

The Effect of Place in MT

Marc'Aurelio Ranzato

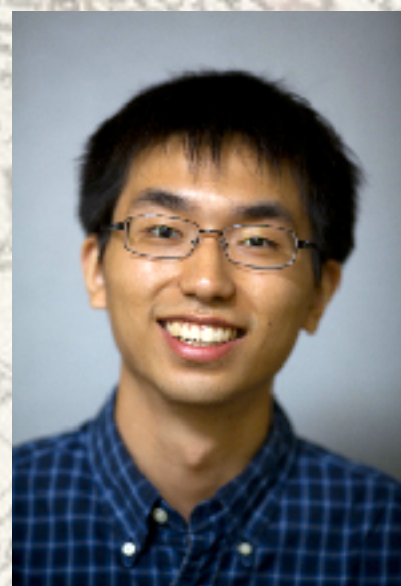
Facebook AI Research - NYC

ranzato@fb.com

joint work with:



Peng-Jen Chen



Jiajun Shen



Matt Le



Junxian He



Myle Ott



Jiatao Gu



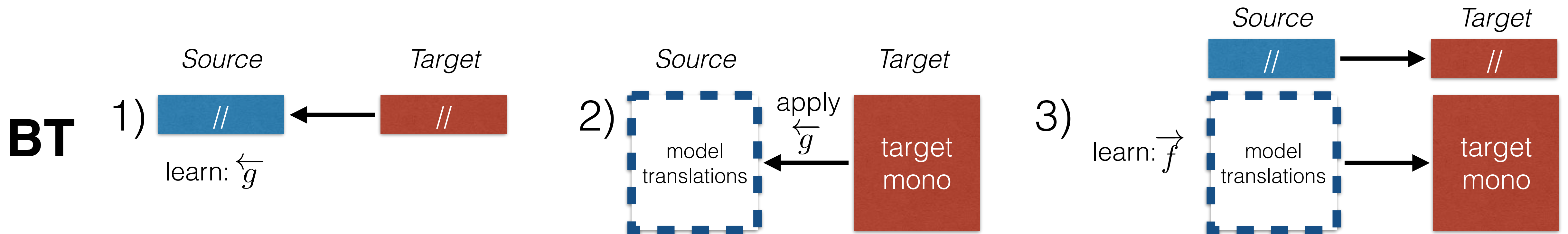
Michael Auli

WeCNLP 6 September 2019

Machine Translation - AD 2019

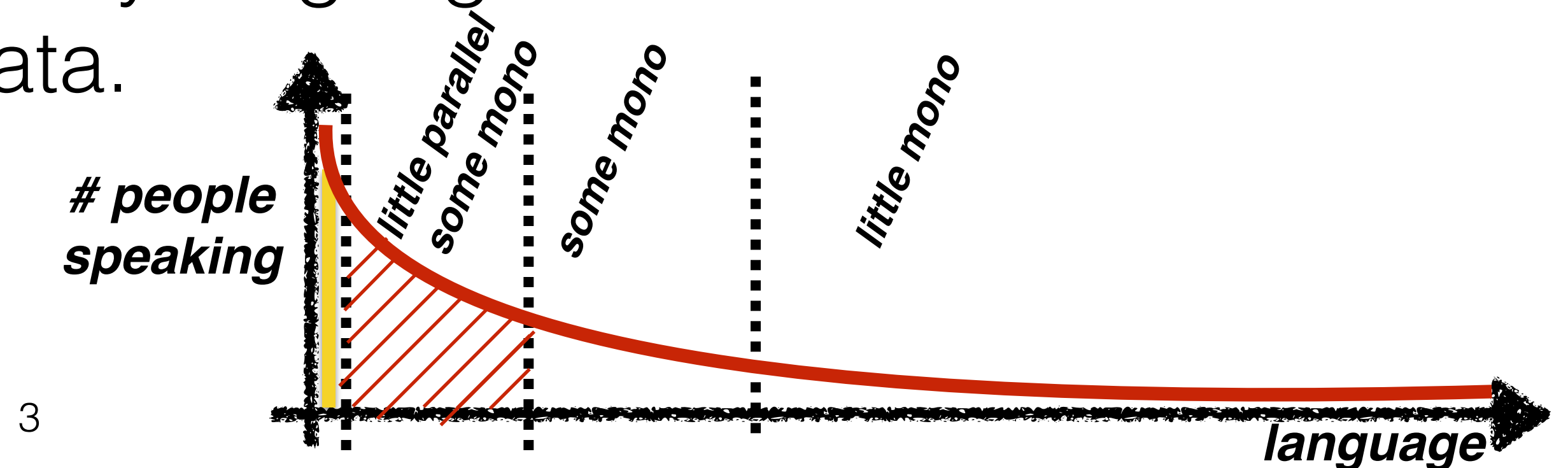
- Human level performance on some languages in some domains under some conditions...
- Key factors of success:
 - big data (parallel & monolingual), big models & big compute
 - better architectures like transformer
 - **back-translation** (BT)

Attention is all you need. Vaswani et al. NIPS 2017
Improving NMT with monolingual data. Sennrich et al. ACL 2015



Machine Translation - AD 2019

- There are 6000+ languages in the world. Very few enjoy large parallel resources.
- **MT for low resource languages.** In addition to back-translation, also:
 - initialization
 - multi-lingual training
 - noisy parallel data from ParaCrawl
- Problem: the tail is very long... for many languages these methods are not applicable because of lack of data.



Setting

- Data:
 - a small parallel dataset ($\sim 10K$).
 - some monolingual data on both target and source language ($\sim 1M$).
- Method:
 - Back-translation (BT).

The Power of BT

Simulating low-resource MT with a high resource language:
using *EuroParl* data with 20K parallel sentences and 100K monolingual target sentences.

<i>EuroParl Fr—>En</i>	
only parallel data	30.4 BLEU
parallel data + BT	33.8 BLEU

+3.4 BLEU!



A Worrisome Finding

BT sometimes yields very mild improvements.

Example #1

FB public posts $En \rightarrow My$

only parallel data	15.2 BLEU
parallel data + BT	15.3 BLEU

+0.1 BLEU!



Example #2

FLORES evaluation set $En \rightarrow Ne$

only parallel data	4.3 BLEU
parallel data + BT	6.8 BLEU

+2.5 BLEU!



Sports



- *football*
- *baseball*
- *basketball*



- *soccer*
- *cricket*
- *rowing*



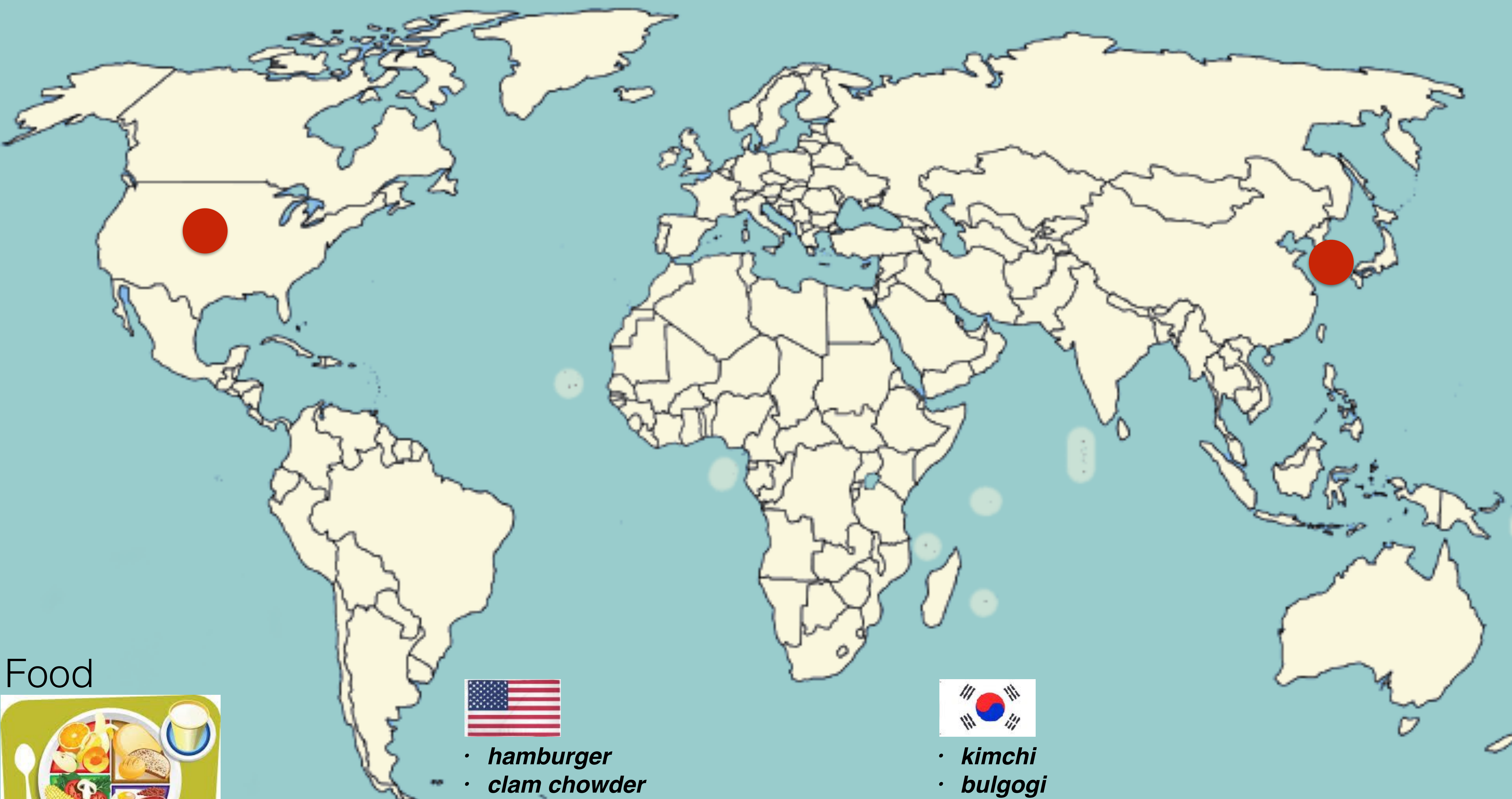
News



- *hurricane*
- *trade war*
- *elections*



- *lake preservation*
- *earthquake*
- *heritage*



Food



- ***hamburger***
- ***clam chowder***
- ***apple pie***

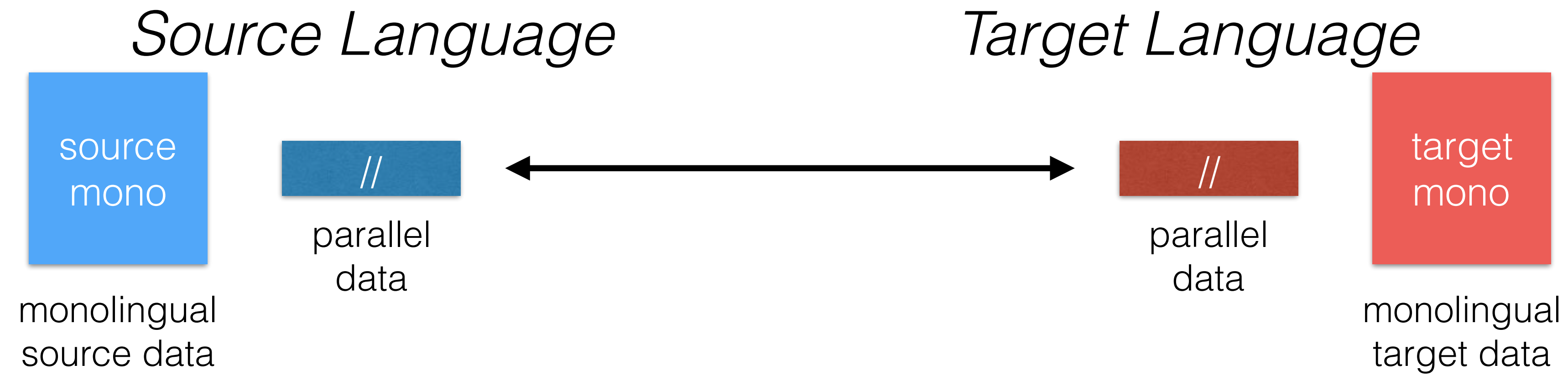


- ***kimchi***
- ***bulgogi***
- ***bibimbap***

The Place Effect

- Def.: Content produced in blogs, social networks, news outlets, etc. varies with the geographic location.
- The place effect is even more pronounced in low resource MT, where source & target geographic locations are typically farther apart and cultures have more distinct traits.
- The place effect makes the MT problem even harder, because of **source/target domain mismatch**.

Source / Target Domain Mismatch



Source / Target Domain Mismatch

Source Language

Target Language

*Source
Domain*



source //

source
mono

translationese target

*Target
Domain*



translationese source

target //

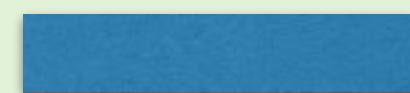
target
mono

STDM Cripples BT

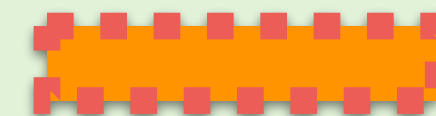
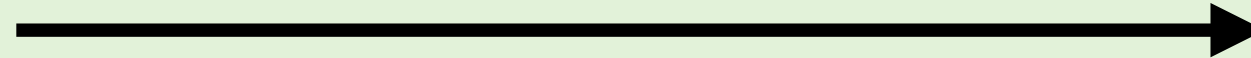
Source Language

Target Language

Source Domain

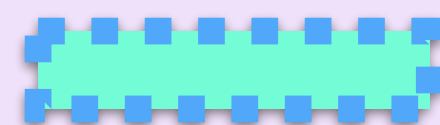


x

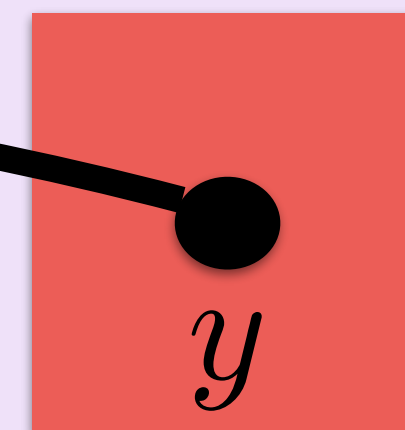


- **backward model trained with mixed-domain data**
- **back-translated data is out-of-domain**

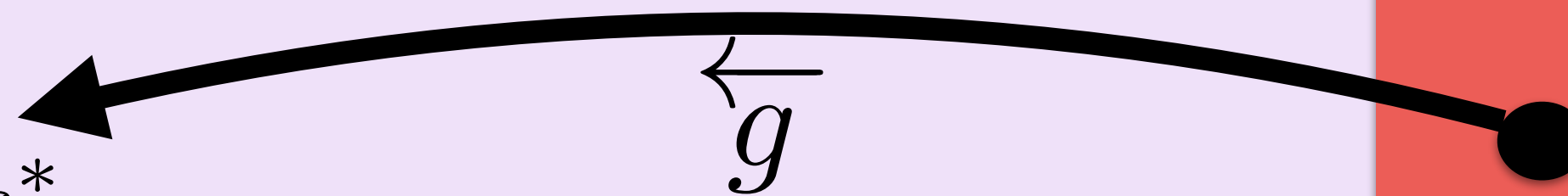
Target Domain



x^*



y



\overleftarrow{g}

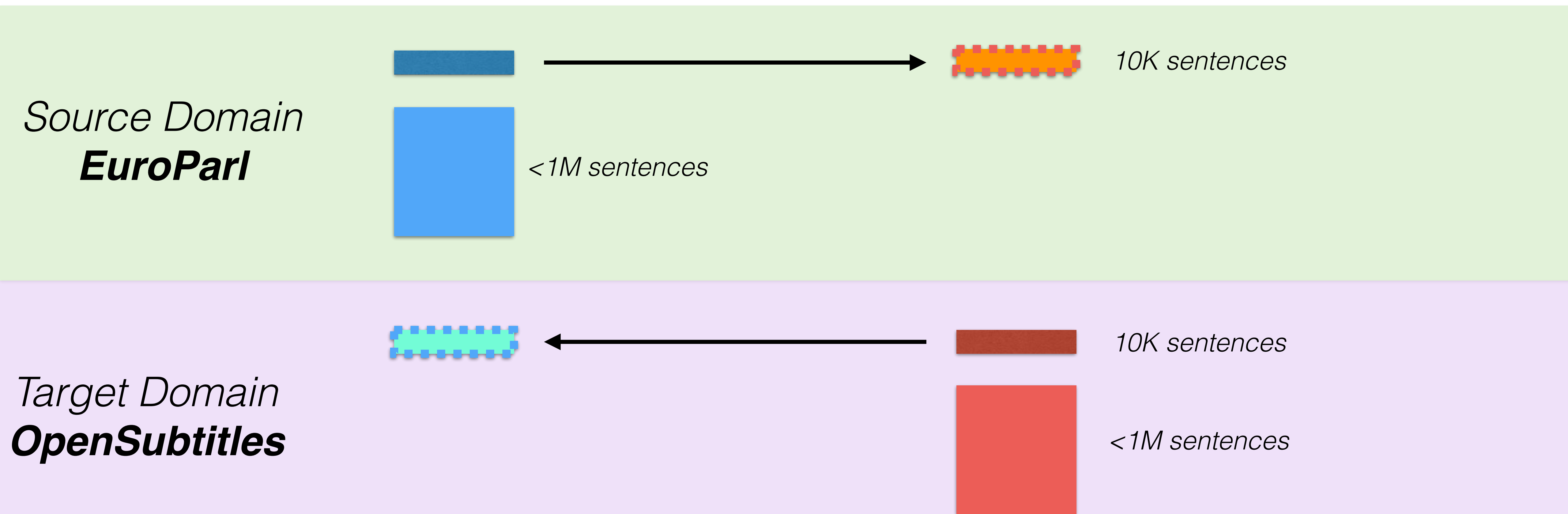
Questions

- Is it true that BT is less effective when there is STDN?
- What other baselines shall we consider when there is STDN?
- Is out-of-domain data worth using when there is STDN?
- What are general best practices when there is STDN?
- How to study STDN in a controlled setting?

Controlled Setting

Source Language: Fr

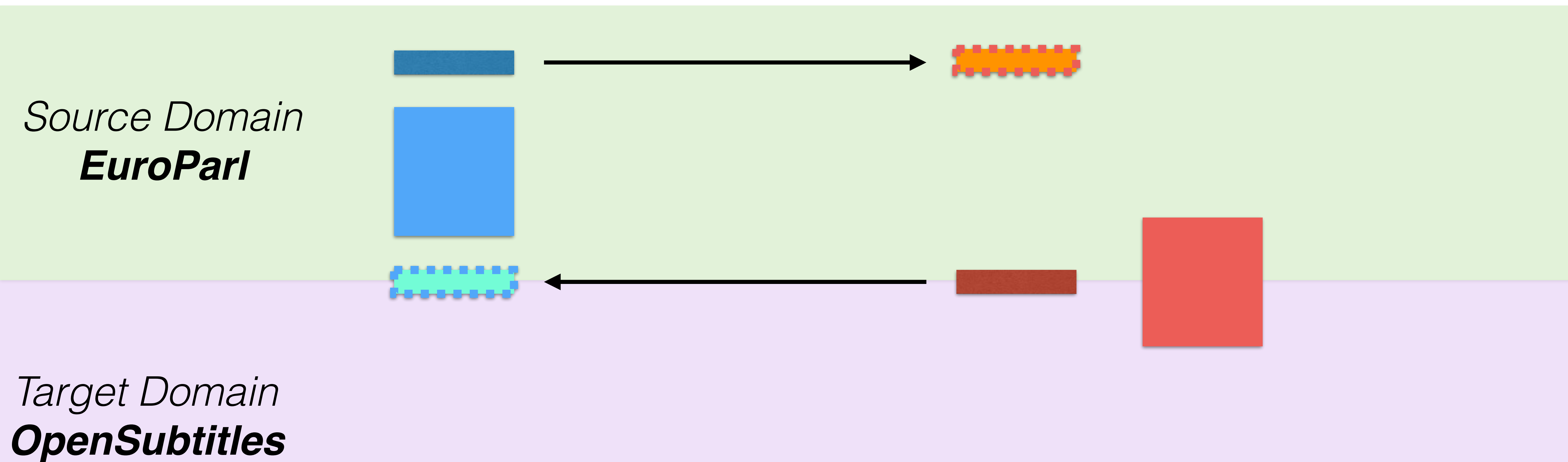
Target Language: En



Controlled Setting

Source Language: Fr

Target Language: En

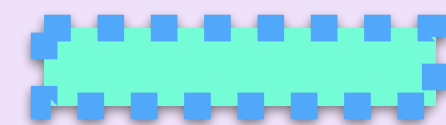
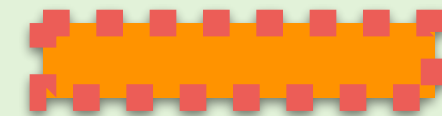
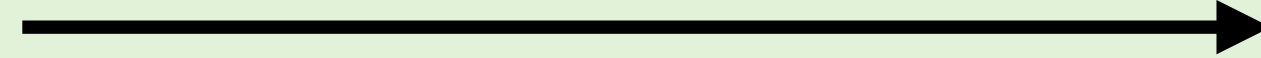
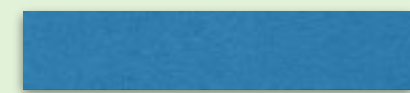


Controlled Setting

Source Language: Fr

Target Language: En

*Source Domain
EuroParl*



Target Domain

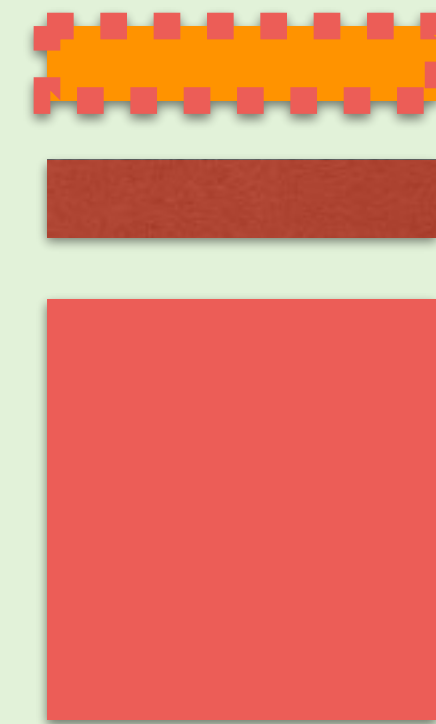
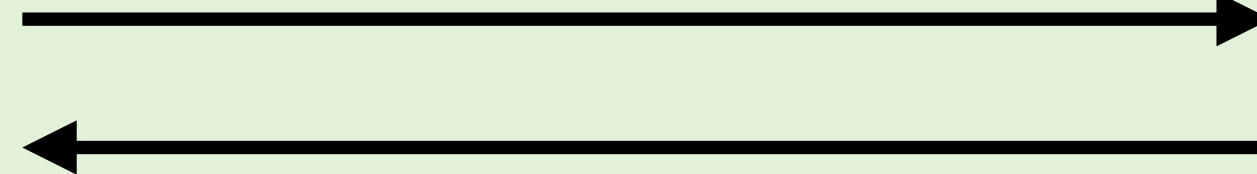
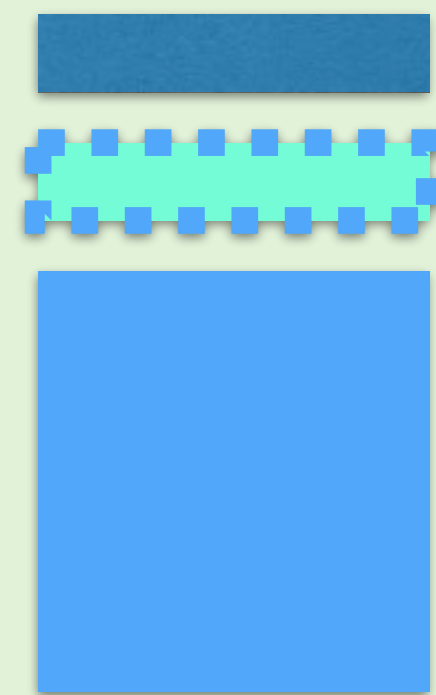
$\alpha \text{ EuroParl} + (1 - \alpha) \text{ OpenSubtitles}$
 $\alpha = 0$

Controlled Setting

Source Language: Fr

Target Language: En

*Source Domain
EuroParl*



Target Domain

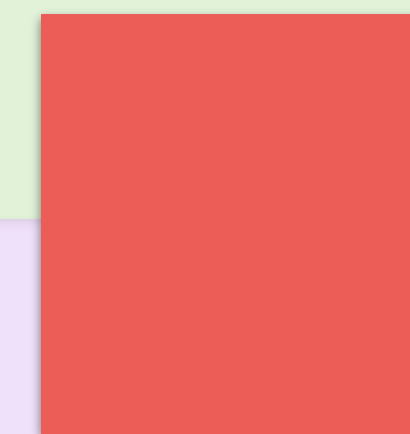
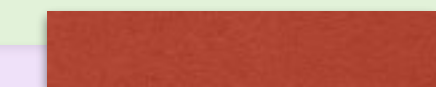
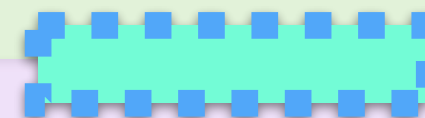
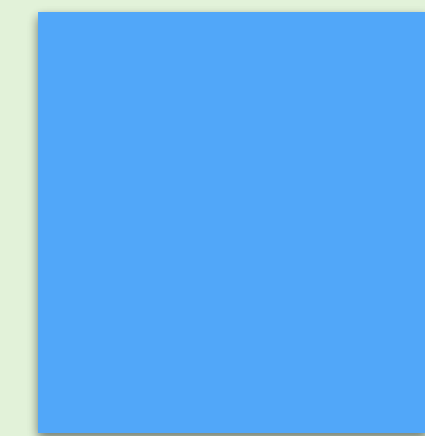
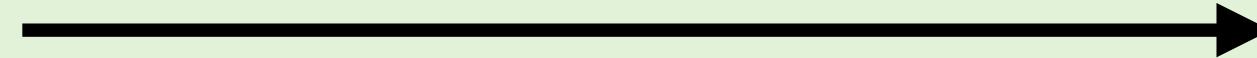
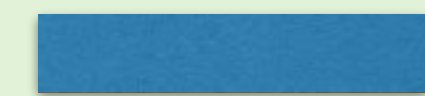
$$\alpha \text{ EuroParl} + (1 - \alpha) \text{ OpenSubtitles}$$
$$\alpha = 1$$

Controlled Setting

Source Language: Fr

Target Language: En

*Source Domain
EuroParl*

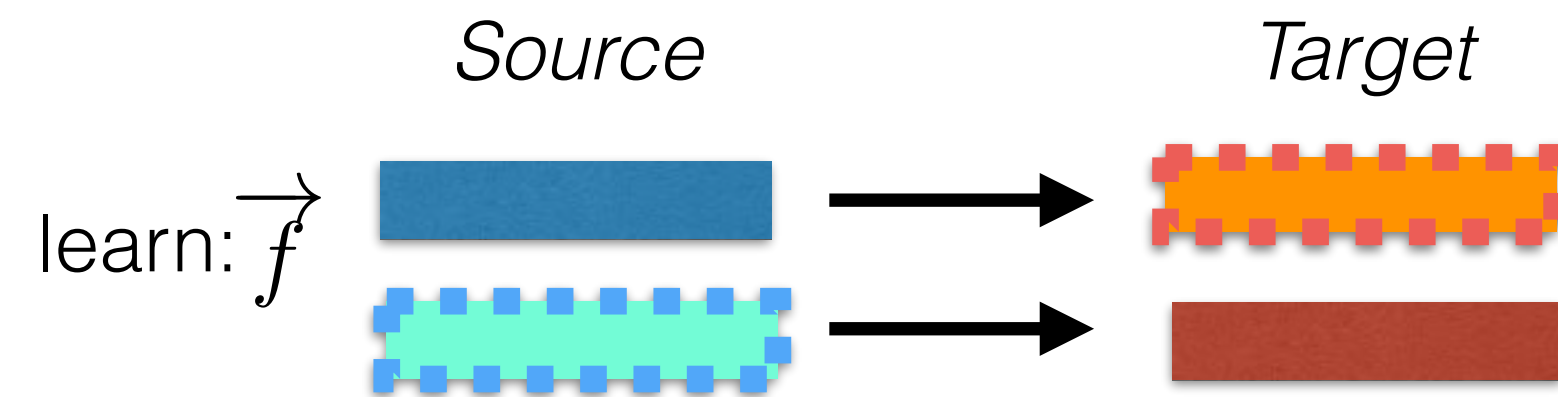


Target Domain

$\alpha \text{ EuroParl} + (1 - \alpha) \text{ OpenSubtitles}$
intermediate value of α

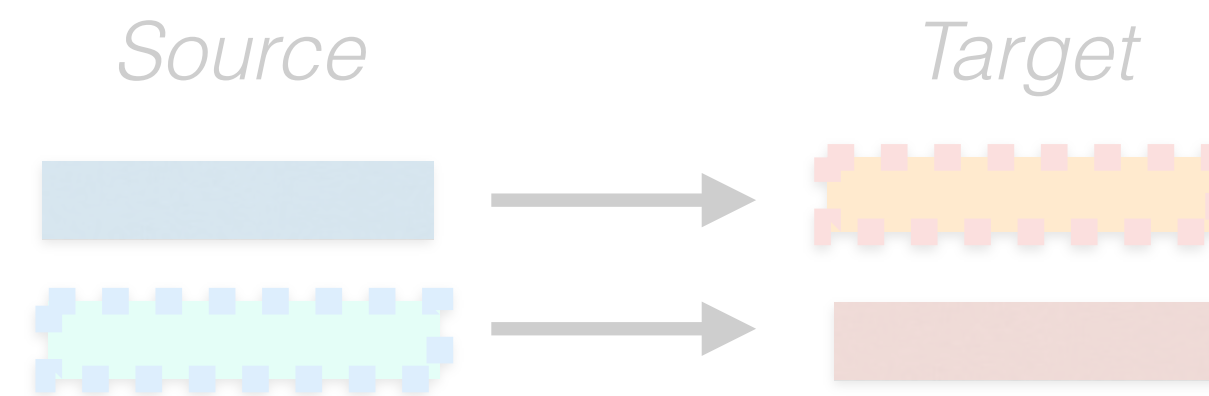
Baseline Approaches

- Bitext only:

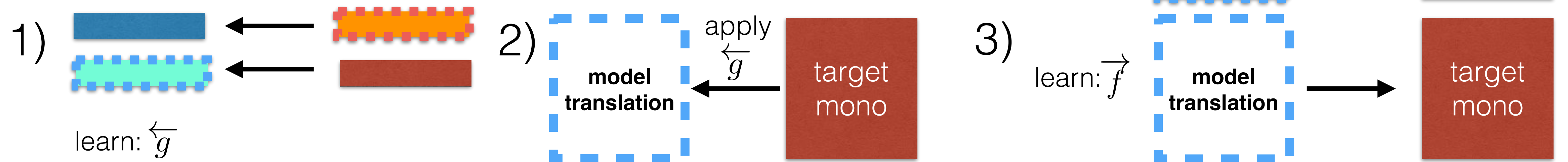


Baseline Approaches

- Bitext only:



- Back-Translation:

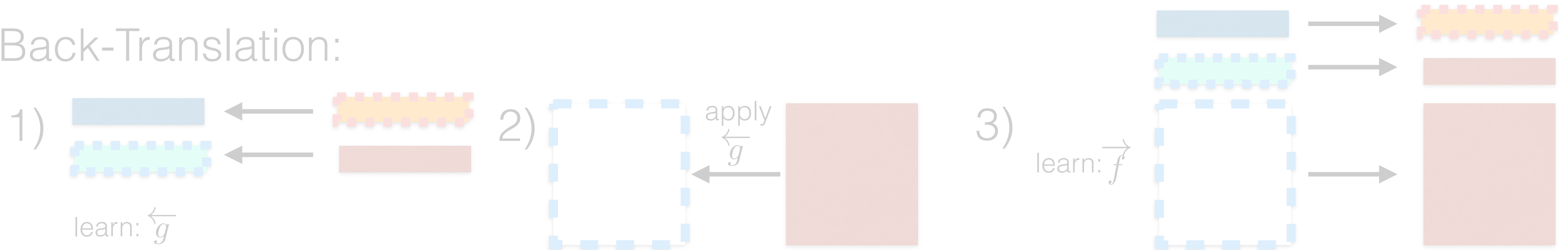


Baseline Approaches

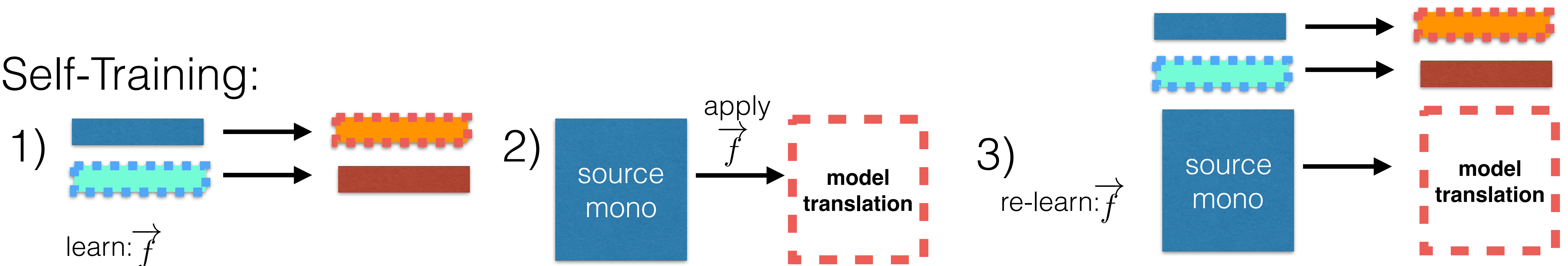
- Bitext only:

Unsupervised word sense disambiguation ... Yarowski ACL 1995
Using monolingual source-language data to improve MT performance. Ueffering IWSLT 2006
Exploiting source-side monolingual data in NMT. Zhang et al. EMNLP 2016

- Back-Translation:

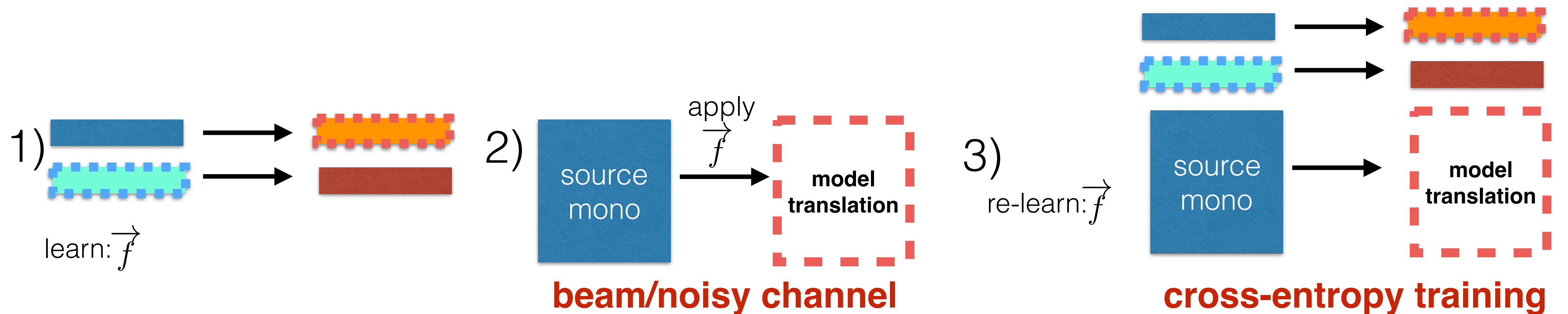


- Self-Training:



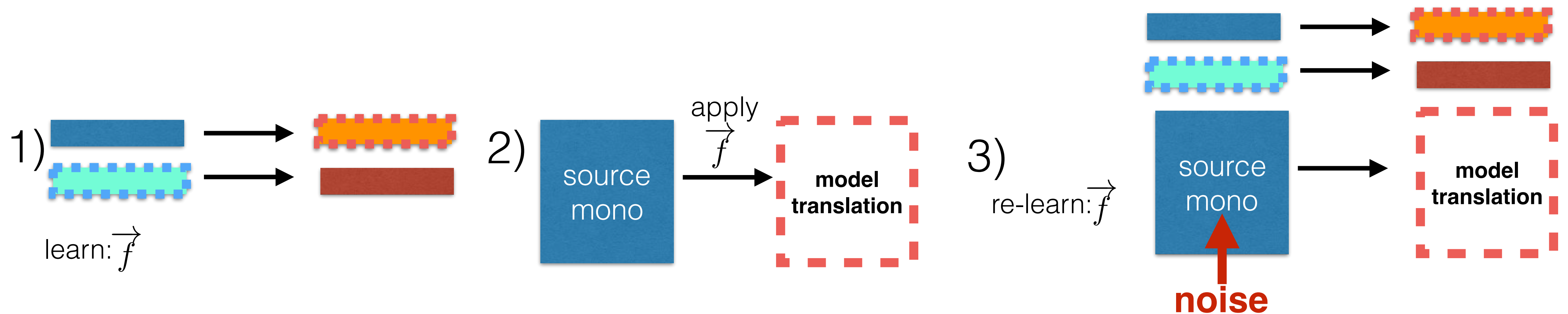
Why Self-Training May Work

- The model learns the decoding process.



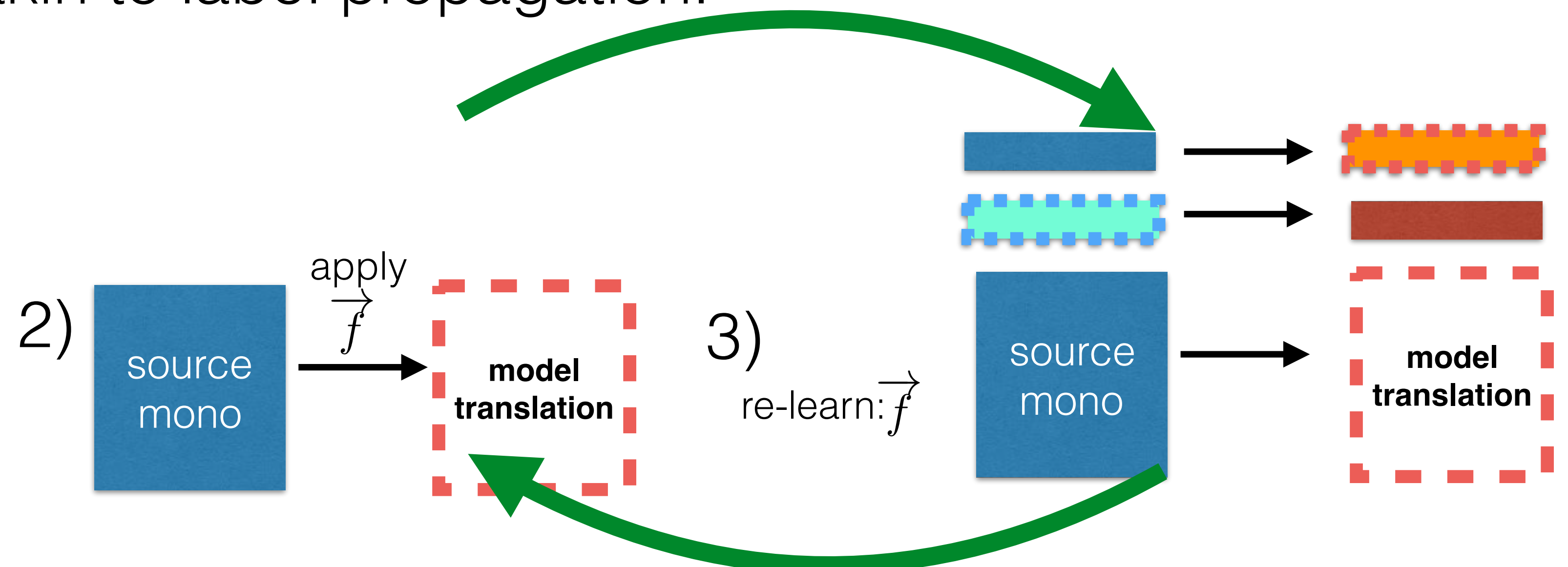
Why Self-Training May Work







- The model learns the decoding process.
- Noise helps mapping similar inputs to the same target.



Why Self-Training May Work

- The model learns the decoding process.
- Noise helps mapping similar inputs to the same target.
- Iterative ST is akin to label propagation.



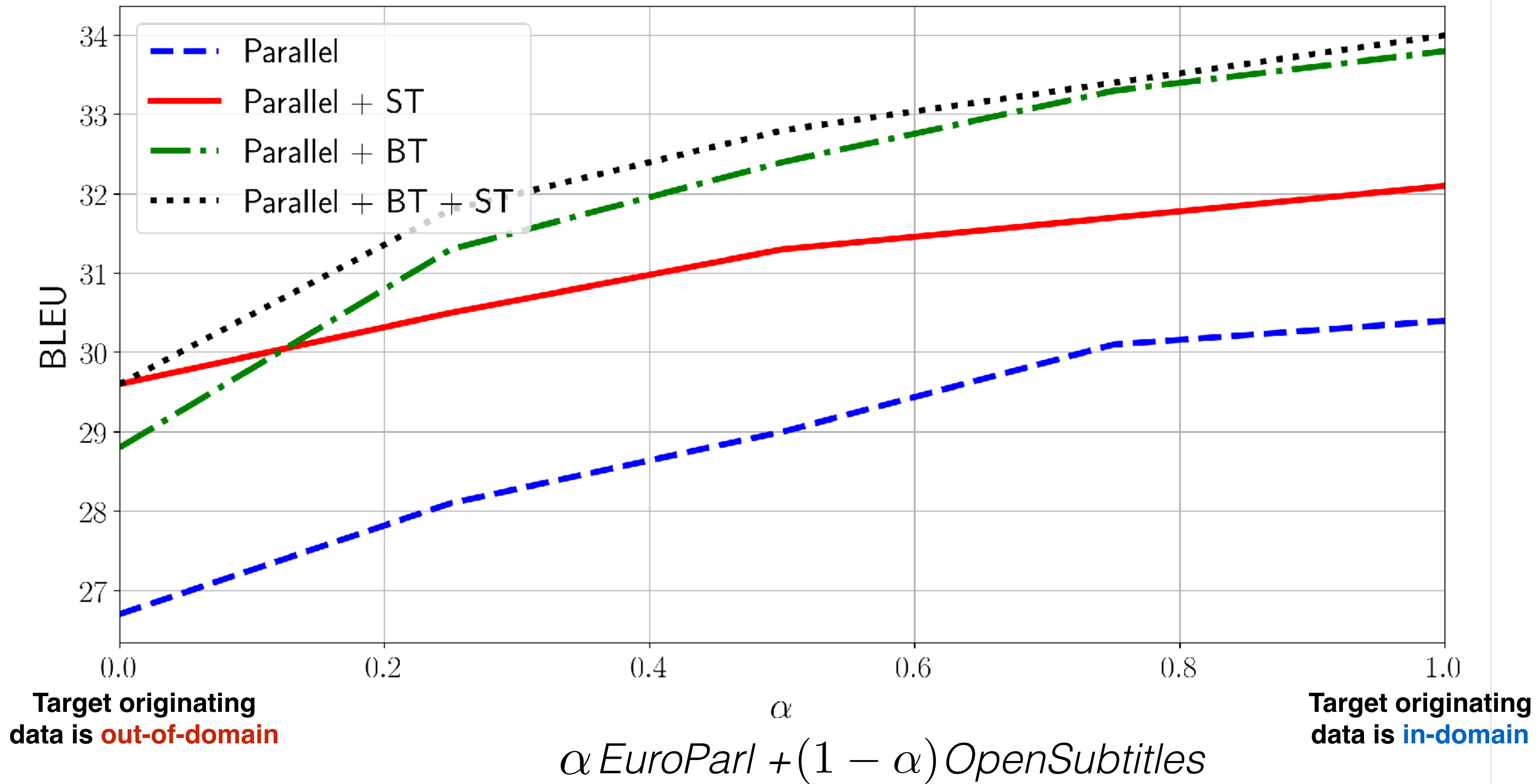
	<i>target mono data</i>	<i>source mono data</i>	<i>in-domain mono data</i>
• Back-Translation:			
• Self-Training:			

Q.: Is it better to have clean targets but out-of-domain data, or noisy targets but in-domain data?

Q.: What's the effect of amount of parallel/monolingual data?

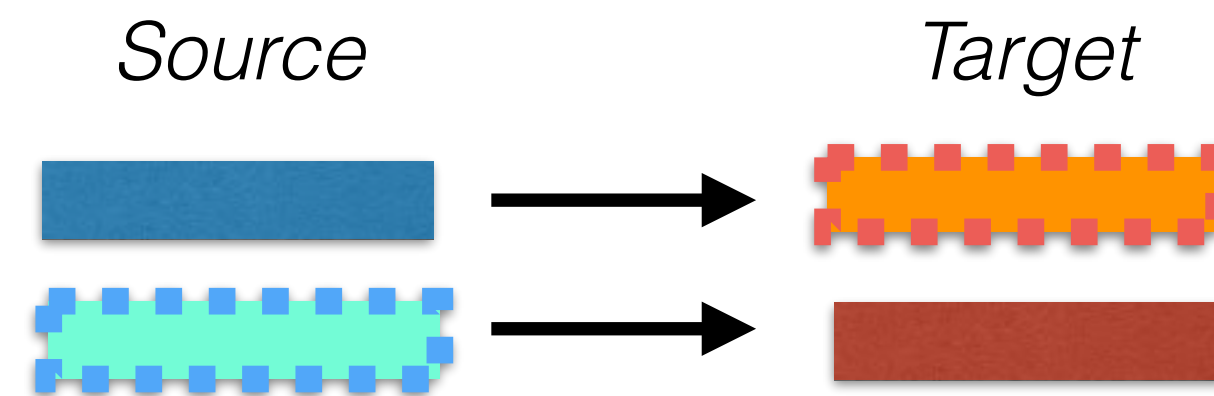
Q.: What's the effect of the quality of the model forward model when training with ST?

Varying Domain of Target Originating Data

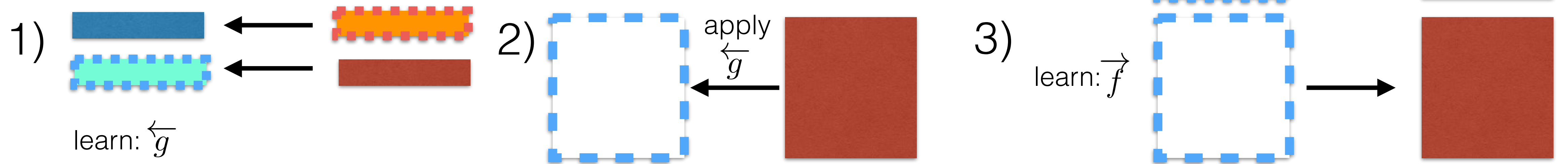


Baseline Approaches

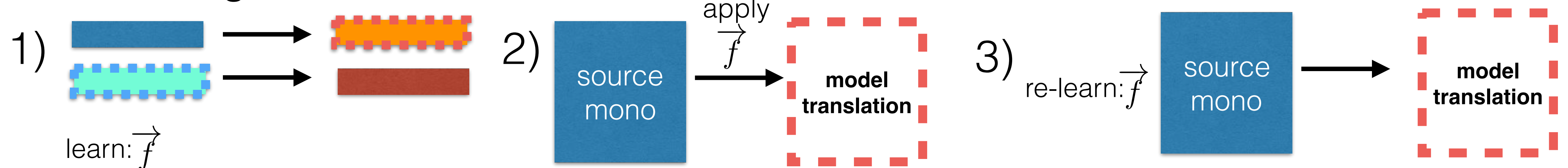
- Bitext only:



- Back-Translation:

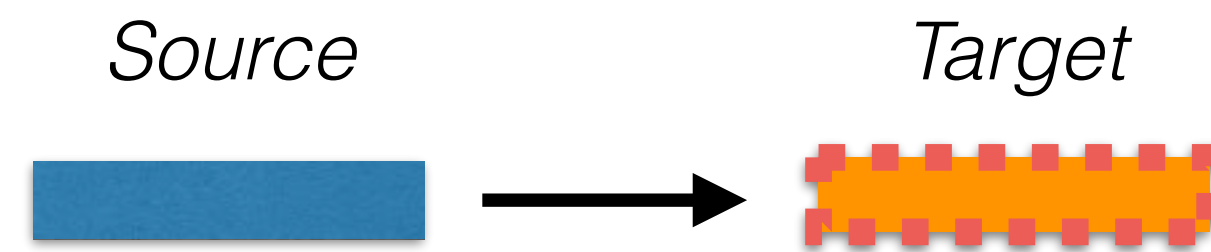


- Self-Training:

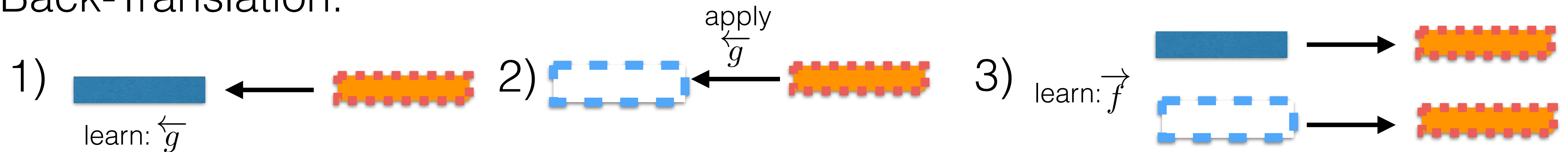


Baseline Approaches: Only In-Domain Data

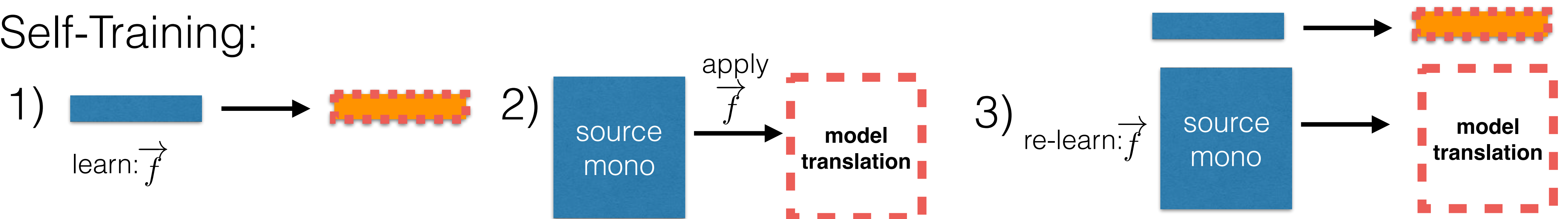
- Bitext only:



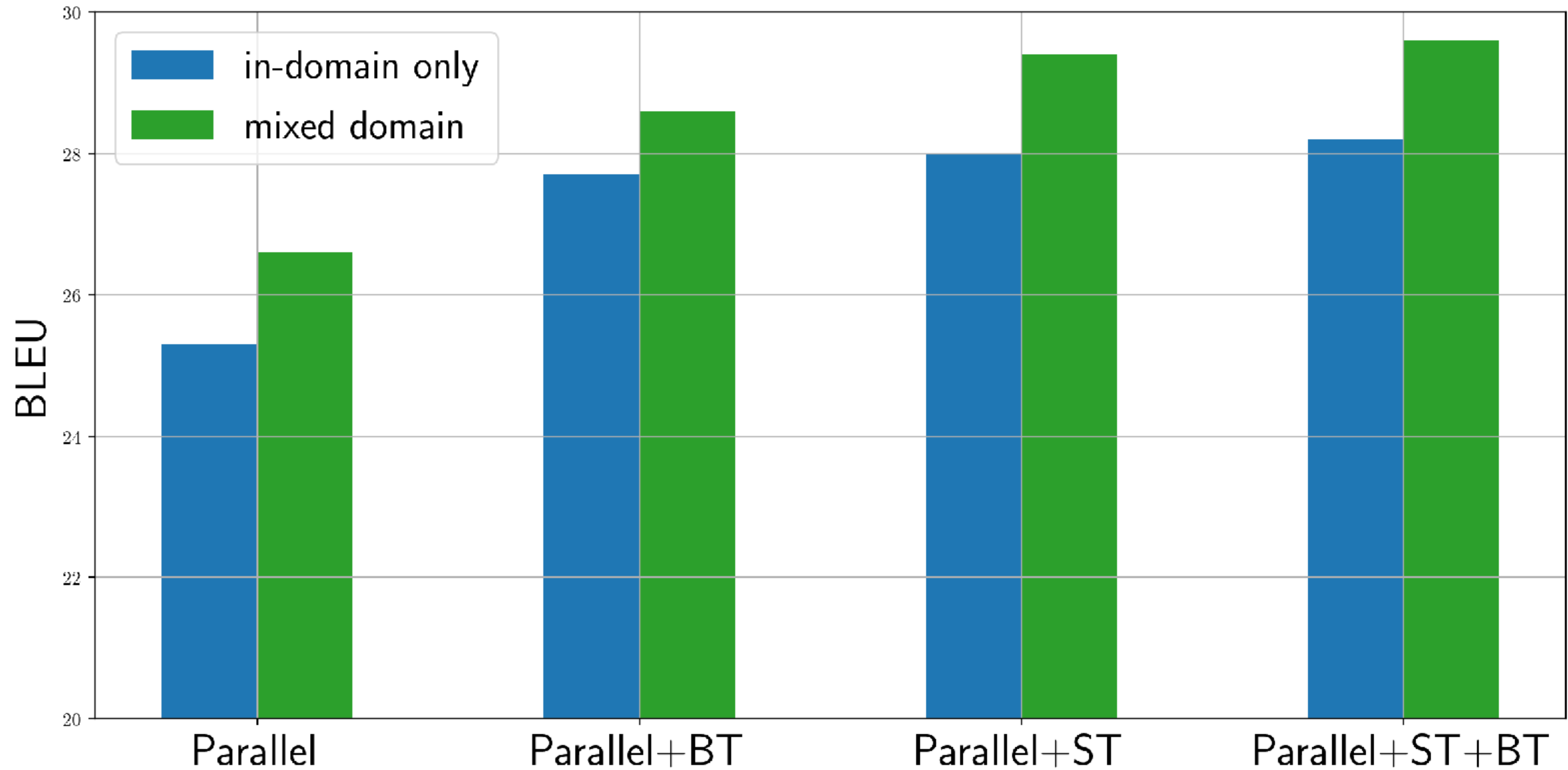
- Back-Translation:



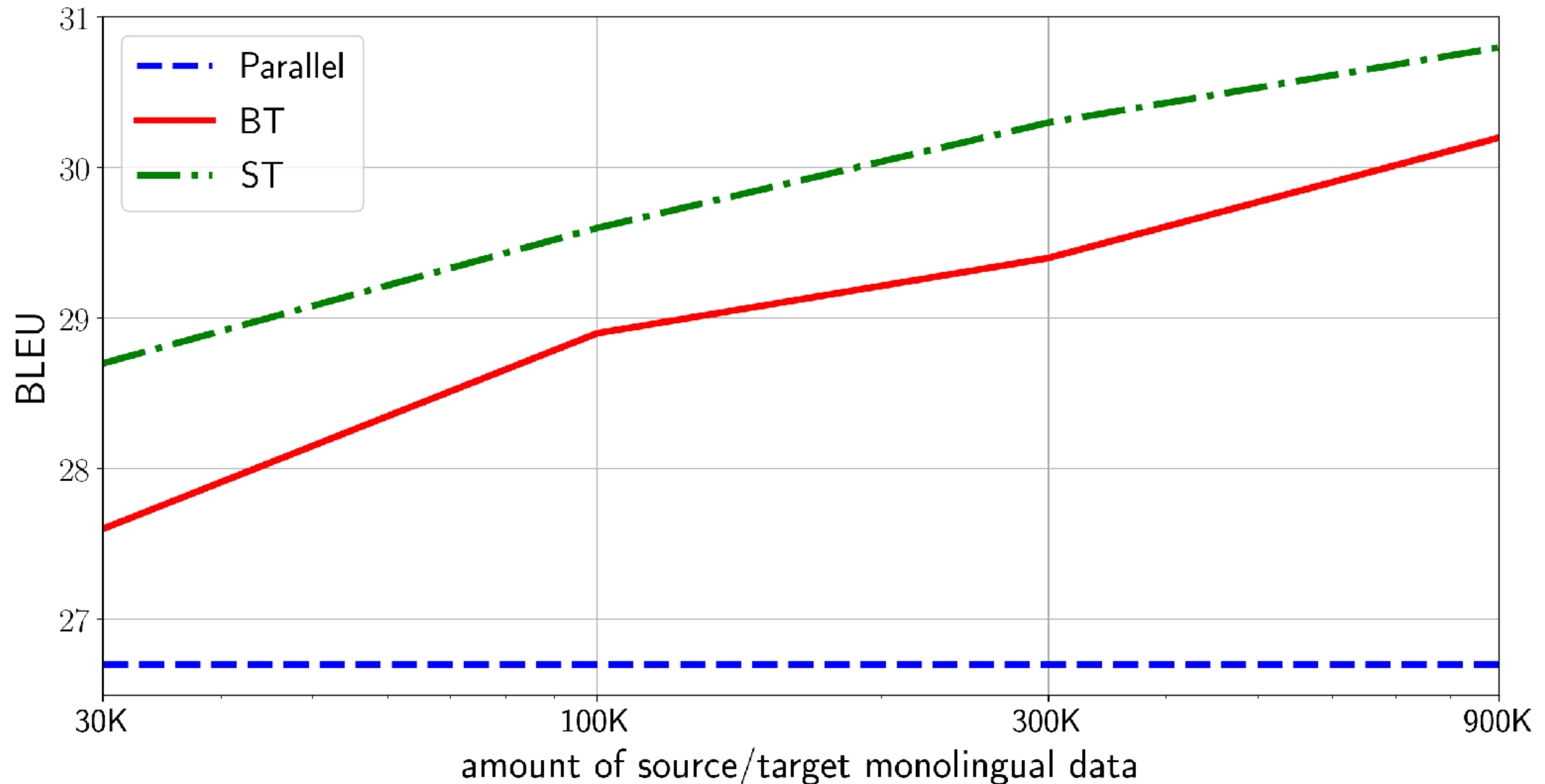
- Self-Training:



In-Domain Only VS. Mixed Domain



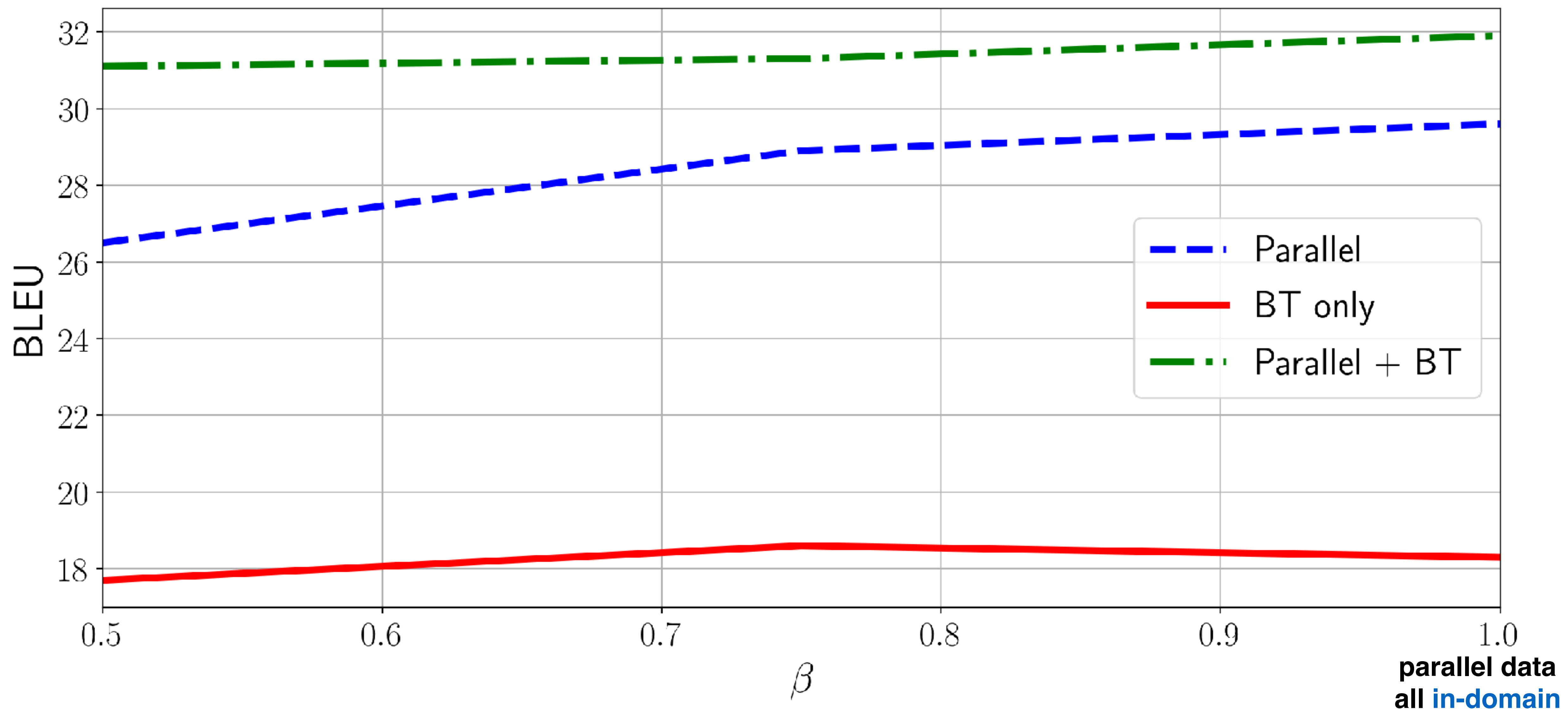
Varying Amount of Monolingual Data ($\alpha = 0$)



How To Construct Parallel Datasets

- Suppose that:
 - we are interested in forward translation only.
 - we can only translate 20K sentences in total.
 - target data is out-of-domain ($\alpha = 0$).
 - we have 100K target monolingual sentences for BT.
- Is it better to translate 20K sentences all originating from the source? or have some originating from the target as well?

$$\text{Parallel Dataset} = \underset{\text{source originating data}}{\beta \text{ EuroParl}} + (1 - \beta) \underset{\text{target originating data}}{\text{OpenSubtitles}}$$



Parallel Dataset = β EuroParl + $(1 - \beta)$ OpenSubtitles
Target Mono Dataset = OpenSubtitles

A Real Case-Study: English-Burmese

ဗဟိုစာမျက်နှာ

ဆွေးနွေးချက်

ဖတ်ရန်


ရင်းမြစ်ကို ကြည့်ရန်

ရာဇဝင်ကြည့်ရန်

ဝီကီပီးဒီးယား တွင် ရှာဖွေရန်

🔍



ဗဟိုစာမျက်နှာ



ဝီကီပီးဒီးယားမှ ကြိုဆိုပါသည်။

မည်သူမဆို ကြည့်ရှုပြင်ဆင်နိုင်သော အခမဲ့လွတ်လပ်စွယ်စုံကျမ်း ဖြစ်ပါသည်။


အကြောင်းအရာပေါင်း ၄၄၈၁၄ ခုကို မြန်မာဘာသာဖြင့် ဖတ်ရှုနိုင်ပါသည်။



- အနုပညာ
- အတ္ထုပ္ပတ္တိ
- ပထဝီဝင်

- သမိုင်း
- သင်္ချာ
- သိပ္ပံ

- လူမှုရေး
- နည်းပညာ
- မုခ်ဦးအားလုံး



အထူးအကြောင်းအရာ

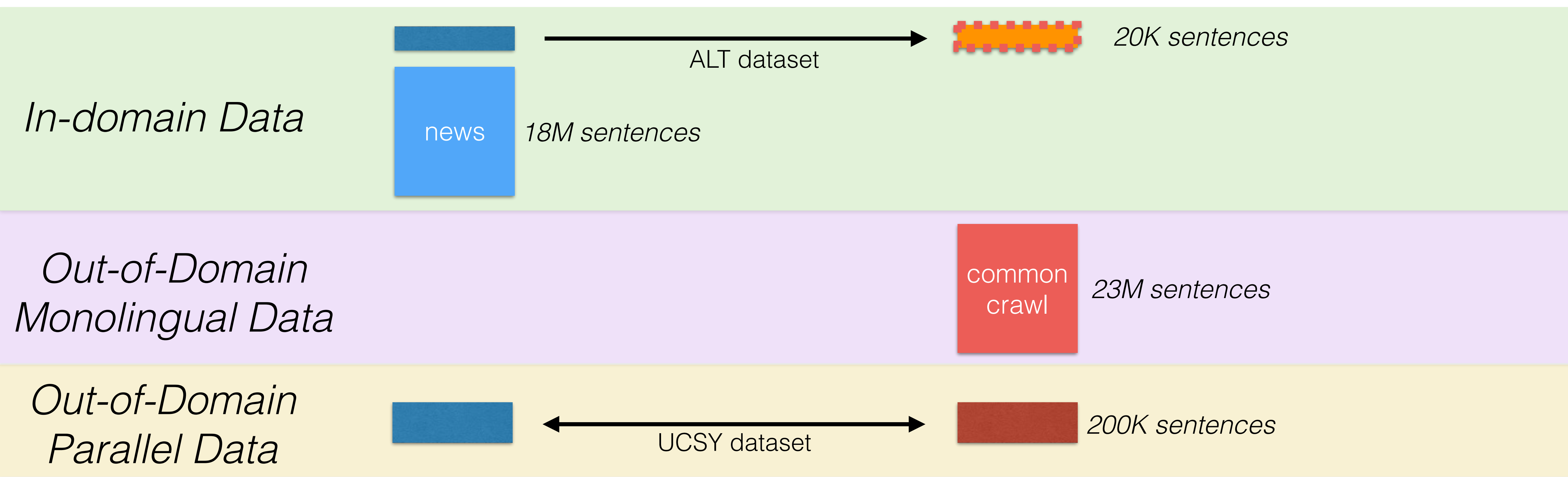
သိန္နီနယ် သည် ယခင် ရှမ်းပဒေသရာဇ်ပြည်နယ်များတွင် ပါဝင်ခဲ့သော ပြည်နယ်တစ်ခု ဖြစ်သည်။ သိန္နီနယ်ကို ခရစ် ၁၈၈၈ ခုနှစ် အင်္ဂလိပ်တို့ ဝင်ရောက်သိမ်းပိုက်ပြီးနောက်မှ မြောက်သိန္နီနယ် (သိန္နီနယ်)နှင့် တောင်သိန္နီနယ်(မိုင်းရယ်နယ်)ဟု ခွဲခြားအုပ်ချုပ်ခဲ့သည်။ ရှေးအခါသမယက သိန္နီနယ်ကြီးသည် အစိတ်စိတ်ကွဲပြားခြင်းမရှိဘဲ ရှမ်းပြည်နယ်တဝှမ်းလုံးတွင် အကျယ်ပြန့်ဆုံး အာဏာအလွှမ်းမိုးဆုံးသော နယ်ကြီးဖြစ်ခဲ့သည်။ သို့သော် မြန်မာဘုရင်များ ဝင်ရောက်တိုက်ခိုက် သိမ်းပိုက်ပြီးသည့်နောက် အုပ်ချုပ်ရေးဝါဒအရ ရာထူးလူသည့်နယ်ရှင် ဇော်ဘွားတို့ကြောင့် သိန္နီနယ်ကြီးသည် ငါးနယ်အထိ အစိတ်စိတ် ကွဲပြားခဲ့လေသည်။ ထိုအတွင်း ဇော်ဘွားအချင်းချင်း စိတ်ဝမ်းကွဲကာ တစ်ဦးနှင့်တစ်ဦးတိုက်ခိုက်၍ ဆိုင်ရာ နယ်ပယ်များကို အုပ်ချုပ်ကြသည်။ နောက်ဆုံးခရစ် ၁၈၈၈ ခုနှစ်၊ အင်္ဂလိပ်တို့ဝင်ရောက်လာမှ အထက်ပါ အတိုင်းနှစ်နယ်ခွဲ၍ အုပ်ချုပ်ခဲ့သည်။ နှစ်နယ်ခွဲ၍ အုပ်ချုပ်စက ခွန်ဆိုင်တုံဟမ်းအား မြောက်သိန္နီနယ်အတွက် ဇော်ဘွားအဖြစ်လည်းကောင်း၊ ဆိုင်နော်ဖ၏သား နော်မိုင်းအား တောင်သိန္နီနယ် ဇော်ဘွားအဖြစ်လည်းကောင်း၊ ခန့်အပ်ခဲ့လေသည်။ ၁၉၂၅ ခုနှစ်တွင် မြောက်သိန္နီနယ်ကို စပ်ဟုံဖက ဇော်ဘွားအဖြစ် ဆောင်ရွက်ခဲ့လေသည်။

A Real Case-Study: English-Burmese

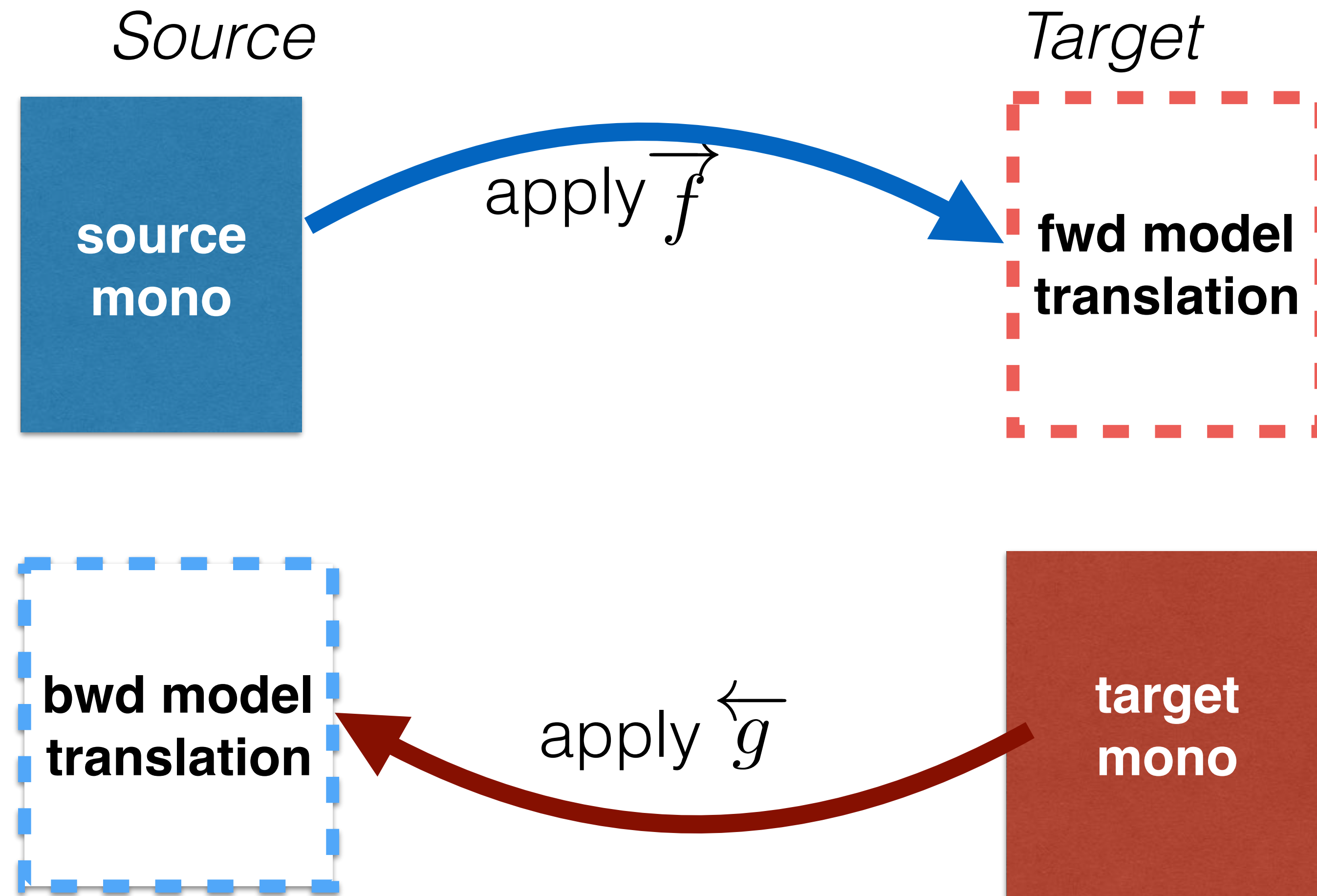
Workshop on Asian Translation @ EMNLP 2019: En-My

Source Language: En

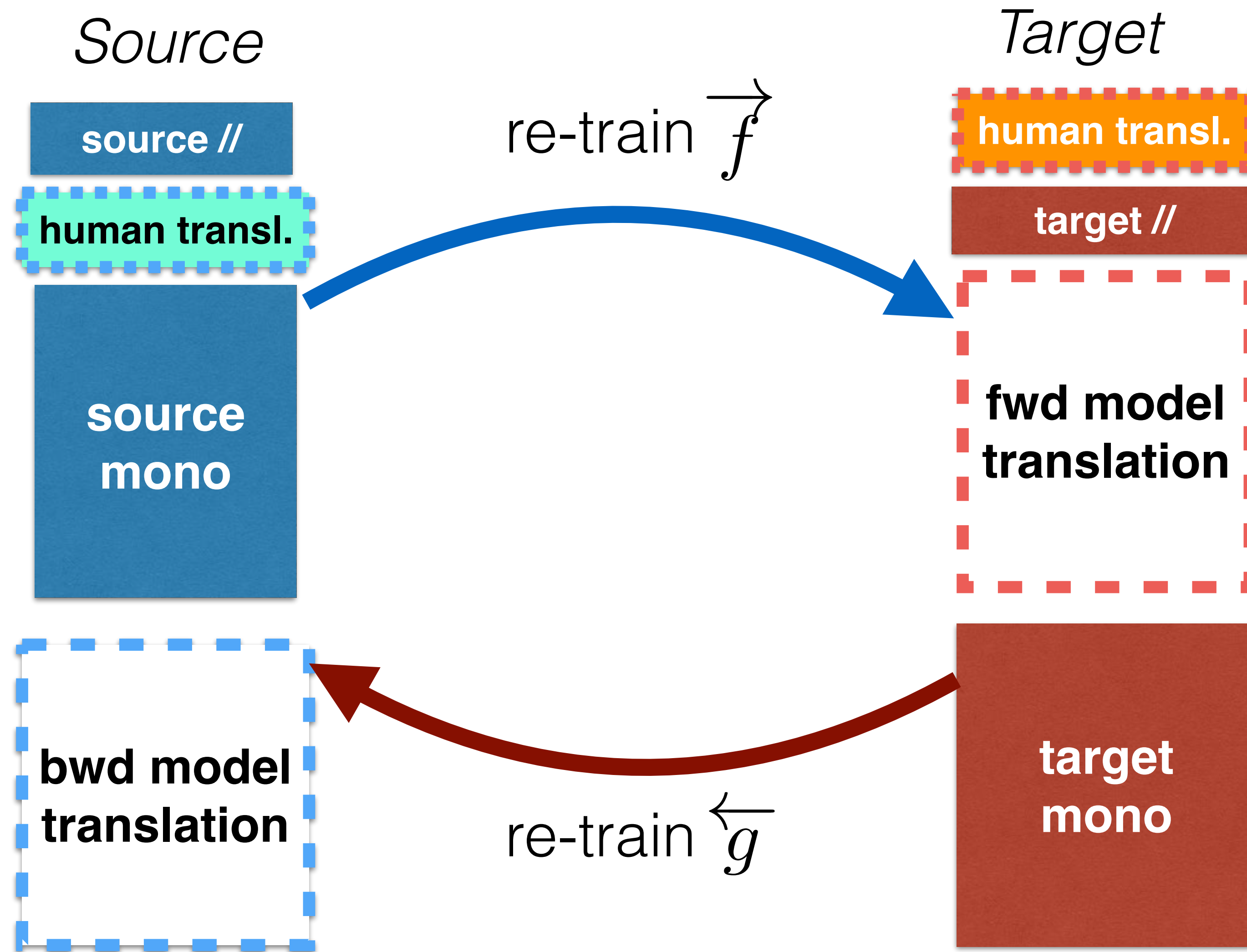
Target Language: My



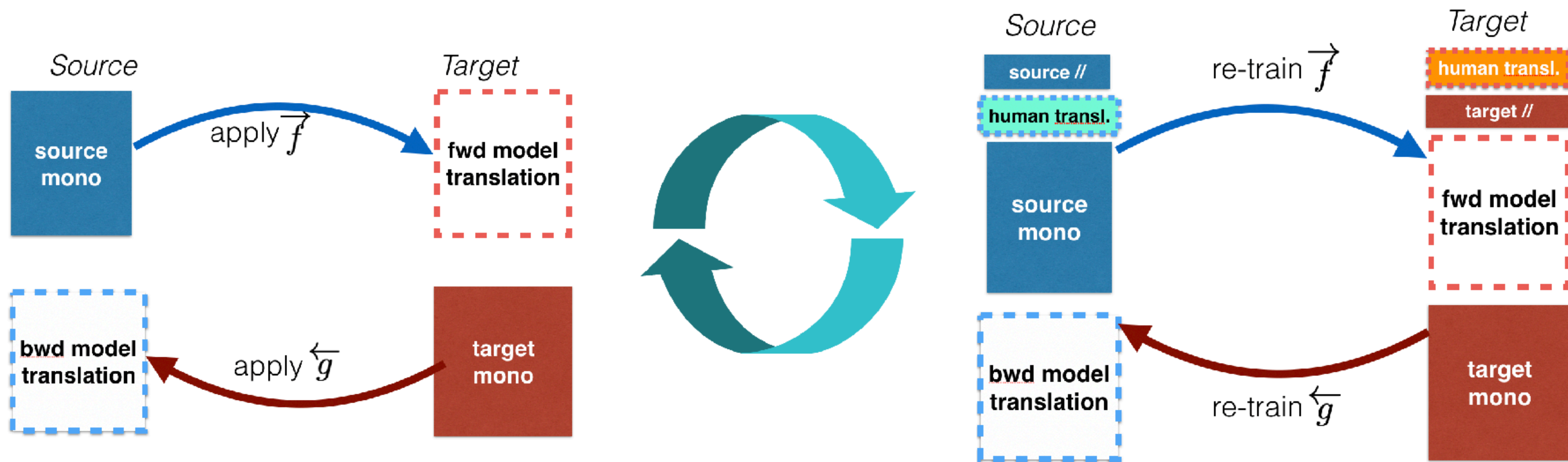
Iterative BT+ST



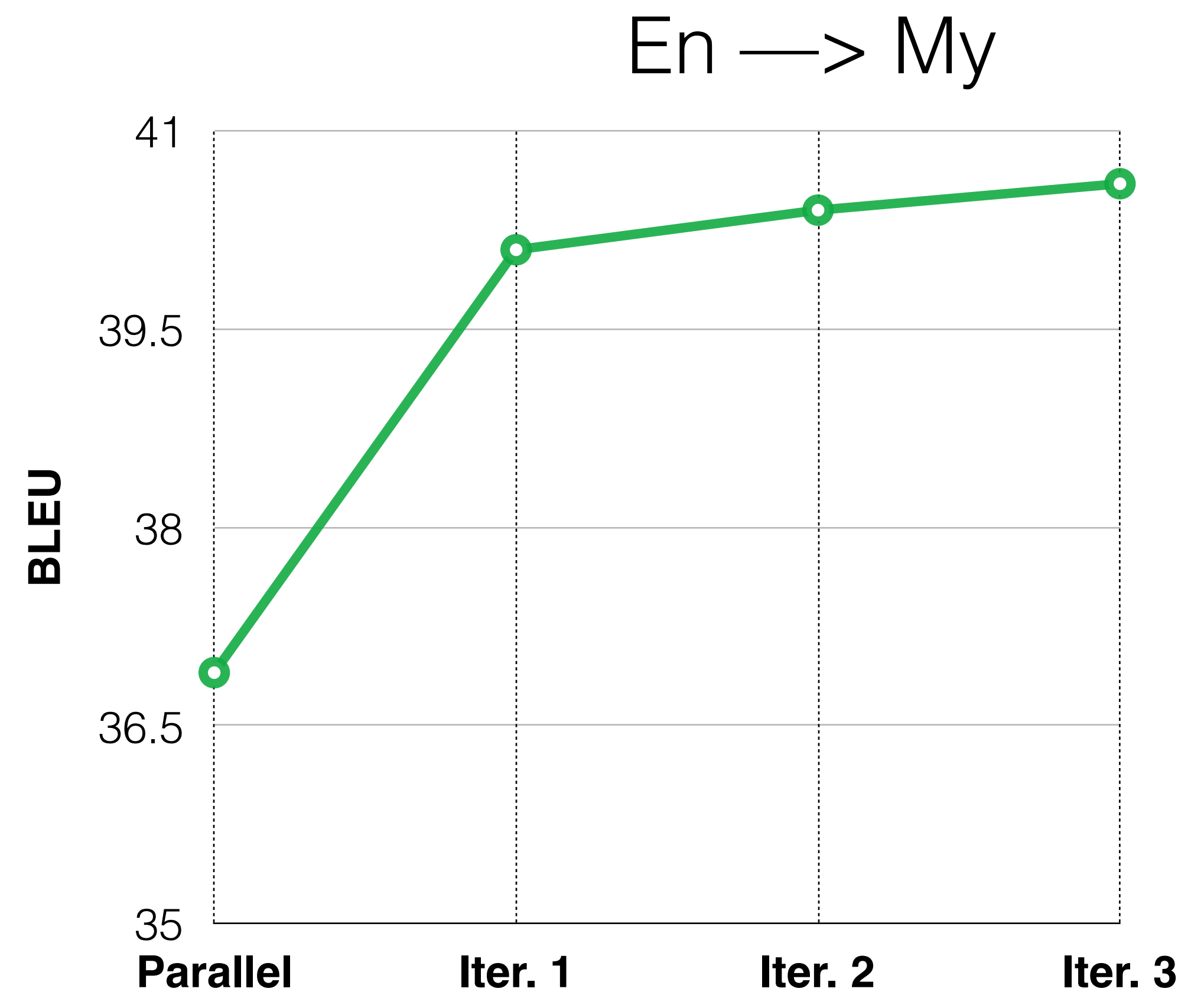
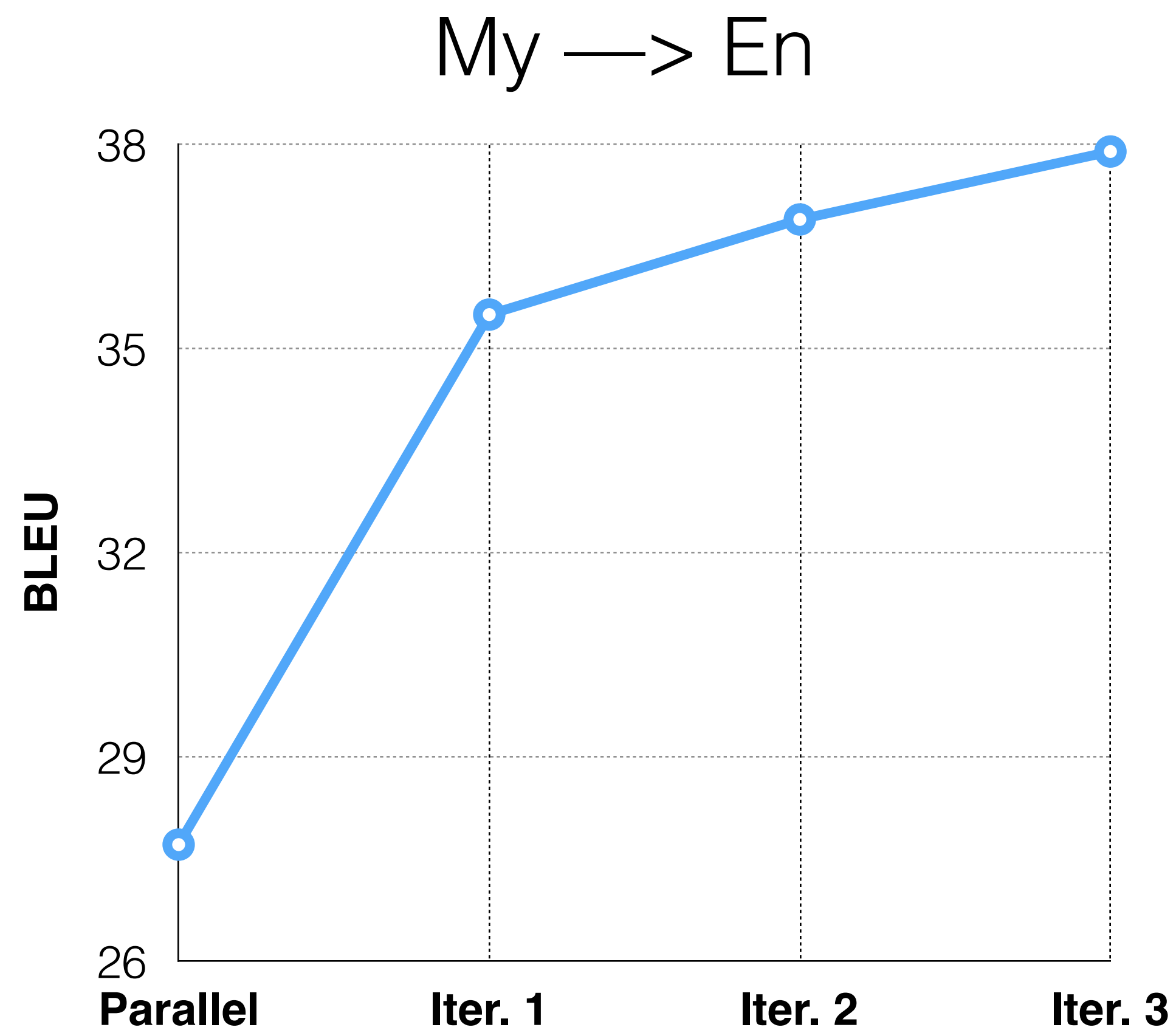
Iterative BT+ST



Iterative BT+ST

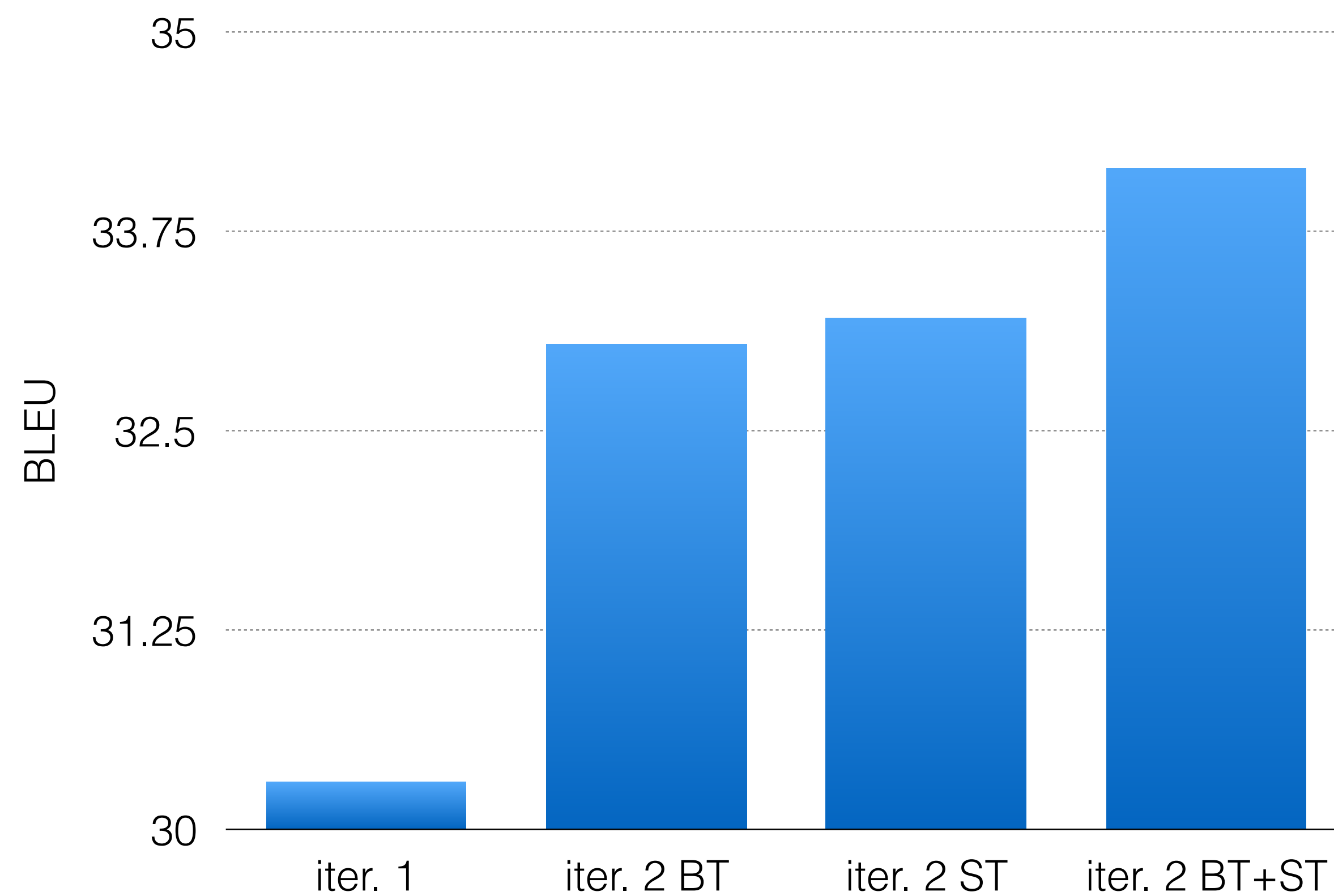


Results: Iterative ST+BT



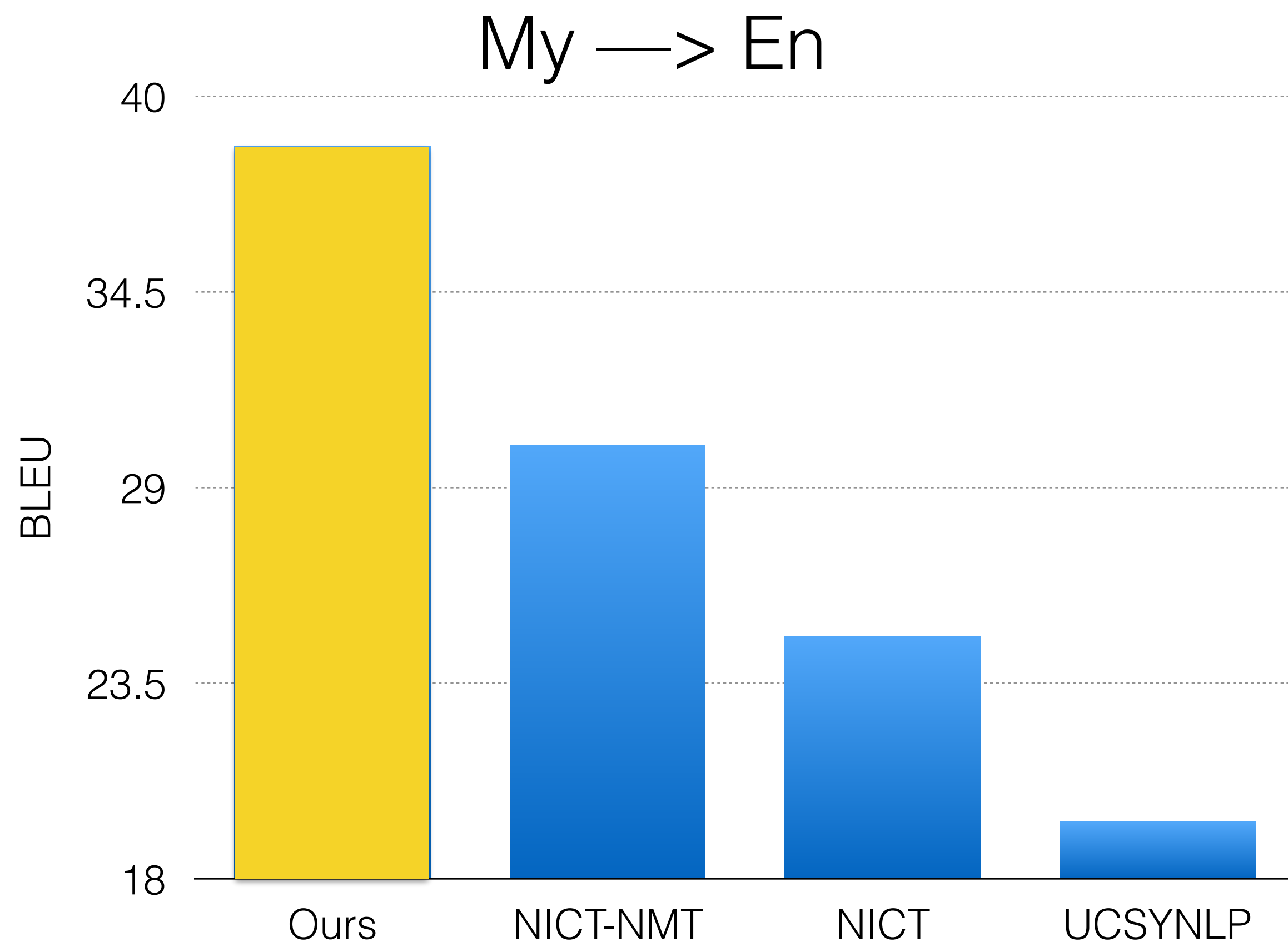
Results: BT vs ST vs BT+ST

My \rightarrow En, iter. 2

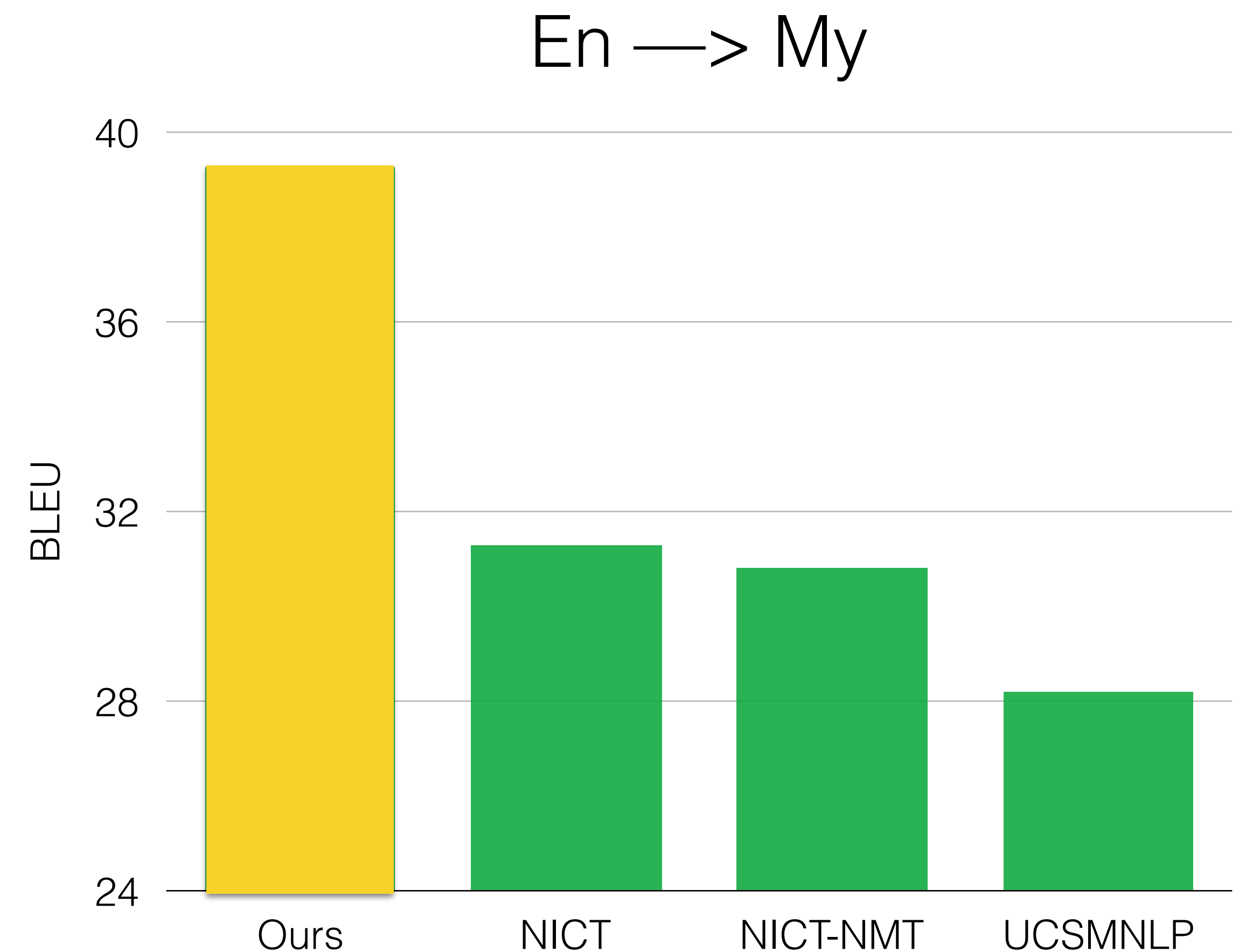


Final Results of 2019 Competition

+8 BLEU compared to second best



<http://lotus.kuee.kyoto-u.ac.jp/WAT/evaluation/list.php?t=70&o=4>



<http://lotus.kuee.kyoto-u.ac.jp/WAT/evaluation/list.php?t=71&o=9>

Conclusion

- The effect of **place** in MT is significant for low resource language pairs.
- Locality of topics is responsible for **source / target domain mismatch**. This decreases the effectiveness of back-translation.
- STDMM can be easily simulated using public benchmarks.
- Self-training works well when there is little monolingual data on the target side and when there is extreme source / target domain mismatch.
- Self-training is complementary to back-translation. They can be combined in an iterative manner.

Open Research Questions

- Are there ways to measure STDMM?
- What are good methods to cope with STDMM?
- In general, how to adapt to the desired domain given little in-domain data (either parallel or just monolingual)?

THANK YOU