

Simultaneous Translation: Recent Advances and Remaining Challenges



Liang Huang

Baidu Research (USA) and Oregon State University



Consecutive vs. Simultaneous Interpretation

consecutive interpretation
multiplicative latency (x2)



simultaneous interpretation
additive latency (+3 secs)

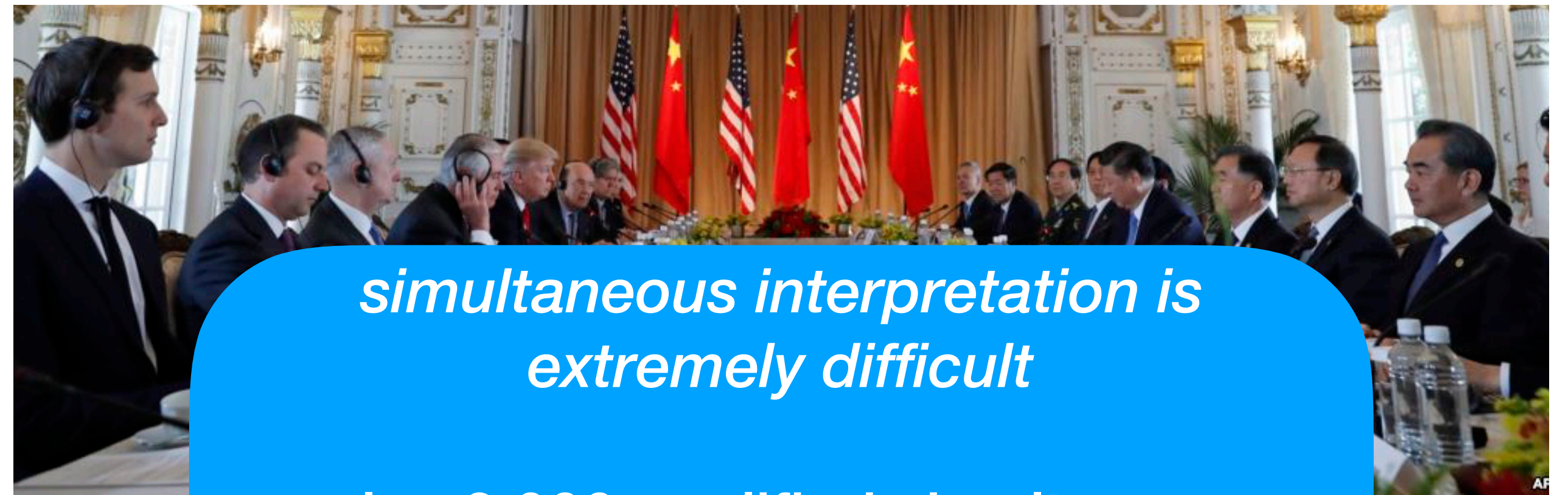


Consecutive vs. Simultaneous Interpretation

consecutive interpretation
multiplicative latency (x2)



simultaneous interpretation
additive latency (+3 secs)



simultaneous interpretation is extremely difficult

only ~3,000 qualified simultaneous interpreters world-wide (AICC)

each interpreter can only sustain for at most 15-20 minutes

the best interpreters can only cover ~60% of the source material

Simultaneous Interpreters: Strategies & Limitations

- anticipation, summarization, generalization, etc...
- and they inevitably make (quite a bit of) mistakes
- “human-level” *quality*: much lower than normal translation
- “human-level” *latency*: very short: 2~4 secs (actually higher latency *hurts* quality...)

我们支持 uh... 玻利维亚大使 和 俄罗斯大使 刚才 所做的立场
we support uh... Bolivia envoy & Russia envoy just-now made position

We support the position of Bolivia & Russia

Simultaneous Interpreters: Strategies & Limitations

- anticipation, summarization, generalization, etc...
- and they inevitably make (quite a bit of) mistakes
- “human-level” *quality*: much lower than normal translation
- “human-level” *latency*: very short: 2~4 secs (actually higher latency *hurts* quality...)

我们支持 uh... 玻利维亚 大使 和 俄罗斯 大使 刚才 所做的 立场
we support uh... Bolivia envoy & Russia envoy just-now made position
We support the position of Bolivia & Russia

Simultaneous Interpreters: Strategies & Limitations

- anticipation, summarization, generalization, etc...
- and they inevitably make (quite a bit of) mistakes
- “human-level” *quality*: much lower than normal translation
- “human-level” *latency*: very short: 2~4 secs (actually higher latency *hurts* quality...)

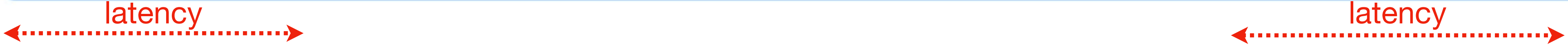
我们支持 uh... 玻利维亚 大使 和 俄罗斯 大使 刚才 所做的 立场
we support uh... Bolivia envoy & Russia envoy just-now made position
We support the position of Bolivia & Russia

uh... 我们 认为 安理会 ah... 没有 必要
uh... we think sec. council uh... no need

Simultaneous Interpreters: Strategies & Limitations

- anticipation, summarization, generalization, etc...
- and they inevitably make (quite a bit of) mistakes
- “human-level” *quality*: much lower than normal translation
- “human-level” *latency*: very short: 2~4 secs (actually higher latency *hurts* quality...)

我们支持 uh... 玻利维亚 大使 和 俄罗斯 大使 刚才 所做的 立场
we support uh... Bolivia envoy & Russia envoy just-now made position
We support the position of Bolivia & Russia

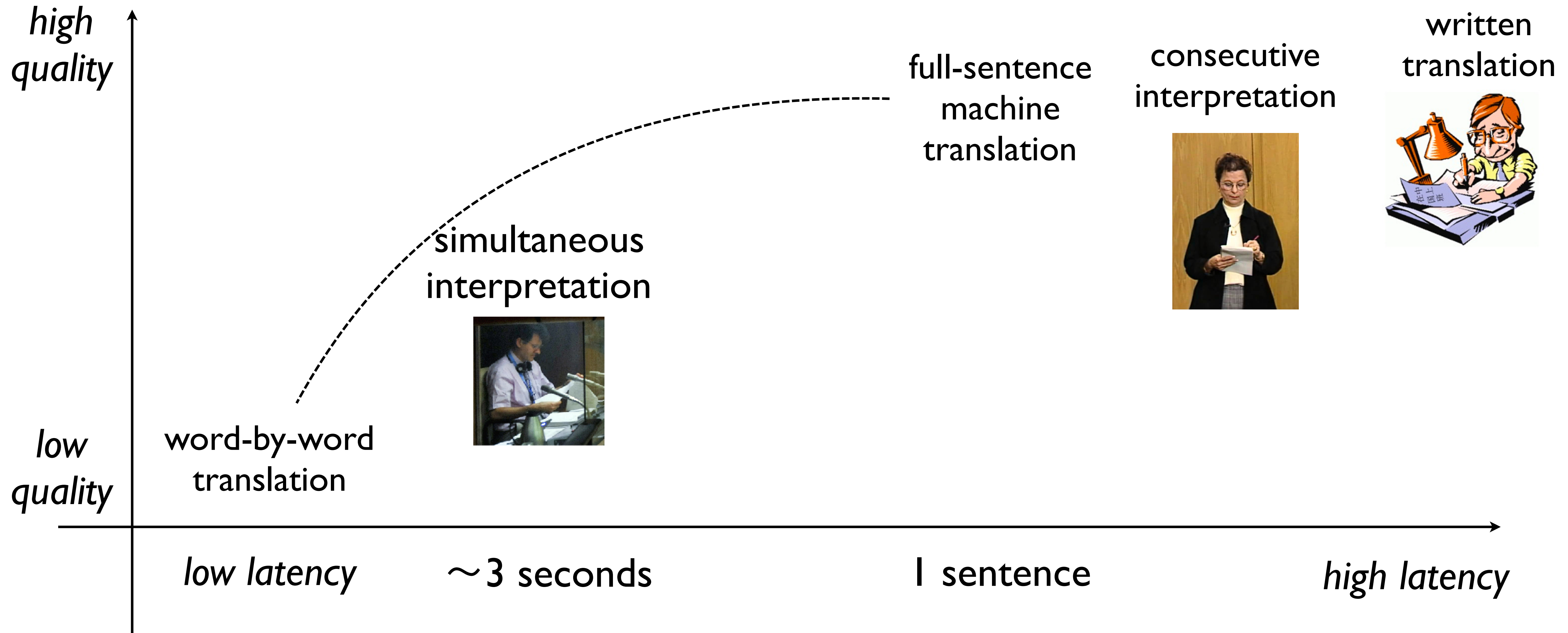


uh... 我们认为 安理会 ah... 没有 必要
uh... we think sec. council uh... no need

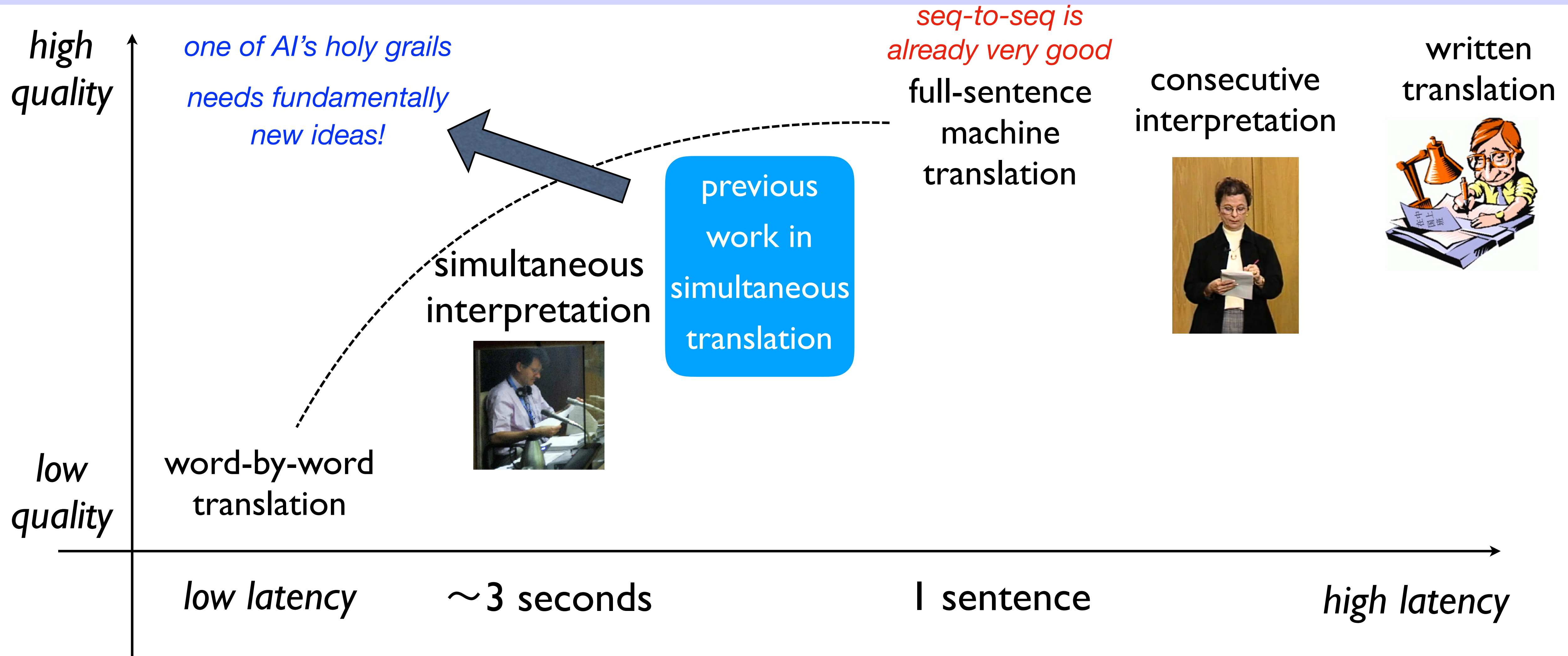


from United Nations Proceedings Speech Corpus (LDC2014S08, Chay et al, 2014)

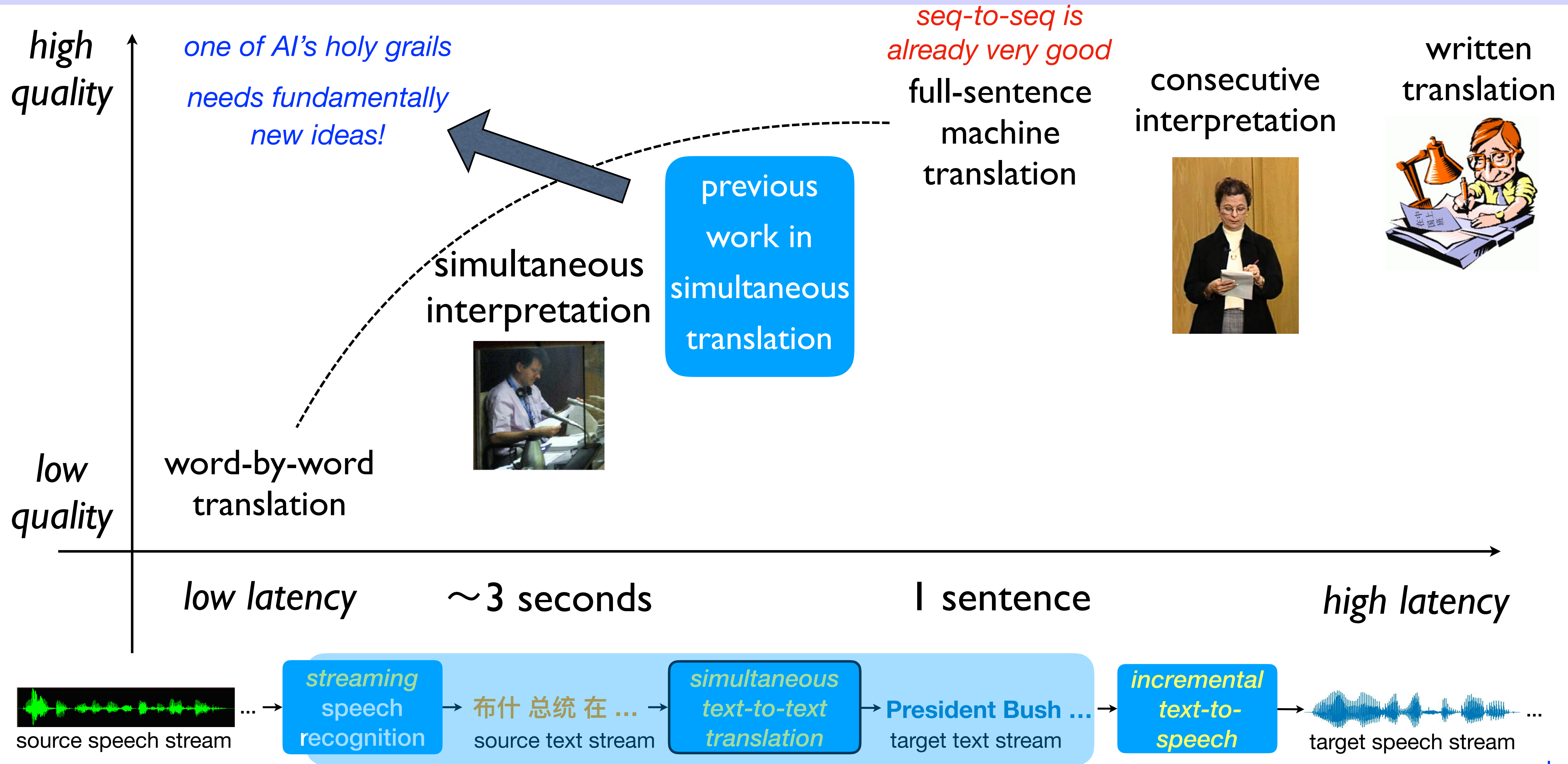
Tradeoff between Latency and Quality



Tradeoff between Latency and Quality



Tradeoff between Latency and Quality



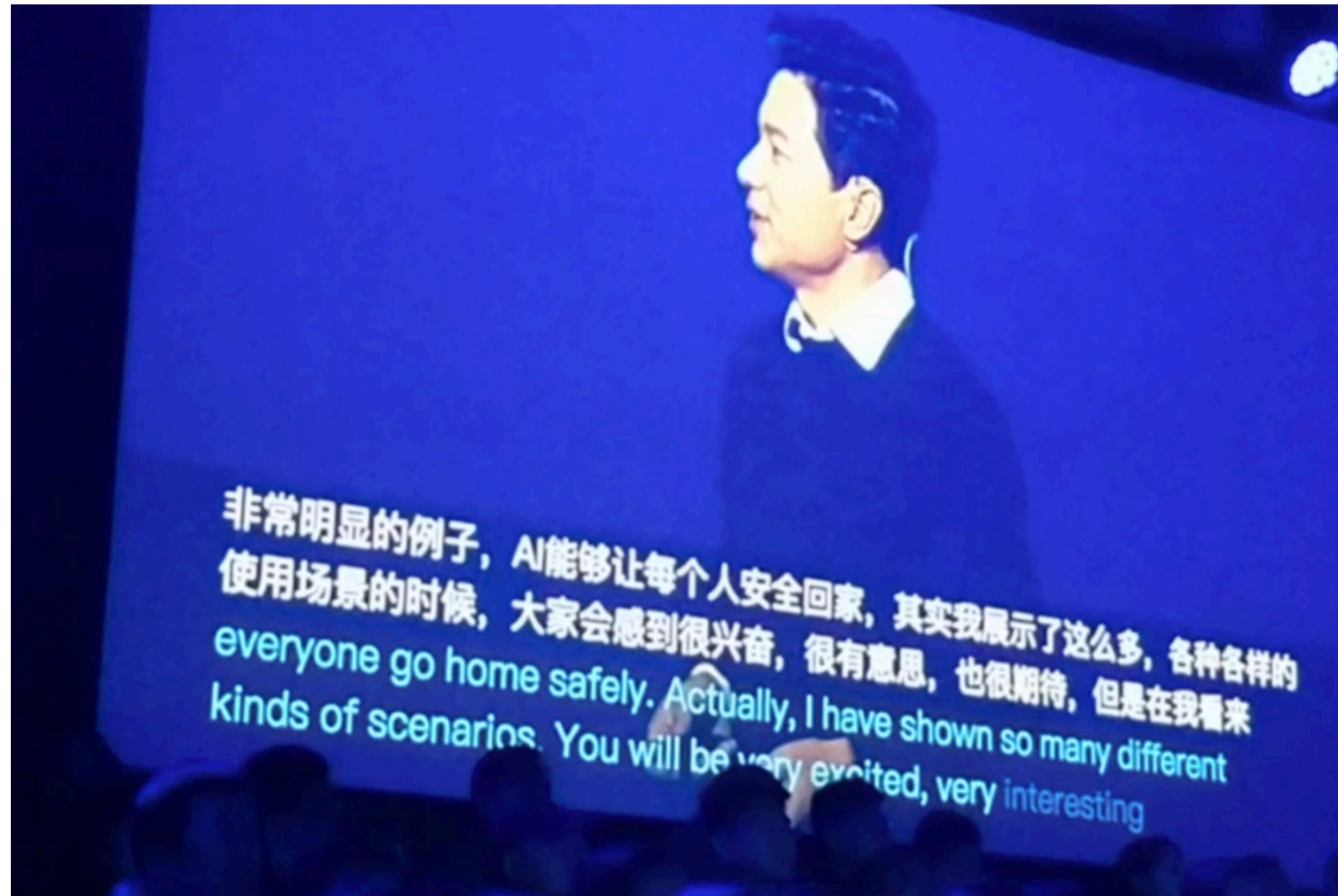
Outline

- Background on Simultaneous Interpretation
- Part I: Our Breakthrough in 2018
 - Prefix-to-Prefix Framework, Integrated Anticipation, Controllable Latency
 - New Latency Metric
 - Demos and Examples
- Part II: Towards Flexible (Adaptive) Translation Policies
- Part III: Remaining Challenges

Our Breakthrough in 2018

Baidu World Conference, Nov. 2017

full-sentence translation (latency: 10+ secs)



our
work

Baidu World Conference, Nov. 2018

low-latency simultaneous translation (latency: ~3 secs)



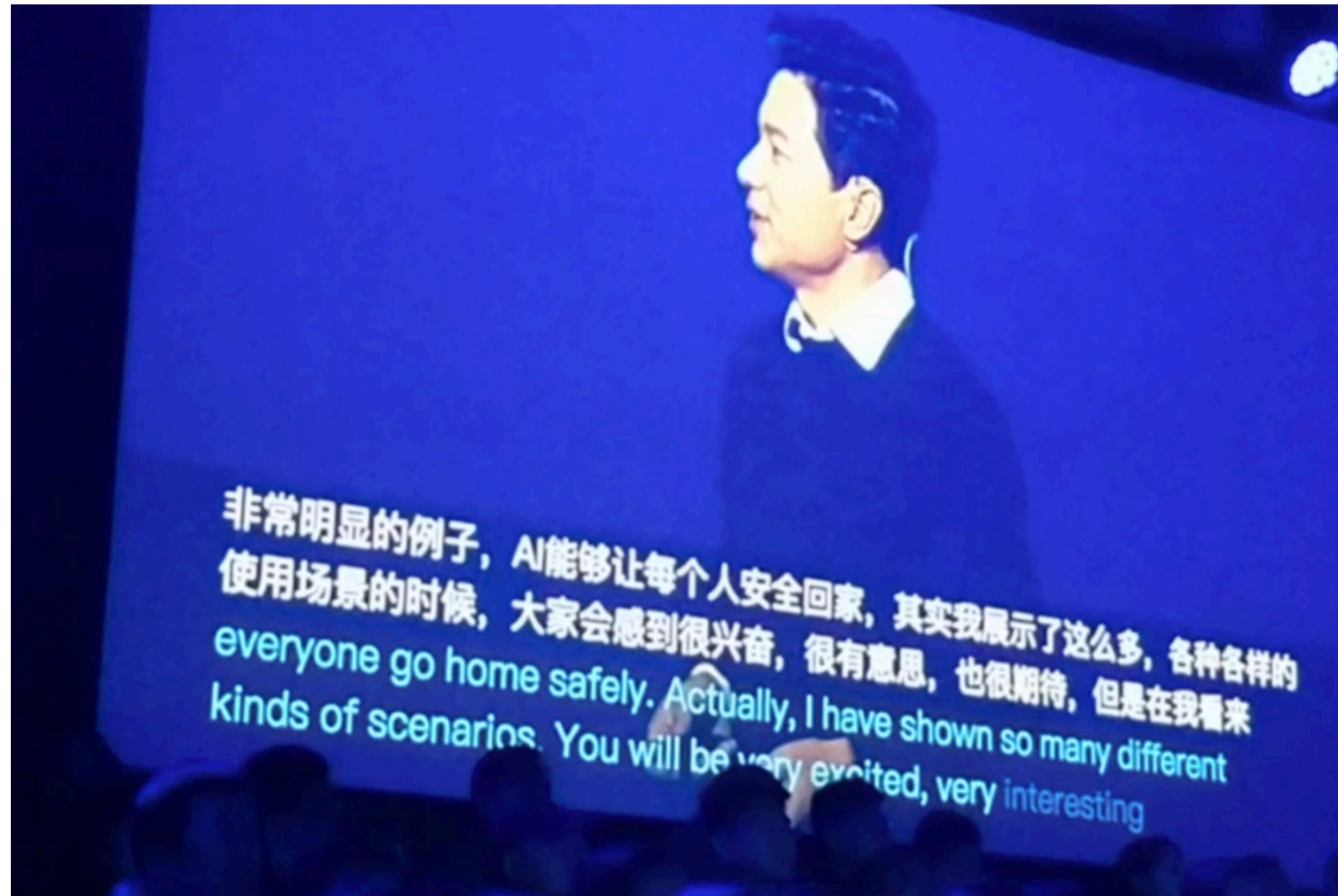
Media coverage:



Our Breakthrough in 2018

Baidu World Conference, Nov. 2017

full-sentence translation (latency: 10+ secs)



our
work

Baidu World Conference, Nov. 2018

low-latency simultaneous translation (latency: ~3 secs)



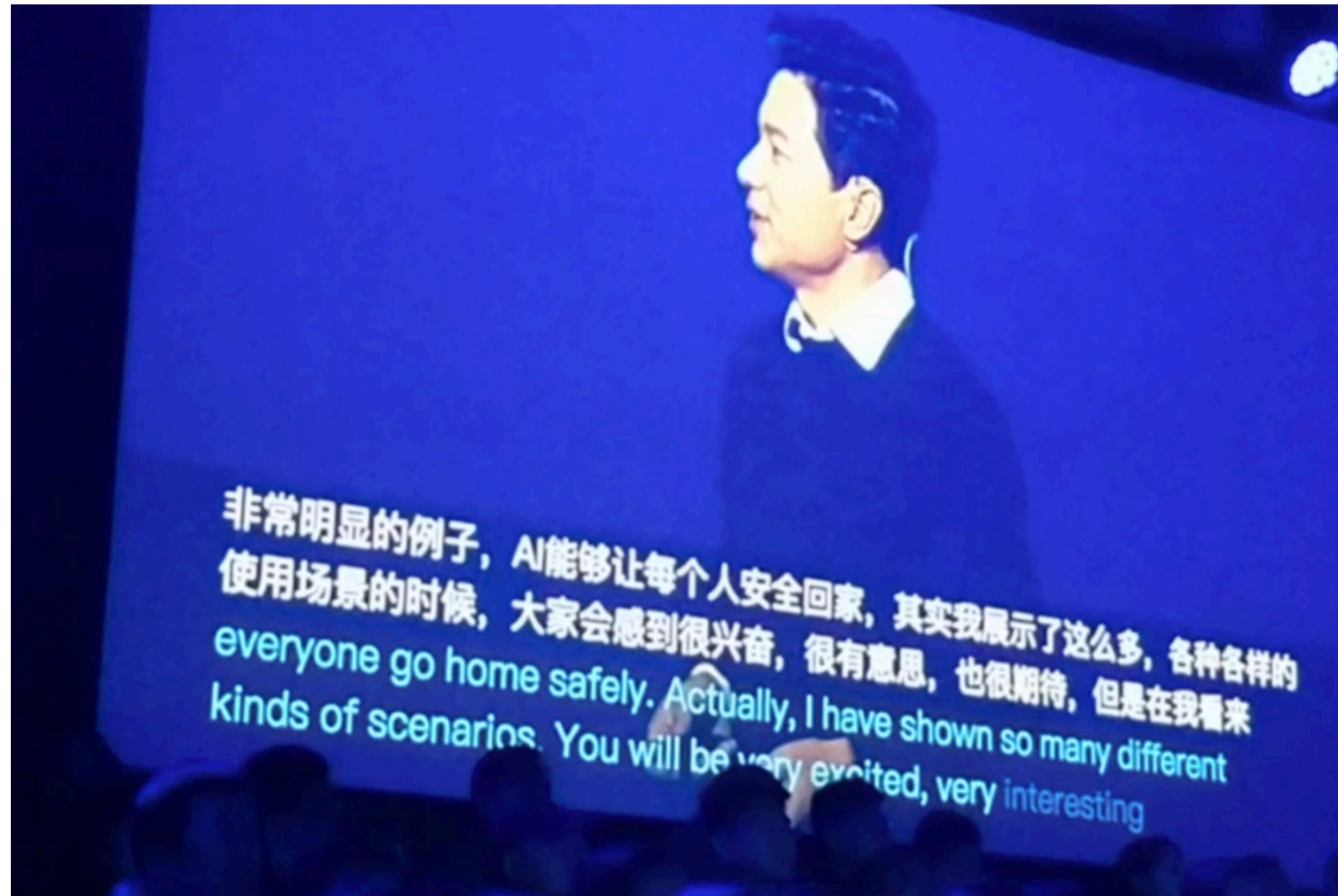
Media coverage:



Our Breakthrough in 2018

Baidu World Conference, Nov. 2017

full-sentence translation (latency: 10+ secs)



our
work

Baidu World Conference, Nov. 2018

low-latency simultaneous translation (latency: ~3 secs)



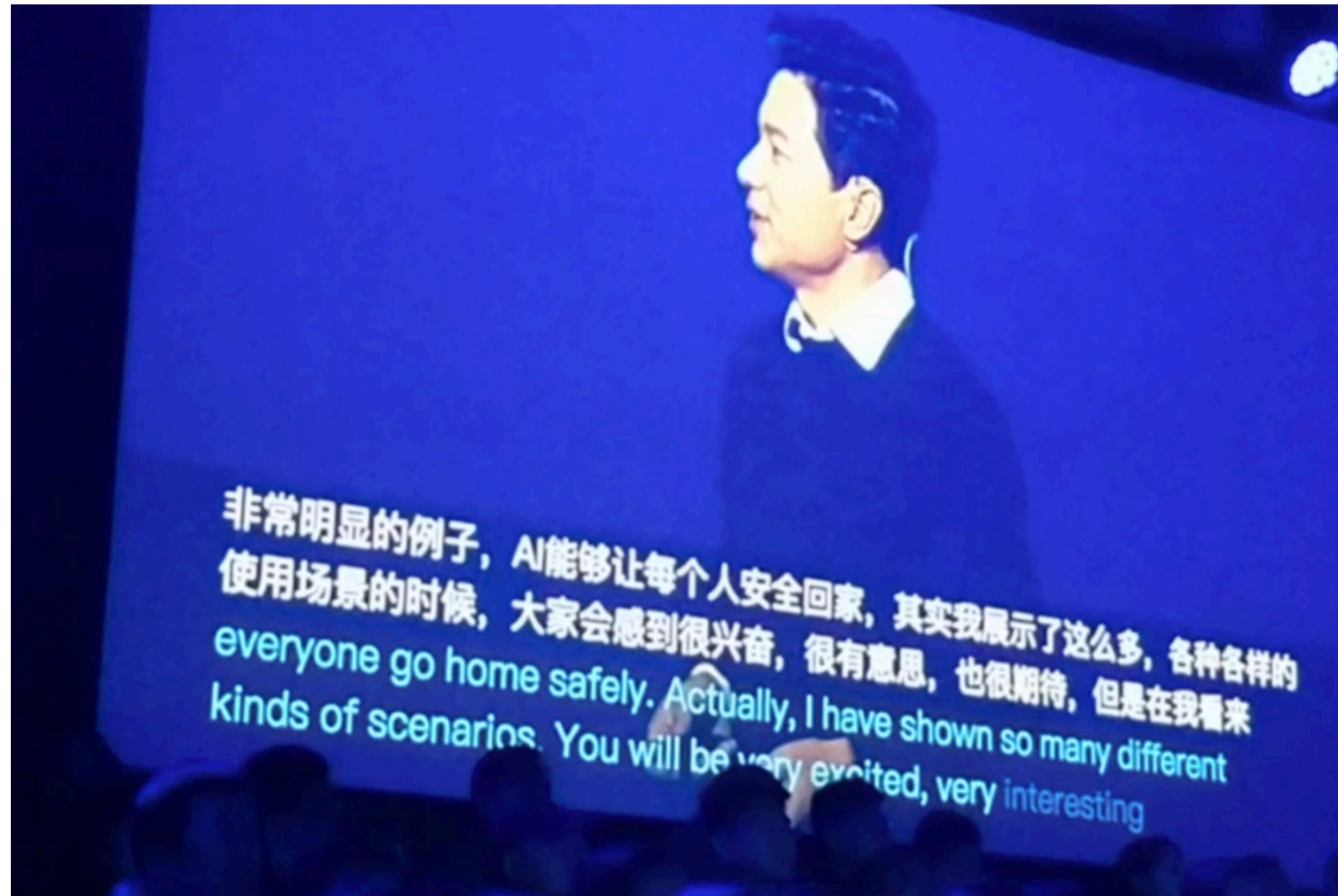
Media coverage:



Our Breakthrough in 2018

Baidu World Conference, Nov. 2017

full-sentence translation (latency: 10+ secs)



our
work

Baidu World Conference, Nov. 2018

low-latency simultaneous translation (latency: ~3 secs)



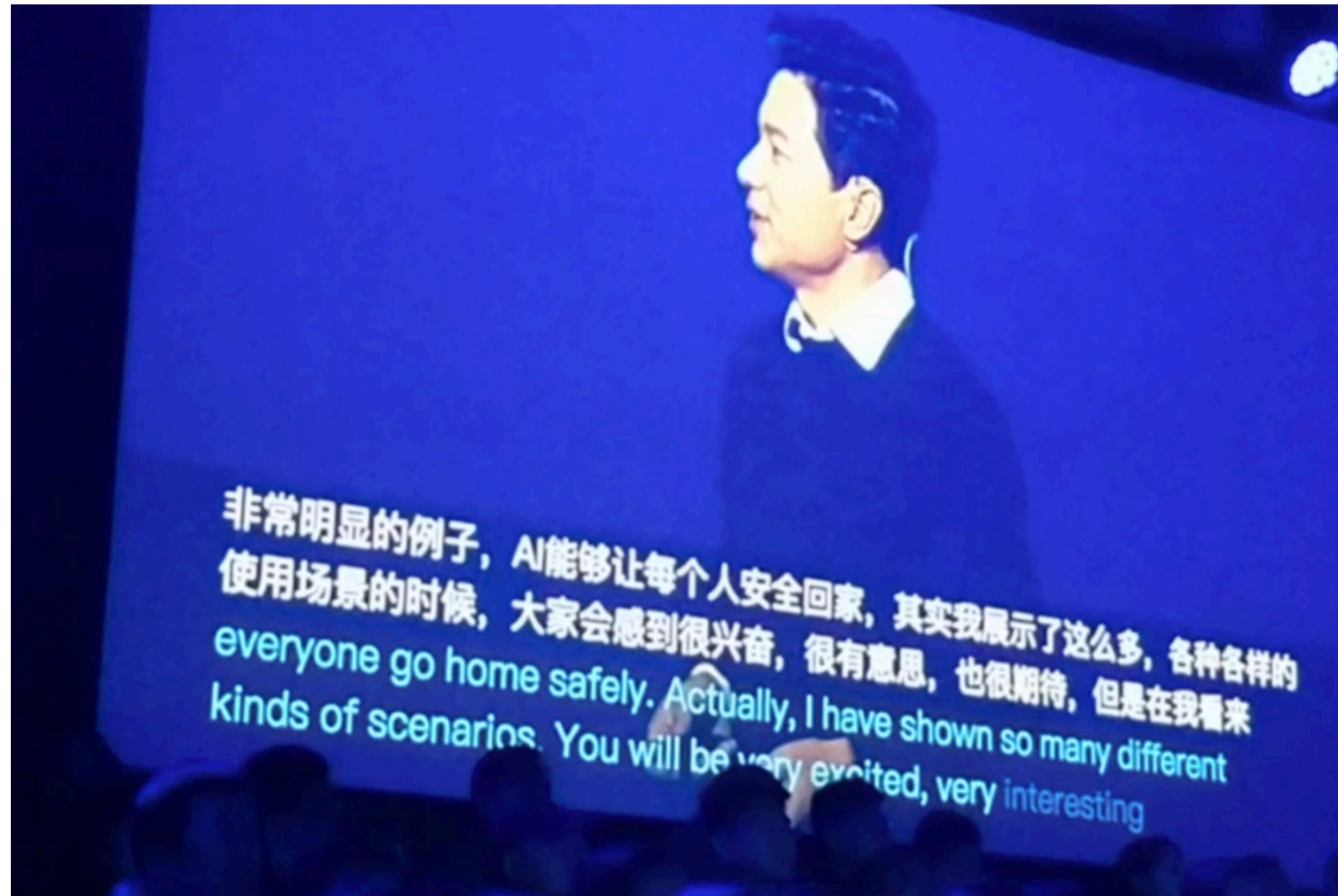
Media coverage:



Our Breakthrough in 2018

Baidu World Conference, Nov. 2017

full-sentence translation (latency: 10+ secs)



our
work

Baidu World Conference, Nov. 2018

low-latency simultaneous translation (latency: ~3 secs)



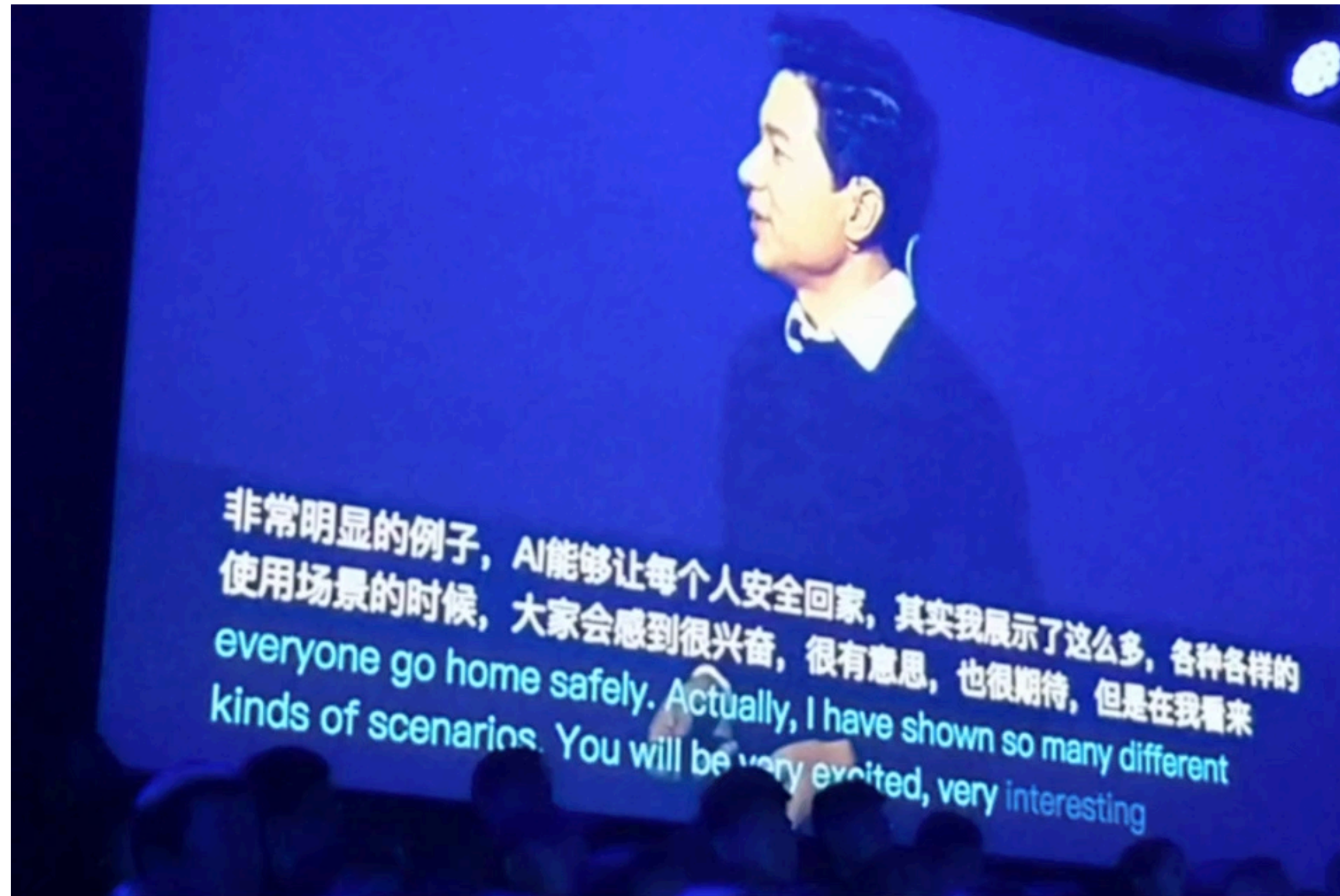
Media coverage:



Our Breakthrough in 2018

Baidu World Conference, Nov. 2017

full-sentence translation (latency: 10+ secs)



our
work

Baidu World Conference, Nov. 2018

low-latency simultaneous translation (latency: ~3 secs)



request



Haifeng Wang



Zhongjun He



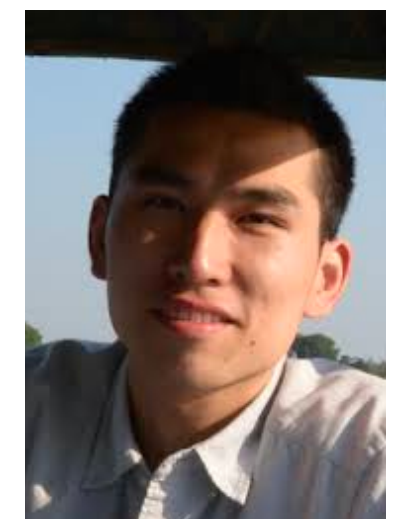
Hao Xiong



Mingbo Ma



Kaibo Liu

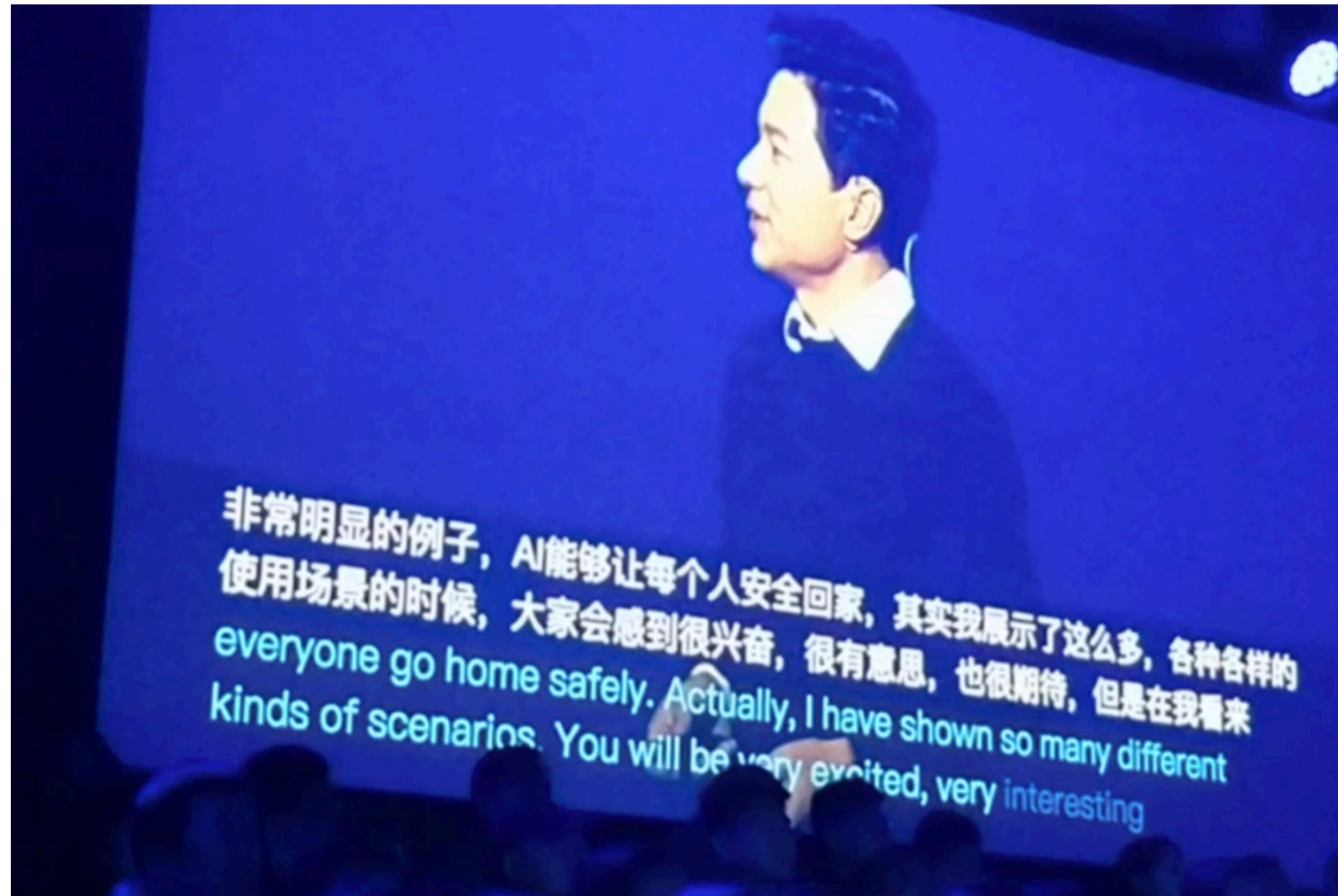


Renjie Zheng

Our Breakthrough in 2018

Baidu World Conference, Nov. 2017

full-sentence translation (latency: 10+ secs)



our work

Baidu World Conference, Nov. 2018

low-latency simultaneous translation (latency: ~3 secs)



request



Haifeng Wang



Zhongjun He



Hao Xiong



Ken Church

I really need low-latency simultaneous translation!



Mingbo Ma



Kaibo Liu



Renjie Zheng

Main Challenge: Word Order Difference

- e.g. translate from Subj-Obj-Verb (Japanese, German) to Subj-Verb-Obj (English)
 - German is underlyingly SOV, and Chinese is a mix of SVO and SOV
 - human simultaneous interpreters routinely “anticipate” (e.g., predicting German verb)

ich bin mit dem Zug nach Ulm **gefahren**

I am with the train to Ulm **traveled**

Grissom et al, 2014

I (..... *waiting*.....) **traveled** by train to Ulm

Main Challenge: Word Order Difference

- e.g. translate from Subj-Obj-Verb (Japanese, German) to Subj-Verb-Obj (English)
 - German is underlyingly SOV, and Chinese is a mix of SVO and SOV
 - human simultaneous interpreters routinely “anticipate” (e.g., predicting German verb)

ich bin mit dem Zug nach Ulm **gefahren**

I am with the train to Ulm **traveled**

Grissom et al, 2014

I (..... *waiting*.....) **traveled** by train to Ulm

<i>Bùshí</i>	<i>zǒngtǒng</i>	<i>zài</i>	<i>Mòsīkē</i>	<i>yǔ</i>	<i>Éluósī</i>	<i>zǒngtǒng</i>	<i>Pǔjīng</i>	<i>huìwù</i>
布什	总统	在	莫斯科	与	俄罗斯	总统	普京	会晤
<i>Bush</i>	<i>President</i>	<i>in</i>	<i>Moscow</i>	<i>with</i>	<i>Russian</i>	<i>President</i>	<i>Putin</i>	<i>meet</i>

President Bush **meets** with Russian President Putin in Moscow

Main Challenge: Word Order Difference

- e.g. translate from Subj-Obj-Verb (Japanese, German) to Subj-Verb-Obj (English)
 - German is underlyingly SOV, and Chinese is a mix of SVO and SOV
 - human simultaneous interpreters routinely “anticipate” (e.g., predicting German verb)

ich bin mit dem Zug nach Ulm **gefahren**

I am with the train to Ulm **traveled**

Grissom et al, 2014

I (..... *waiting*.....) **traveled** by train to Ulm

<i>Bùshí</i>	<i>zǒngtǒng</i>	<i>zài</i>	<i>Mòsīkē</i>	<i>yǔ</i>	<i>Éluósī</i>	<i>zǒngtǒng</i>	<i>Pǔjīng</i>	<i>huìwù</i>
布什	总统	在	莫斯科	与	俄罗斯	总统	普京	会晤
<i>Bush</i>	<i>President</i>	<i>in</i>	<i>Moscow</i>	<i>with</i>	<i>Russian</i>	<i>President</i>	<i>Putin</i>	<i>meet</i>

President Bush **meets** with Russian President Putin in Moscow

non-anticipative: President Bush (..... *waiting*.....) **meets** with Russian ...

Main Challenge: Word Order Difference

- e.g. translate from Subj-Obj-Verb (Japanese, German) to Subj-Verb-Obj (English)
- German is underlyingly SOV, and Chinese is a mix of SVO and SOV
- human simultaneous interpreters routinely “anticipate” (e.g., predicting German verb)

ich bin mit dem Zug nach Ulm **gefahren**

I am with the train to Ulm **traveled**

Grissom et al, 2014

I (..... *waiting*.....) **traveled** by train to Ulm

<i>Bùshí</i>	<i>zǒngtǒng</i>	<i>zài</i>	<i>Mòsīkē</i>	<i>yǔ</i>	<i>Éluósī</i>	<i>zǒngtǒng</i>	<i>Pǔjīng</i>	<i>huìwù</i>
布什	总统	在	莫斯科	与	俄罗斯	总统	普京	会晤
<i>Bush</i>	<i>President</i>	<i>in</i>	<i>Moscow</i>	<i>with</i>	<i>Russian</i>	<i>President</i>	<i>Putin</i>	<i>meet</i>

President Bush **meets** with Russian President Putin in Moscow

non-anticipative: President Bush (..... *waiting*.....) **meets** with Russian ...

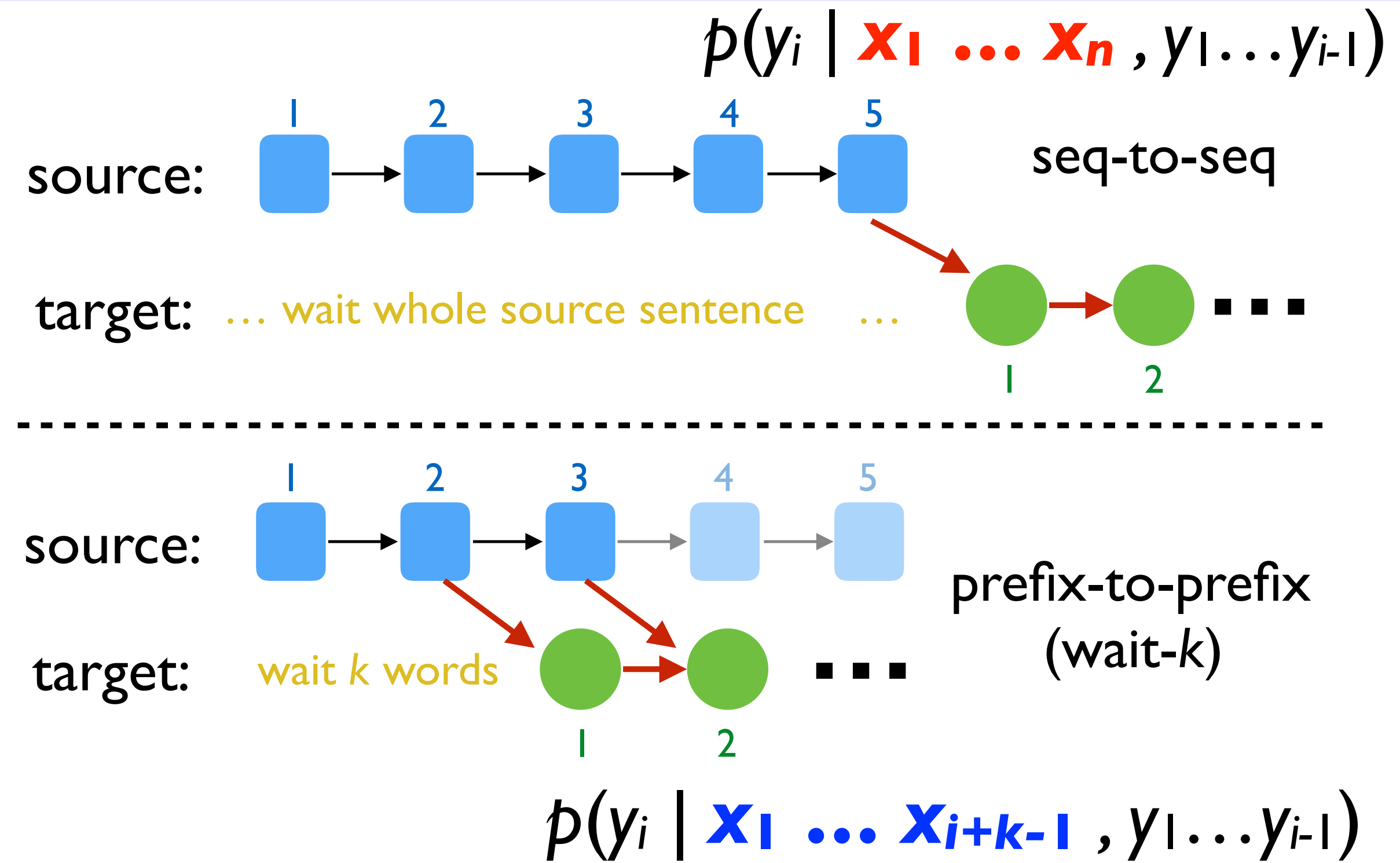
anticipative: President Bush **meets** with Russian President Putin in Moscow

Previous Solutions

- industrial systems
 - almost all “real-time” translation systems use full-sentence translation
 - some systems “repeatedly retranslate”, but constantly changing translations is annoying to the users and can’t be used for speech-to-speech translation
- academic papers (just to sample a few)
 - explicit prediction of German verbs (Grissom et al, 2014)
 - reinforcement learning (Gu et al, 2017) to decide READ or WRITE
 - segment-based (Bangalore et al, 2012; Fujita et al, 2013; Oda et al, 2014)
 - these efforts (a) use full-sentence translation model; (b) can’t ensure a given latency

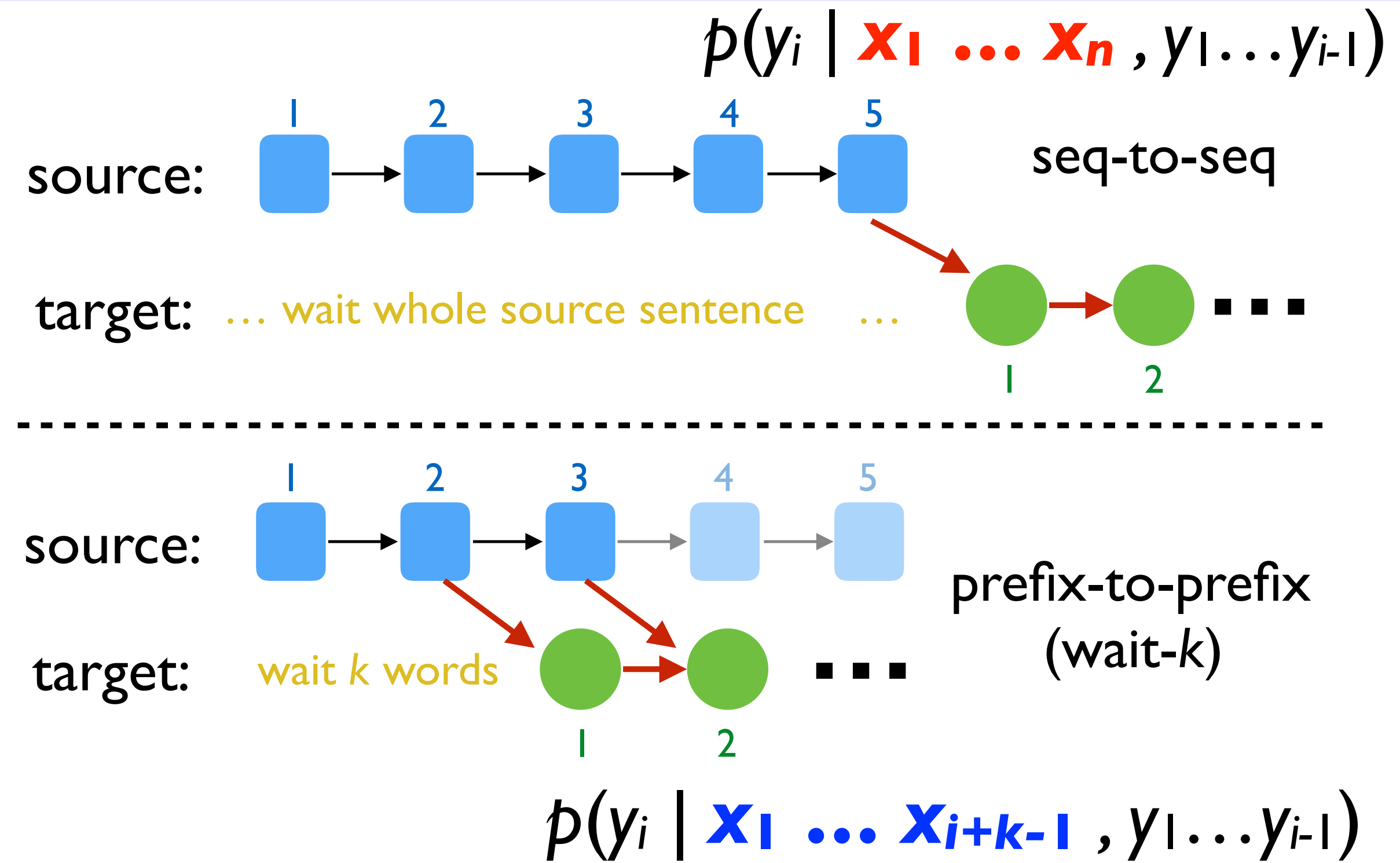
Our Idea: Prefix-to-Prefix, *not* Seq-to-Seq

- standard **seq-to-seq** is only suitable for conventional full-sentence MT
- we propose **prefix-to-prefix framework** tailed to tasks with simultaneity
- special case: **wait- k policy**: translation is always k words behind source sentence
- decoding this way \Rightarrow **controllable latency**
- training this way \Rightarrow **implicit anticipation on the target-side**



Our Idea: Prefix-to-Prefix, *not* Seq-to-Seq

- standard **seq-to-seq** is only suitable for conventional full-sentence MT
- we propose **prefix-to-prefix framework** tailed to tasks with simultaneity
- special case: **wait- k policy**: translation is always k words behind source sentence
- decoding this way \Rightarrow **controllable latency**
- training this way \Rightarrow **implicit anticipation on the target-side**

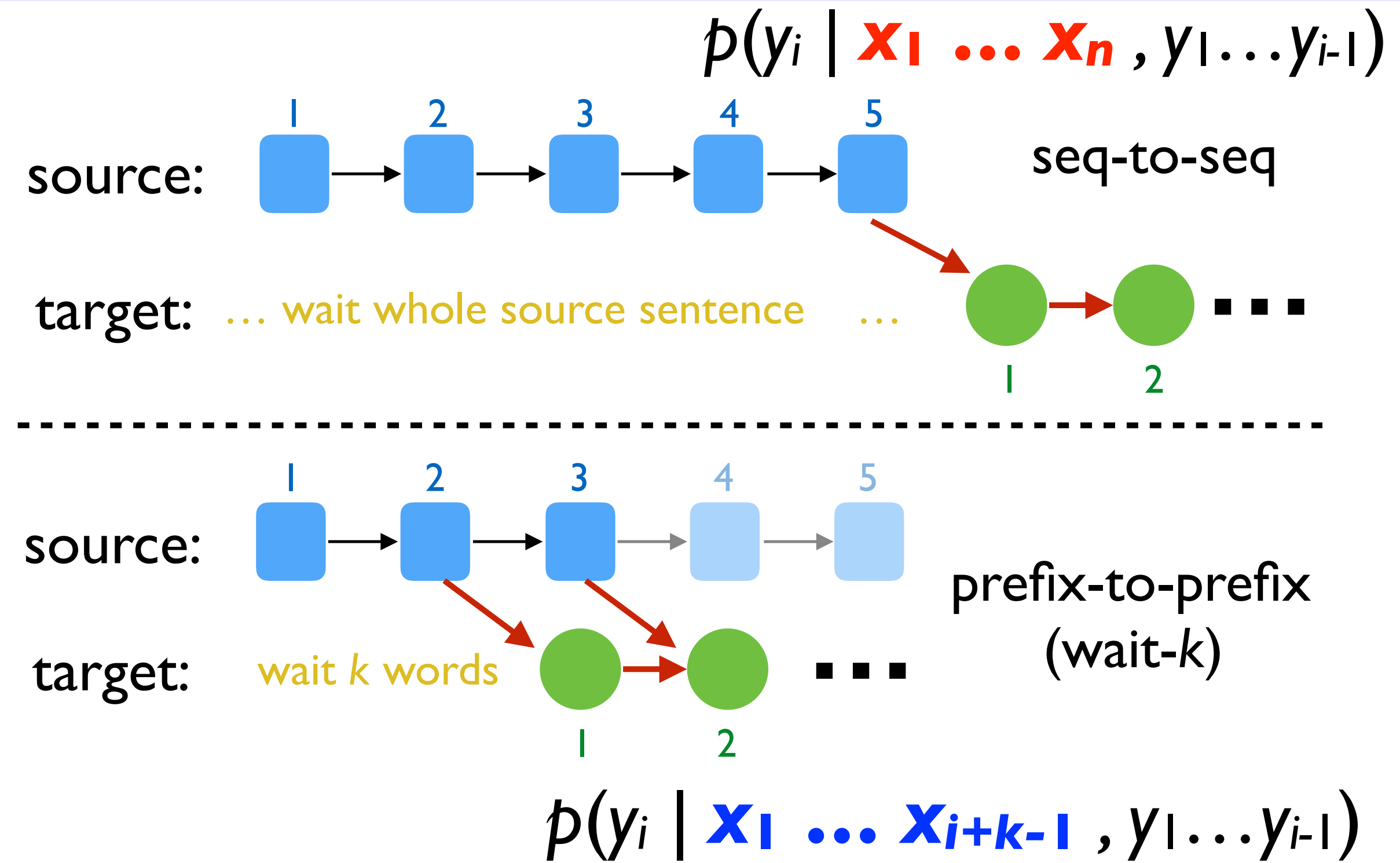


Bùshí	zǒngtǒng	zài	Mòsīkē
布什	总统	在	莫斯科
Bush	President	in	Moscow

wait 2 President Bush meets

Our Idea: Prefix-to-Prefix, *not* Seq-to-Seq

- standard **seq-to-seq** is only suitable for conventional full-sentence MT
- we propose **prefix-to-prefix framework** tailed to tasks with simultaneity
- special case: **wait- k policy**: translation is always k words behind source sentence
- decoding this way => **controllable latency**
- training this way => **implicit anticipation on the target-side**



Bùshí	zǒngtǒng	zài	Mòsīkē	yǔ	Éluósī	zǒngtǒng	Pǔjīng	huìwù
布什	总统	在	莫斯科	与	俄罗斯	总统	普京	会晤
Bush	President	in	Moscow	with	Russian	President	Putin	meet

wait 2 President Bush **meets** with Russian President Putin in Moscow

More General Prefix-to-Prefix

- seq-to-seq (given full source sent)

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

- prefix-to-prefix (given source prefix)

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

$g(\cdot)$ is a monotonic non-decreasing function

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

	Bush	Pres.	in	Moscow	with	Putin	meet
	布什	总统	在	莫斯科	与	普京	会晤
	President						
	Bush						
$t=3$	meets	$g(3) = 4$					
	with						
	Putin						
	in						
	Moscow						

this general framework can be used for other tasks such as incremental parsing and incremental text-to-speech

Research Demo

江泽民对法国总统的来华

jiang zemin expressed his appreciation

Research Demo

江泽民对法国总统的来华

jiang zemin expressed his appreciation

Research Demo

江泽民对法国总统的来华

jiang zemin expressed his appreciation

jiāng zémín duì fǎ guó zǒngtǒng de

江泽民对法国总统的

jiang zemin to French President's

láihuá fǎngwèn

来华访问

to-China visit

biǎoshì gǎnxiè

表示感谢。

express gratitude

jiang zemin expressed his appreciation for the visit by french president .

Research Demo

江泽民对法国总统的来华

jiang zemin expressed his appreciation

jiāng zé mǐn duì fǎ guó zǒng tǒng de

江泽民对法国总统的

jiang zemin to French President's

lái huá fǎng wèn

来华访问

to-China visit

biǎo shì gǎn xiè

表示感谢。

express gratitude

jiang zemin expressed his appreciation for the visit by french president .

我们支持 uh... 玻利维亚大使 和 俄罗斯大使 刚才 所做的立场

we support uh... Bolivia envoy & Russia envoy just-now made position

We support

the position

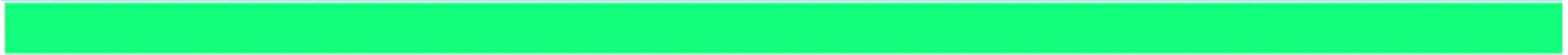
of Bolivia & Russia

Latency-Accuracy Tradeoff

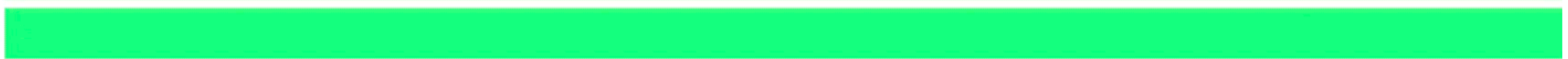
Chinese input:



Pinyin:



Word-by-Word
Translation:



Simultaneous
Translation (wait 3):



Simultaneous
Translation (wait 5):



Baseline
Translation (greedy):



Baseline
Translation (beam 5):



Latency-Accuracy Tradeoff

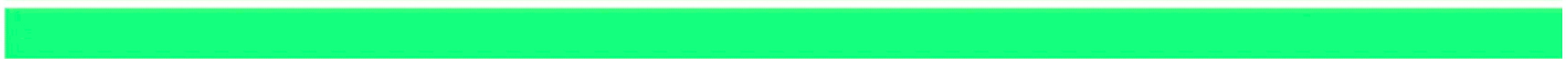
Chinese input:



Pinyin:



Word-by-Word
Translation:



Simultaneous
Translation (wait 3):



Simultaneous
Translation (wait 5):



Baseline
Translation (greedy):



Baseline
Translation (beam 5):



Deployment Demo



Deployment Demo



German=>English Anticipation Example

German source:

doch während man sich im kongress **nicht** auf ein vorgehen **einigen kann** , warten mehrere bundesstaaten nicht länger .
but while they self in congress not on one action agree can wait several states not longer

English translation (simultaneous, wait 3):

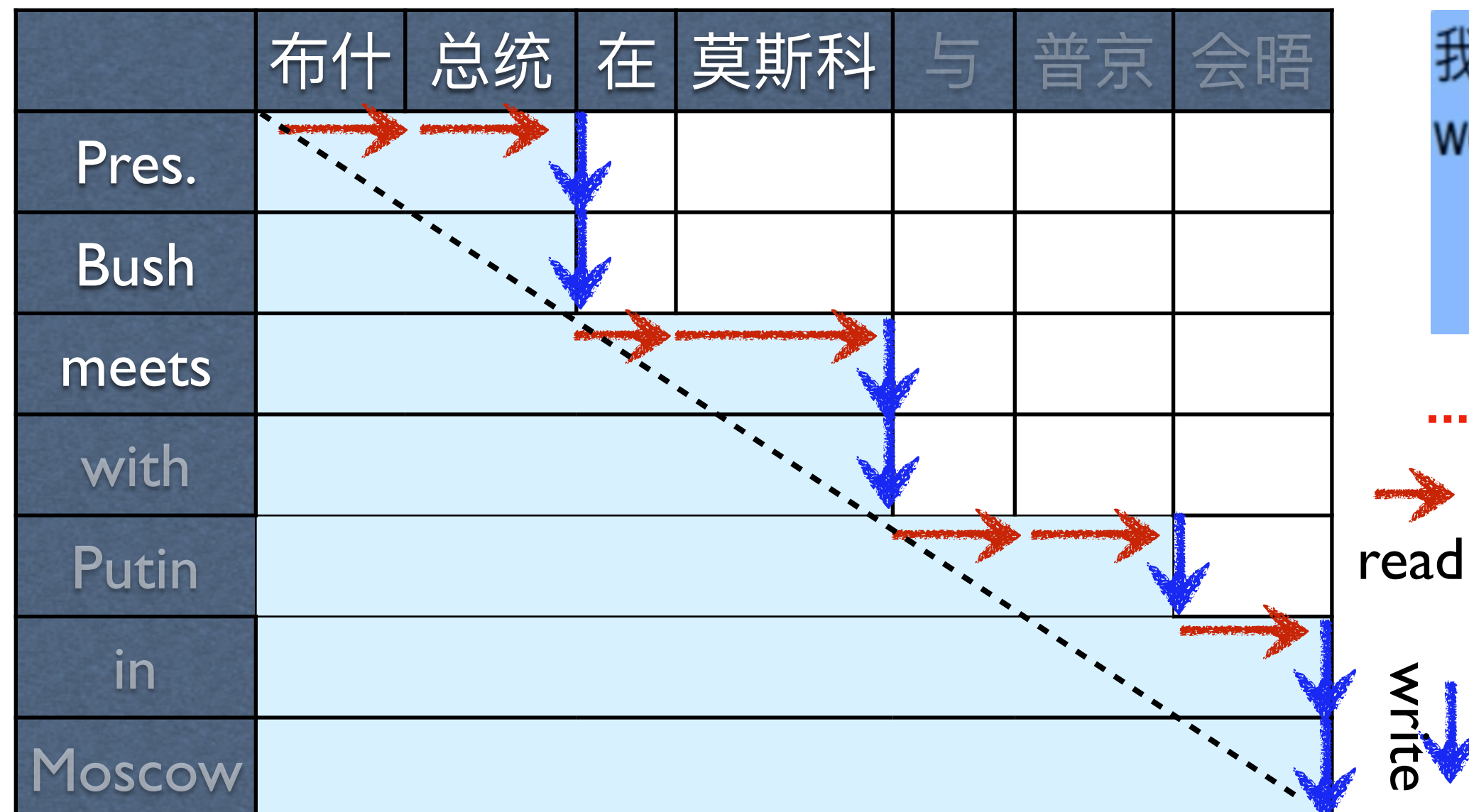
but , while congress **does not agree** on a course of action , several states no longer wait .

English translation (full-sentence baseline):

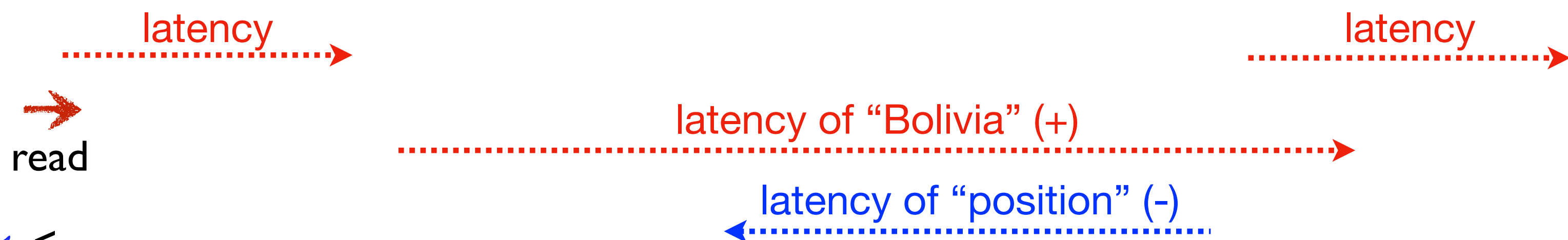
but , while congressional action **can not** be **agreed** , several states are no longer waiting .

New Latency Metric: Average Lagging

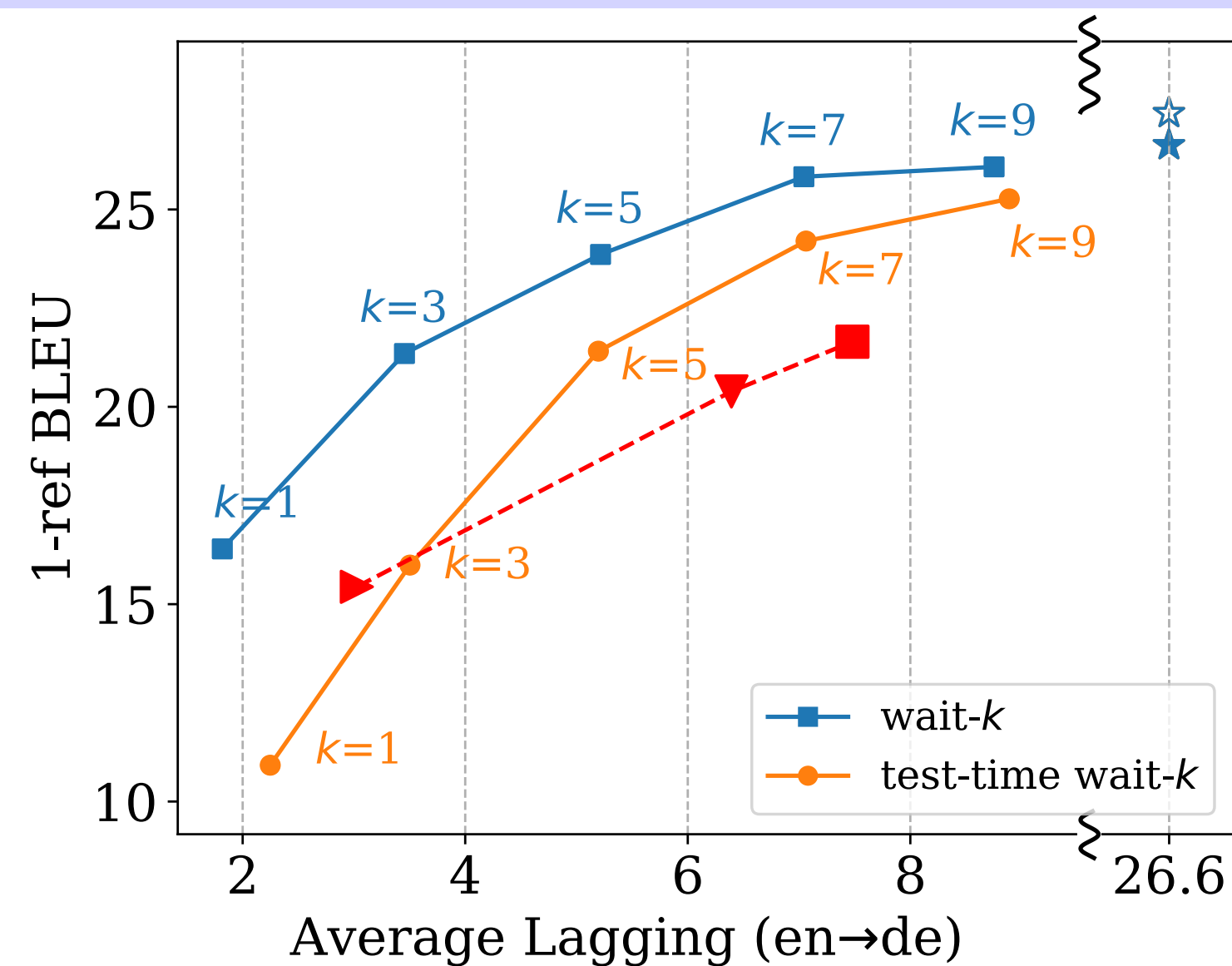
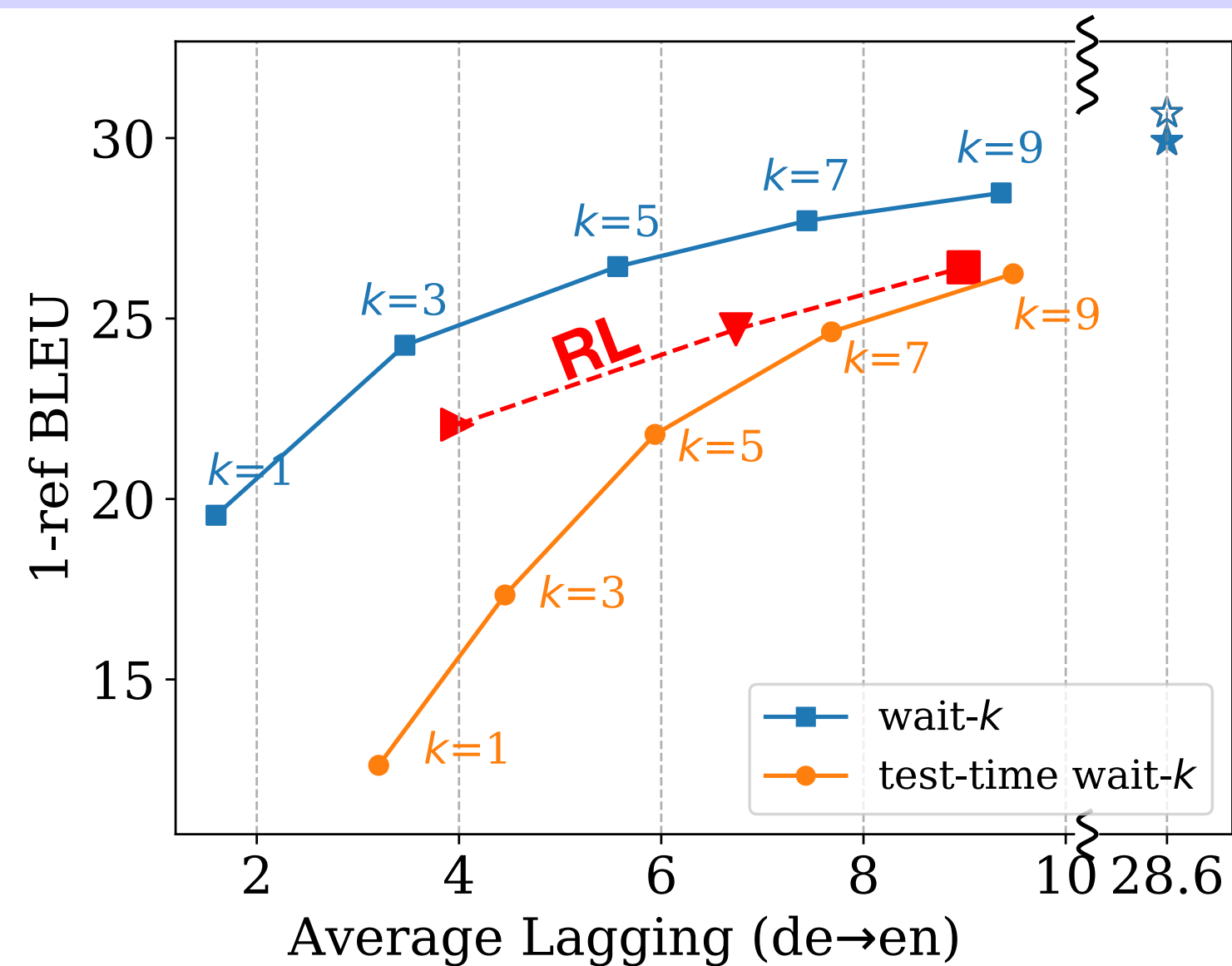
- previous metrics: CW (consecutive wait) and AP (average proportion)
 - they do *not* directly measure the level of “lagging behind” (Gu et al '17; Cho & Esipova '16)
- our metric, *Average Lagging (AL)*, measures on average how many source words the translation lags behind the source speech; ideally, $AL(\text{wait-}k) \approx k$
- closely related to “**ear-voice span**” (EVS) in the interpretation literature



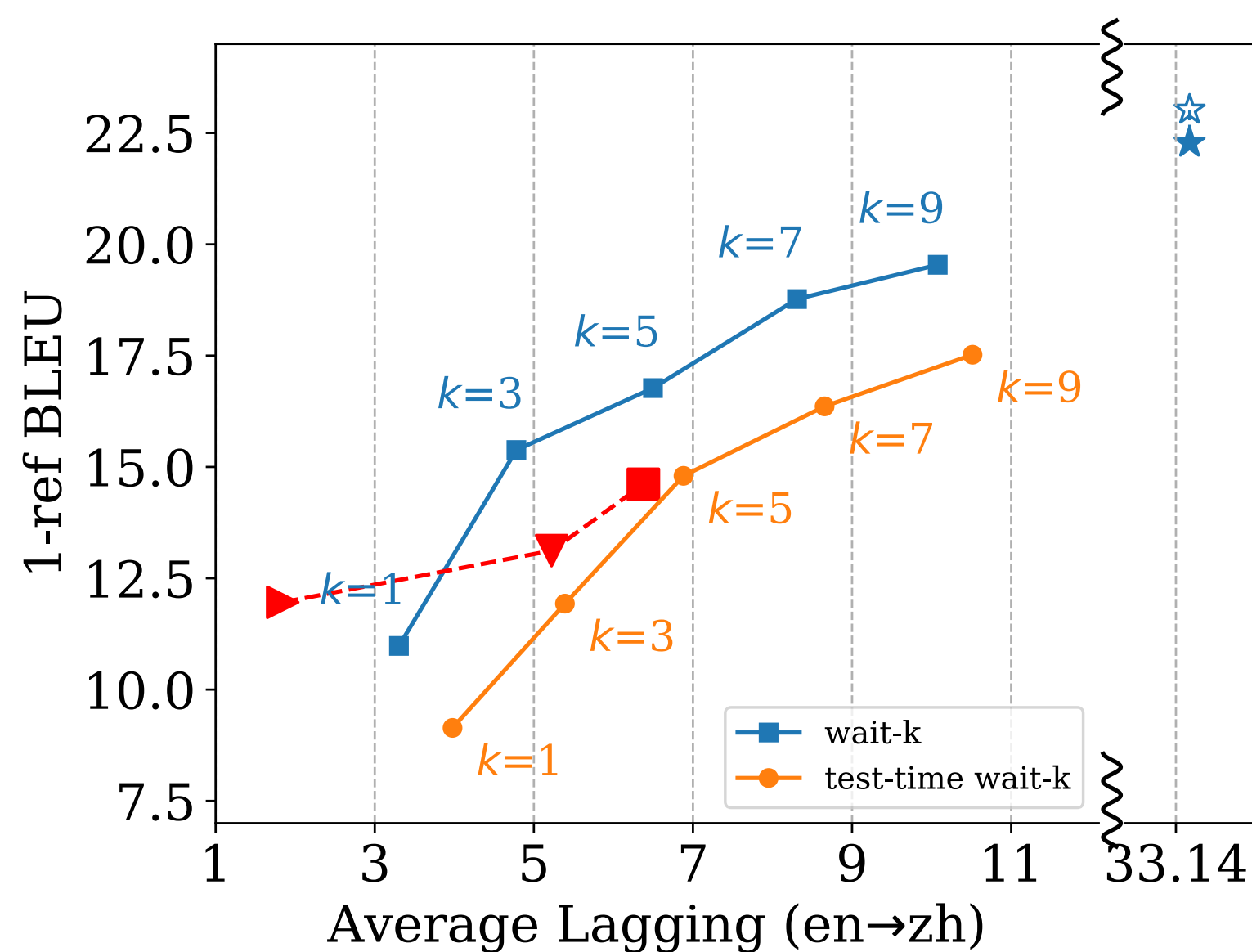
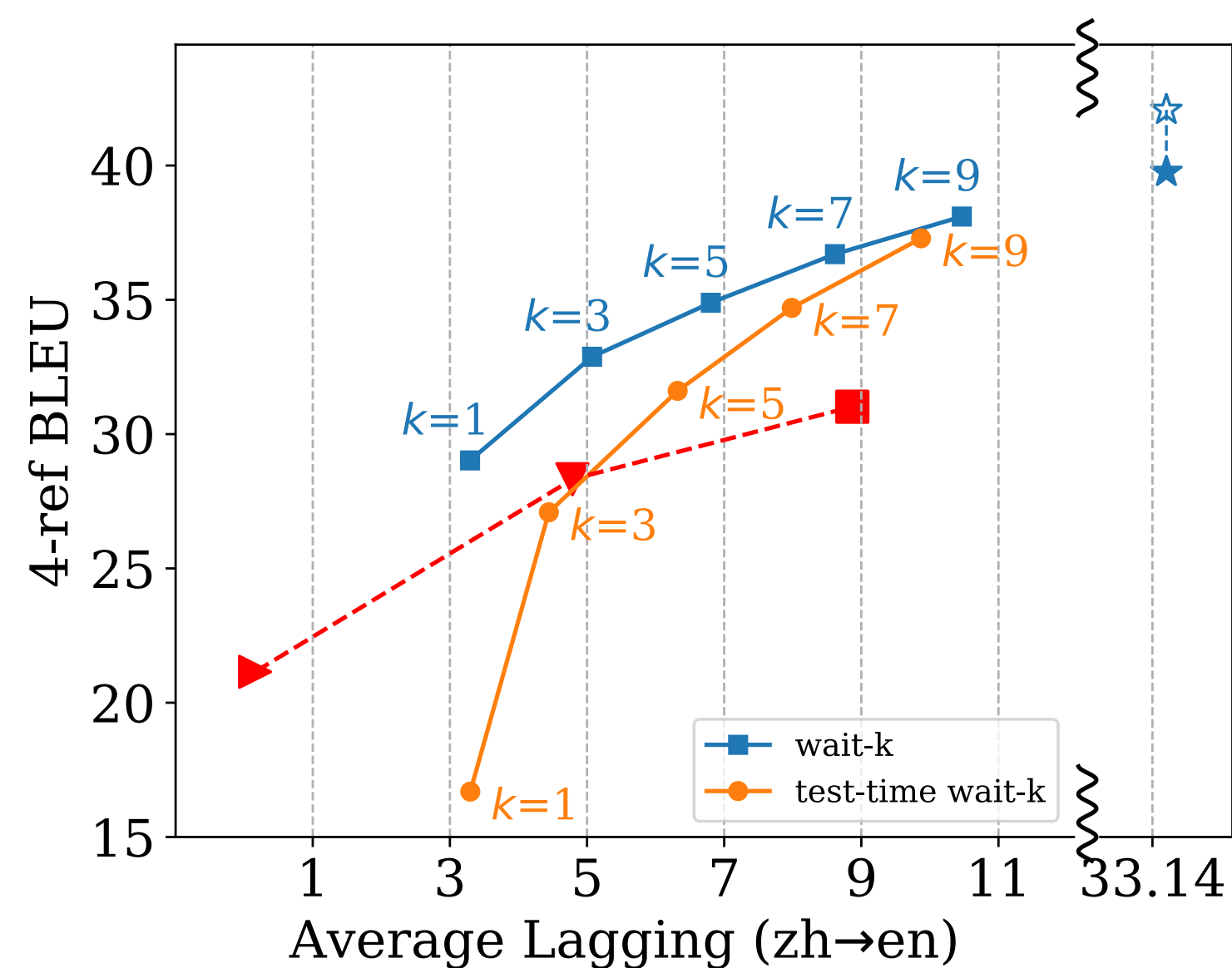
我们支持 uh... 玻利维亚大使 和 俄罗斯大使 刚才 所做的立场
 we support uh... Bolivia envoy & Russia envoy just-now made position
 We support the position of Bolivia & Russia



Experiments (de↔en & zh↔en)



RL: our adaptation of Gu et al (2017) on the same Transformer codebase, trained with CW=2, 5, 8.



Summary of Innovations in 2018

- prefix-to-prefix framework tailed to simultaneity (incremental on both sides)
 - first genuinely simultaneous translation model (rather than full-sentence model)
 - decoding like this => controllable latency
 - training like this => implicit anticipation on the target side
- very easy to train and scalable — minor changes to most neural MT codebase
- prefix-to-prefix is very general; can be used in other tasks with simultaneity
- a new latency metric (AL) that resembles “ear-voice span” in interpretation

Media coverage:



Part II: Towards Adaptive Translation Policies

Part II: Towards Adaptive Translation Policies

	<i>fixed-latency policies</i>	<i>adaptive policies</i>
<i>full-sentence MT model</i>	Dalvi et al. (2018); test-time wait-<i>k</i> (Ma et al. 2018)	Grissom et al. (2014); Cho & Esipova (2016); Satija & Pineau (2016); Gu et al. (2017); Alinejad et al (2018); ...
<i>simultaneous MT model (our invention)</i>	wait-<i>k</i> (Ma et al. 2018)	Arivazhagan et al. (ACL 2019) Zheng et al. (ACL 2019)

Limitations of Fixed-Latency (wait- k) Policy

- can be too aggressive (**anticipation errors**) with small k (too fast)

Limitations of Fixed-Latency (wait- k) Policy

- can be too aggressive (**anticipation errors**) with small k (too fast)

input	wǒ	shàng	wèi	dédào	yǒuguān	bùmén
	我	尚	未	得到	有关	部门
	<i>I</i>	<i>yet</i>	<i>not</i>	<i>receive</i>	<i>relevant</i>	<i>department</i>
wait-1 (AL=1.4)	I	have	not	received	relevant	

Limitations of Fixed-Latency (wait- k) Policy

- can be too aggressive (**anticipation errors**) with small k (too fast)

input	wǒ	shàng	wèi	dédào	yǒuguān	bùmén	
	我	尚	未	得到	有关	部门	
	<i>I</i>	<i>yet</i>	<i>not</i>	<i>receive</i>	<i>relevant</i>	<i>department</i>	
wait-1 (AL=1.4)	I	have	not	received	relevant	documents	

Limitations of Fixed-Latency (wait- k) Policy

- can be too aggressive (**anticipation errors**) with small k (too fast)

input	wǒ 我	shàng 尚	wèi 未	dédào 得到	yǒuguān 有关	bùmén 部门	de 的	huíyìng 回应
	<i>I</i>	<i>yet</i>	<i>not</i>	<i>receive</i>	<i>relevant</i>	<i>department</i>	's	<i>response</i>
wait-1 (AL=1.4)	I	have	not	received	relevant	documents	from	relevant departments

Limitations of Fixed-Latency (wait- k) Policy

- can be too aggressive (**anticipation errors**) with small k (too fast)
- can also be too conservative with large k (too slow)

input	wǒ 我 <i>I</i>	shàng 尚 <i>yet</i>	wèi 未 <i>not</i>	dédào 得到 <i>receive</i>	yǒuguān 有关 <i>relevant</i>	bùmén 部门 <i>department</i>	de 的 <i>'s</i>	huíying 回应 <i>response</i>
wait-1 (AL=1.4)	I	have	not	received	relevant	documents	from	relevant departments
wait-4 (AL=4.0)		have	not	received	response	from	relevant departments	

Limitations of Fixed-Latency (wait- k) Policy

- can be too aggressive (**anticipation errors**) with small k (too fast)
- can also be too conservative with large k (too slow)

input	wǒ 我 <i>I</i>	shàng 尚 <i>yet</i>	wèi 未 <i>not</i>	dédào 得到 <i>receive</i>	yǒuguān 有关 <i>relevant</i>	bùmén 部门 <i>department</i>	de 的 <i>'s</i>	huíying 回应 <i>response</i>
wait-1 (AL=1.4)	I	have	not	received	relevant	documents	from	relevant departments
wait-4 (AL=4.0)				I	have	not	received	response from relevant departments
adaptive (AL=1.8)	I	<u>have not</u>	received					response from relevant departments

Previous Work on Adaptive Policy

- READ and WRITE actions

Bùshí
布什
Bush

zǒngtǒng
总统
President

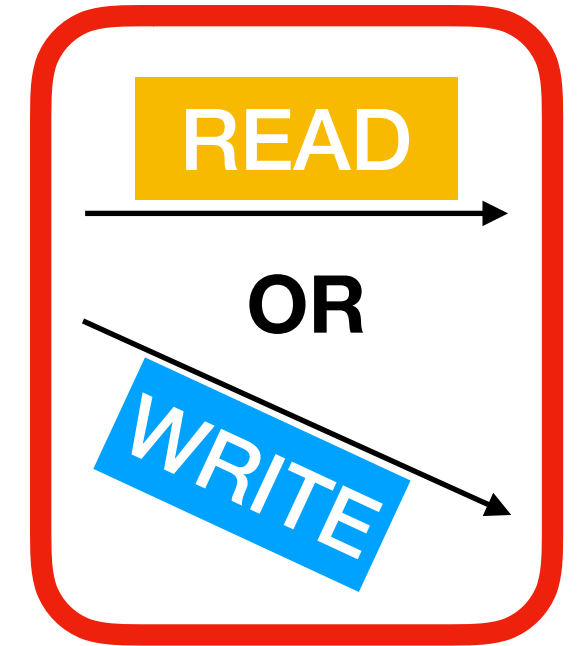
zài
在
in

Mòskē
莫斯科
Moscow

President

Bush

Action

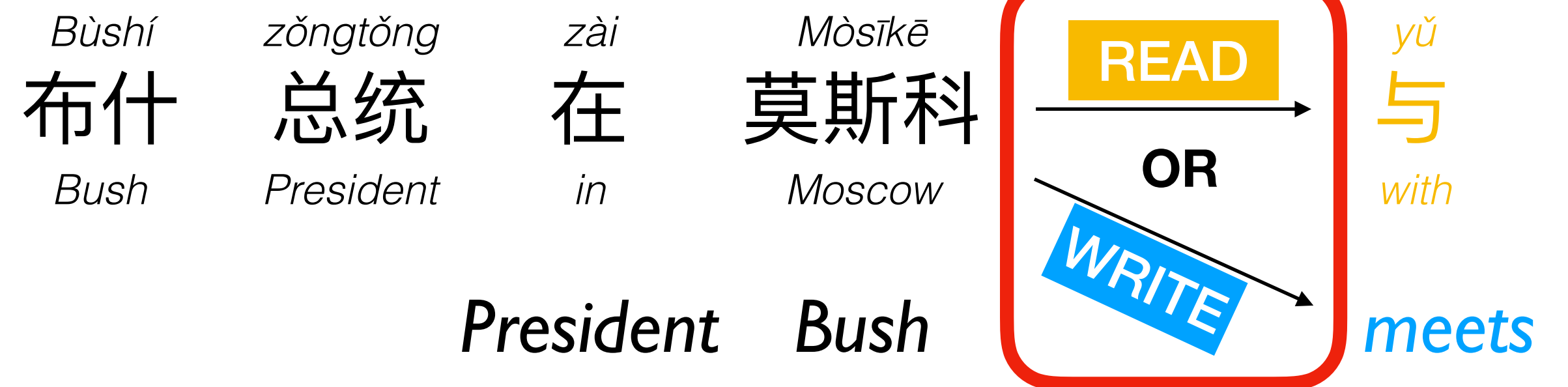


yǔ
与
with

meets

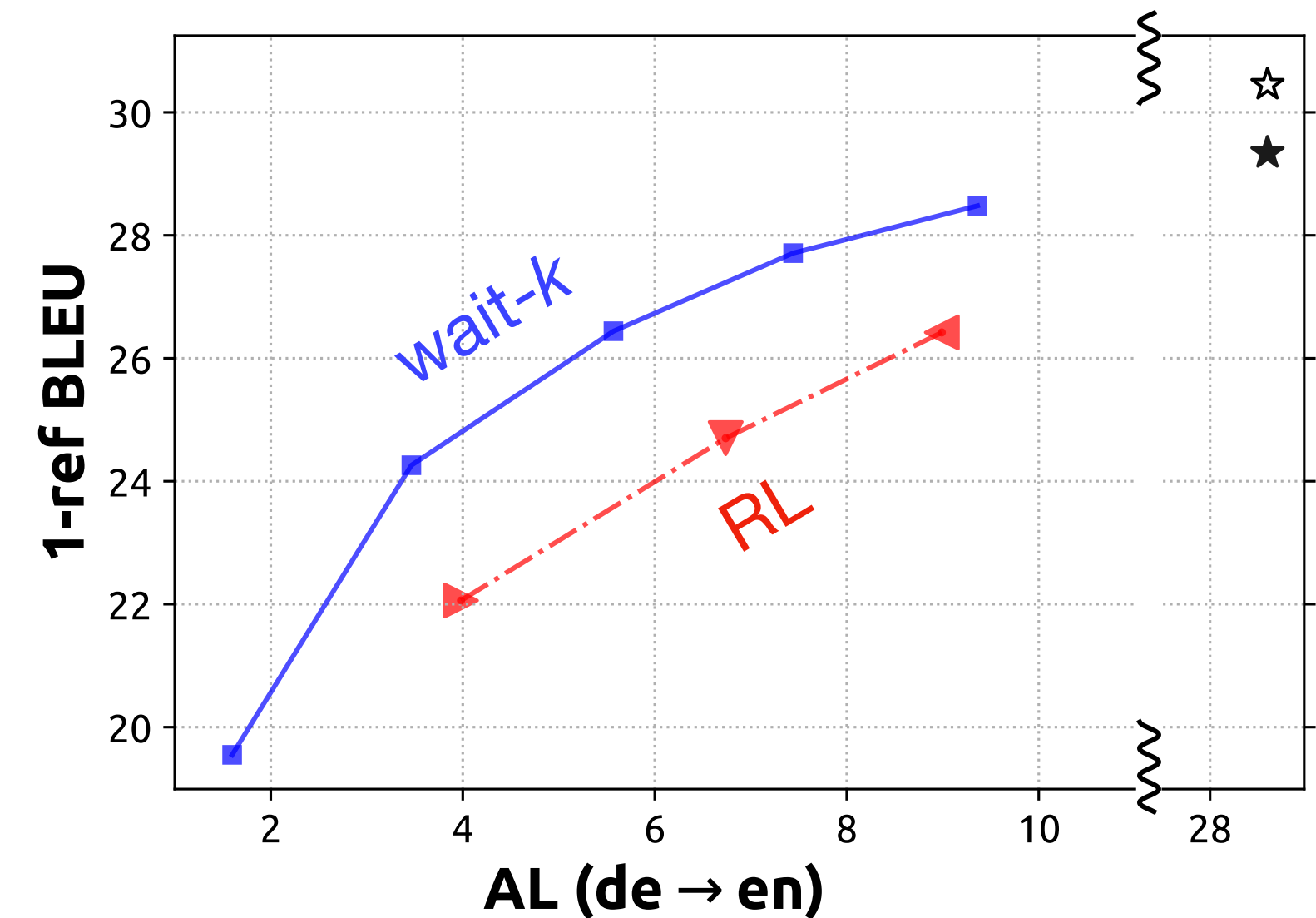
Previous Work on Adaptive Policy

- READ and WRITE actions



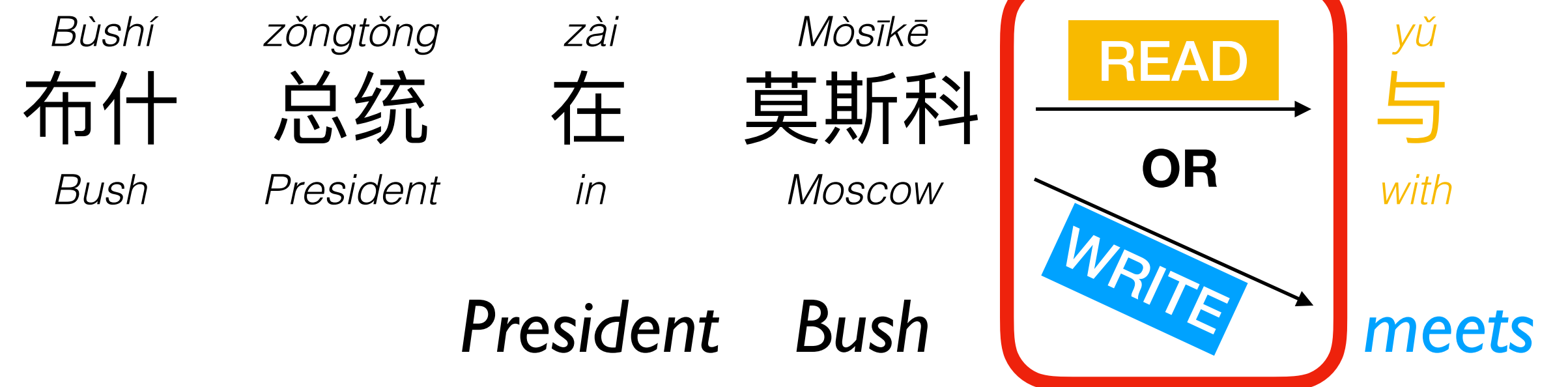
- sequential decision making → reinforcement learning (Gu et al. 2017)

- unstable training (randomness in exploration)
- complicated (two models trained in two stages)
- worse performance (than wait-k model)



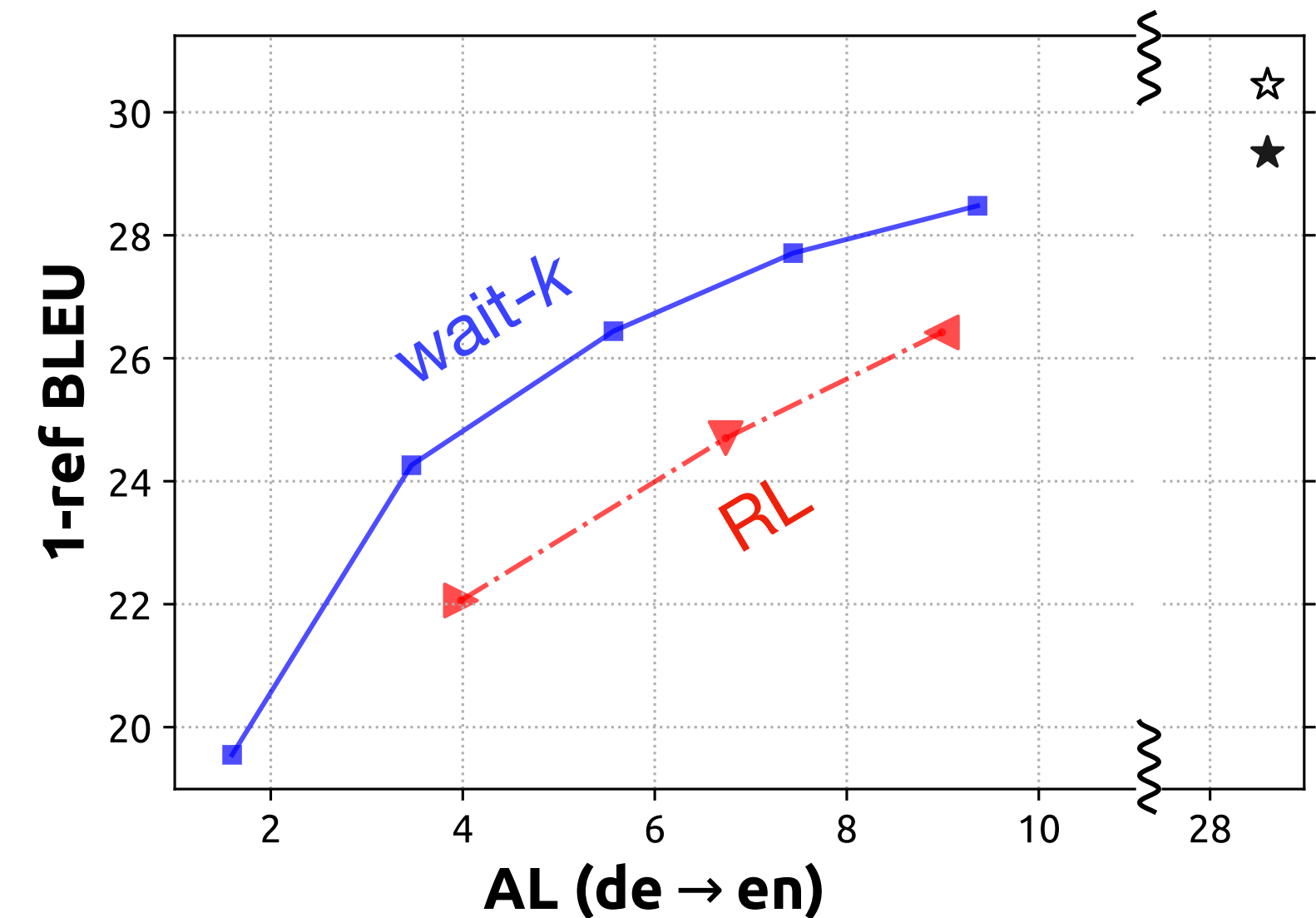
Previous Work on Adaptive Policy

- READ and WRITE actions



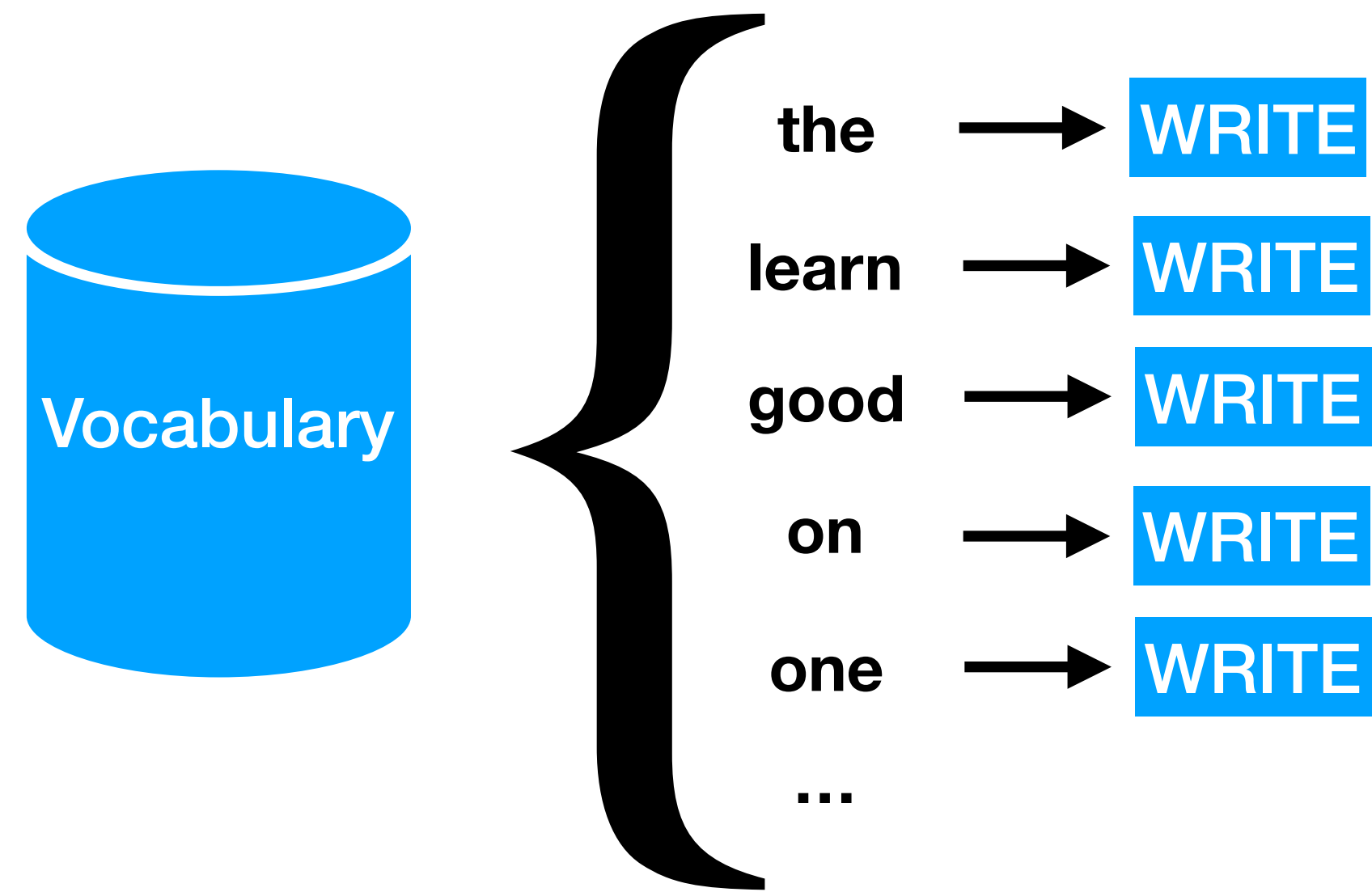
- sequential decision making → reinforcement learning (Gu et al. 2017)

- unstable training (randomness in exploration)
- complicated (two models trained in two stages)
- worse performance (than wait-k model)

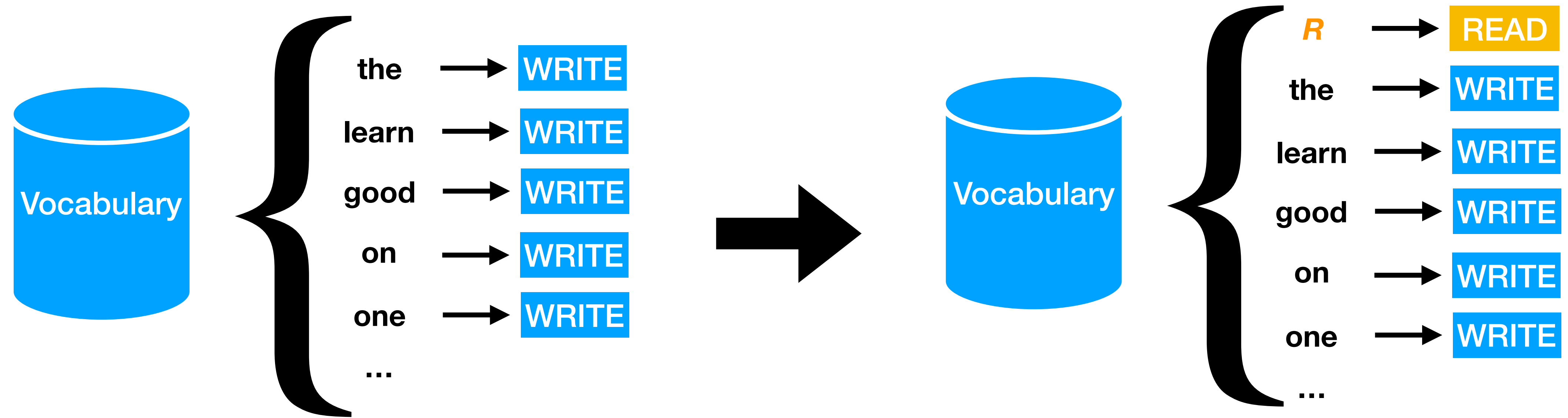


- can we learn a better model with adaptive policy via simpler methods ?

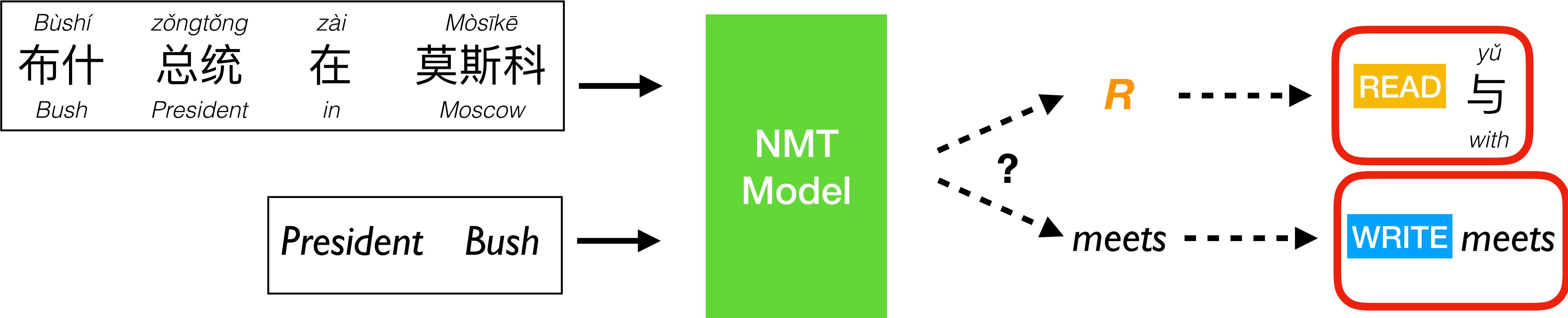
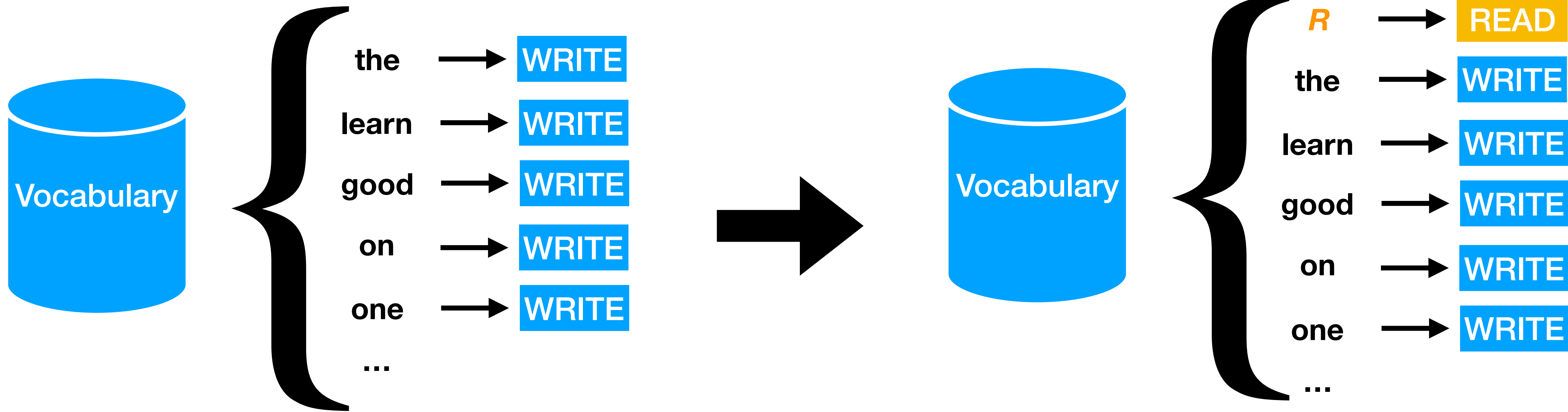
Our Idea: Single Model, with READ as a Word



Our Idea: Single Model, with READ as a Word



Our Idea: Single Model, with READ as a Word



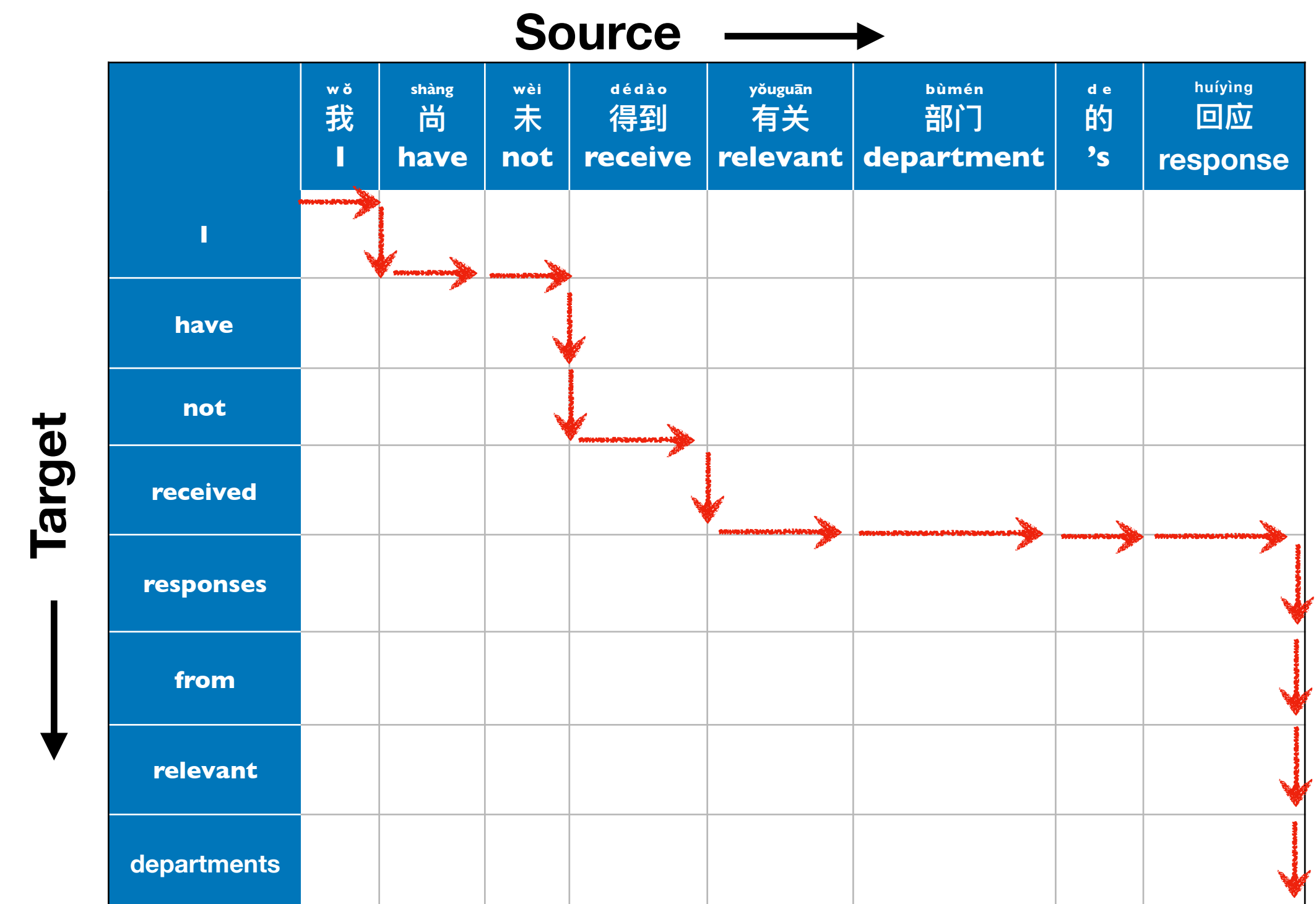
Learn a Single Model via Imitation Learning

- imitation learning
 - learn to imitate a given expert policy



Learn a Single Model via Imitation Learning

- imitation learning
 - learn to imitate a given expert policy
- basic ideas
 - merge two models into one
 - add read action into target vocabulary
 - end-to-end training
 - design an expert policy to use imitation learning



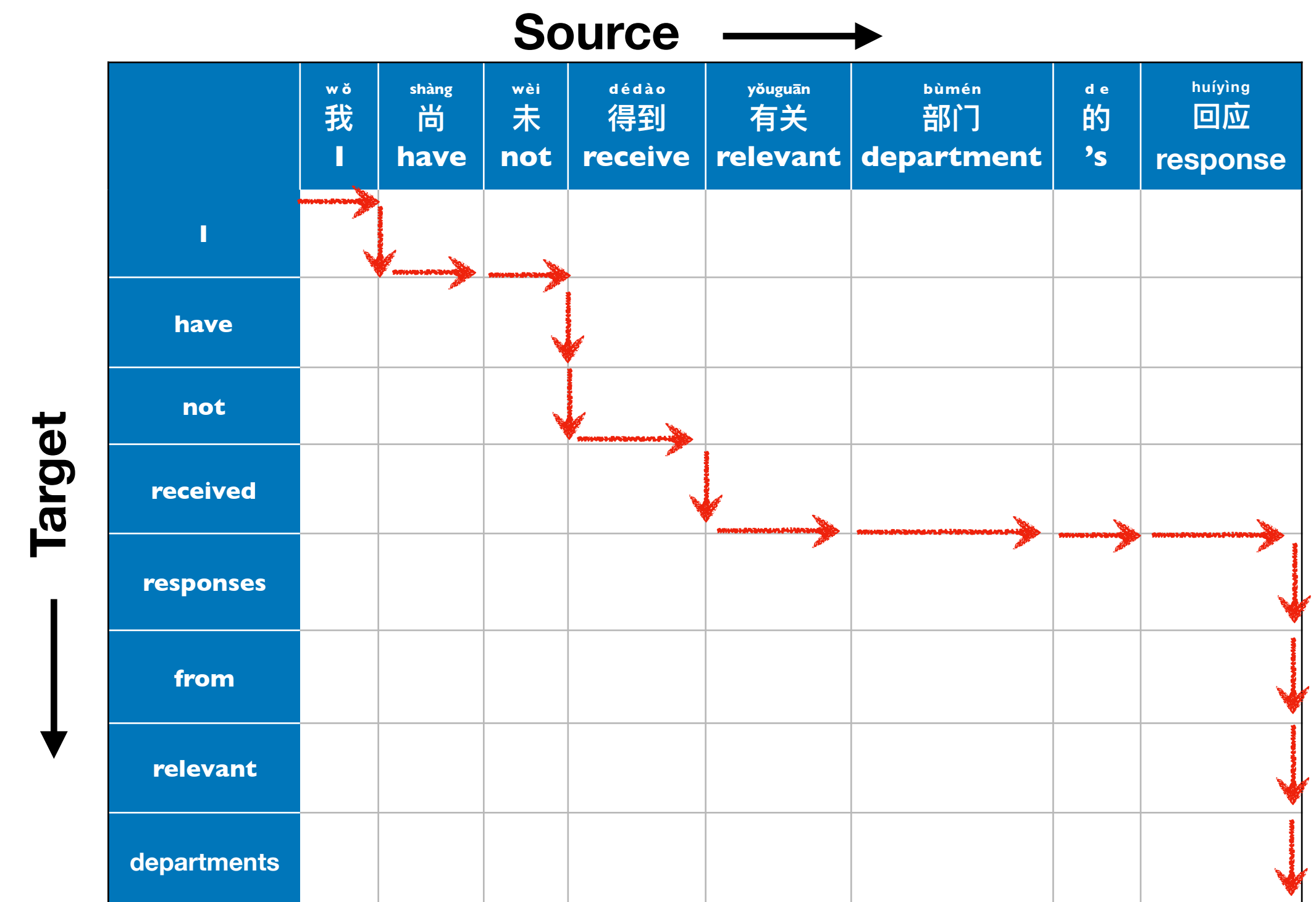
Learn a Single Model via Imitation Learning

- imitation learning
- learn to imitate a given expert policy



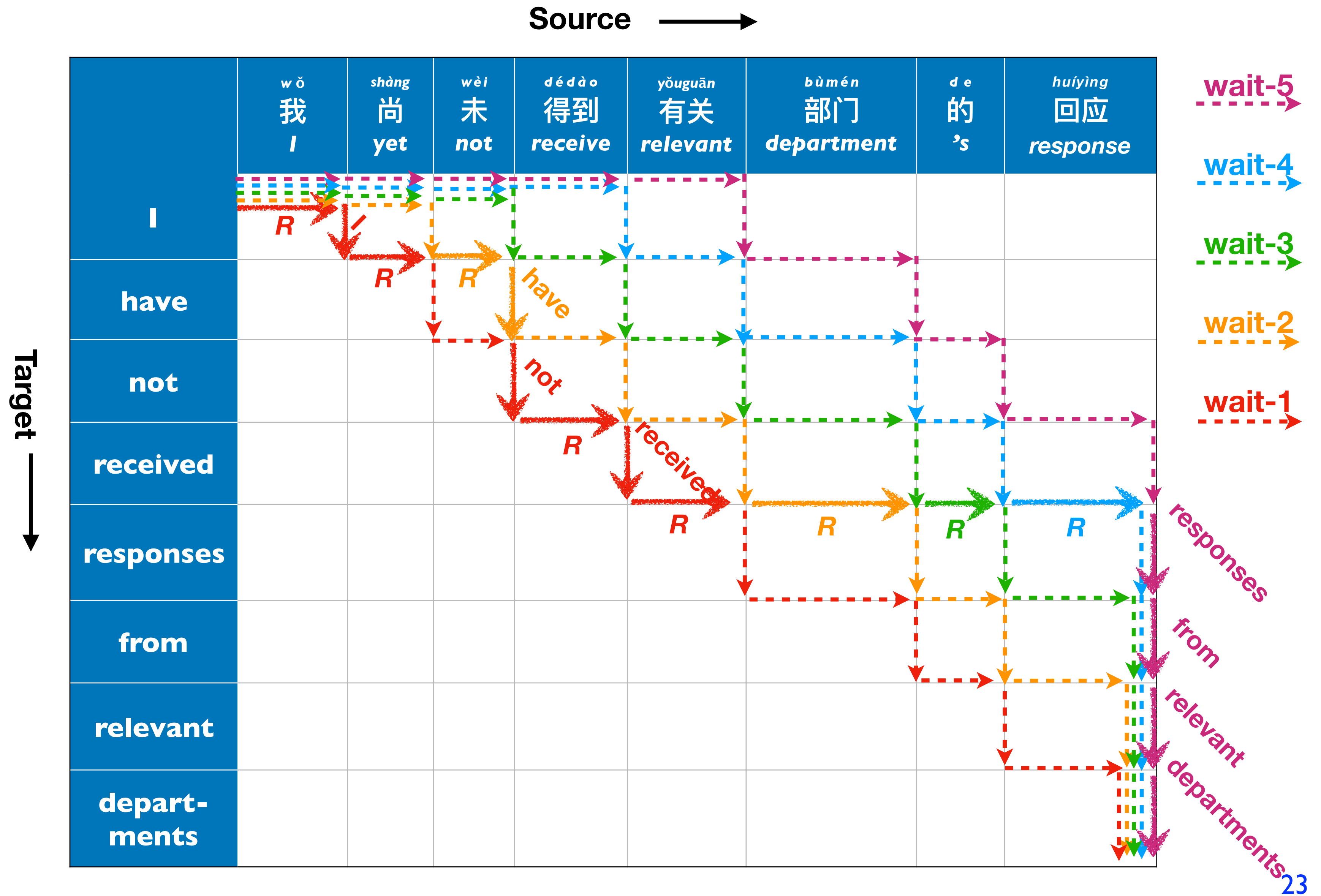
*for more details
come to my short talk tomorrow*

- basic ideas
- merge two models into one
 - add read action into target vocabulary
- end-to-end training
 - design an expert policy to use imitation learning



Another Much Simpler Idea

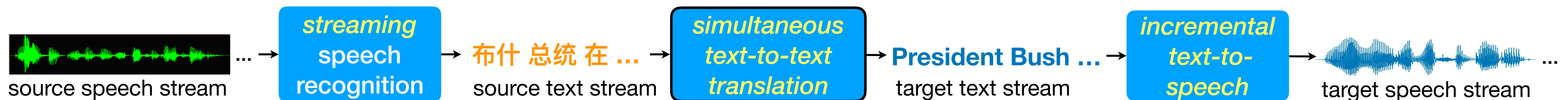
- on-the-fly decide READ or WRITE
- depending on $p(y_i | \dots)$
- if not confident enough, READ
 - switch to wait-(k+1) (more conservative)
- otherwise WRITE
 - switch to wait-(k-1) (more aggressive)



Part III: Remaining Challenges

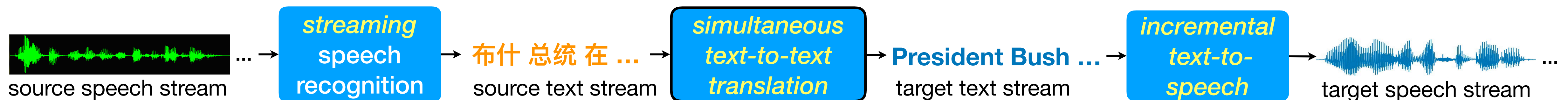
Part III: Remaining Challenges

- Speech Recognition-related
 - coping with ASR noise, esp. homophones
 - code switching
 - sentence breaking
 - prosody lost in translation
 - directly speech-to-speech without text-to-text?



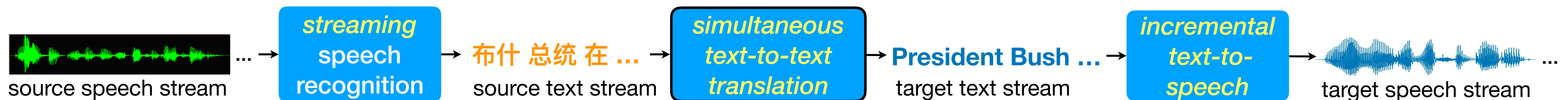
Part III: Remaining Challenges

- Speech Recognition-related
 - coping with ASR noise, esp. homophones
 - code switching
 - sentence breaking
 - prosody lost in translation
 - directly speech-to-speech without text-to-text?
- Incremental Text-to-Speech Synthesis (TTS)



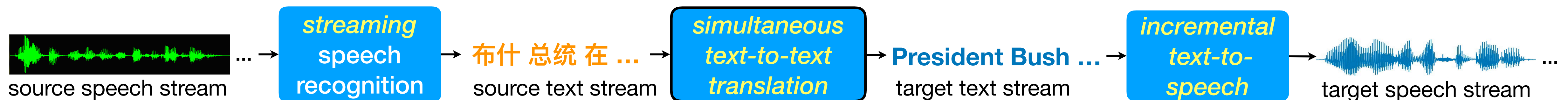
Part III: Remaining Challenges

- Speech Recognition-related
 - coping with ASR noise, esp. homophones
 - code switching
 - sentence breaking
 - prosody lost in translation
 - directly speech-to-speech without text-to-text?
- Incremental Text-to-Speech Synthesis (TTS)
- Better Dataset for Training



Part III: Remaining Challenges

- Speech Recognition-related
 - coping with ASR noise, esp. homophones
 - code switching
 - sentence breaking
 - prosody lost in translation
 - directly speech-to-speech without text-to-text?
- Incremental Text-to-Speech Synthesis (TTS)
- Better Dataset for Training
- Detecting and Fixing Mistakes (esp. anticipation errors)



Coping with ASR noise

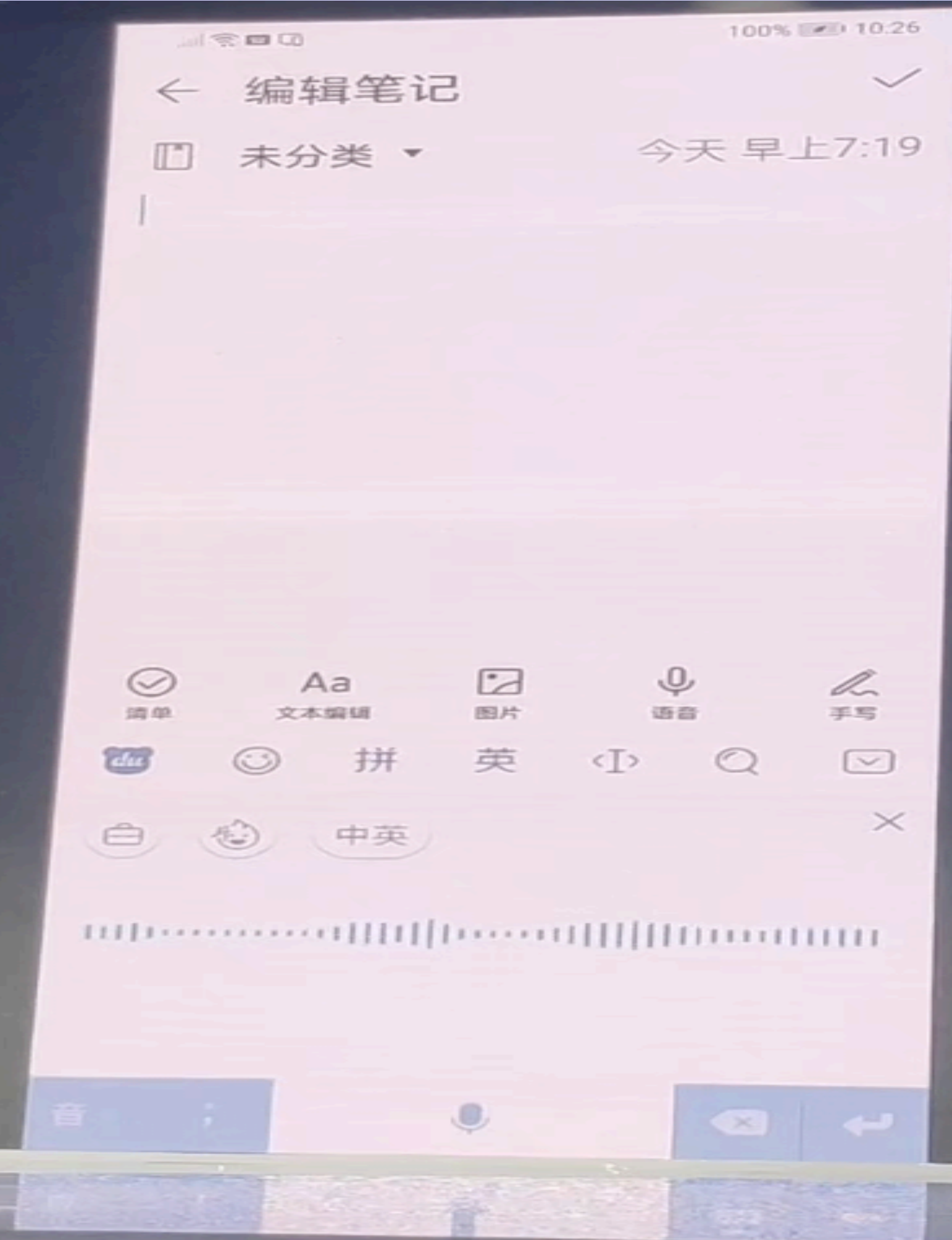
- neural MT is fragile, and automatic speech recognition (ASR) output is noisy
- our work (Liu et al, ACL 2019): robust neural MT using phonetic information

Clean Input	目前已发现有109人死亡, 另有57人获救	yǒu 有
Output of Transformer	at present, 109 people have been found dead and 57 have been rescued	have
Noisy Input	目前已发现又109人死亡, 另有57人获救	yòu 又
Output of Transformer	the hpv has been found dead so far and 57 have been saved	again
Output of Our Method	so far, 109 people have been found dead and 57 others have been rescued	

Table 1: The translation results on Mandarin sentences without and with homophone noises. The word ‘有’ (yǒu, “have”) in clean input is replaced by one of its homophone, ‘又’ (yòu, “again”), to form a noisy input. This seemingly minor change completely fools the Transformer to generate something irrelevant (“hpv”). Our method, by contrast, is very robust to homophone noises thanks to phonetic information.

Baidu ASR's Code-Switching Capabilities

Baidu AI Create, July 2019

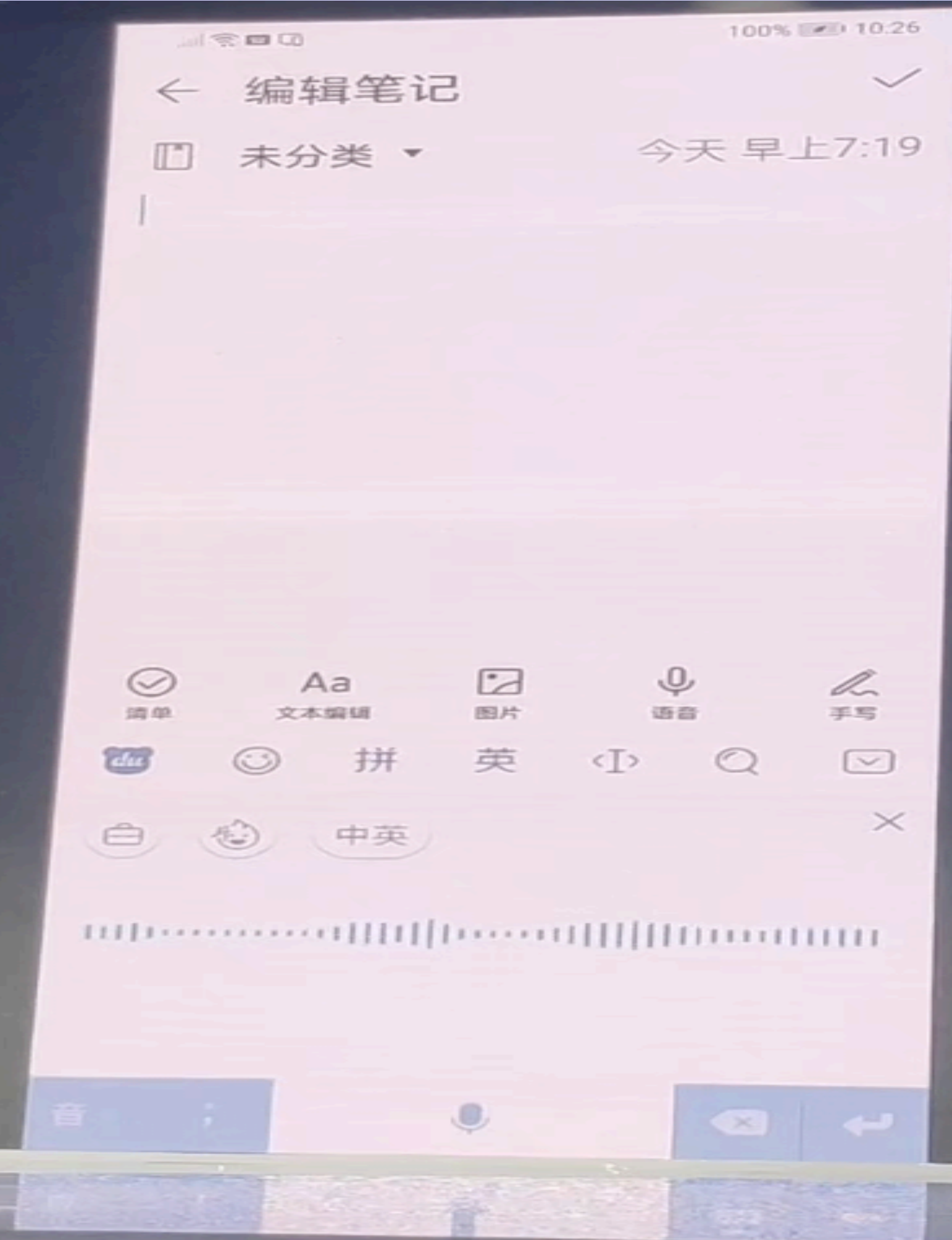


百度输入法
中英文自由说

- Baidu ASR is awesome at code-switching (English terms in Chinese speech)

Baidu ASR's Code-Switching Capabilities

Baidu AI Create, July 2019



百度输入法
中英文自由说

- Baidu ASR is awesome at code-switching (English terms in Chinese speech)

Baidu ASR's Code-Switching Capabilities

Baidu AI Create, July 2019

百度输入法
中英文自由说

- Baidu ASR is awesome at code-switching (English terms in Chinese speech)

Better Dataset for Training Simultaneous Translation

- standard parallel text is not made for simultaneous translation
 - involves too many “unnecessary long-distance reorderings”
- simultaneous interpretation corpora is not ideal training data either
 - contains too many mistakes, speech repairs, and compressions
- again, our goal is short latency (like human simultaneous interpretation) and good quality (like human written translation)

uh... 我们 认为 安理会 ah... 没有 必要 介入
uh... we think sec. council uh... no need intervene

We believe that it is (not) necessary for the security council to get involved

Better Dataset for Training Simultaneous Translation

- standard parallel text is not made for simultaneous translation
 - involves too many “unnecessary long-distance reorderings”
- simultaneous interpretation corpora is not ideal training data either
 - contains too many mistakes, speech repairs, and compressions
- again, our goal is short latency (like human simultaneous interpretation) and good quality (like human written translation)

uh... 我们 认为 安理会 ah... 没有 必要 介入
uh... we think sec. council uh... no need intervene

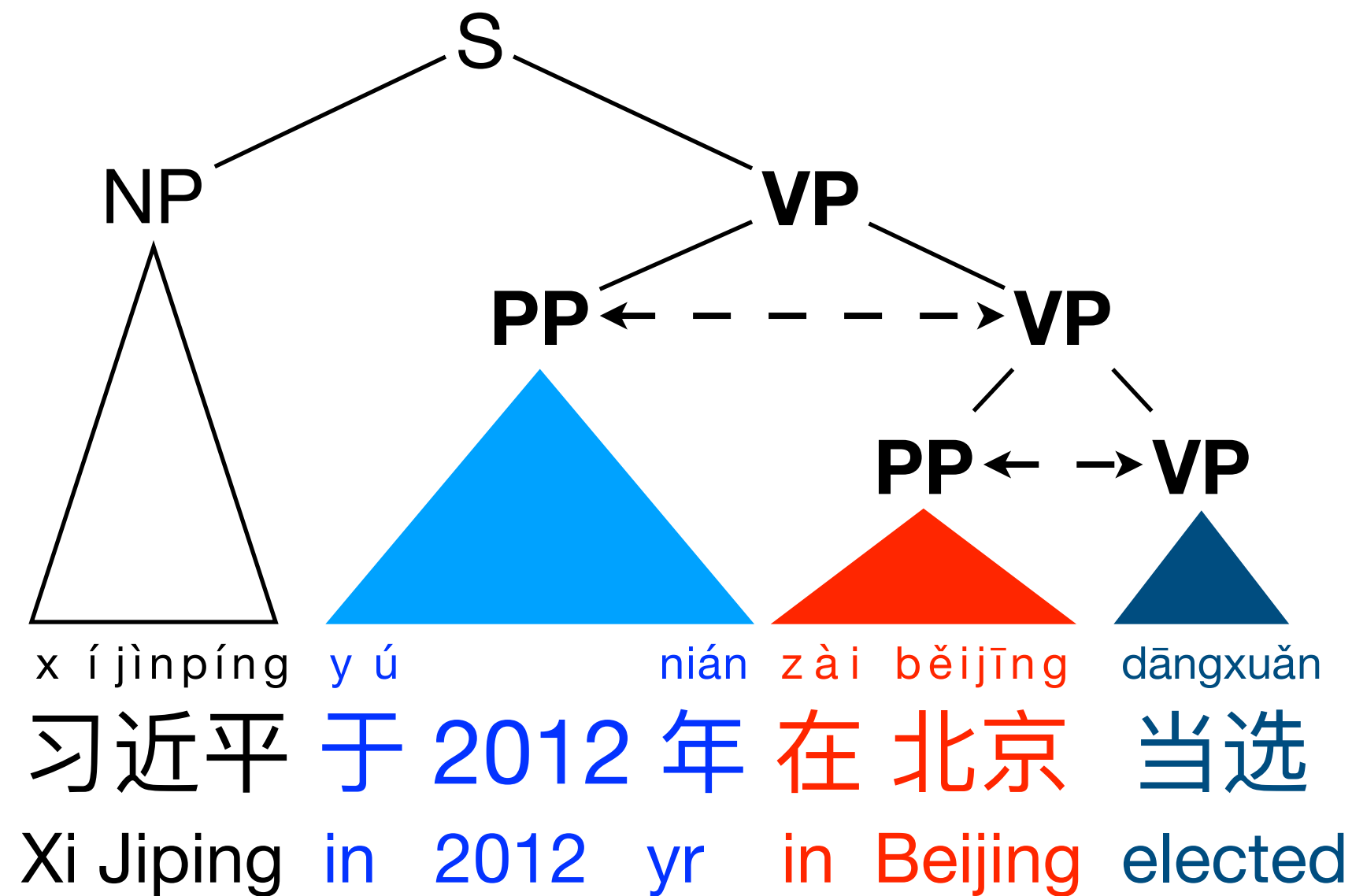
We believe that it is (not) necessary for the security council to get involved

Better Dataset for Training Simultaneous Translation

- idea: rephrase target side of parallel text to remove unnecessary reorderings

mandatory reordering

(Chinese) PP VP => (English) VP PP

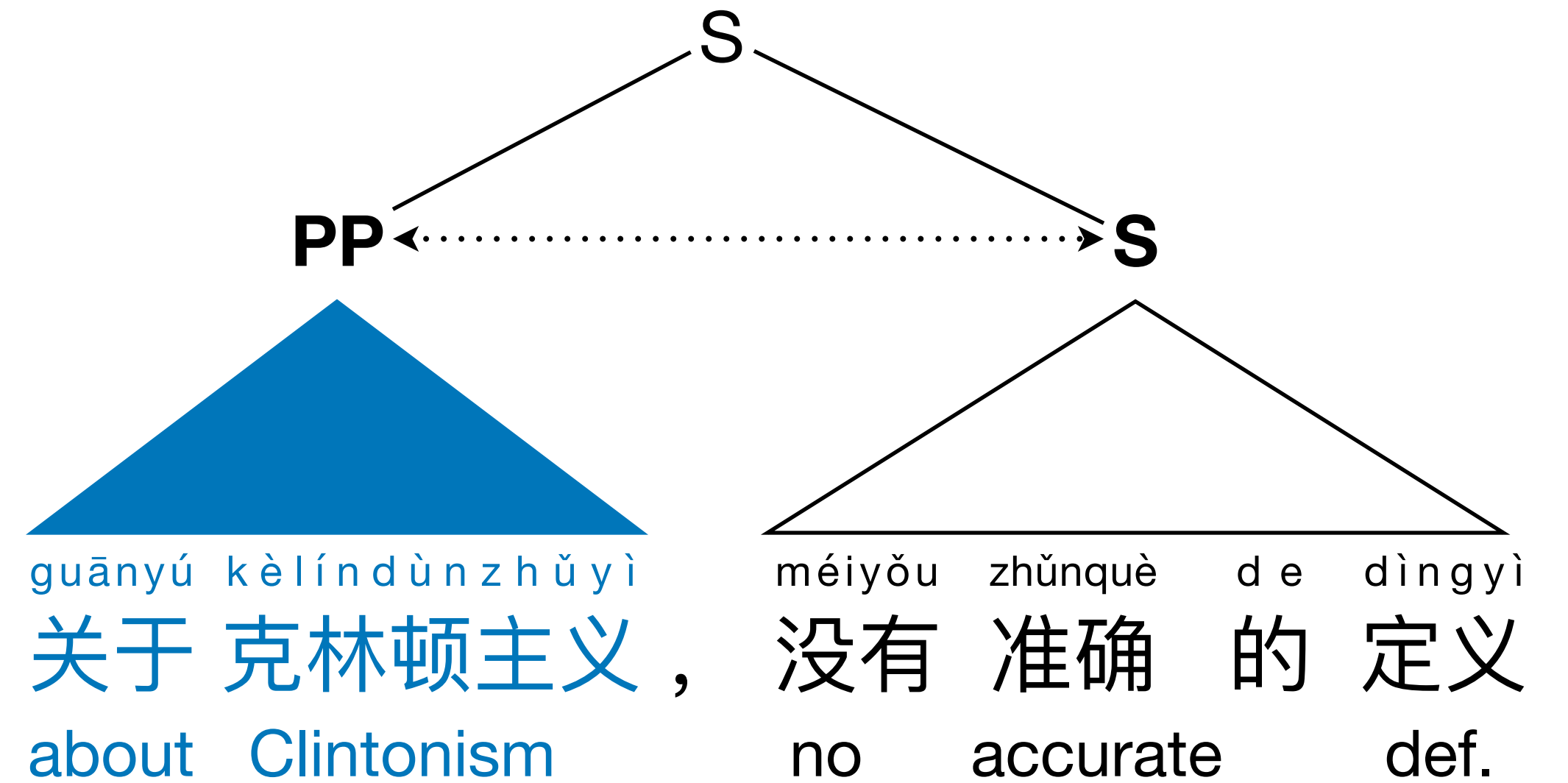


reference translation

“Xi Jinping was elected in Beijing in 2012”

optional reordering

(Chinese) PP S => (English) PP S or S PP



“There is no accurate definition of Clintonism.”

ideal => About Clintonism, there is no accurate definition.

Detecting and Fixing Mistakes

- idea: use a slower policy to verify the current policy's output along the way

	1	2	3	4	5	6	7
	<i>Bùshí</i>	<i>zǒngtǒng</i>	<i>zài</i>	Mòsīkē	<i>cānjiā</i>	<i>zhòngyào</i>	<i>fēnghuì</i>
	布什	总统	在	莫斯科	参加	重要	峰会
	Bush	president	in	Moscow	attend	important	summit
wait-2 decode			president	bush	met		
wait-3 check				president	bush	met	
				0.9	0.8	0.0	
revision						I mean <i>attended</i>	
wait-3 continue						an	important summit in moscow

The point of this talk is to “抛砖引玉”, i.e.,
to stimulate interests in this long-standing problem.



AI同传 ^热 百度wifi翻译机 人工翻译 | 插件下载 APP下载

检测到中文 ▾



英语 ▾

翻 译

人工翻译

抛砖引玉



Throw away a brick in order to get a gem

抛砖引玉 [pāo zhuān yǐn yù] throw away a brick in order to





非常感谢您 来听我的演讲

Thank you very much for listening to my speech

非常感谢您 来听我的演讲

Thank you very much for listening to my speech

Code (will be) available at <https://nlp.baidu.com/paddlenlp>
using <https://github.com/PaddlePaddle> framework

(it supports both static & dynamic graphs)

(the code for robust decoding with ASR noise is already available)

Two Posters after the coffee break (10:30), Session 4A (#4 & #6)

Short Talk tomorrow, Session 8D (17:13, CAVANIGLIA)

