EN-DE)

Trained on 4.5M sentence pairs (WMT 15)



基于语义单元（ Meaningful Unit ）的自适应策略

人类同传策略：顺句驱动、语义单元



如何定义语义单元？

- 尽量短，以缩短时延
- 含有较充分的信息（不依赖于后文），保证翻译质量

Meaningful Unit ?

wo zai kan

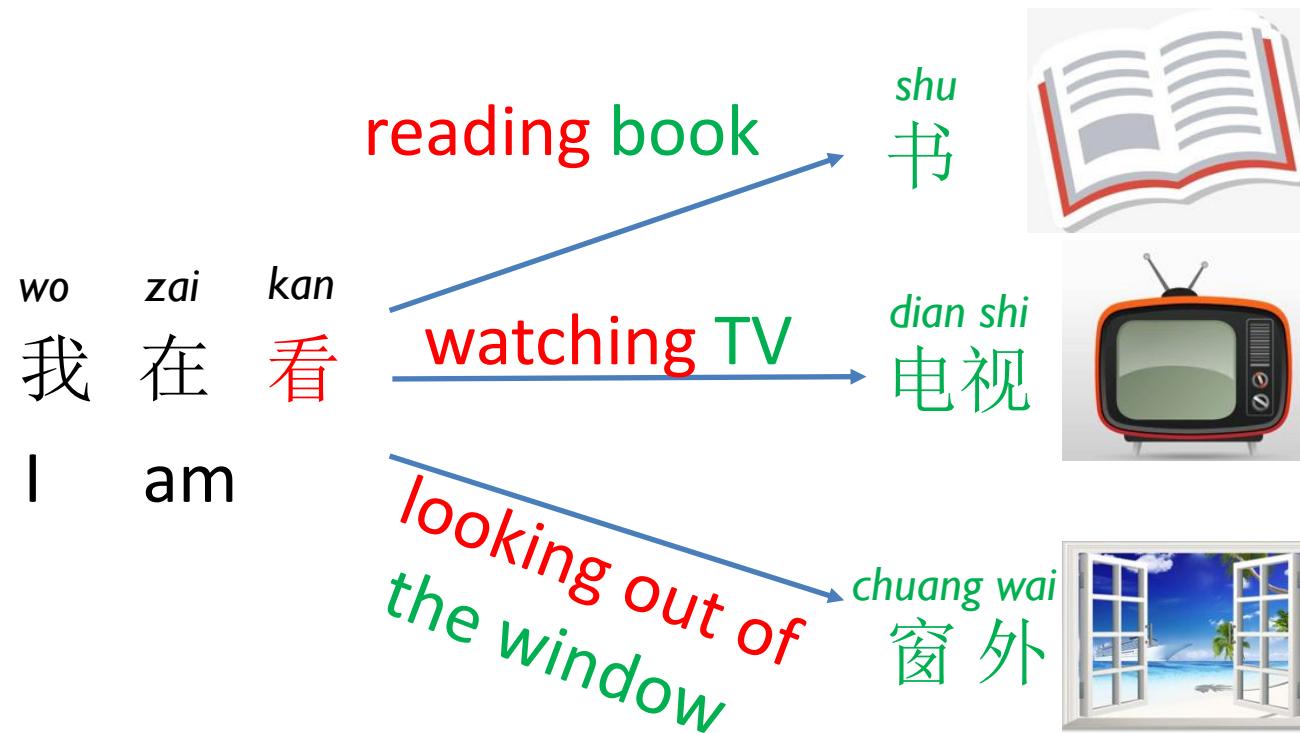
我 在 看

I am



如何定义语义单元？

- 尽量短，以缩短时延
- 含有较充分的信息（不依赖于后文），保证翻译质量



语义单元识别

- 分类问题，单词是否是语义单元边界？
- 基于预训练模型的边界识别 Pre-training & Fine-tuning

语义单元识别

history

buffer

上午

十点

我

去了

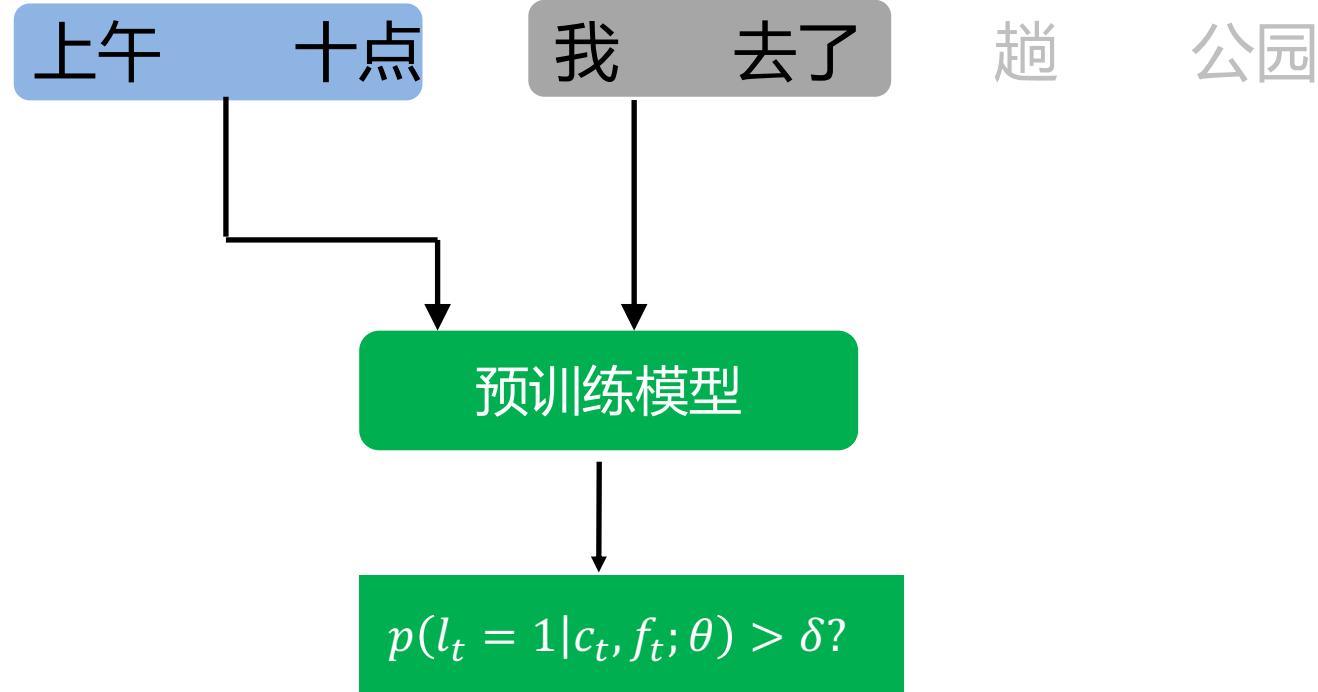
趟

公园

预训练模型

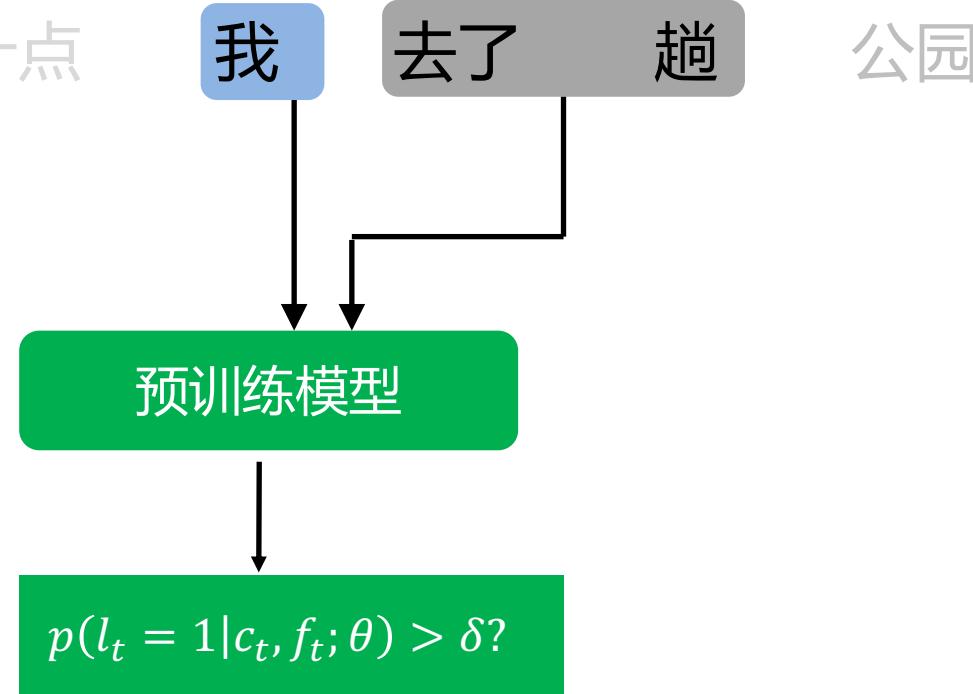
$$p(l_t = 1 | c_t, f_t; \theta) > \delta?$$

语义单元识别



语义单元识别

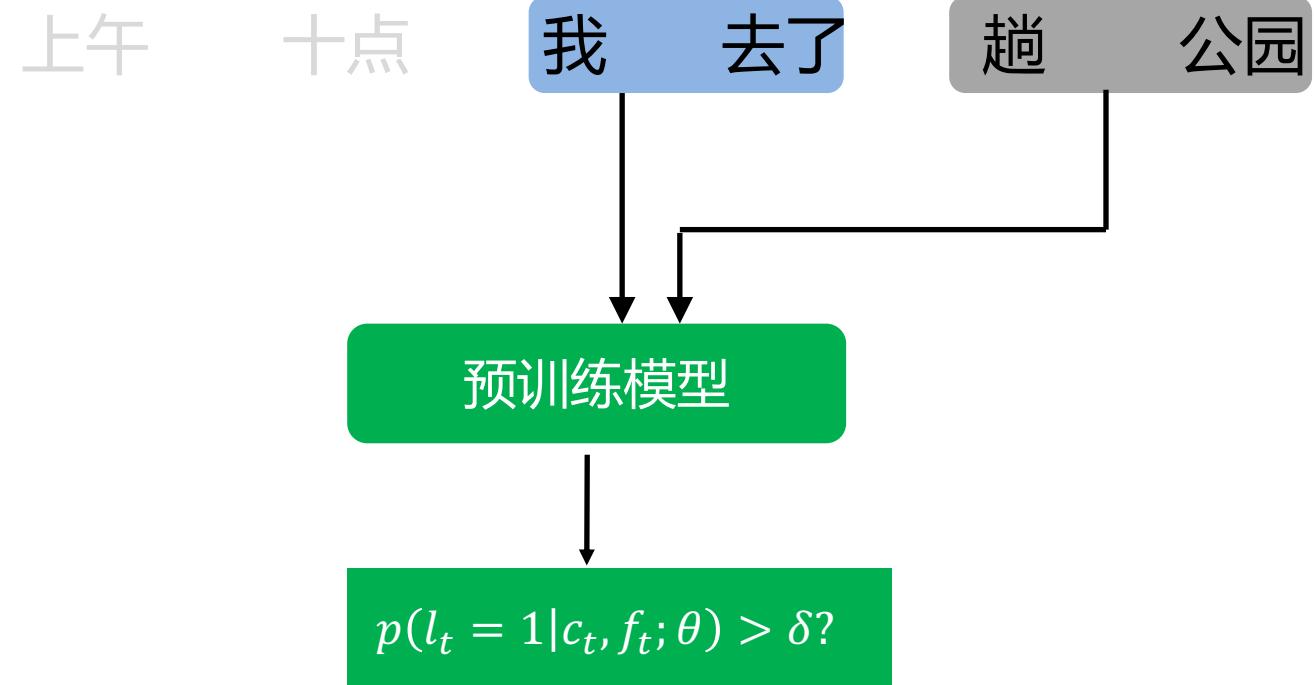
上午 十点 我 去了 趟 公园



上午 十点

At 10:00 a.m.

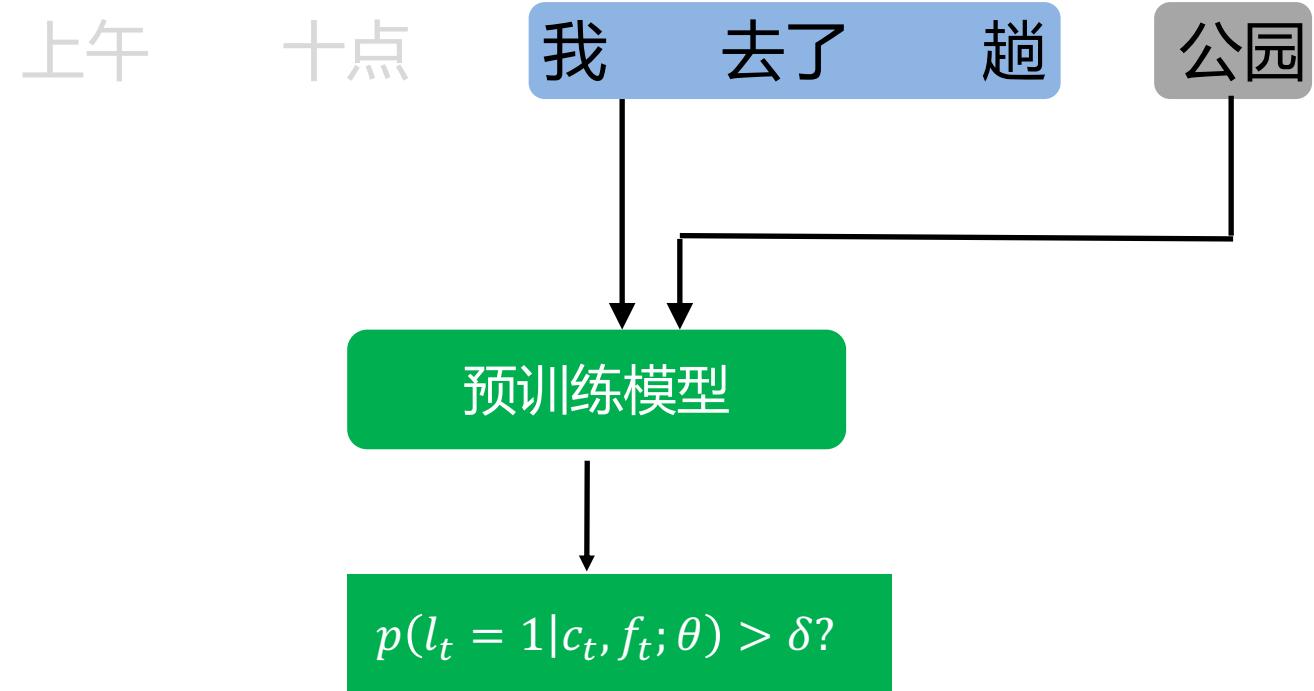
语义单元识别



上午 十点

At 10:00 a.m.

语义单元识别



上午 十点

At 10:00 a.m.

语义单元识别

上午 十点 我 去了 趟 公园

预训练模型

$$p(l_t = 1 | c_t, f_t; \theta) > \delta?$$

上午 十点

我 去了 趟

At 10:00 a.m.



I went to

语义单元识别

上午 十点 我 去了 趟 公园

预训练模型



$$p(l_t = 1 | c_t, f_t; \theta) > \delta?$$

上午 十点

我 去了 趟

公园

At 10:00 a.m.

I went to

the park

构造训练数据

- source prefix whose translation is also a prefix of the full-sentence translation

<i>Source</i>	shàngwǔ	10	diǎn	wǒ	qùle	tàng	gōngyuán
	上午	10	点	我	去了	趟	公园
<i>full sentence translation</i>	At 10 a.m., I went to the park.						

构造训练数据

- source prefix whose translation is also a prefix of the full-sentence translation

<i>Source</i>	shàngwǔ	10	diǎn	wǒ	qùle	tàng	gōngyuán
	上午	10	点	我	去了	趟	公园
<i>full sentence translation</i>	At 10 a.m., I went to						the park.
$M'_{nmt}(x_{\leq 1})$	Morning						

构造训练数据

- source prefix whose translation is also a prefix of the full-sentence translation

<i>Source</i>	shàngwǔ	10	diǎn	wǒ	qùle	tàng	gōngyuán
	上午	10	点	我	去了	趟	公园
<i>full sentence translation</i>	At 10 a.m.,			I went to			the park.
$M'_{nmt}(x_{\leq 1})$	Morning						
$M'_{nmt}(x_{\leq 2})$	Morning 10						

构造训练数据

- source prefix whose translation is also a prefix of the full-sentence translation

<i>Source</i>	shàngwǔ	10	diǎn	wǒ	qùle	tàng	gōngyuán
	上午	10	点	我	去了	趟	公园
<i>full sentence translation</i>	At 10 a.m., I went to the park.						
$M'_{nmt}(x_{\leq 1})$	Morning						
$M'_{nmt}(x_{\leq 2})$	Morning	10					
$M'_{nmt}(x_{\leq 3})$	At 10 a.m.						

构造训练数据

- source prefix whose translation is also a prefix of the full-sentence translation

Source	shàngwǔ	10	diǎn	wǒ	qùle	tàng	gōngyuán
	上午	10	点	我	去了	趟	公园
full sentence translation	At 10 a.m., I went to the park.						
$M'_{nmt}(x_{\leq 1})$	Morning						
$M'_{nmt}(x_{\leq 2})$	Morning	10					
$M'_{nmt}(x_{\leq 3})$	At 10 a.m.						
$M'_{nmt}(x_{\leq 4})$	At 10 a.m.			me			
$M'_{nmt}(x_{\leq 5})$	At 10 a.m.		I went there				
$M'_{nmt}(x_{\leq 6})$	At 10 a.m.		I went to				

构造训练数据

- source prefix whose translation is also a prefix of the full-sentence translation

<i>Source</i>	shàngwǔ	10	diǎn	wǒ	qùle	tàng	gōngyuán	
<i>full sentence translation</i>	上午	10	点	我	去了	趟	公园	
$M'_{nmt}(x_{\leq 1})$	Morning							
$M'_{nmt}(x_{\leq 2})$	Morning	10						
$M'_{nmt}(x_{\leq 3})$	At 10 a.m.							
$M'_{nmt}(x_{\leq 4})$	At 10 a.m.			me				
$M'_{nmt}(x_{\leq 5})$	At 10 a.m.			I went there				
$M'_{nmt}(x_{\leq 6})$	At 10 a.m.			I went to				
$M'_{nmt}(x_{\leq 7})$	At 10 a.m.			I went to		the park		
<i>Extracted MUs</i>	shàngwǔ 10 diǎn			wǒ qùle tàng		gōngyuán		

构造训练数据

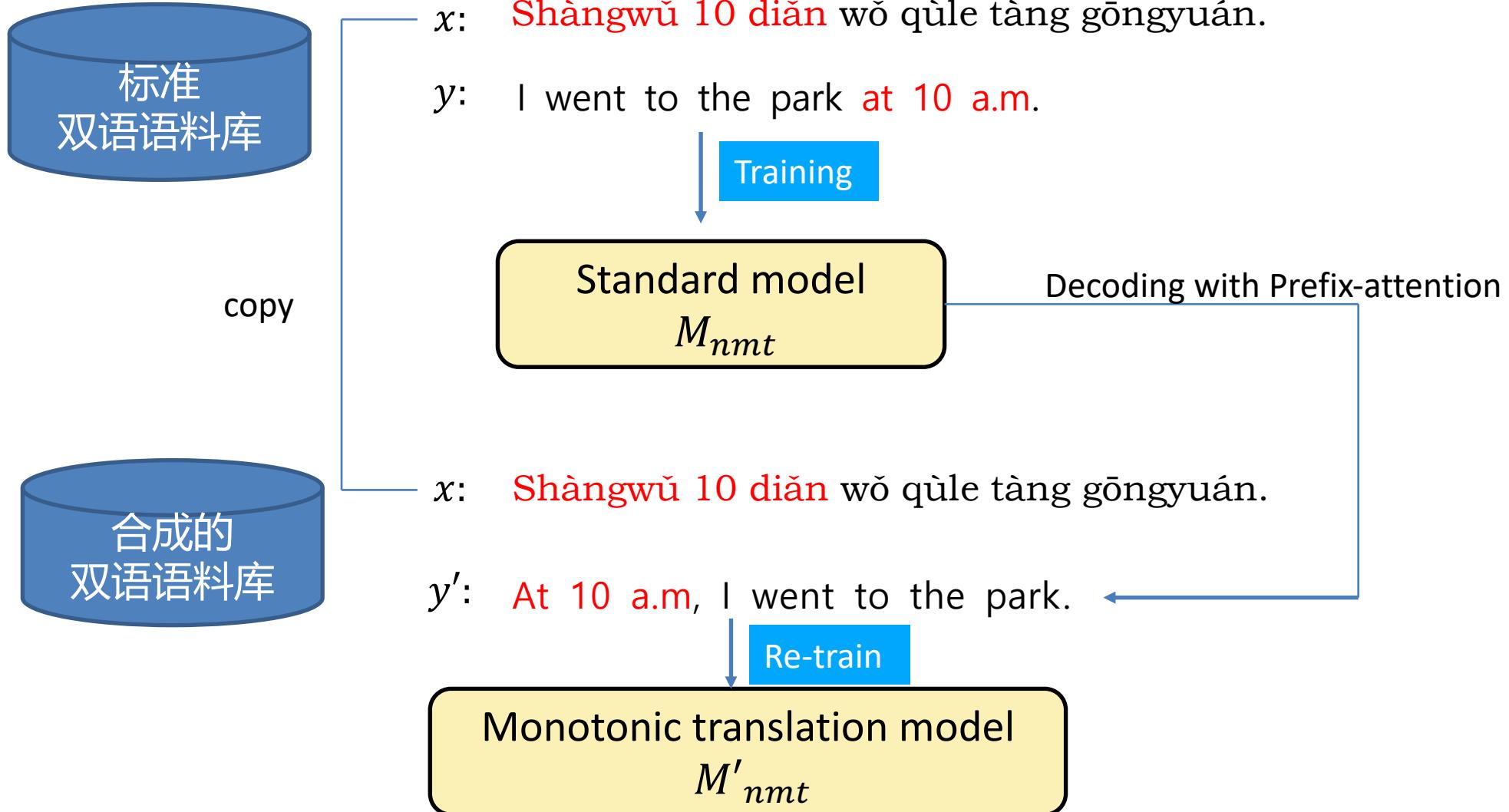
- source prefix whose translation is also a prefix of the full-sentence translation

<i>Source</i>	shàngwǔ	10	diǎn	wǒ	qùle	tàng	gōngyuán
	上午	10	点	我	去了	趟	公园
<i>full sentence translation</i>	At 10 a.m., I went to the park.						
<i>Extracted MUs</i>	shàngwǔ 10 diǎn			wǒ qùle tàng			gōngyuán

- Long distance reorderings in full sentence translation generate long MUs

<i>Source</i>	shàngwǔ	10	diǎn	wǒ	qùle	tàng	gōngyuán
	上午	10	点	我	去了	趟	公园
<i>full sentence translation with long reorderings</i>	I went to the park at 10 a.m.						
<i>Extracted MUs</i>	shàngwǔ	10	diǎn	wǒ	qùle	tàng	gōngyuán

训练数据构建

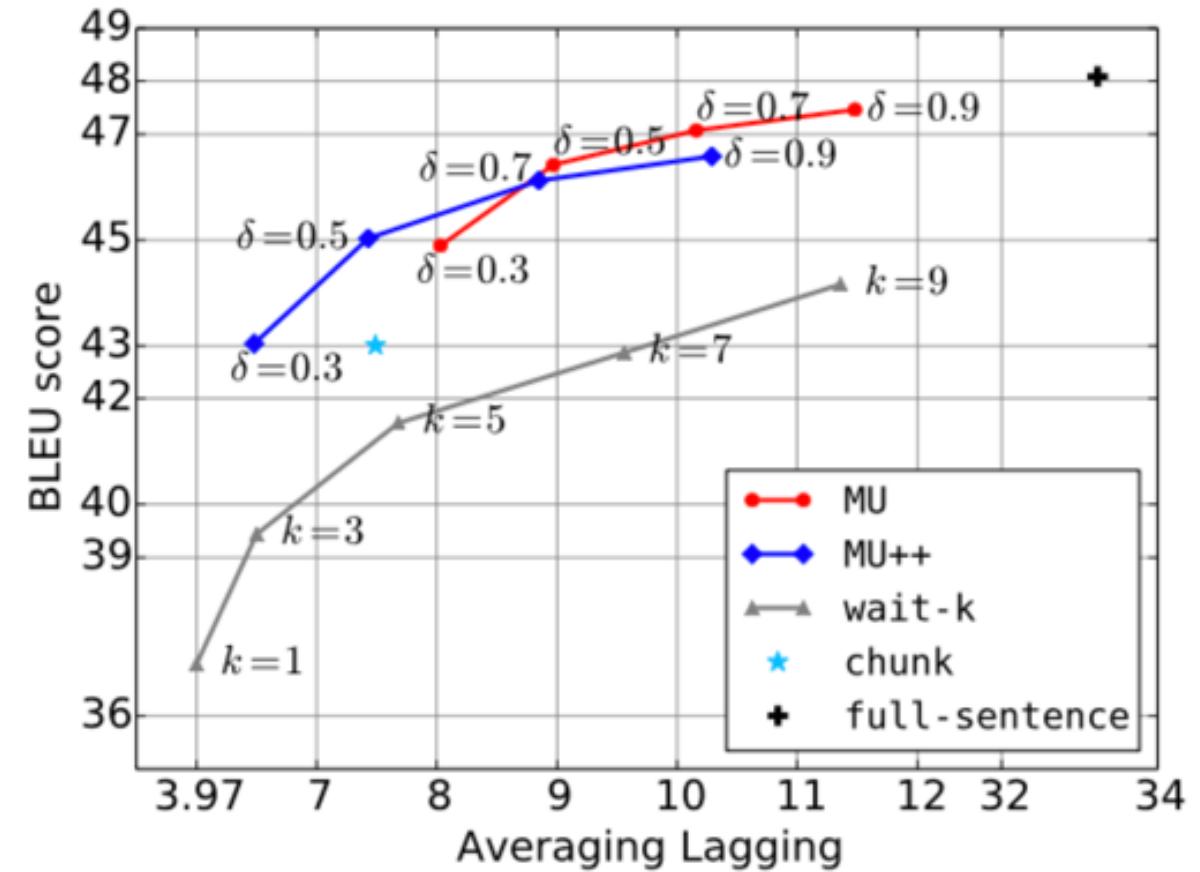


训练数据构建

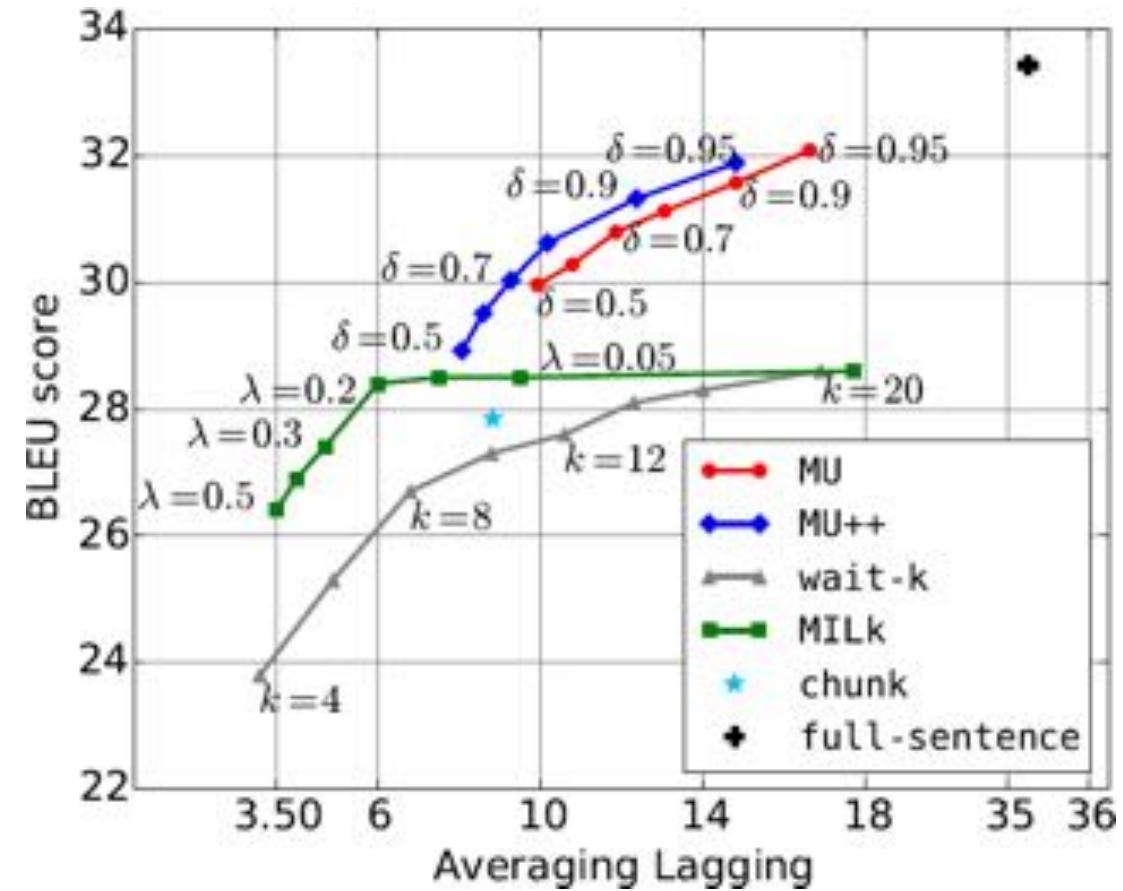
shàngwǔ 10 **diǎn** | wǒ qùle **tàng** | gōngyuán

t	history	future words	MU label
1	shàngwǔ	10 diǎn	0
2	shàngwǔ 10	diǎn wǒ	0
3	shàngwǔ 10 diǎn	wǒ qùle	1
4	shàngwǔ 10 diǎn wǒ	qùle tàng	0
5	shàngwǔ 10 diǎn wǒ qùle	tàng gōngyuán	0
6	shàngwǔ 10 diǎn wǒ qùle tàng	gōngyuán	1
7

实验结果



中英 (NIST 测试集)



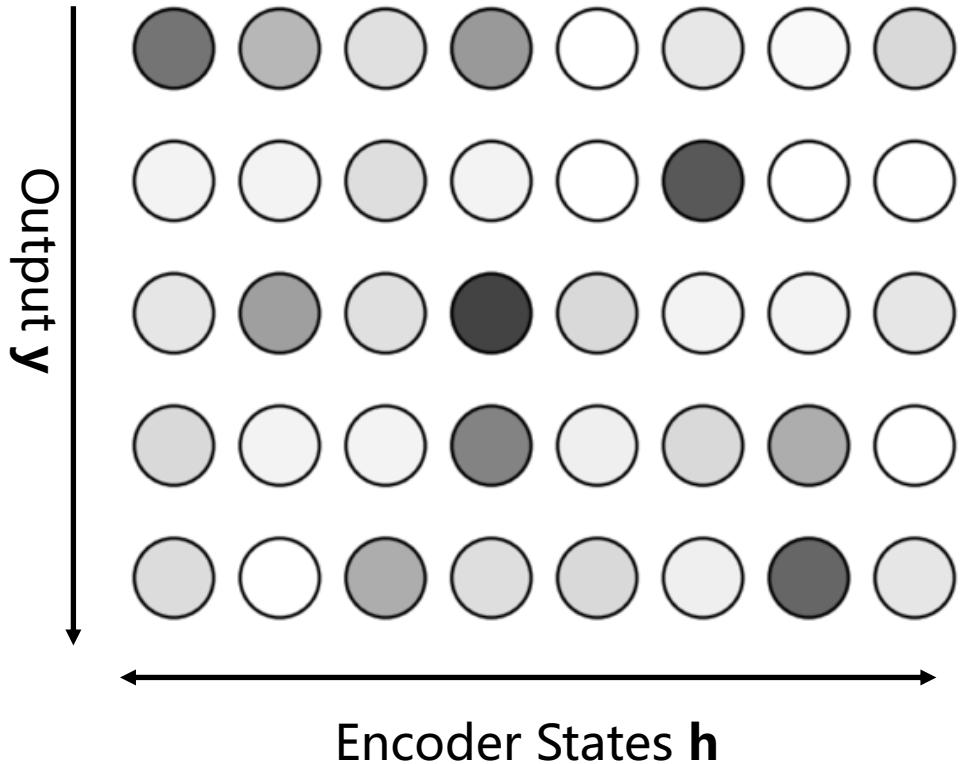
德英 (WMT 测试集)

级联模型 – 同传策略

- 固定策略 (Fixed Policy) : 读入长度 (翻译单元) 固定
 - Static read-write (Dalvi et al., 2018)
 - STACL (Ma et al., 2018), etc.
- 自适应策略 (Adaptive Policy) : 动态调整翻译单元长度
 - Rule-based (Cho et al., 2016)
 - RL-based (Gu et al., 2017)
 - Supervised policy (Zheng et al., 2019)
 - Meaningful unit (Zhang et al., 2020)
 - MILk (Arivazhagan et al., 2019)
 - Multihead monotonic attention (Ma et al., 2020), etc.

Softmax Attention V.S. Monotonic Attention (单调注意力)

Softmax Attention



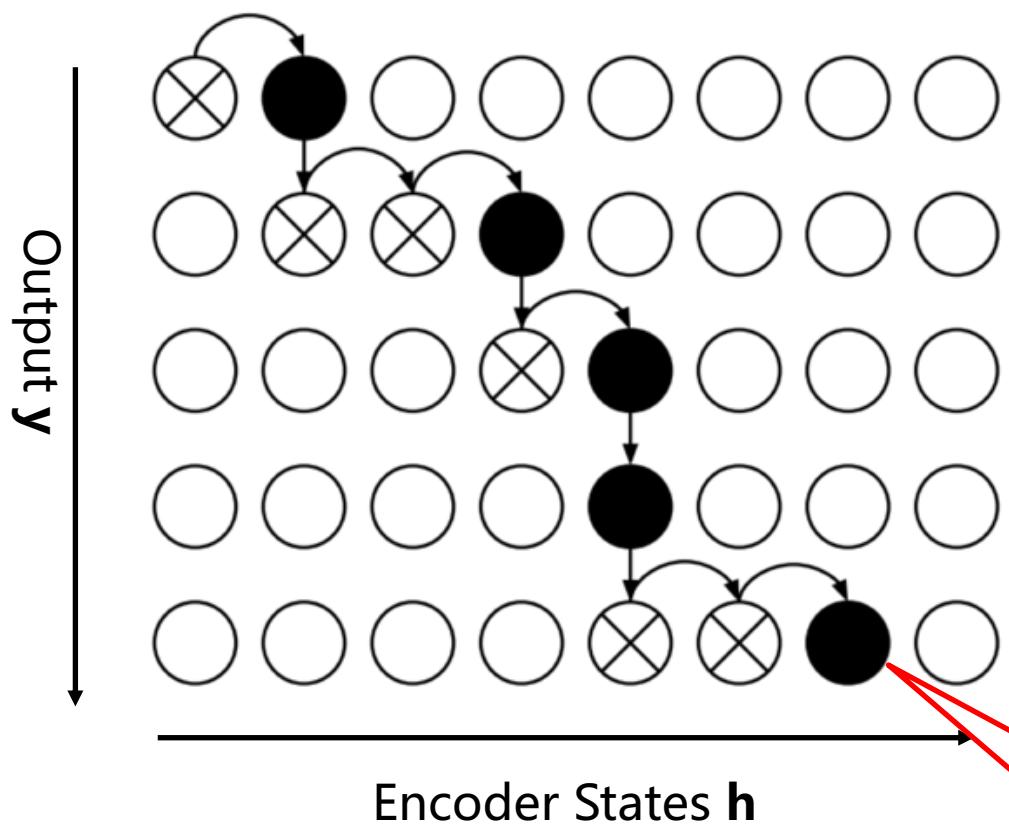
$$e_{i,j} = \text{FeedForward}(s_{i-1}, h_j)$$

$$\alpha_{i,j} = \frac{\exp(e_{i,j})}{\sum_{k=1}^T \exp(e_{i,k})}$$

$$c_i = \sum_{j=1}^{|\mathbf{x}|} \alpha_{i,j} h_j$$

Softmax Attention V.S. Monotonic Attention (单调注意力)

Monotonic Attention



$$e_{i,j} = a(s_{i-1}, h_j)$$

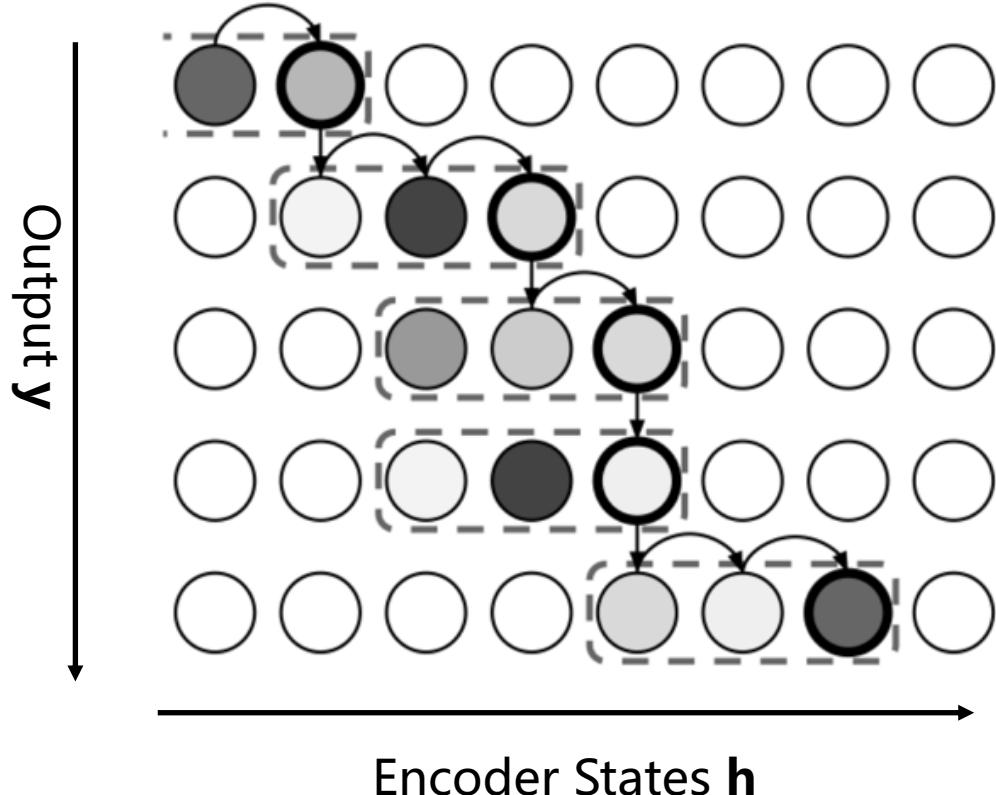
$$p_{i,j} = \sigma(e_{i,j})$$

$$z_{i,j} \sim \text{Bernoulli}(p_{i,j})$$

$$z_{i,j} = \begin{cases} 0 & \otimes \\ 1 & \bullet \end{cases} \quad \begin{array}{l} \text{Read} \\ \text{Write} \end{array}$$

Hard Attention

Monotonic Chunkwise Attention

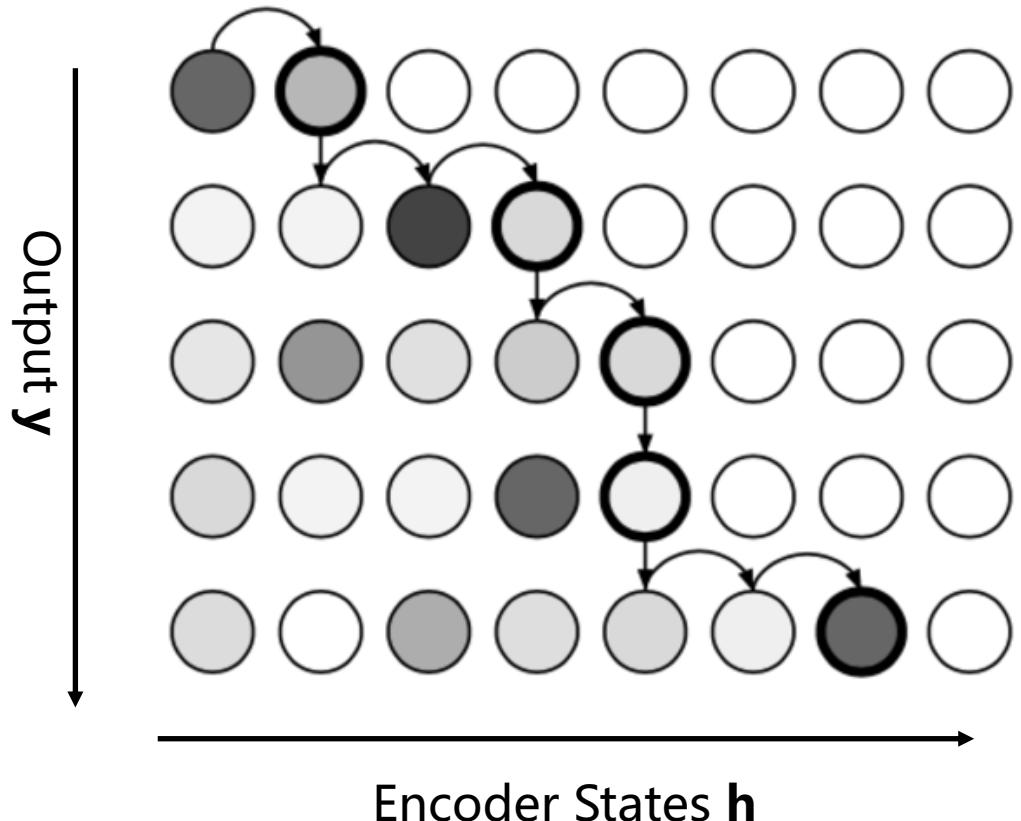


$$v = t_i - w + 1$$

$$u_{i,k} = \text{ChunkEnergy}(s_{i-1}, h_k), k \in \{v, v+1, \dots, t_i\}$$

$$c_i = \sum_{k=v}^{t_i} \frac{\exp(u_{i,k})}{\sum_{l=v}^{t_i} \exp(u_{i,l})} h_k$$

Monotonic Infinite Lookback Attention



$$u_{i,k} = \text{SoftmaxEnergy}(h_k, s_{i-1})$$

$$c_i = \sum_{j=1}^{t_i} \frac{\exp(u_{i,j})}{\sum_{l=1}^{t_i} \exp(u_{i,l})} h_j$$

提纲

- 同声传译
- 主要挑战
- 主流方法
 - 级联模型
 - 端到端模型
- 开放数据
- 产品及应用
- 总结及展望

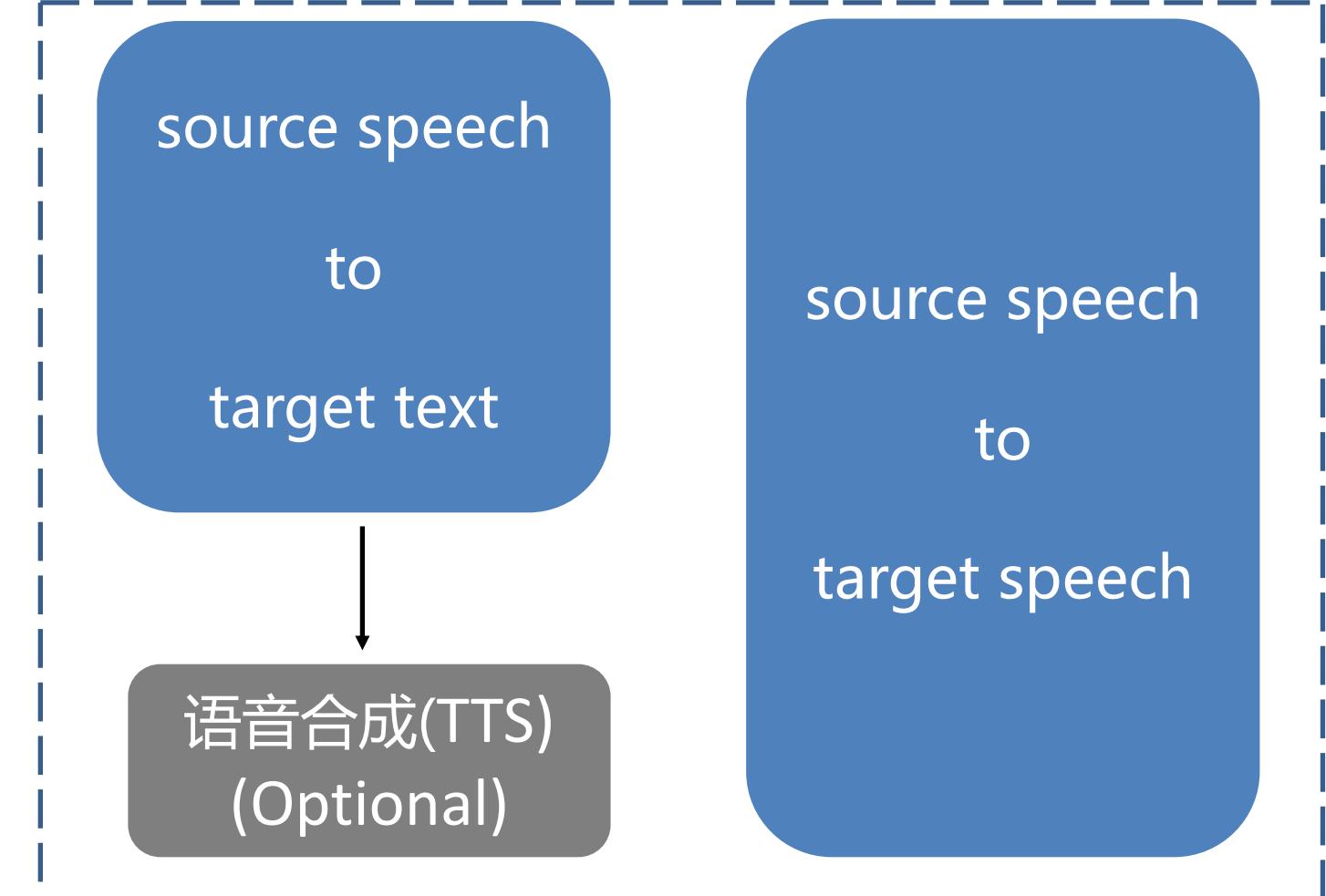
端到端模型

优点

- 缓解错误传递
- 降低时间延迟

挑战

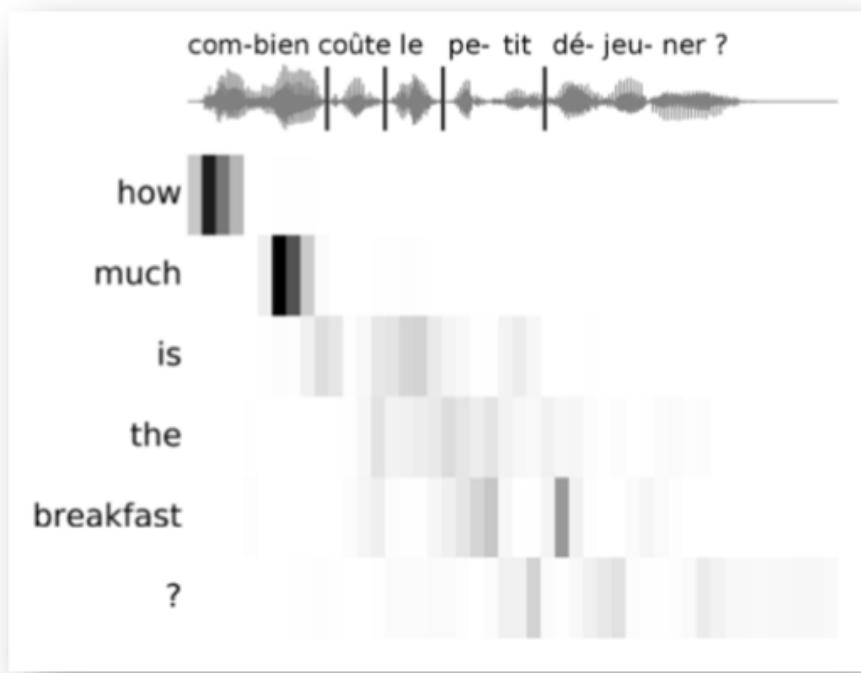
- 缺乏训练数据
- 异构知识共享



语音到文本端到端模型

French→English BTEC corpus

训练数据仅2万句对



源语言语音信号



编码器

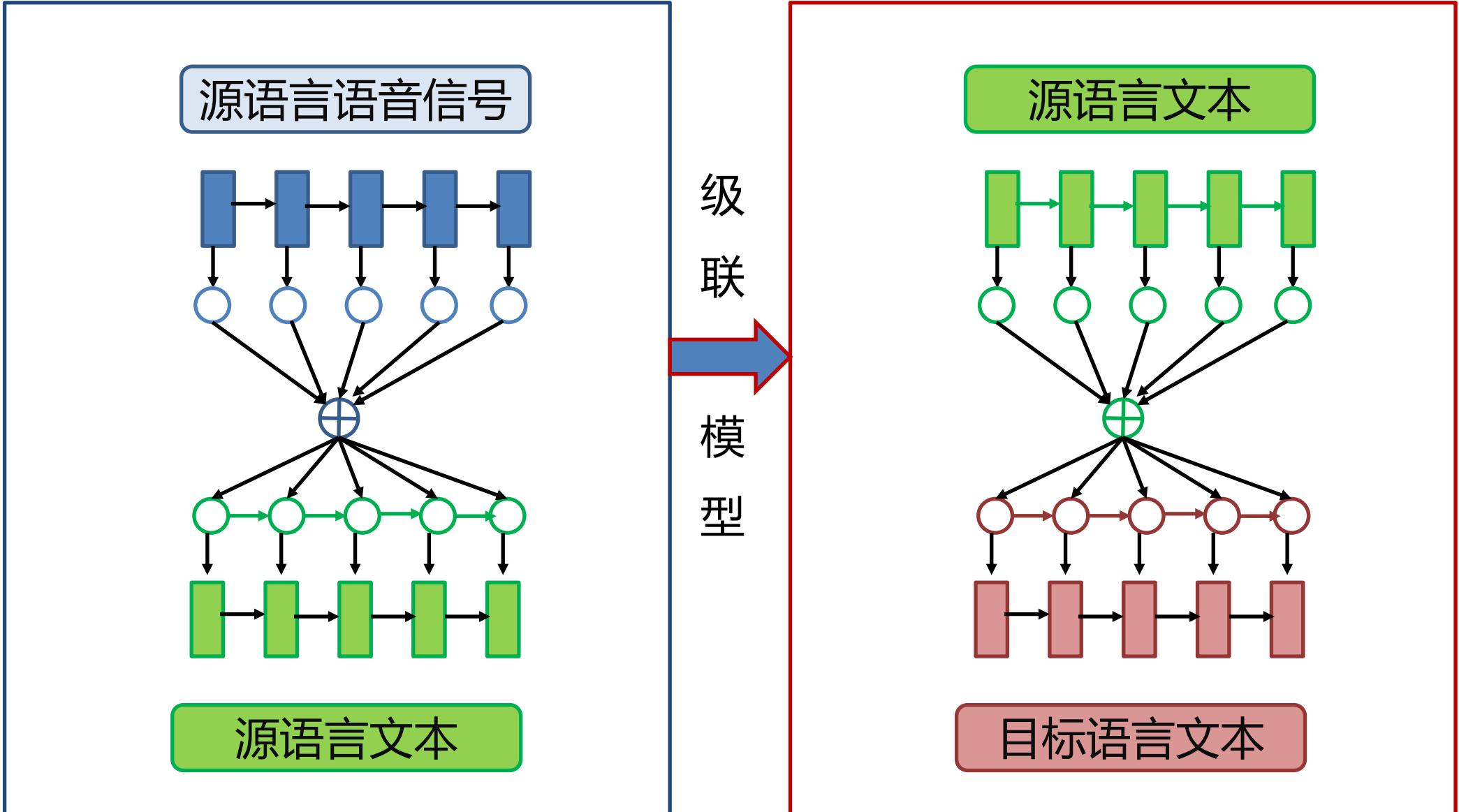


解码器

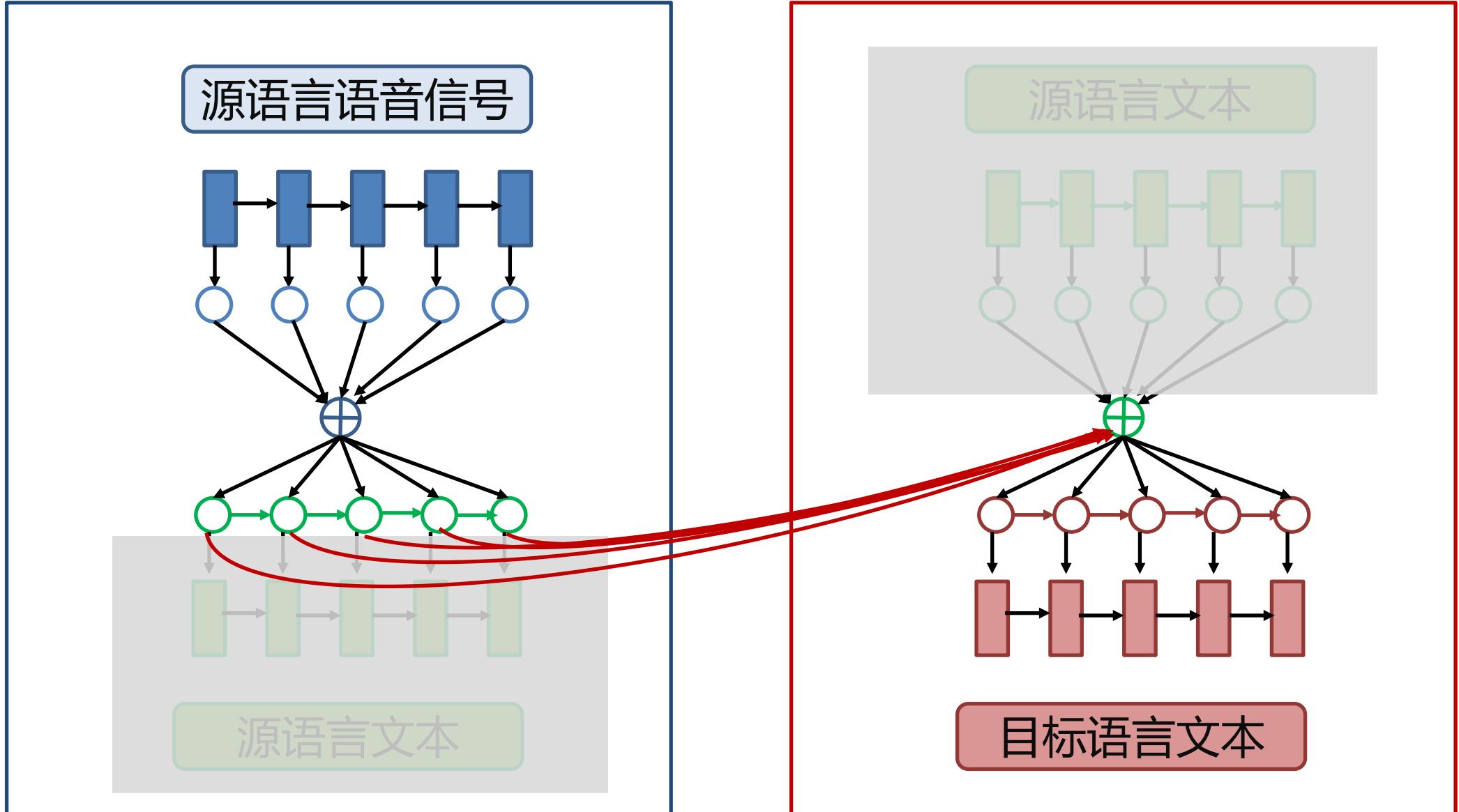


目标语言文本

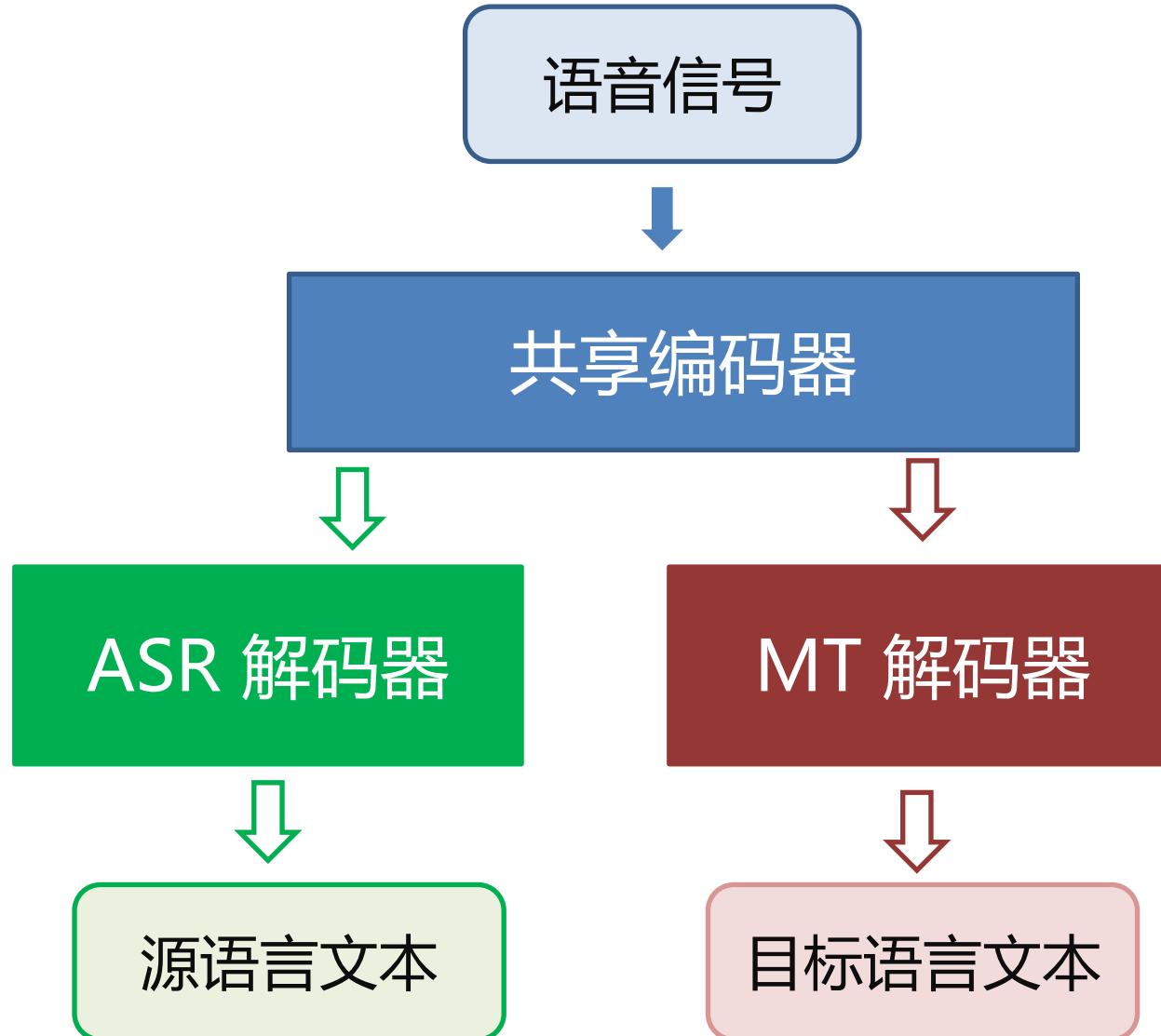
两阶段 (two-stage) 模型



两阶段 (two-stage) 模型



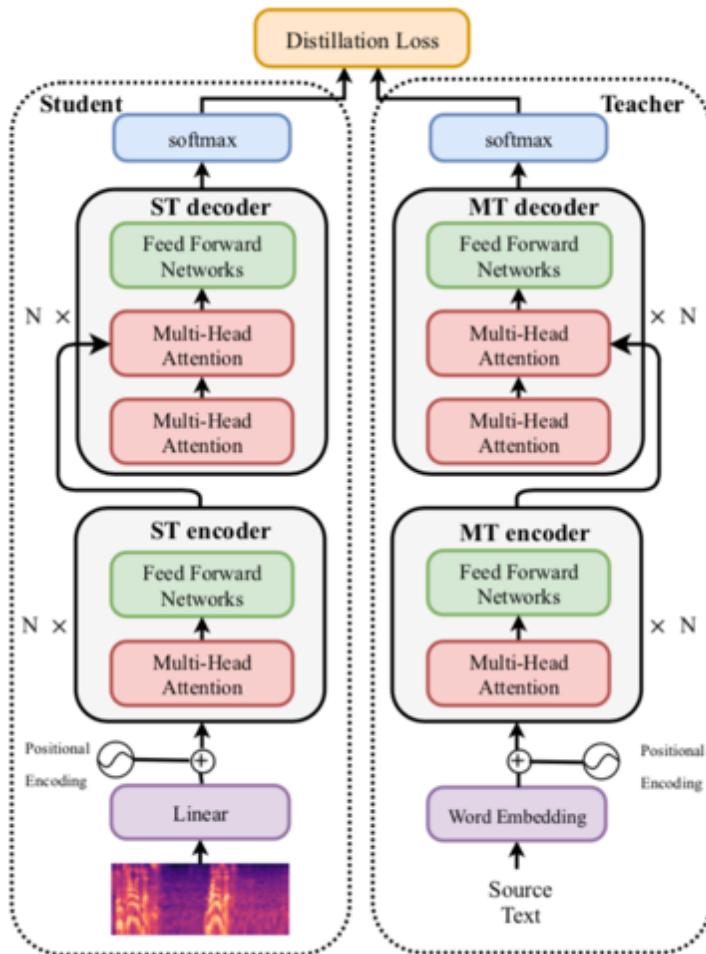
多任务学习



基于知识蒸馏的模型

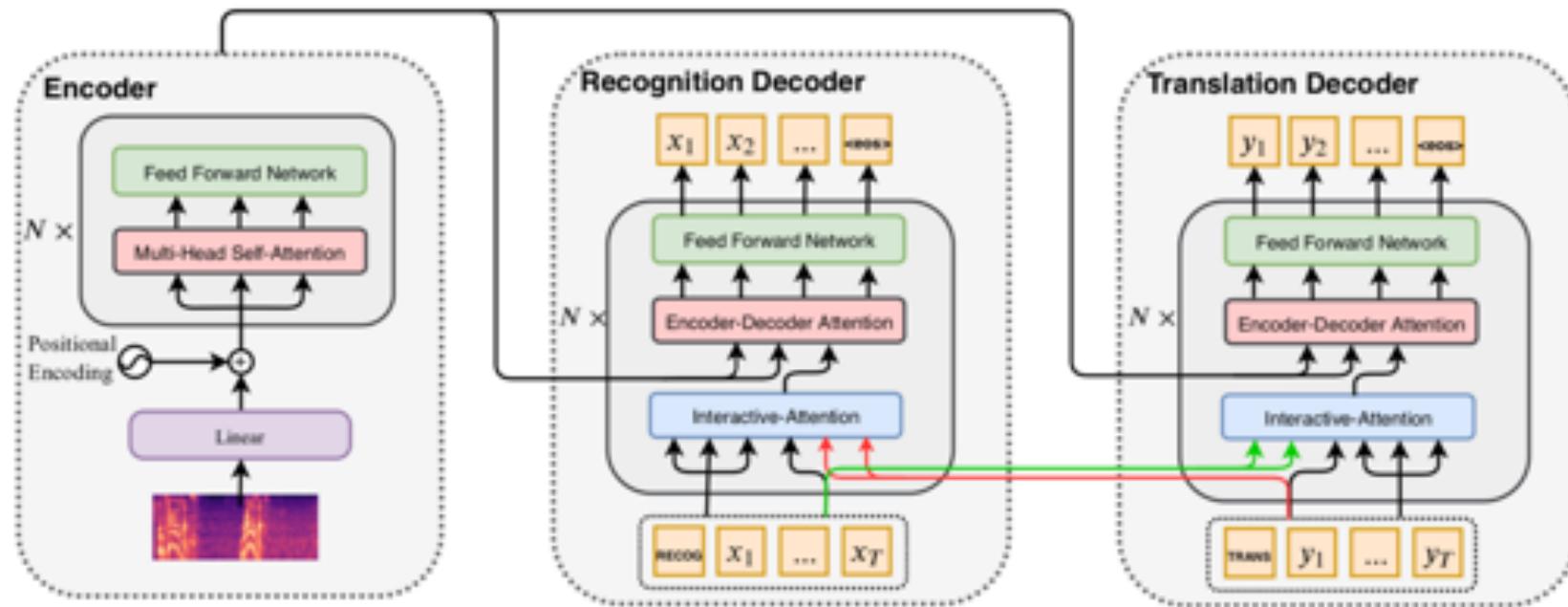
$$L_{\text{ALL}}(D; \theta; \theta_T) = (1 - \lambda)L_{\text{ST}}(D; \theta) + \lambda L_{\text{KD}}(D; \theta, \theta_T)$$

Student
Speech-Text
Translation Model



Teacher
Text-Text Translation
Model

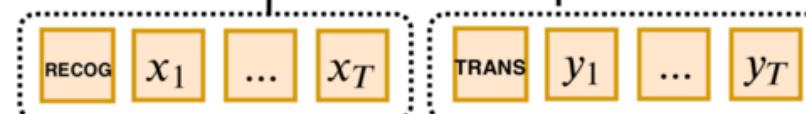
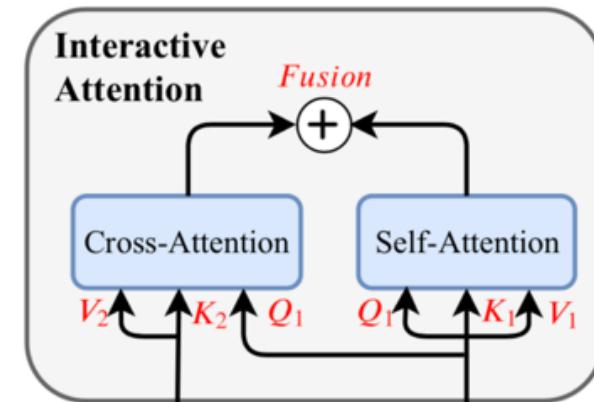
语音识别和翻译同步交互解码



$$\mathbf{H}_{\text{final}} = \mathbf{H}_{\text{self}} + \lambda * \mathbf{H}_{\text{cross}}$$

$$\mathbf{H}_{\text{self}} = \text{Attention}(\mathbf{Q}_1, \mathbf{K}_1, \mathbf{V}_1)$$

$$\mathbf{H}_{\text{cross}} = \text{Attention}(\mathbf{Q}_1, \mathbf{K}_2, \mathbf{V}_2)$$



实验

- Training set: TED corpus
 - Audio: English Speech
 - Translation pairs: English to German (De), French (Fr), Chinese (Zh), and Japanese (Ja)
- Dev and Test sets
 - IWSLT tst2014 / IWSLT tst2015

Model	En-De		En-Fr		En-Zh		En-Ja	
	WER(↓)	BLEU(↑)	WER(↓)	BLEU(↑)	WER(↓)	BLEU(↑)	WER(↓)	BLEU(↑)
Text MT	/	22.19	/	30.68	/	25.01	/	22.93
Pipeline	16.19	19.50	14.20	26.62	14.20	21.52	14.21	20.87
E2E	16.19	16.07	14.20	27.63	14.20	19.15	14.21	16.59
Multi-task	15.20	18.08	13.04	28.71	13.43	20.60	14.01	18.73
Two-stage	15.18	19.08	13.34	30.08	13.55	20.99	14.12	19.32
Interactive	14.76	19.82	12.58	29.79	13.38	21.68	13.91	20.06

提纲

- 同声传译
- 主要挑战
- 主流方法
 - 级联模型
 - 端到端模型
- 数据、鲁棒性模型及评价方法
- 产品及应用
- 总结及展望



开放数据集

同传数据集

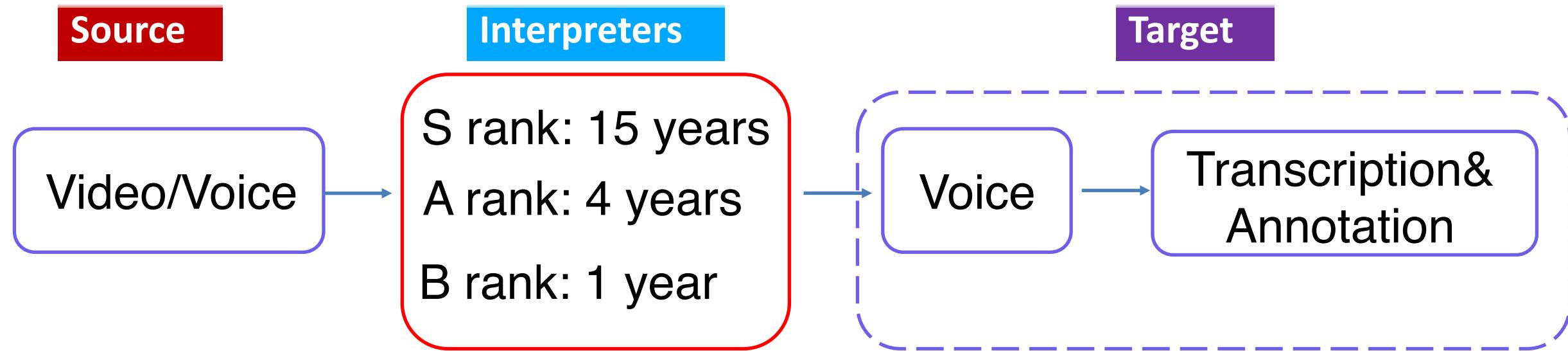
<i>Speech Translation</i>	Languages	Hours
F-C (2013)	Es→En	38
KIT-Disfluency (2014)	De→En	13
BTEC (2016)	En→Fr	17
MSLT V1.0 (2016)	En↔Fr/De	23
	En→Zh/Jp	6
MSLT V1.1 (2017)	Zh→En	5
	Jp →En	9
Travel (2017)	Am→En	8
Aug-LibriSpeech (2018)	En→Fr	236
MuST-C (2019)	En→8 Euro langs	3617
Europarl-ST (2020)	9 Euro langs	1642
Covost (2020a; 2020b)	En↔21 langs	2880
<i>Simultaneous Translation</i>	Languages	Hours
CIAIR (2004)	En↔Jp	182
EPPS (2009)	En↔Es	217
Simul-Trans (2014)	En↔Jp	22
BSTC (ours)	Zh→En	68

NAIST Corpus

Language	English-Japanese / Japanese-English
Domain	Academic Lectures,, News, General
Source Lang. Material	TED, CNN, CSJ(corpus of Spontaneous Japanese), NHK
Total Words	22 hours (387K words of transcribed data)
Link:	https://ahcweb01.naist.jp/resource/stc/

NAIST Corpus

Data	Domain	Format	Lang	Number	Minutes (avg.)	Words (avg.)
TED (S rank)	Lectures	Video	English	46	558 (12.1)	98,034 (2,131)
TED (A, B rank)	Lectures	Video	English	34	415 (12.2)	70,228 (2,066)
CSJ	Lectures	Voice	Japanese	30	326 (10.9)	85,042 (2,835)
CNN	News	Voice	English	8	27 (3.4)	4,639 (580)
NHK	News	Voice	Japanese	10	16 (1.8)	4,121 (412)



NAIST Corpus

Example of a transcript in English and Japanese

	Start-End Time
ID	0001 - 00:20:393 - 00:25:725
Content	<p>So I'm going to present, first of all, the background of my research and purpose of it</p> <p>0002 - 00:26:236 - 00:27:858 and also analytical methods.</p> <p>0003 - 00:28:397 - 00:30:828 Then (F ah) talk about my experiment.</p>
Discourse tags	<p>0001 - 00:44:107 - 00:45:043 本日は<H></p> <p>0002 - 00:45:552 - 00:49:206 みなさまに(F え)難しい話題についてお話ししたいと思います。</p> <p>0003 - 00:49:995 - 00:52:792 (F え)みなさんにとっても意外と身近な話題です。</p>

Example of comparing the translation and simultaneous interpretation

Source	but this understates the seriousness of this particular problem because it doesn't show the thickness of the ice	
Reference (translator)	しかし / もっと深刻な / 問題 / というのは / 実は / 氷河の厚さなのです <i>but / more serious / problem / is / in fact / the thickness of the ice</i>	
Reference (S rank)	しかし / これ本当は / もっと深刻で / 氷の厚さまでは / 見せてないんですね <i>but / this is really / more serious and / the thickness of the ice / it isn't shown</i>	15 years
Reference (A rank)	この / 本当に / 問題に / なっているのは / 氷の厚さです <i>this / real / problem / becoming is / the thickness of the ice</i>	4 years
Reference (B rank)	この / 問題は <i>this / problem is</i>	1 year

CIAIR Simultaneous Interpretation Corpus

Language	English Japanese
Domain	Monologue Speech : economics, history, culture, etc. Dialogue Speech : Travel conversation (airports and hotels)
Total Length	Monologue Speech (Speaker) : 21.5 hours Dialogue Speech (Speaker) : 56 hours
Link:	http://shachi.org/resources/3270

CIAIR Simultaneous Interpretation Corpus

Time	05'106"-09'158"	09'558"-13'654"			
Lecture's Utterance	The theme for this speech is going to be	east coast America versus West coast America			
Interpreter's utterance		次の テーマですが これが アメリカの東海岸対 西海岸			
Time	06'232"- 07'103"	07'344"- 08'387"	09'048"- 09'555"	10'199"- 12'568"	12'900"- 14'471"

CIAIR Simultaneous Interpretation Corpus

Monologue Speech		No. of words	No. of utterance	Recording time (min)
Speaker	English	90249	8422	695
	Japanese	84278	6529	597
	Total	174527	14951	1292
Interpreter	E-J	266050	25507	1639
	J-E	127991	16083	1255
	Total	394041	41590	2904
Sum Total		568568	56541	4196

CIAIR Simultaneous Interpretation Corpus

Dialogue Speech		No. of words	No. of utterance	Recording time (min)
Speaker	English	107850	14223	1678
	Japanese	106258	16485	1678
	Total	214108	30708	3356
Interpreter	E-J	116776	15286	1678
	J-E	91743	13719	1678
	Total	208519	29005	3356
Sum Total		422627	59713	6712

EPIC: European Parliament Interpreting Corpus

Language	English, Italian, Spanish
Domain	public domain
Source Lang.	Europe by Satellite (EbS) TV channel
Material	
Total Words	357 speeches (18 hours, 177K words)
Link:	https://corpora.dipintra.it

EPIC: European Parliament Interpreting Corpus

	sub-corpus	n. of speeches	total word count	% of EPIC
Original Speeches (en, it, es)	Org-en	81	42705	24%
	Org-it	17	6765	3.8%
	Org-es	21	14468	8.2%
Simultaneously Interpreted Speeches	Int-it-en	17	6708	3.8%
	Int-es-en	21	12995	7.3%
	Int-en-it	81	35765	20.1%
	Int-es-it	21	12833	7.2%
	Int-en-es	81	38435	21.6%
	Int-it-es	17	7073	4%
TOTAL		357	177748	100%

MuST-C: a Multilingual Speech Translation Corpus

Language	English – De, Es, Fr, It, Ni, Pt, Ro, Ru
Domain	public domain, business, science, entertainment, etc.
Source Lang. Material	TED Talks
Total Words	385 ~ 504 hours per language
Link:	https://ict.fbk.eu/must-c/

MuST-C: a Multilingual Speech Translation Corpus

Tgt	#Talk	#Sent	Hours	src w	tgt w
De	2,093	234K	408	4.3M	4.0M
Es	2,564	270K	504	5.3M	5.1M
Fr	2,510	280K	492	5.2M	5.4M
It	2,374	258K	465	4.9M	4.6M
Nl	2,267	253K	442	4.7M	4.3M
Pt	2,050	211K	385	4.0M	3.8M
Ro	2,216	240K	432	4.6M	4.3M
Ru	2,498	270K	489	5.1M	4.3M

BSTC: Baidu Simultaneous Translation Corpus

Language	Chinese-English
Domain	science, technology, economy, culture, art, etc.
Source Lang.	Chinese talks
Material	
Total Words	68 hours (237 talks)
Link:	https://ai.baidu.com/broad/introduction?dataset=bstc

BSTC: Baidu Simultaneous Translation Corpus

	talks	sentences	Characters / words		Hours
			Chinese	English	
Training set	215	37901	1,028,538	524,395	64.71
Dev set	16	956	26,059	13,277	1.58
Test set	6	975	25,832	12,724	1.46

BSTC: Baidu Simultaneous Translation Corpus

Training samples

Field	Content
Audio	
ASR	那么我们今天呢，就希望从一个20年的AI工作者来说，如何从专业的角度去解读一下，我们现在究竟发生了什么事情？他的权势金生。
Transcript	那么我们今天呢就希望，从一个二十年的AI工作者来说，如何从专业的角度去解读一下，我们现在究竟发生了什么事情，它的前世今生。
Translation	So today, as one who has been working on AI for twenty years, I wish I could give you a professional interpretation of what exactly is going on, its origin, history, characteristic, and where it is going.

BSTC: Baidu Simultaneous Translation Corpus

Test Set: 3 interpreters to interpret 6 lectures,
simulating real interpreting scenario

Lectures ID	Domain	Length
1	Art	15'
2	AI	15'
3	Art	19'
4	Story	11'
5	Big Data	14'
6	AI	10'



语音容错翻译

词语、音节联合建模

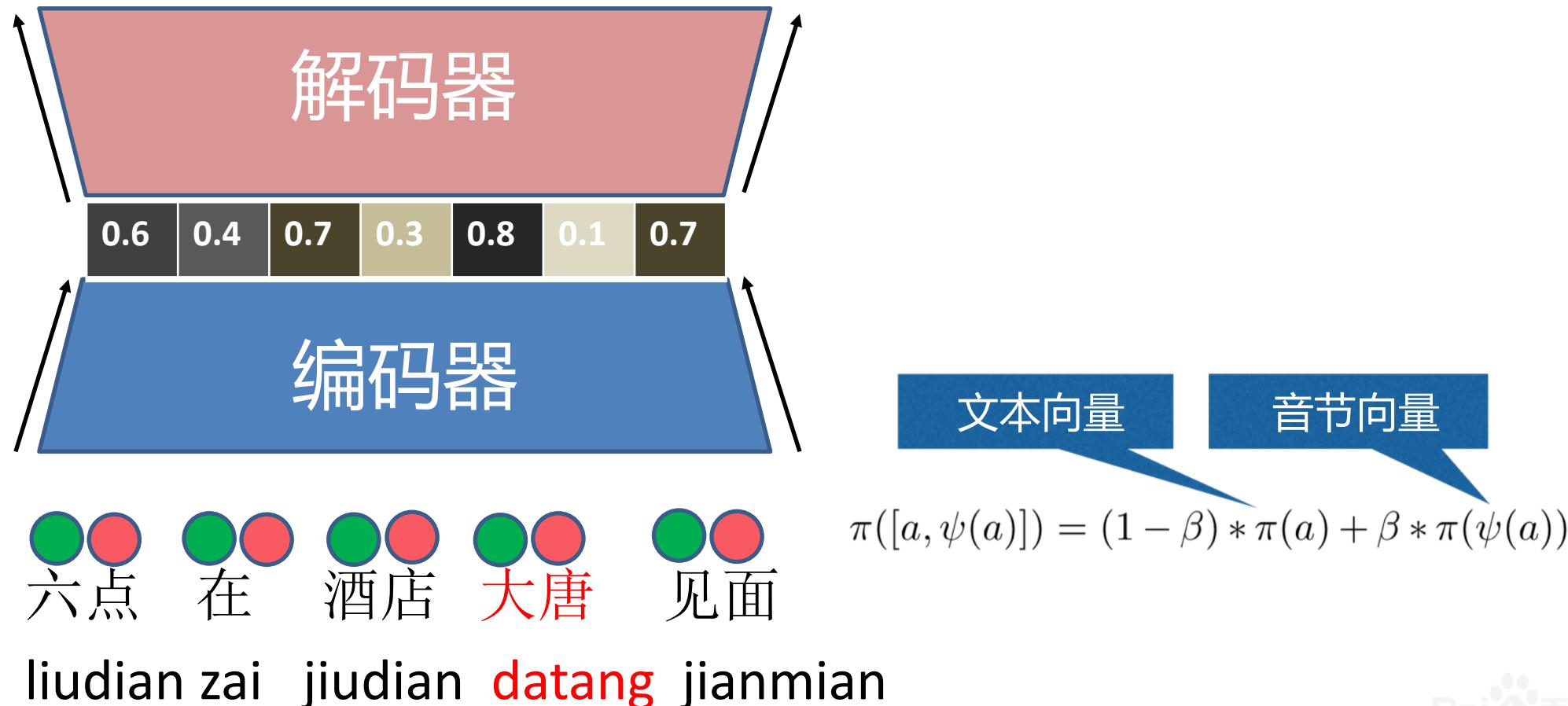
Let' s meet at 6 o'clock in the hotel DaTang



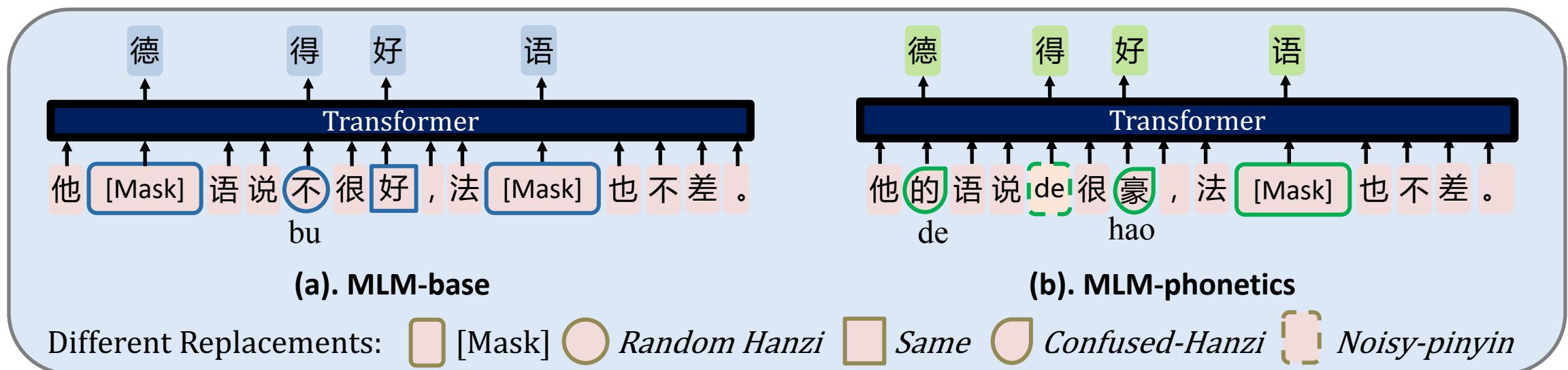
六点 在 酒店 大唐 见面

词语、音节联合建模

Let's meet at 6 o'clock in the hotel **lobby**



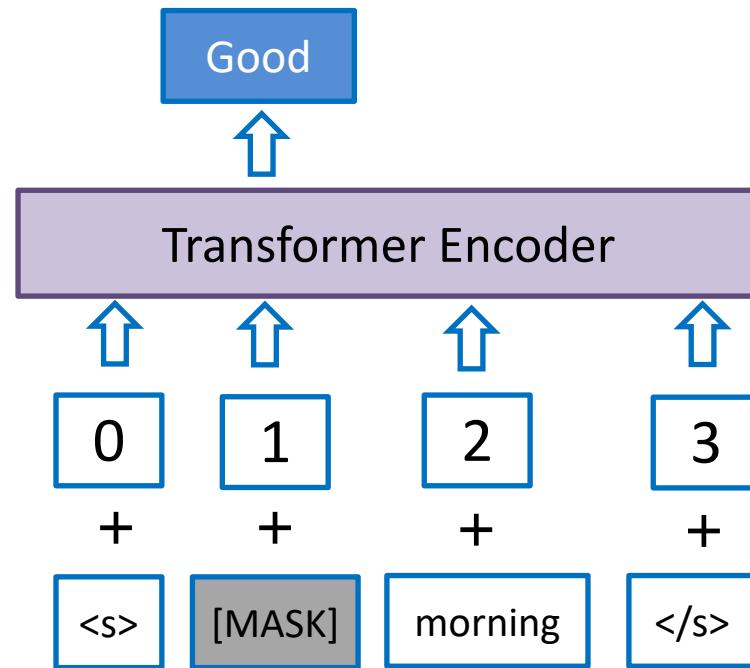
语音识别纠错





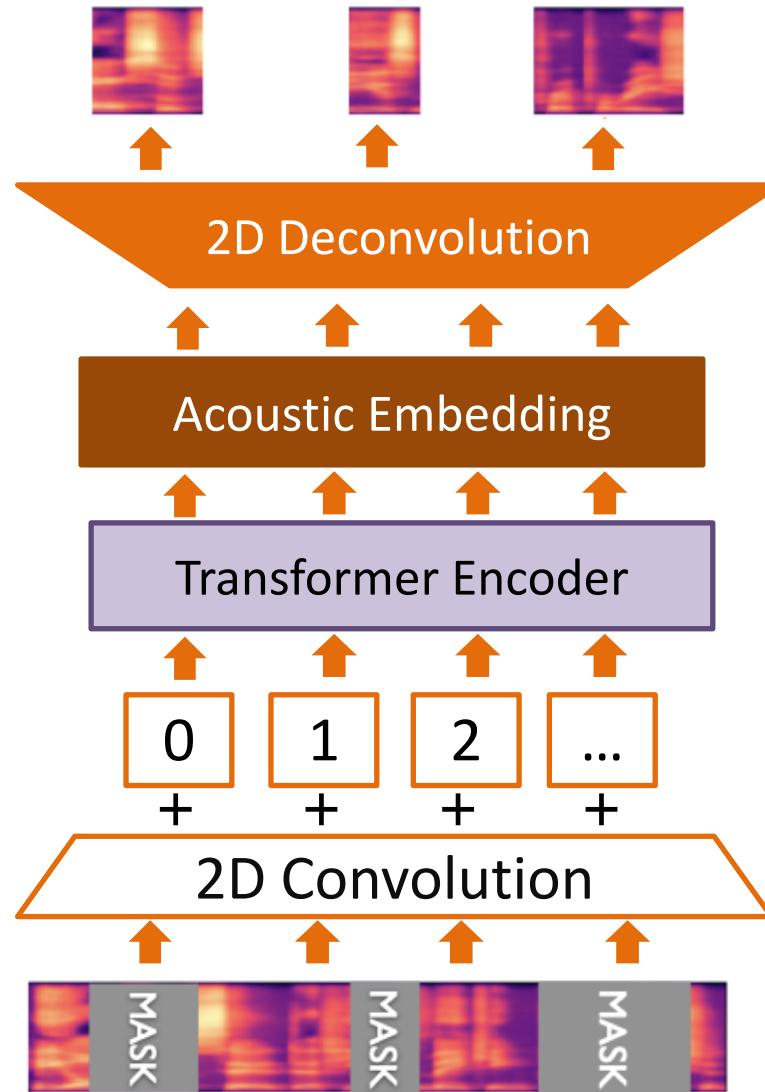
多模态信息融合

语音、语言联合预训练模型

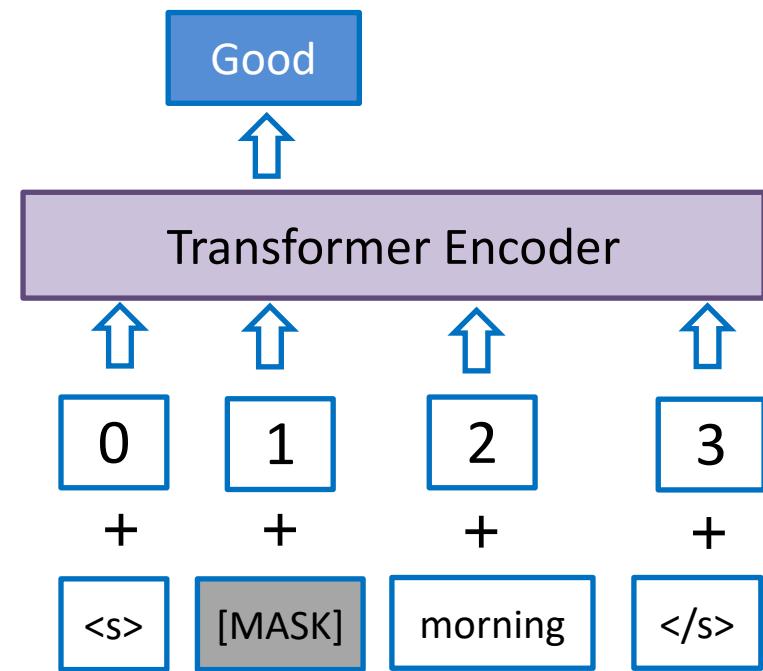


Masked Language Model

语音、语言联合预训练模型

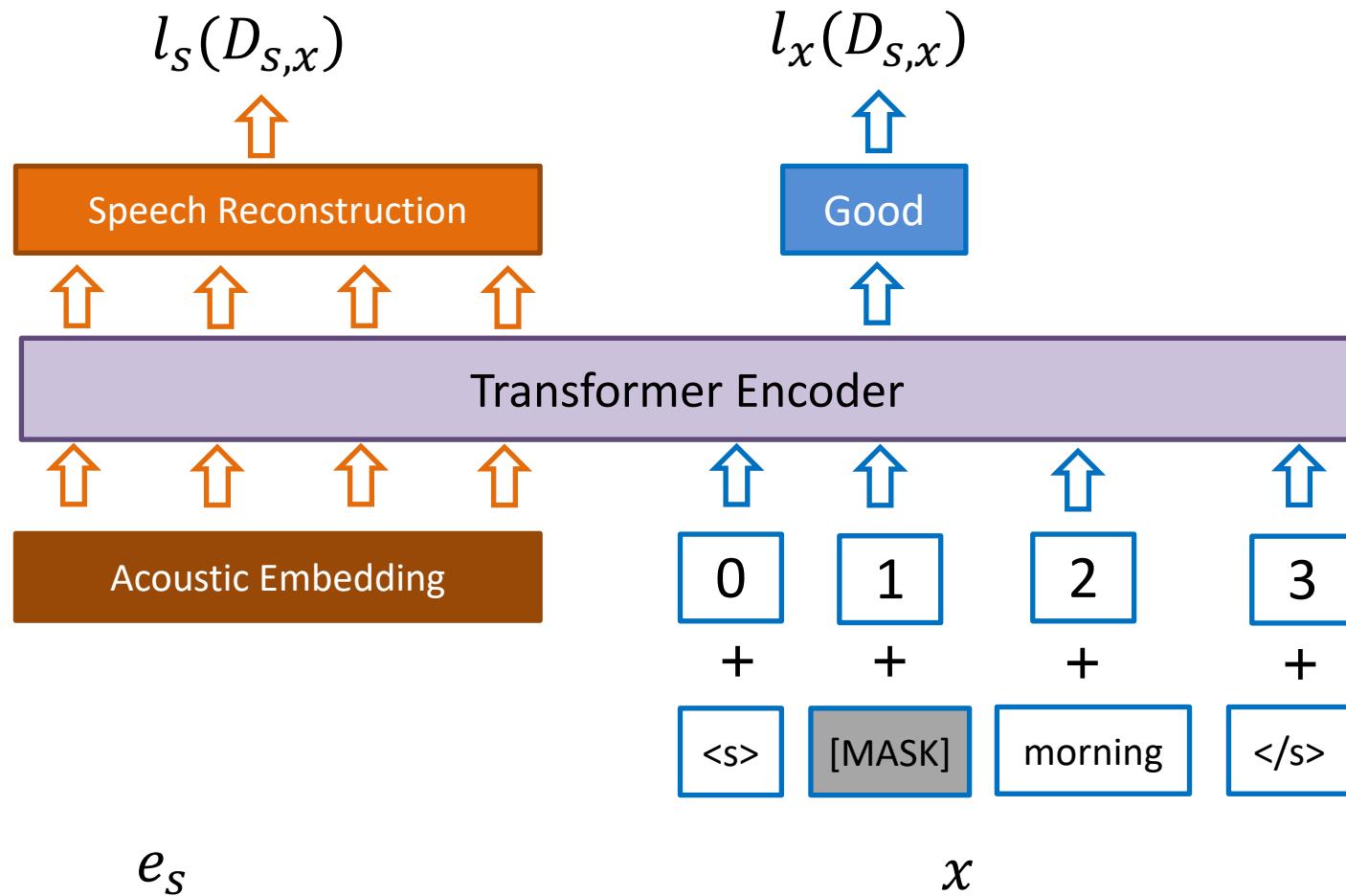


Masked Acoustic Model



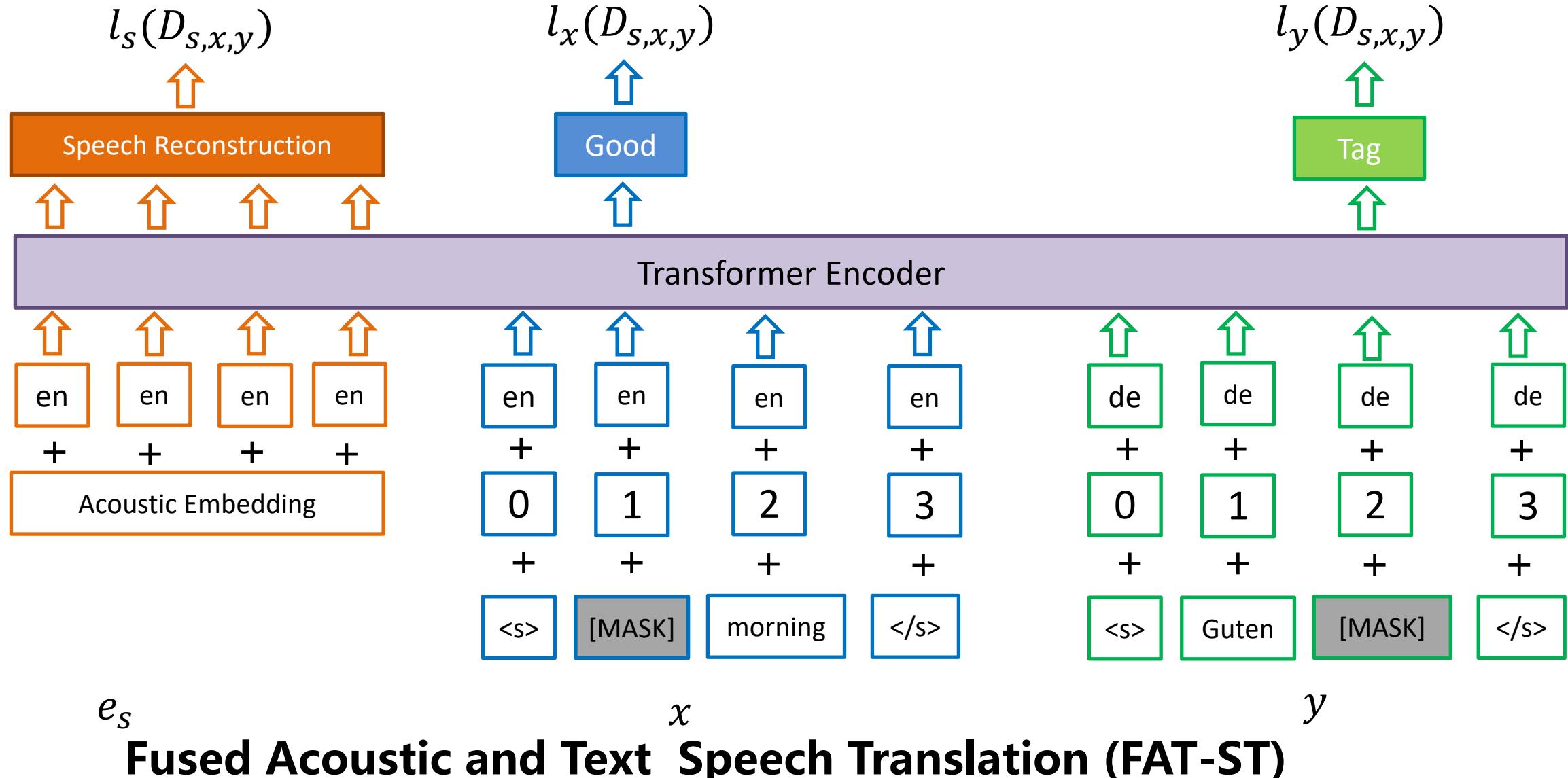
Masked Language Model

语音、语言联合预训练模型



Fused Acoustic and Text Masked Language Model (FAT-MLM)

语音、语言联合预训练模型

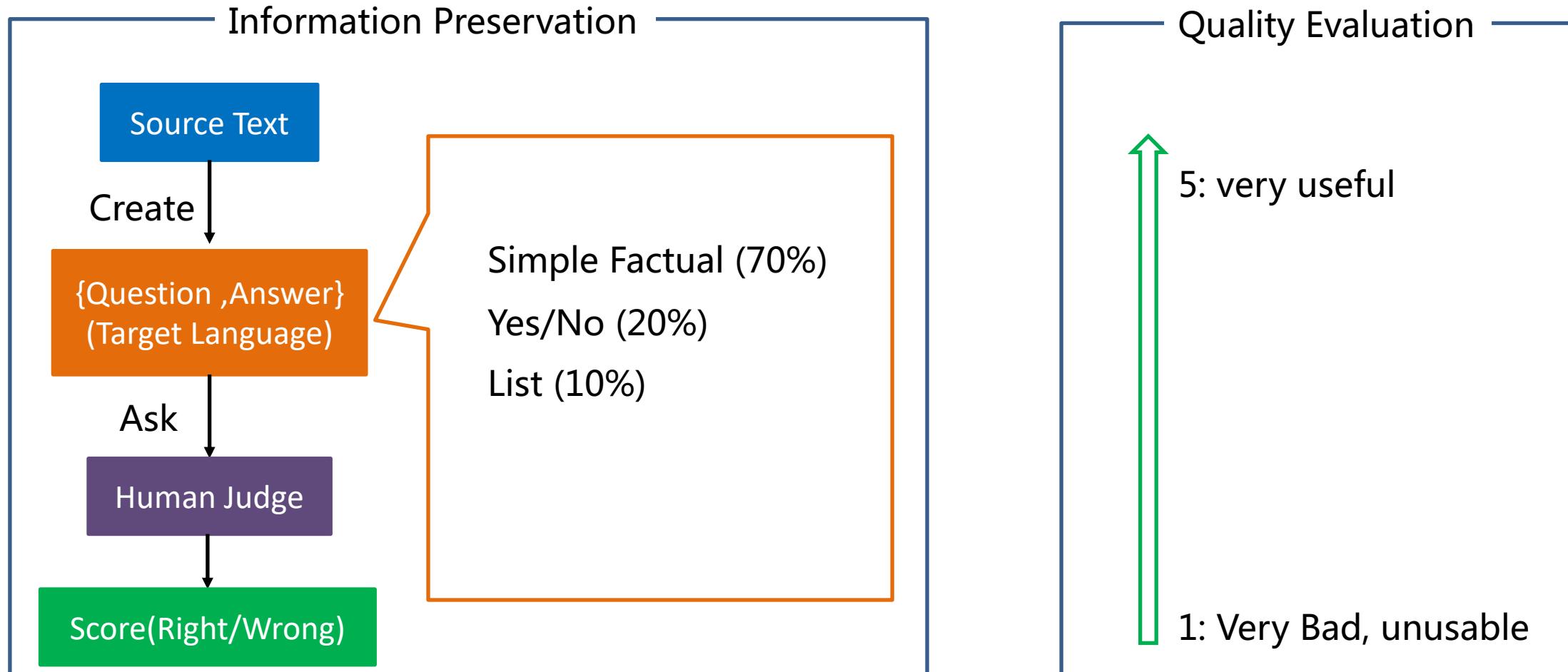




同传评价

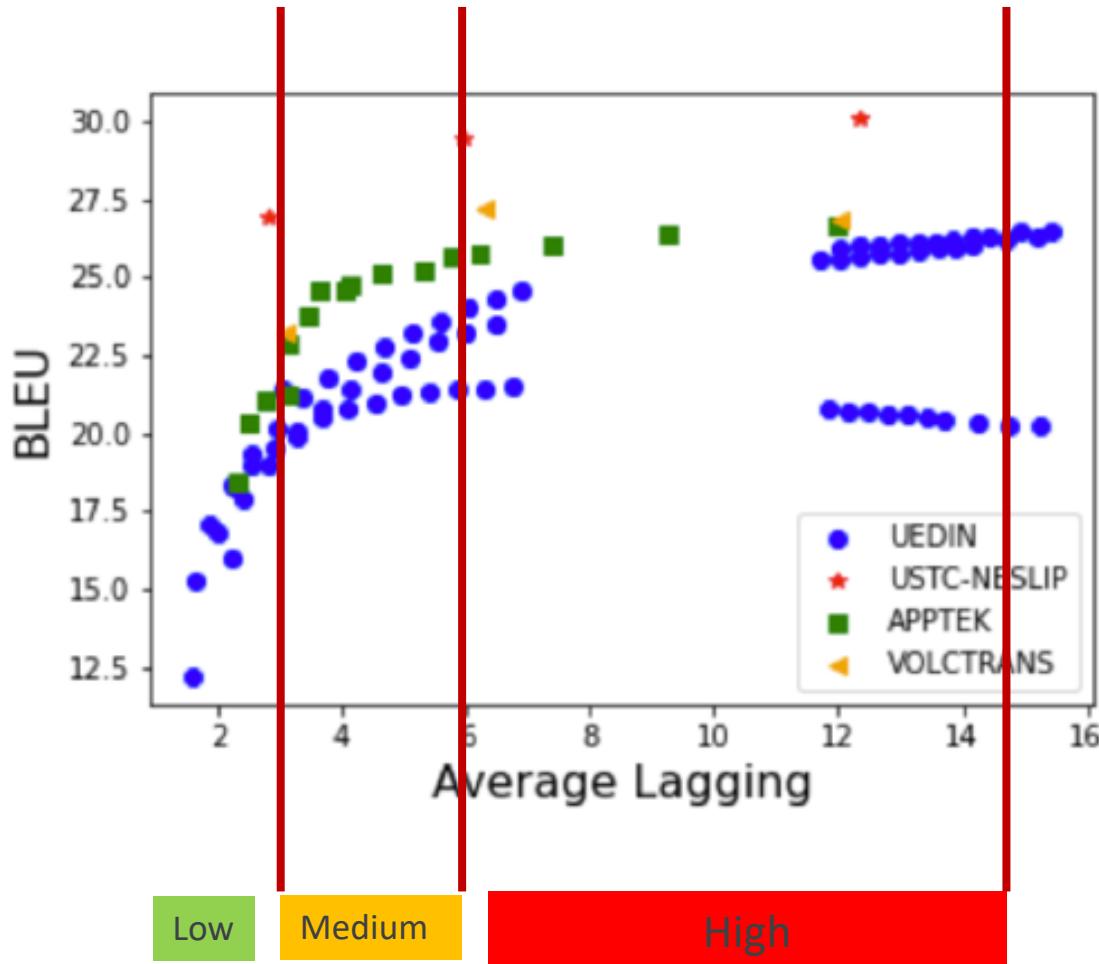
基于阅读理解的端到端评价

针对同传场景，问题和答案设计较为关键，人力成本高、周期长



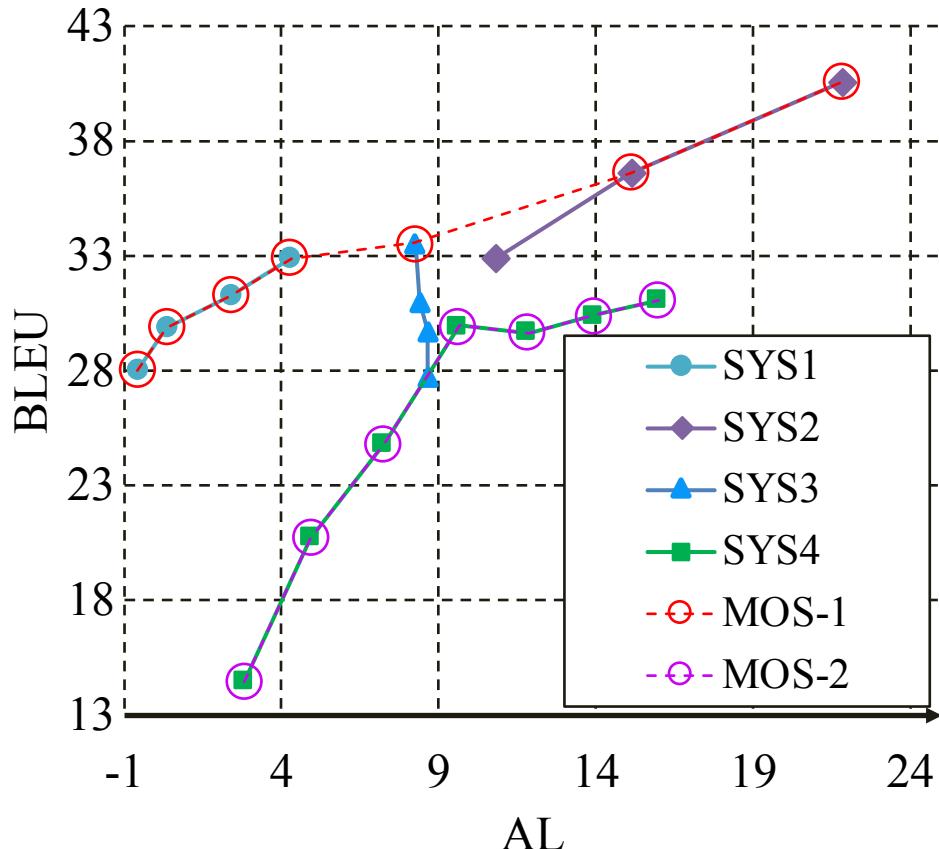
如何兼顾质量与时延？

时延分为3个等级，分别评估翻译质量



如何兼顾质量与时延？

构建凸包集合，统计各系统落在凸包上的点的个数



同传研讨会

AutoSimTrans 2020 **Program** Home Invited Talks Shared Task Call for Papers Organization Important Dates

The 1st Workshop on Automatic Simultaneous Translation Challenges, Recent Advances, and Future Directions

Workshop at [ACL 2020](#), Seattle, July 10, 2020



清华大学
Tsinghua University

AutoSimTrans 2021 **Program** Home Invited Talks Shared Task Call for Papers Organization Past Workshops

The 2nd Workshop on Automatic Simultaneous Translation Challenges, Recent Advances, and Future Directions

Workshop at [NAACL 2021](#), Mexico City, June 10, 2021



HUAWEI



提纲

- 同声传译
- 主要挑战
- 主流方法
 - 级联模型
 - 端到端模型
- 数据、鲁棒性模型及评价方法
- 产品及应用
- 总结及展望

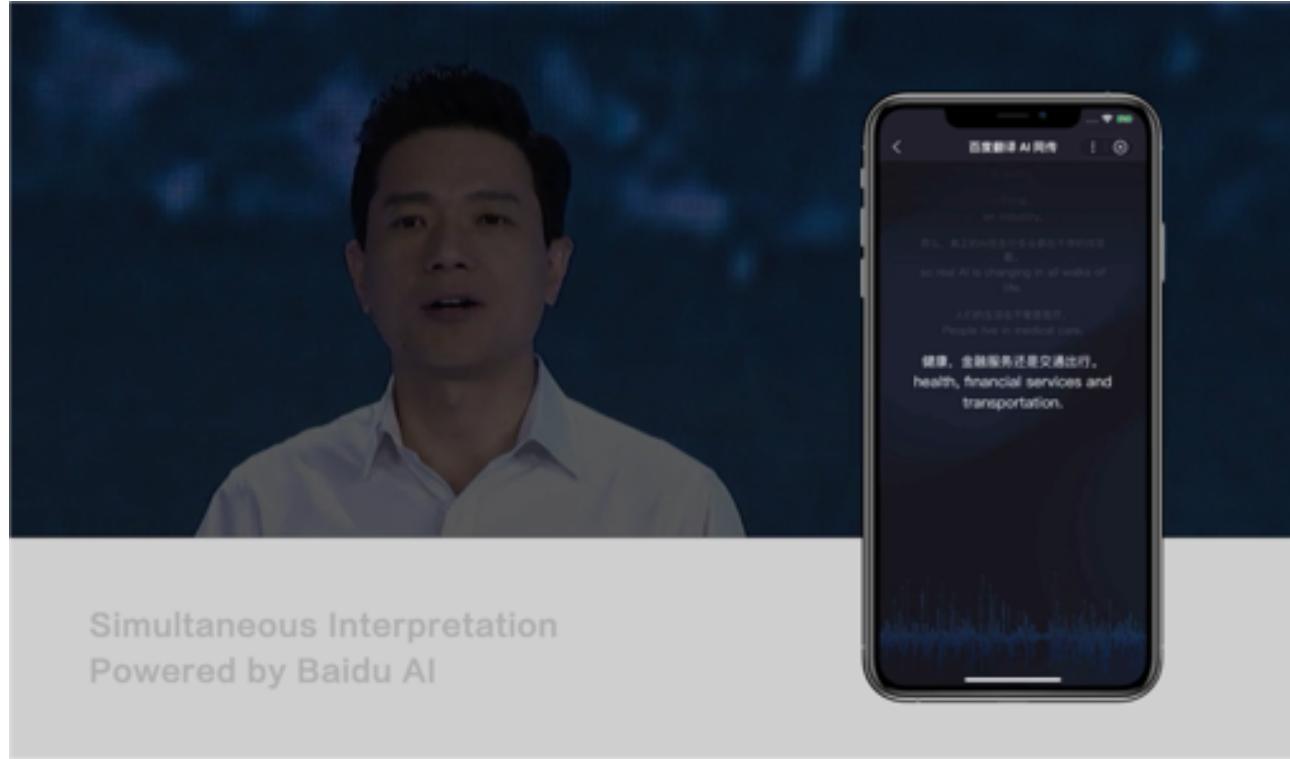
字幕式（语音到文字，看同传）



主流产品形式
部署方便
无需额外硬件

多语言展现困难
分散观众注意力看字幕
放大识别与翻译错误

语音式（语音到语音，听同传）

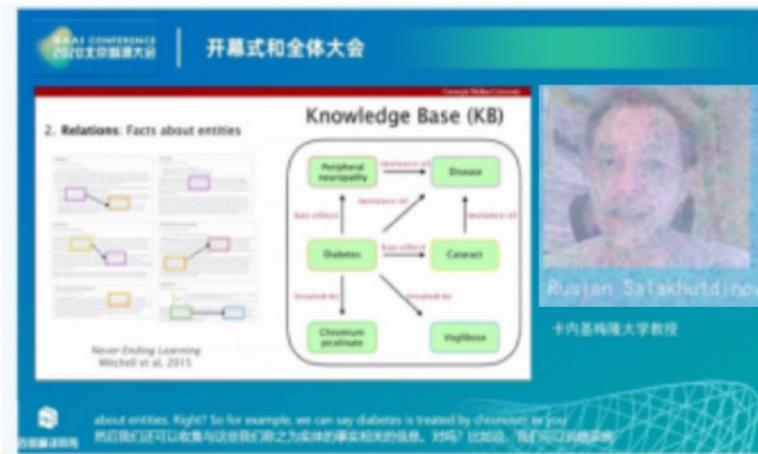


Simultaneous Interpretation
Powered by Baidu AI

多语言扩展方便
沉浸式体验
缓解错误（如同音字）

需要接收装置（如手机）
语音合成，增加时延

同传应用



总结与展望

有一天，当你在人民大会堂和世界各国外友人聚会的时候，你会发现，无论哪个国家的人在台上讲话，与会者都能从耳机里听到自己国家的语言，同时你会发现，在耳机里做翻译的不是人，而是我们的“万能翻译博士”，因为“博士”的语调不像人那样委婉自如。

--《机器翻译浅说》 1964

总结与展望

- **模型**
 - 鲁棒的同传模型
 - 端到端同传模型
 - 融合视觉信息的多模态同传
- **数据**
 - 构建更大规模的数据集
- **评价**
 - 设计符合同传场景的评价指标



谢 谢 !