

CS11-737 Multilingual NLP

# Multilingual Neural Machine Translation Pre-training and Joint Training Strategies

Lei Li

<https://lileicc.github.io/course/11737mnlp23fa/>



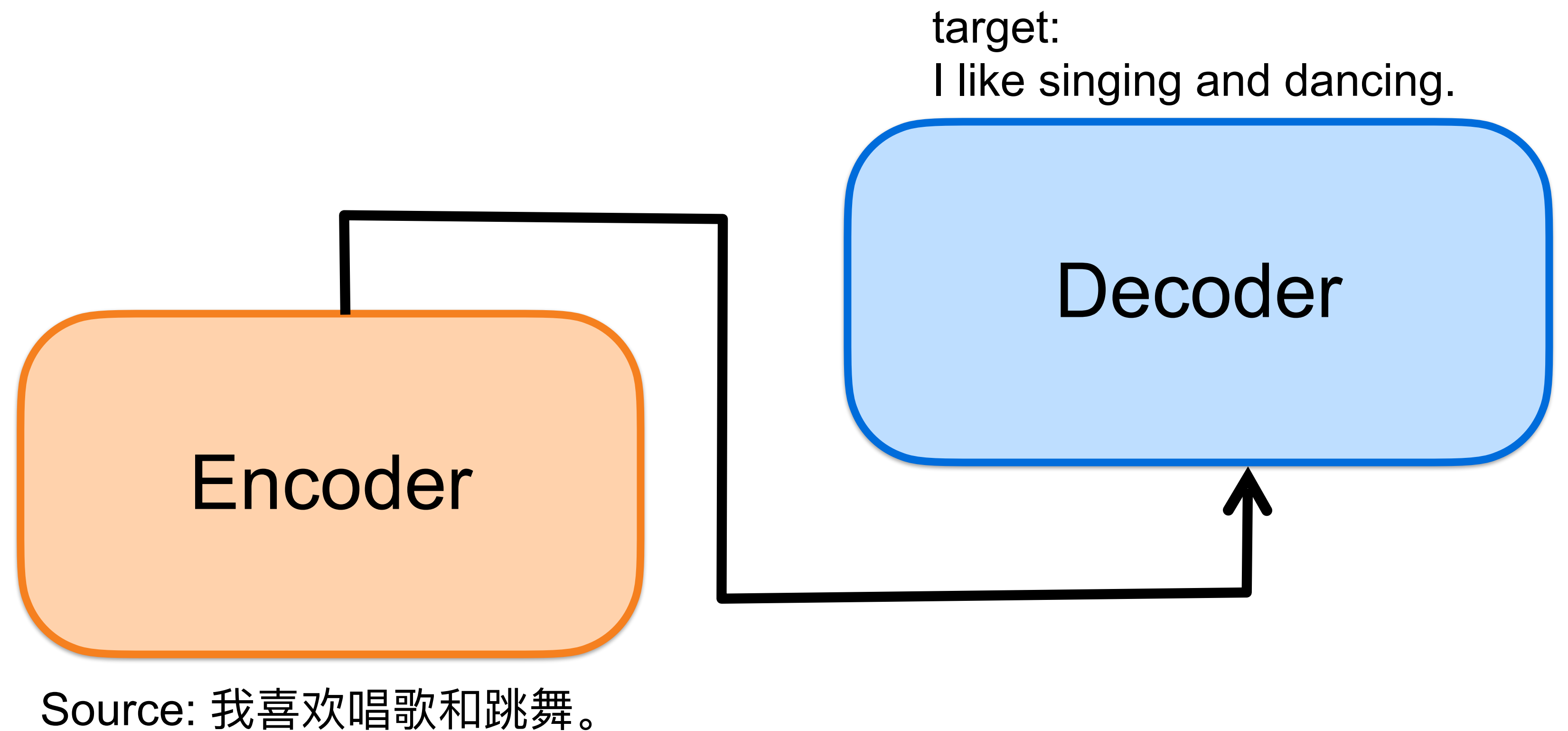
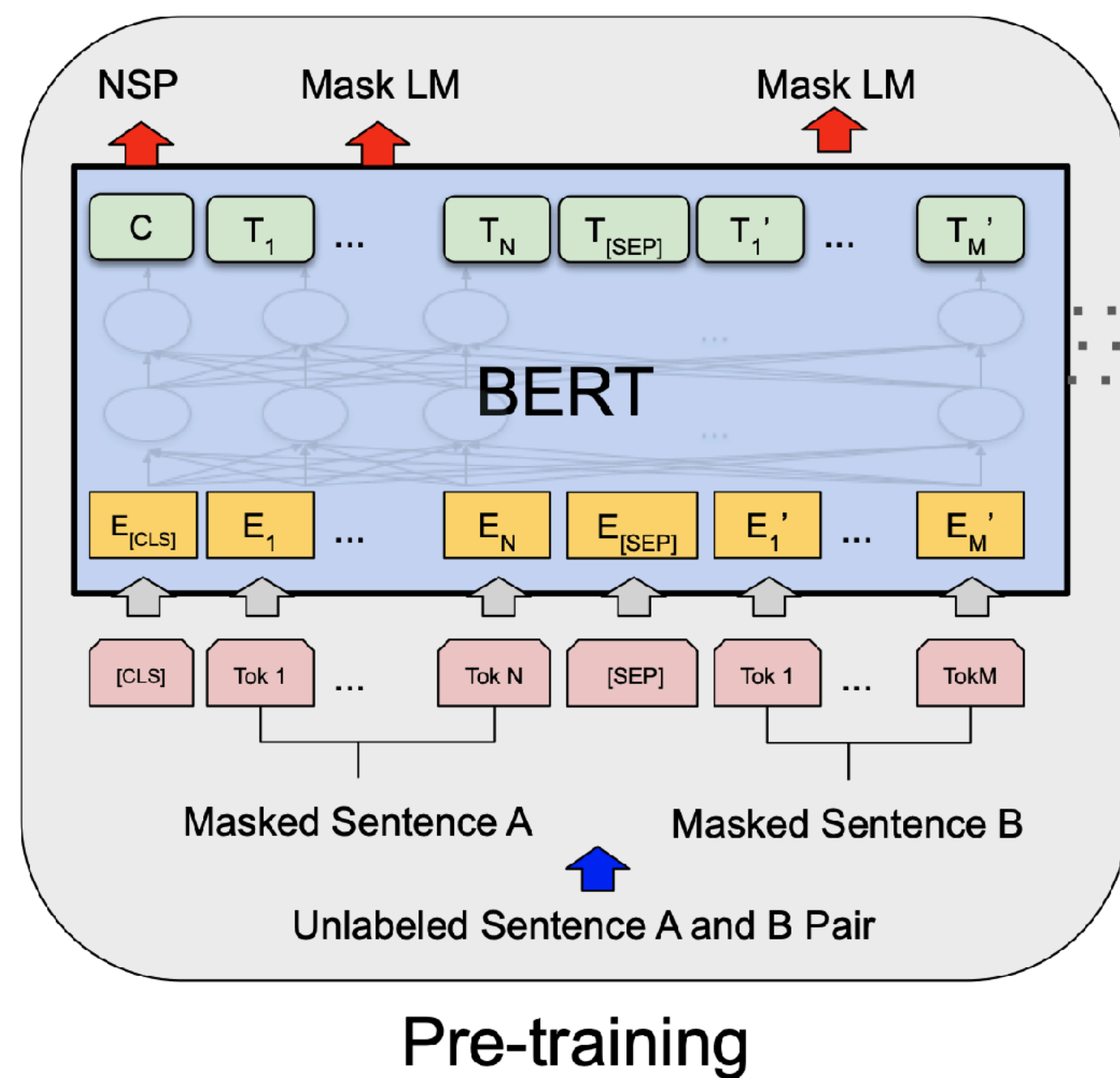
**Carnegie Mellon University**

Language Technologies Institute

# Sequence-to-sequence Pre-training

# Mismatch between Pre-trained LM and MT

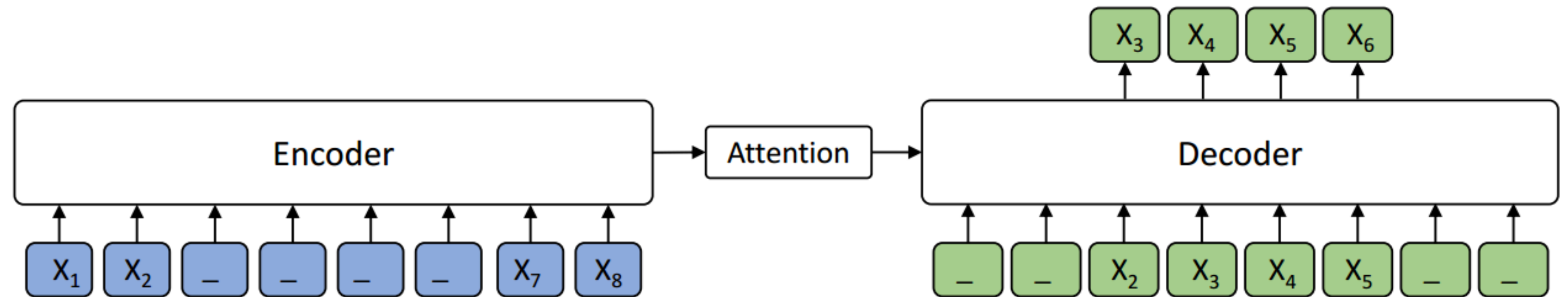
- BERT/GPT pre-training objective is different from MT



# MASS: Pre-train for Sequence to Sequence Generation

- MASS is carefully designed to jointly pre-train the encoder

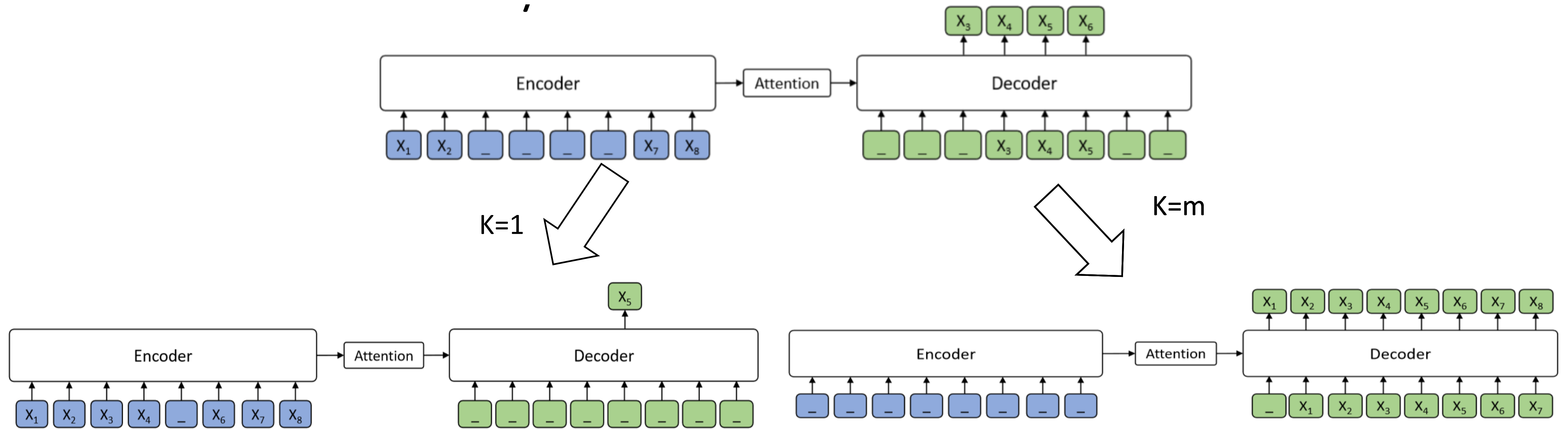
and decoder



- Mask  $k$  consecutive tokens (segment)

- Force the decoder to attend on the source representations, i.e., encoder-decoder attention
- Develop the decoder with the ability of language modeling

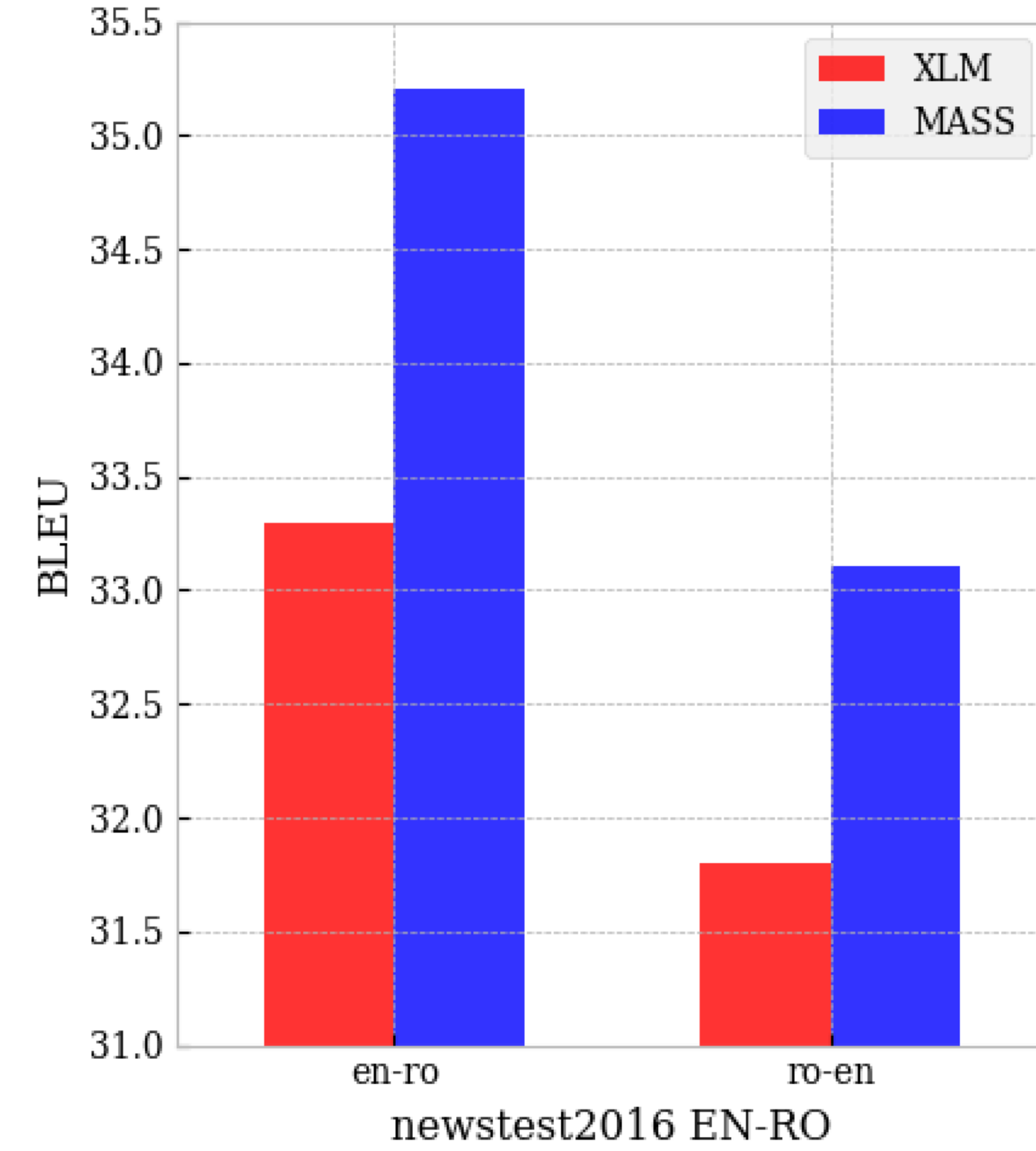
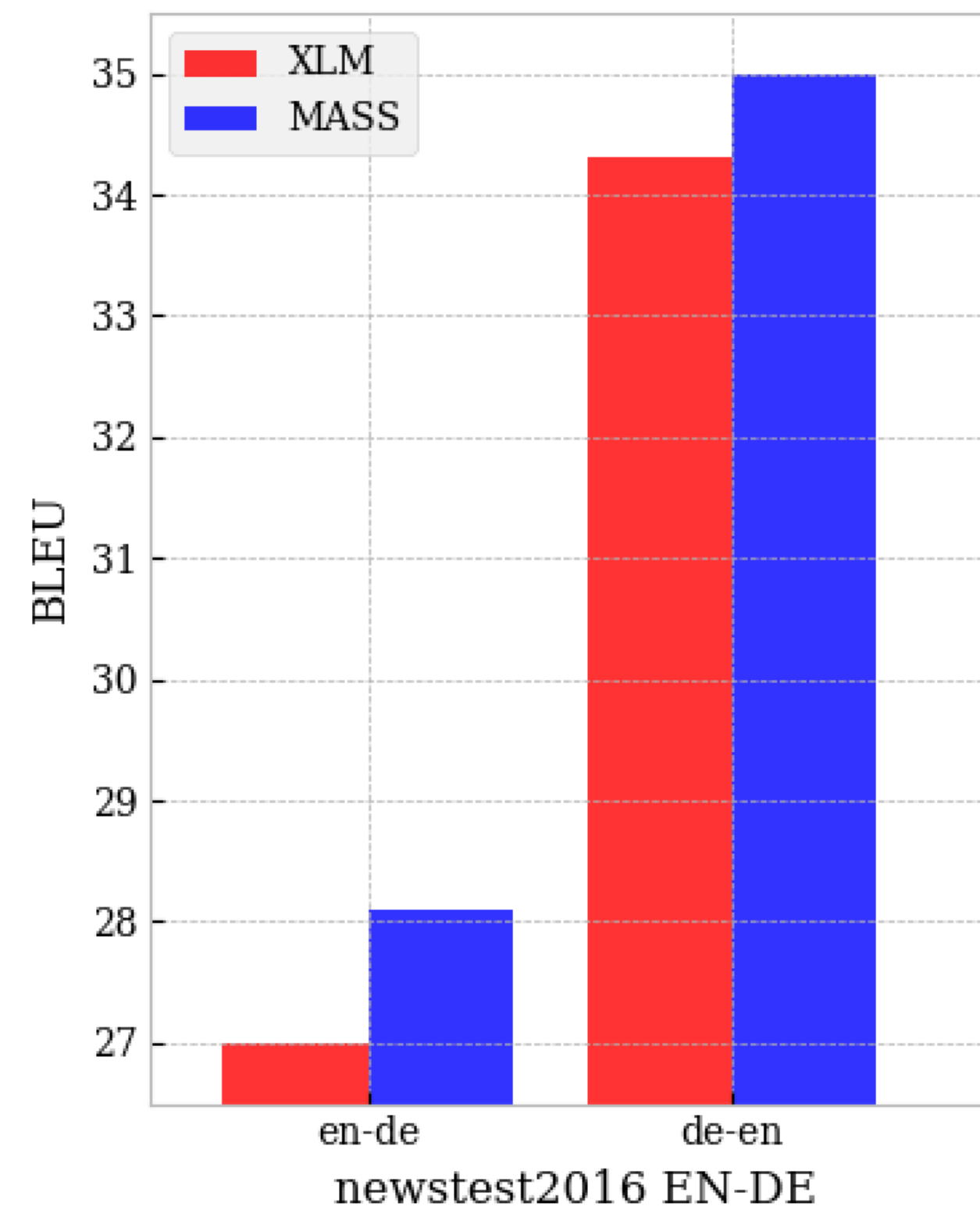
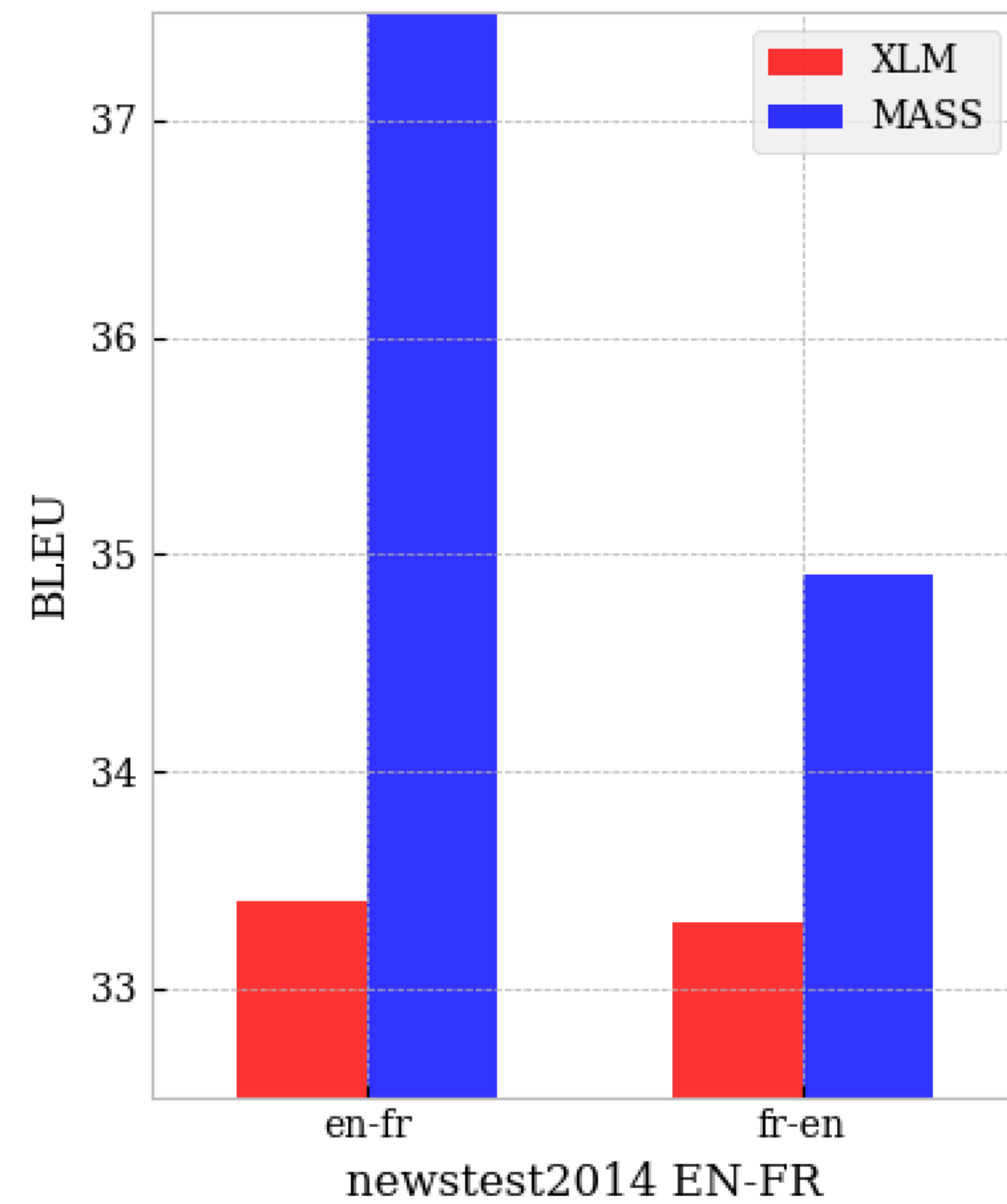
# MASS vs. BERT/GPT



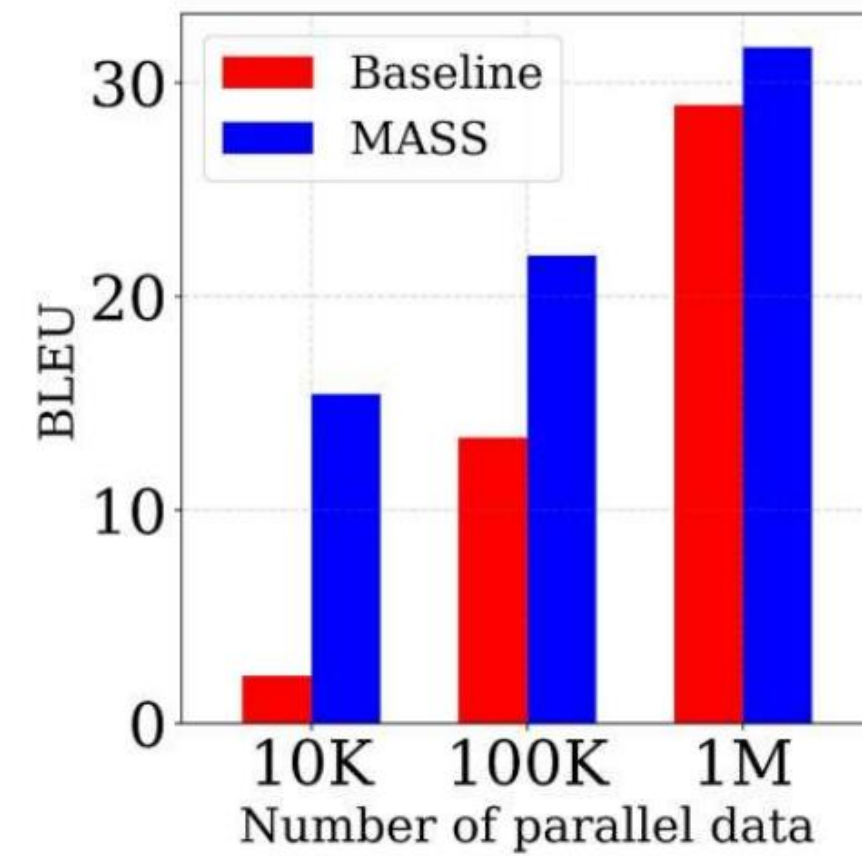
Length	Probability	Model
$k = 1$	$P(x^u   x^{\setminus u}; \theta)$	masked LM in BERT
$k \in [1, m]$	$P(x^{u:v}   x^{\setminus u:v}; \theta)$	MASS

Length	Probability	Model
$k = m$	$P(x^{1:m}   x^{\setminus 1:m}; \theta)$	standard LM in GPT
$k \in [1, m]$	$P(x^{u:v}   x^{\setminus u:v}; \theta)$	MASS

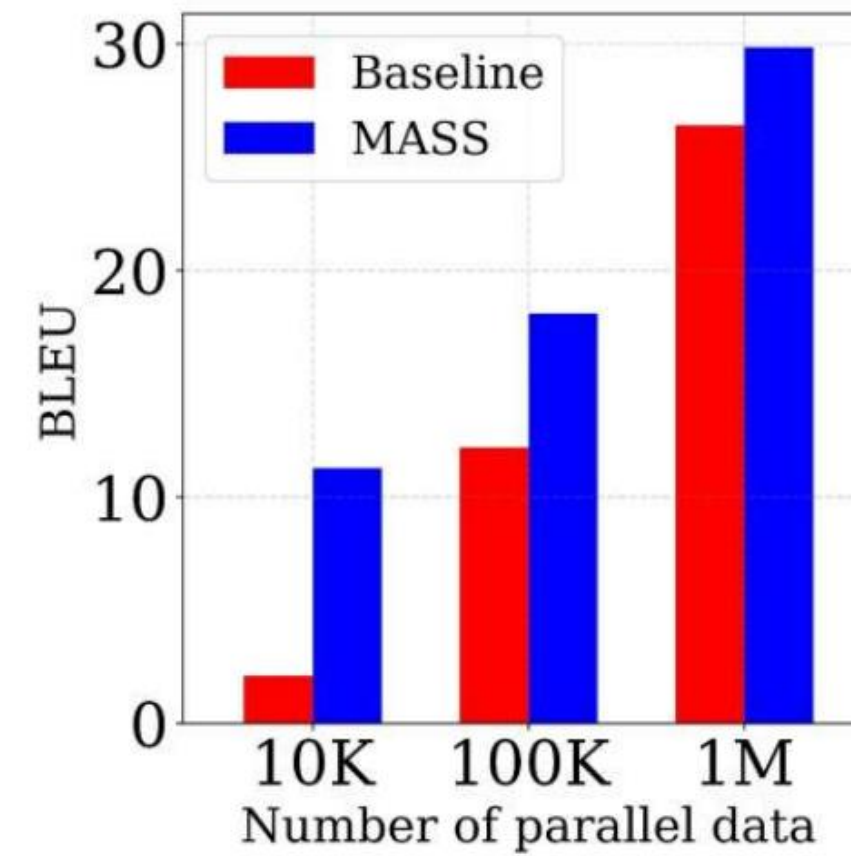
# Unsupervised NMT



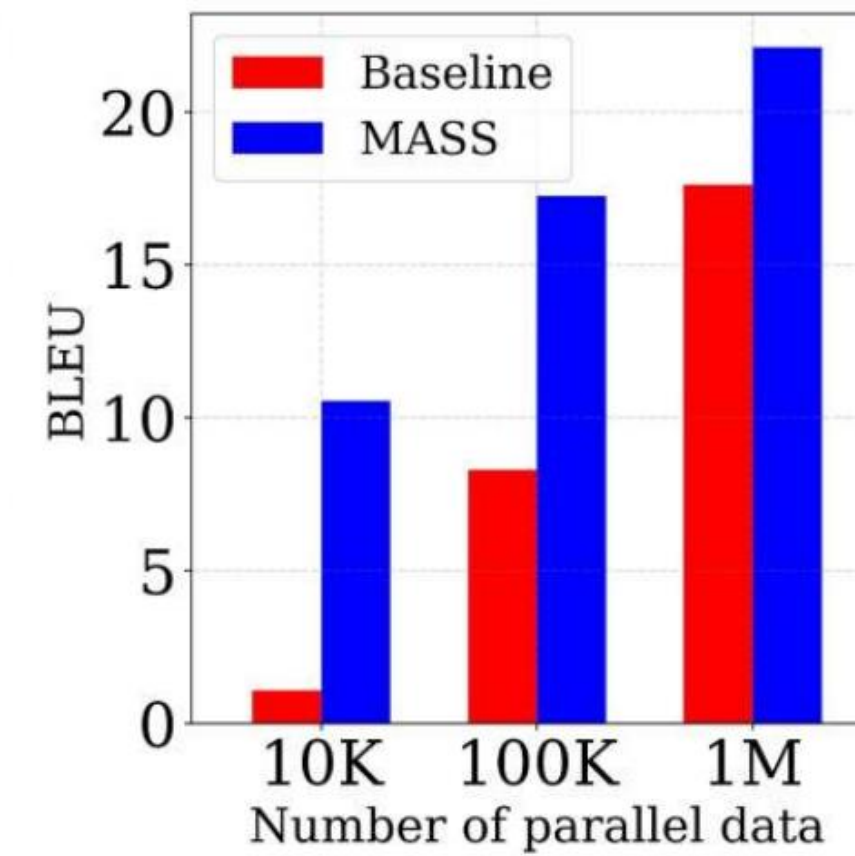
# Low-resource NMT



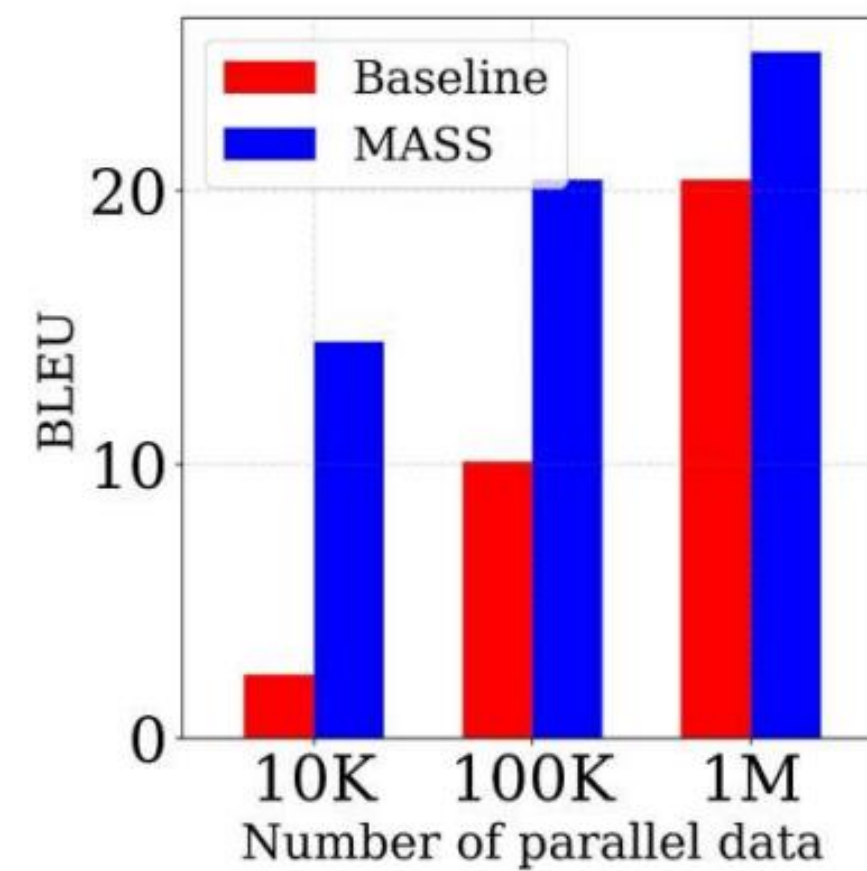
(a) en-fr



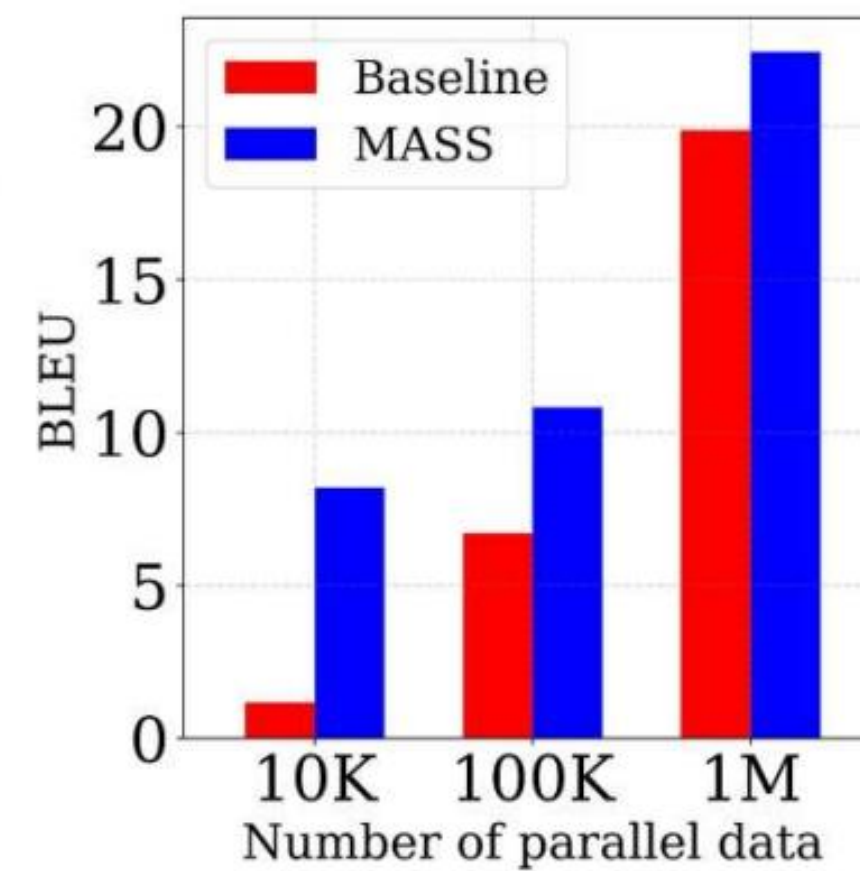
(b) fr-en



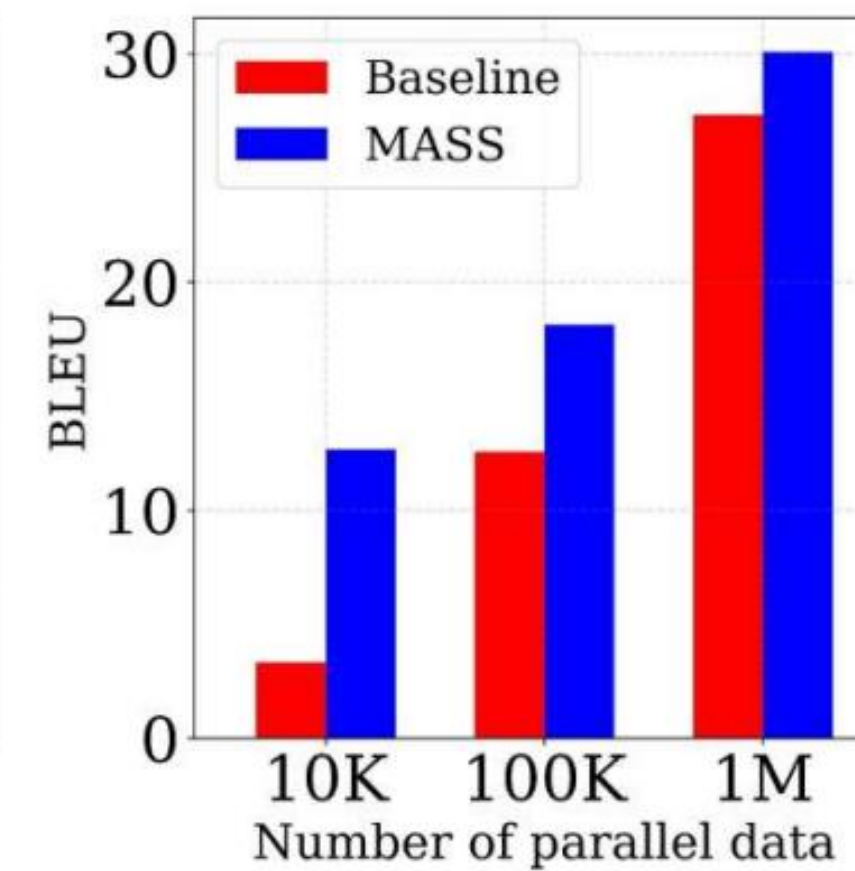
(c) en-de



(d) de-en



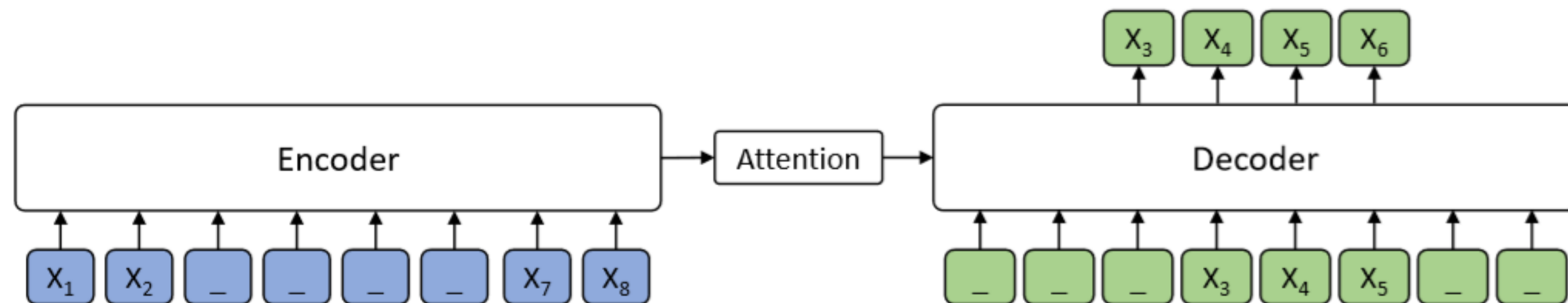
(e) en-ro



(f) ro-en

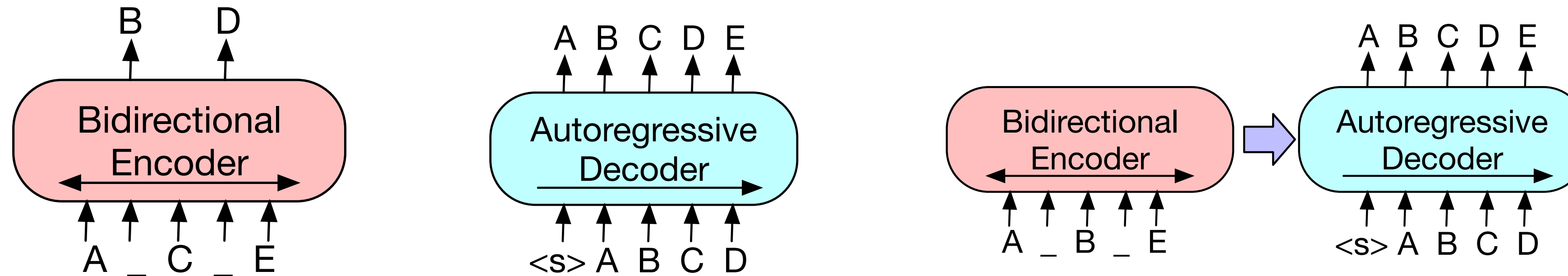
# Summary

- Advantages
  - Unified sequence-to-sequence pretraining which jointly pretrains encoder, decoder and cross attention
  - Achieves improvements on zero-shot / unsupervised NMT
- Limitations
  - No evidence on rich resource NMT
  - Pre-training objective inconsistent with NMT, e.g. [monolingual v.s. multilingual](#)





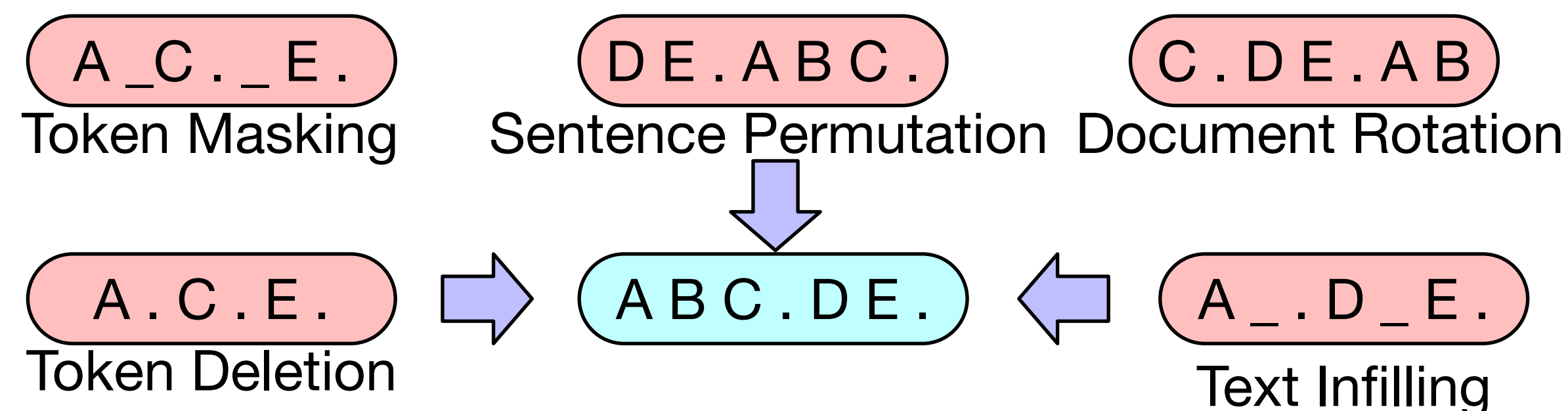
# BART: Denoising Sequence-to-Sequence Pre-training



A schema comparison with BERT, GPT and BART.

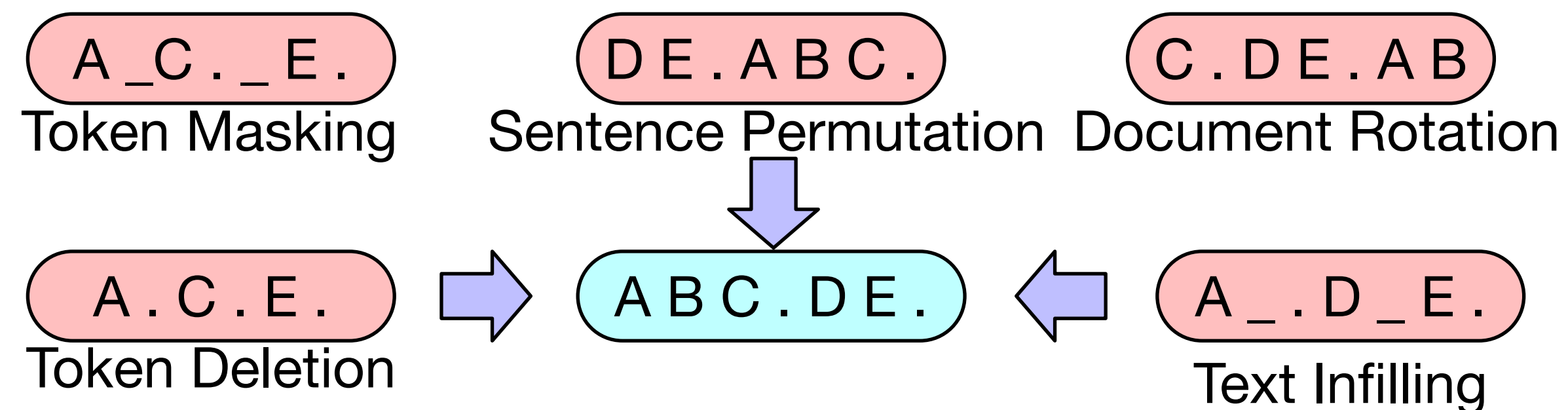
- Standard sequence-to-sequence Transformer architecture
- Trained by corrupting documents and then optimizing a reconstruction loss
- Allows to apply *any* type of document corruption.

# Noising the input



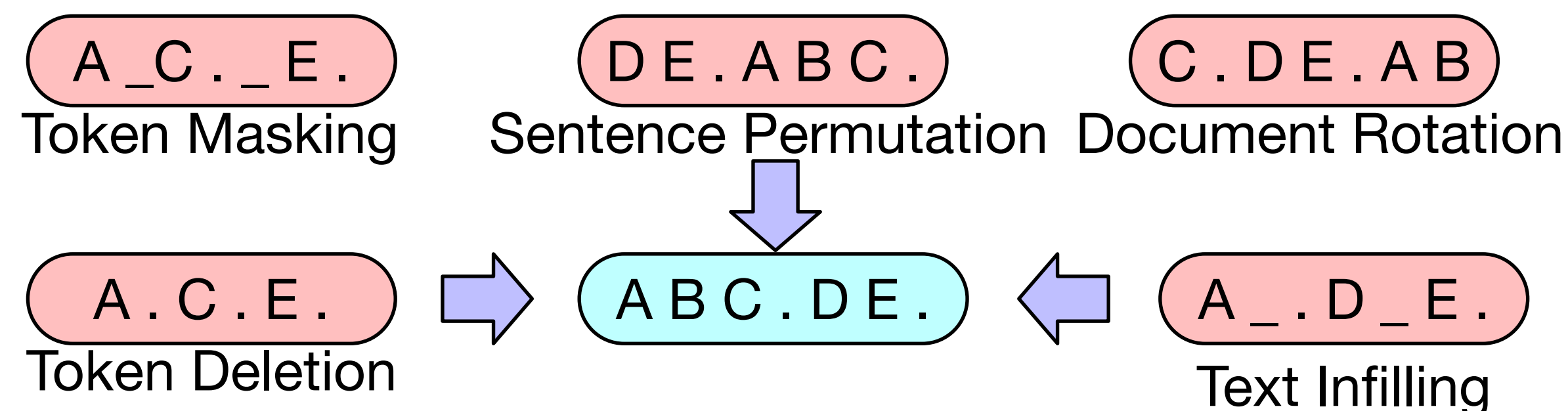
- Token masking: Random tokens are sampled and replaced with [MASK]
- Token deletion: Random tokens are deleted from the input.
- Text infilling: A number of span are sampled. Each span is replaced with [MASK]. 0-length span corresponding the insertion of [MASK].
- Sentence permutation: Sentences are shuffled with random order.
- Document Rotation: A token is chosen uniformly at random, and the document is rotated so that it begins with that token.

# Noising the input



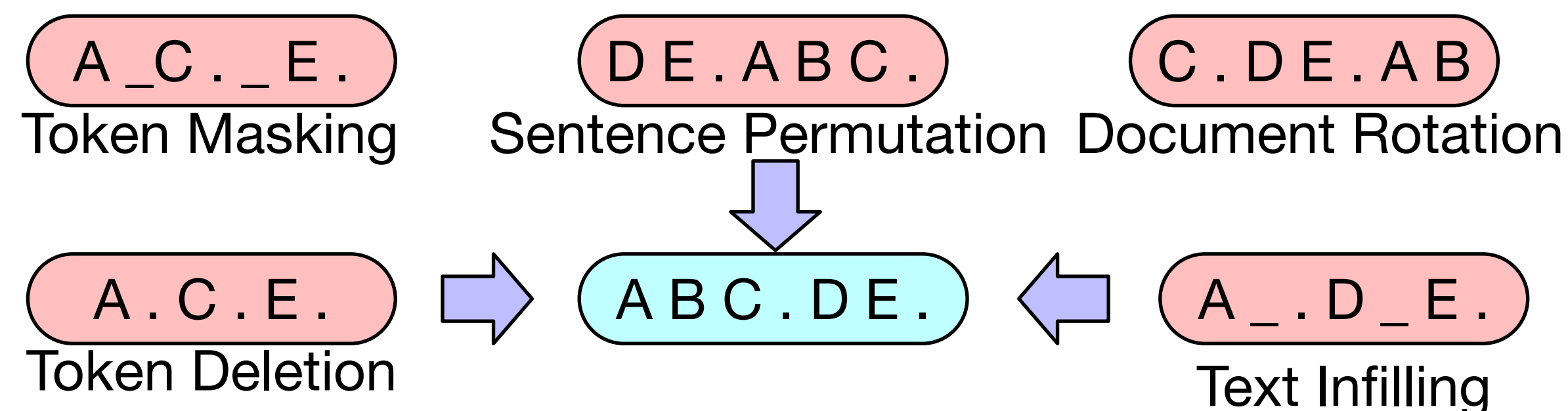
- Token masking: Random tokens are sampled and replaced with [MASK]
- Token deletion: Random tokens are deleted from the input.
- Text infilling: A number of span are sampled. Each span is replaced with [MASK]. 0-length span corresponding the insertion of [MASK].
- Sentence permutation: Sentences are shuffled with random order.
- Document Rotation: A token is chosen uniformly at random, and the document is rotated so that it begins with that token.

# Noising the input



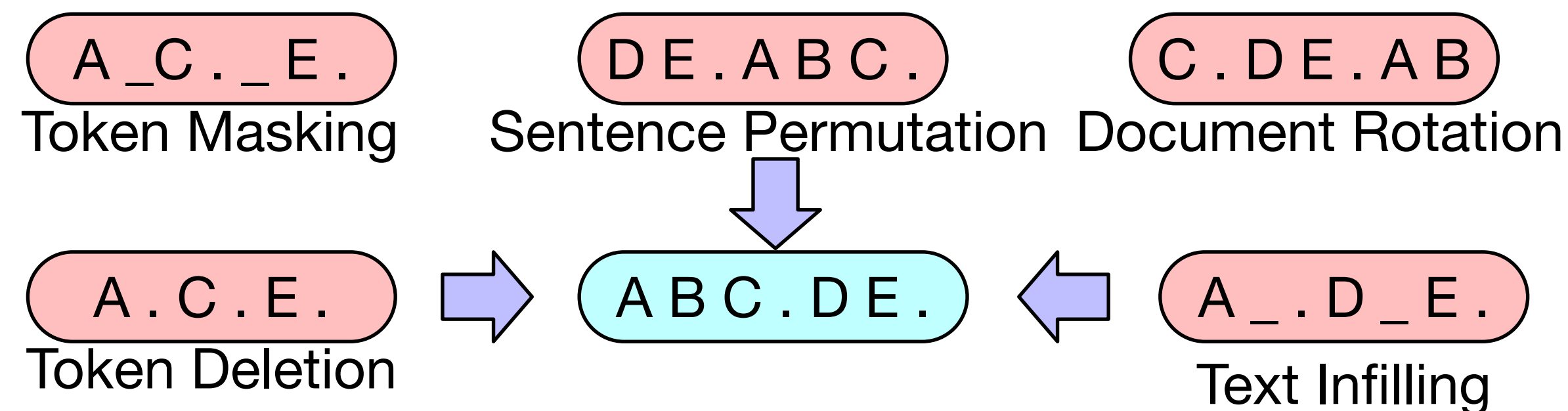
- Token masking: Random tokens are sampled and replaced with [MASK]
- Token deletion: Random tokens are deleted from the input.
- Text infilling: A number of span are sampled. Each span is replaced with [MASK]. 0-length span corresponding the insertion of [MASK].
- Sentence permutation: Sentences are shuffled with random order.
- Document Rotation: A token is chosen uniformly at random, and the document is rotated so that it begins with that token.

# Noising the input



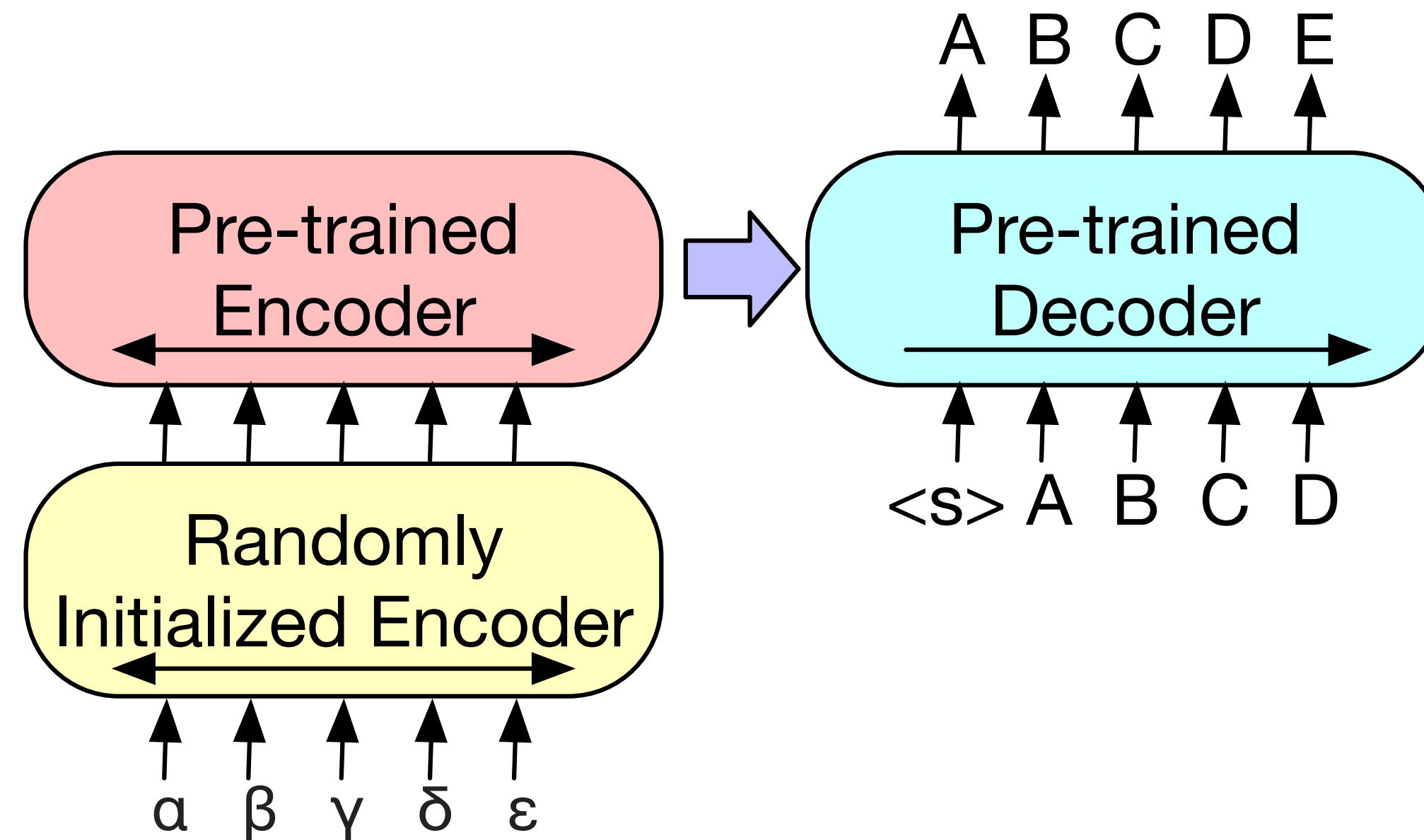
- Token masking: Random tokens are sampled and replaced with [MASK]
- Token deletion: Random tokens are deleted from the input.
- Text infilling: A number of span are sampled. Each span is replaced with [MASK]. 0-length span corresponding the insertion of [MASK].
- Sentence permutation: Sentences are shuffled with random order.
- Document Rotation: A token is chosen uniformly at random, and the document is rotated so that it begins with that token.

# Noising the input



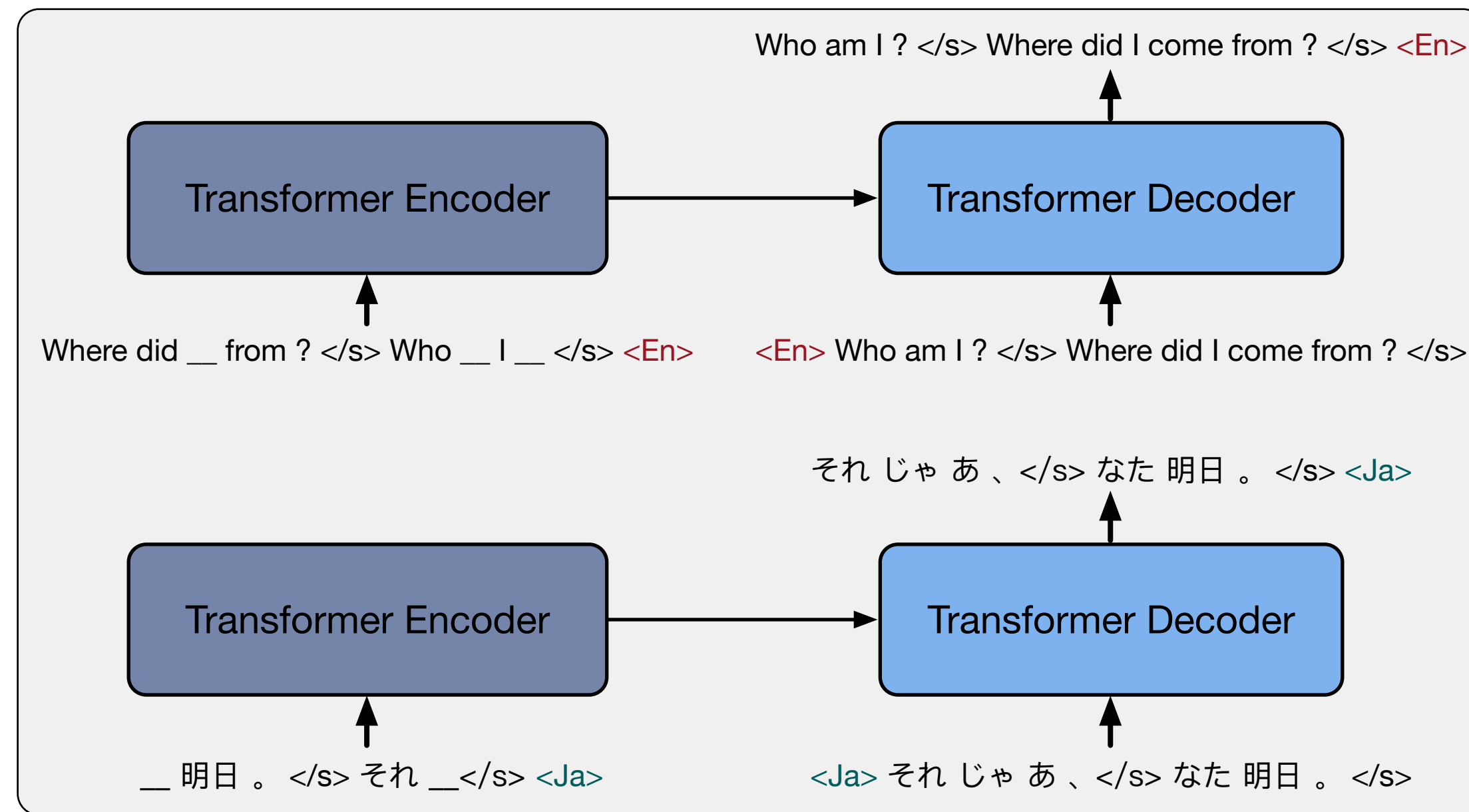
- Token masking: Random tokens are sampled and replaced with [MASK]
- Token deletion: Random tokens are deleted from the input.
- Text infilling: A number of span are sampled. Each span is replaced with [MASK]. 0-length span corresponding the insertion of [MASK].
- Sentence permutation: Sentences are shuffled with random order.
- Document Rotation: A token is chosen uniformly at random, and the document is rotated so that it begins with that token.

# Fine-Tune on Neural Machine Translation

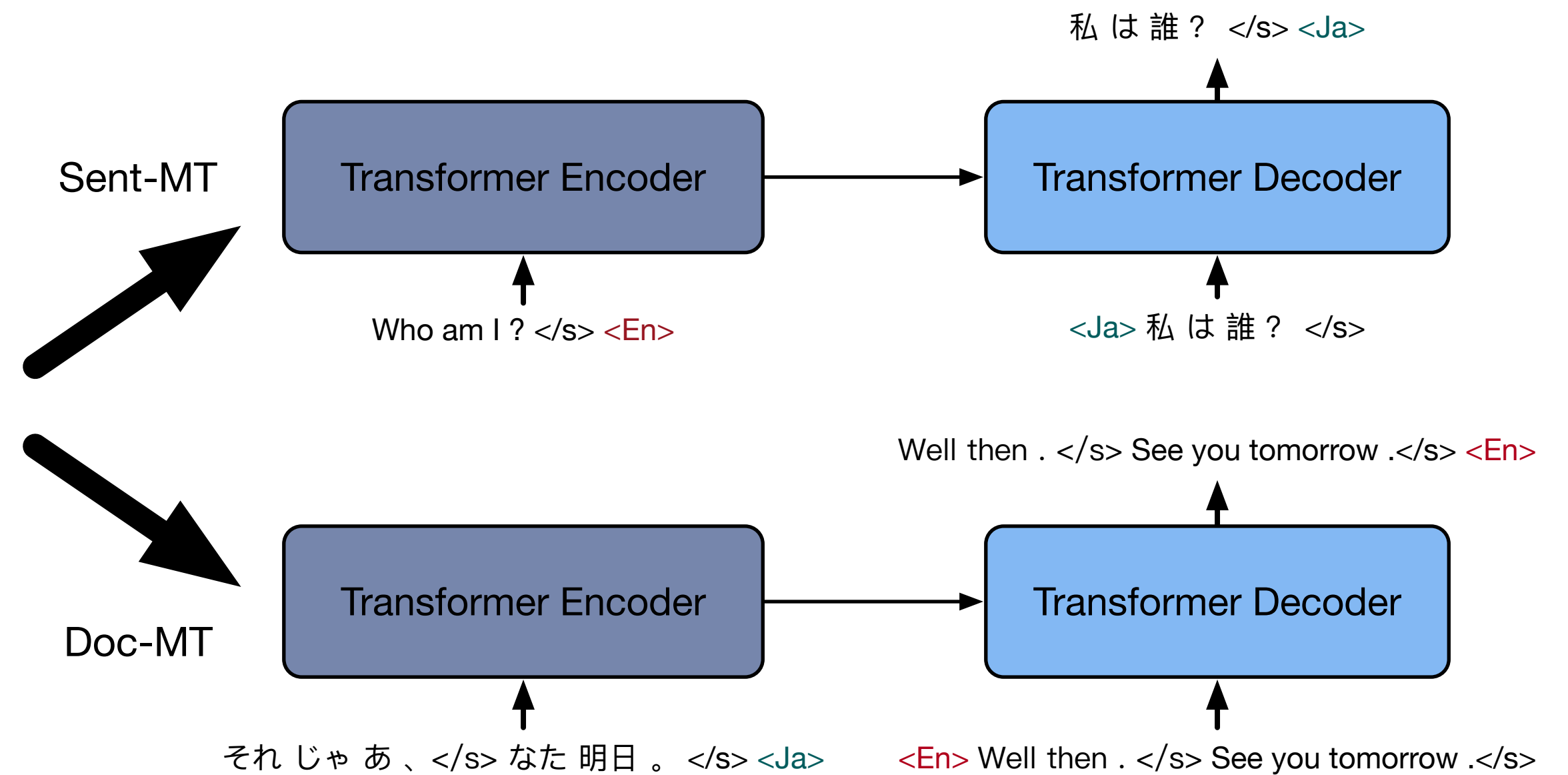


- Replace BART's encoder embedding layer with a new randomly initialized encoder
- The new encoder uses a separate vocabulary from the original BART mode
- First, freeze BART parameters and **only** update the randomly initialized source encoder. Then, jointly tuning with a few steps.

# mBART: Multilingual Denoising Pre-training for Neural Machine Translation



Multilingual Denoising **Pre-Training** (mBART)



**Fine-tuning** on Machine Translation

- Multilingual denoising **pre-training** (25 languages)
  - Sentence permutation
  - Word-span masking
- **Fine-tuning** on MT with special language id



# Dataset

- Data: CC25 corpus

- CC25 includes 25 languages from different families and with varied amounts of text from Common Crawl (CC)
- Rebalanced the corpus by up/down-sampling text

$$\lambda_i = \frac{1}{p_i} \cdot \frac{p_i^\alpha}{\sum_i p_i^\alpha},$$

- Sentence Piece which includes 25,000 subwords
- Noisy function follows BART

Code	Language	Tokens/M	Size/GB
<b>En</b>	English	55608	300.8
<b>Ru</b>	Russian	23408	278.0
<b>Vi</b>	Vietnamese	24757	137.3
<b>Ja</b>	Japanese	530 (*)	69.3
<b>De</b>	German	10297	66.6
<b>Ro</b>	Romanian	10354	61.4
<b>Fr</b>	French	9780	56.8
<b>Fi</b>	Finnish	6730	54.3
<b>Ko</b>	Korean	5644	54.2
<b>Es</b>	Spanish	9374	53.3
<b>Zh</b>	Chinese (Sim)	259 (*)	46.9
<b>It</b>	Italian	4983	30.2
<b>Nl</b>	Dutch	5025	29.3
<b>Ar</b>	Arabic	2869	28.0
<b>Tr</b>	Turkish	2736	20.9
<b>Hi</b>	Hindi	1715	20.2
<b>Cs</b>	Czech	2498	16.3
<b>Lt</b>	Lithuanian	1835	13.7
<b>Lv</b>	Latvian	1198	8.8
<b>Kk</b>	Kazakh	476	6.4
<b>Et</b>	Estonian	843	6.1
<b>Ne</b>	Nepali	237	3.8
<b>Si</b>	Sinhala	243	3.6
<b>Gu</b>	Gujarati	140	1.9
<b>My</b>	Burmese	56	1.6

# mBART: Low-medium translation results

Languages	En-Gu	En-Kk	En-Vi	En-Tr	En-Ja	En-Ko						
Data Source	WMT19	WMT19	IWSLT15	WMT17	IWSLT17	IWSLT17						
Size	10K	91K	133K	207K	223K	230K						
Direction	← →	← →	← →	← →	← →	← →						
Random	0.0	0.0	0.8	0.2	23.6	24.8	12.2	9.5	10.4	12.3	15.3	16.3
mBART25	<b>0.3</b>	<b>0.1</b>	<b>7.4</b>	<b>2.5</b>	<b>36.1</b>	<b>35.4</b>	<b>22.5</b>	<b>17.8</b>	<b>19.1</b>	<b>19.4</b>	<b>24.6</b>	<b>22.6</b>

Languages	En-Nl	En-Ar	En-It	En-My	En-Ne	En-Ro						
Data Source	IWSLT17	IWSLT17	IWSLT17	WAT19	FLoRes	WMT16						
Size	237K	250K	250K	259K	564K	608K						
Direction	← →	← →	← →	← →	← →	← →						
Random	34.6	29.3	27.5	16.9	31.7	28.0	23.3	34.9	7.6	4.3	34.0	34.3
mBART25	<b>43.3</b>	<b>34.8</b>	<b>37.6</b>	<b>21.6</b>	<b>39.8</b>	<b>34.0</b>	<b>28.3</b>	<b>36.9</b>	<b>14.5</b>	<b>7.4</b>	<b>37.8</b>	<b>37.7</b>

Languages	En-Si	En-Hi	En-Et	En-Lt	En-Fi	En-Lv						
Data Source	FLoRes	ITTB	WMT18	WMT19	WMT17	WMT17						
Size	647K	1.56M	1.94M	2.11M	2.66M	4.50M						
Direction	← →	← →	← →	← →	← →	← →						
Random	7.2	1.2	10.9	14.2	22.6	17.9	18.1	12.1	21.8	20.2	15.6	12.9
mBART25	<b>13.7</b>	<b>3.3</b>	<b>23.5</b>	<b>20.8</b>	<b>27.8</b>	<b>21.4</b>	<b>22.4</b>	<b>15.3</b>	<b>28.5</b>	<b>22.4</b>	<b>19.3</b>	<b>15.9</b>

Low resource: more than 6 BLEU. But fails in extremely low-resource setting

# mBART: Low-medium translation results

Languages	En-Gu	En-Kk	En-Vi	En-Tr	En-Ja	En-Ko						
Data Source	WMT19	WMT19	IWSLT15	WMT17	IWSLT17	IWSLT17						
Size	10K	91K	133K	207K	223K	230K						
Direction	← →	← →	← →	← →	← →	← →						
Random	0.0	0.0	0.8	0.2	23.6	24.8	12.2	9.5	10.4	12.3	15.3	16.3
<b>mBART25</b>	<b>0.3</b>	<b>0.1</b>	<b>7.4</b>	<b>2.5</b>	<b>36.1</b>	<b>35.4</b>	<b>22.5</b>	<b>17.8</b>	<b>19.1</b>	<b>19.4</b>	<b>24.6</b>	<b>22.6</b>

Languages	En-Nl	En-Ar	En-It	En-My	En-Ne	En-Ro						
Data Source	IWSLT17	IWSLT17	IWSLT17	WAT19	FLoRes	WMT16						
Size	237K	250K	250K	259K	564K	608K						
Direction	← →	← →	← →	← →	← →	← →						
Random	34.6	29.3	27.5	16.9	31.7	28.0	23.3	34.9	7.6	4.3	34.0	34.3
<b>mBART25</b>	<b>43.3</b>	<b>34.8</b>	<b>37.6</b>	<b>21.6</b>	<b>39.8</b>	<b>34.0</b>	<b>28.3</b>	<b>36.9</b>	<b>14.5</b>	<b>7.4</b>	<b>37.8</b>	<b>37.7</b>

Languages	En-Si	En-Hi	En-Et	En-Lt	En-Fi	En-Lv						
Data Source	FLoRes	ITTB	WMT18	WMT19	WMT17	WMT17						
Size	647K	1.56M	1.94M	2.11M	2.66M	4.50M						
Direction	← →	← →	← →	← →	← →	← →						
Random	7.2	1.2	10.9	14.2	22.6	17.9	18.1	12.1	21.8	20.2	15.6	12.9
<b>mBART25</b>	<b>13.7</b>	<b>3.3</b>	<b>23.5</b>	<b>20.8</b>	<b>27.8</b>	<b>21.4</b>	<b>22.4</b>	<b>15.3</b>	<b>28.5</b>	<b>22.4</b>	<b>19.3</b>	<b>15.9</b>

Low resource: more than 6 BLEU. But fails in extremely low-resource setting

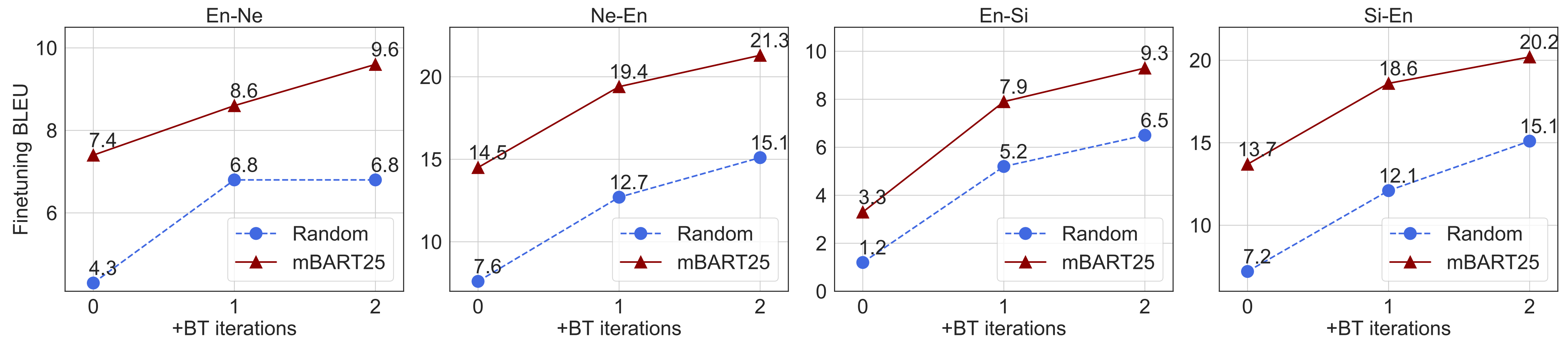
Medium resource: more than 3 BLEU

# mBART: Rich-resource translation

Languages	Cs	Es	Zh	De	Ru	Fr
Size	11M	15M	25M	28M	29M	41M
Random	16.5	33.2	<b>35.0</b>	<b>30.9</b>	<b>31.5</b>	<b>41.4</b>
mBART25	<b>18.0</b>	<b>34.0</b>	33.3	30.5	31.3	41.0

- Pre-training slightly hurts performance when >25M parallel sentence are available.
- When a significant amount of bi-text data is given, supervised training are supposed to wash out the pre-trained weights completely.

# mBART: Pre-training complementary to BT



- Test on low resource FLoRes dataset [Guzmán et al., 2019]
- Use the same monolingual data to generate BT data
- Initializing the model with mBART25 pre-trained parameters improves BLEU scores at each iteration of back translation, resulting in new state-of-the-art results in all four translation directions

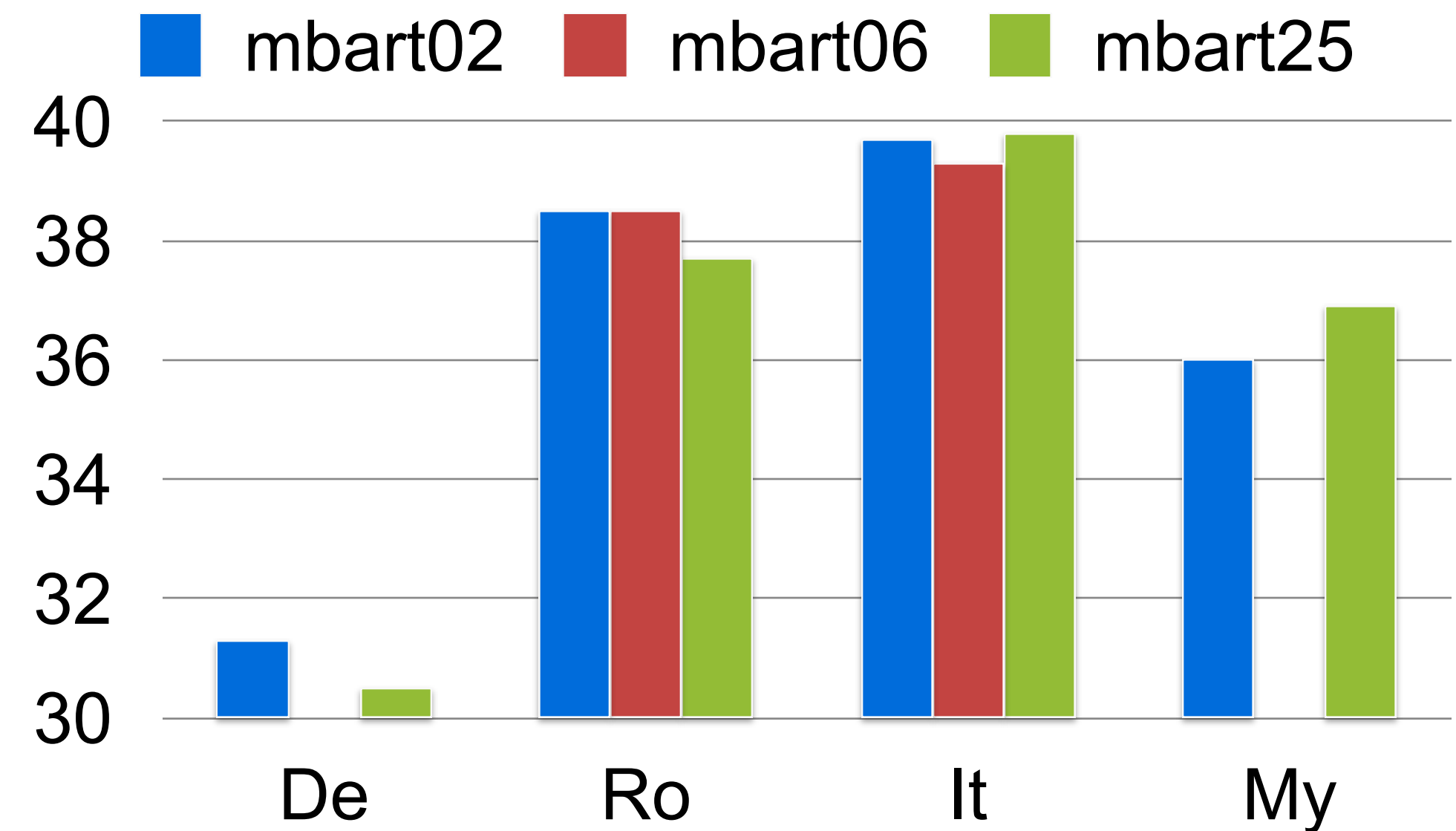
# Is pre-training on multilingual better than on single language?

Model	Pre-training	Fine-tuning		
	Data	En→Ro	Ro→En	+BT
Random	None	34.3	34.0	36.8
XLM (2019)	En Ro	-	35.6	38.5
MASS (2019)	En Ro	-	-	39.1
BART (2019)	En	-	-	38.0
XLM-R (2019)	CC100	35.6	35.8	-
BART-En	En	36.0	35.8	37.4
BART-Ro	Ro	37.6	36.8	38.1
<b>mBART02</b>	En Ro	<b>38.5</b>	<b>38.5</b>	<b>39.9</b>
mBART25	CC25	37.7	37.8	38.8

- BART model trained on the same En and Ro data only. Both have improvements over baselines, while worse than mBART results, indicating pre-training in a multilingual setting is essential.
- Combining BT leads to additional gains, resulting in a new state-of-the-art for Ro-En translation
- mBART02 is better than mBART25. The more seems not the better?

# How many languages should you pre-train on?

Languages	De	Ro	It	My	En
Size/GB	66.6	61.4	30.2	1.6	300.8
mBART02	31.3	38.5	39.7	36.5	
mBART06	-	38.5	39.3	-	
mBART25	30.5	37.7	39.8	36.9	



- Pretraining on more languages helps most when the target language monolingual data is limited
- When monolingual data is plentiful (De, Ro), pre-training on multiple languages slightly hurts the final results (<1 BLEU)

# Analysis: Generalization to unseen languages

	Monolingual	NI-En	En-NI	Ar-En	En-Ar	NI-De	De-NI
<b>Random</b>	None	34.6 (-8.7)	29.3 (-5.5)	27.5 (-10.1)	16.9 (-4.7)	21.3 (-6.4)	20.9 (-5.2)
<b>mBART02</b>	En Ro	41.4 (-2.9)	34.5 (-0.3)	34.9 (-2.7)	21.2 (-0.4)	26.1 (-1.6)	25.4 (-0.7)
<b>mBART06</b>	En Ro Cs It Fr Es	43.1 (-0.2)	34.6 (-0.2)	37.3 (-0.3)	21.1 (-0.5)	26.4 (-1.3)	25.3 (-0.8)
<b>mBART25</b>	All	<b>43.3</b>	<b>34.8</b>	<b>37.6</b>	<b>21.6</b>	<b>27.7</b>	<b>26.1</b>

NI-De and Ar are not included in the pre-training corpus

- mBART can improve performance even with fine tuning for languages that did not appear in the pre-training corpora,
- Pre-training has language universal aspects, especially within the parameters learned at the Transformer layers.
- The more pre-trained languages the better

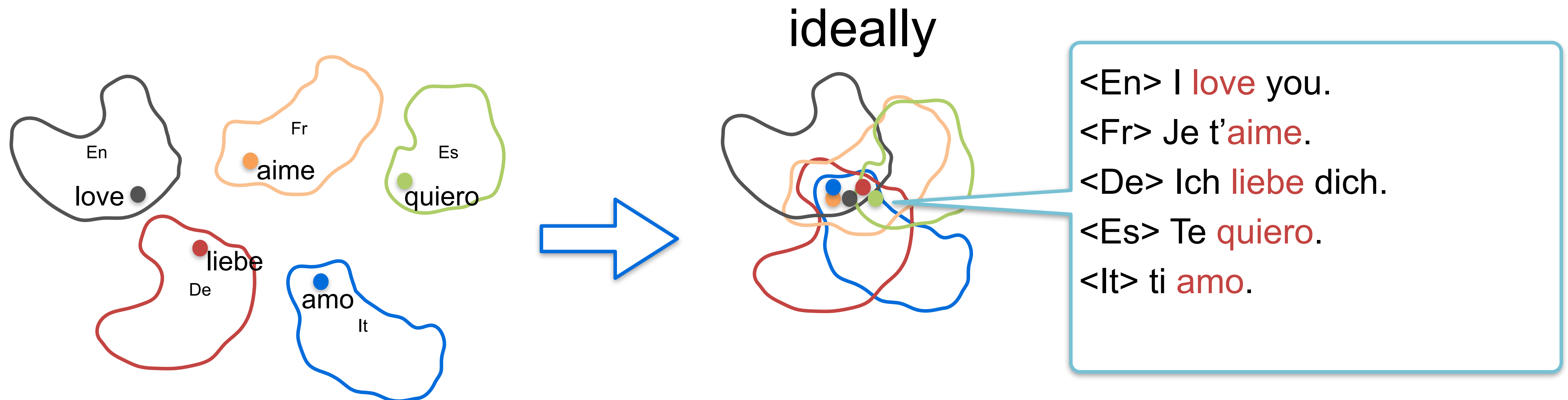


# Multilingual Training

**How can we build a single unified Multilingual MT models with superior performance on all language directions?**

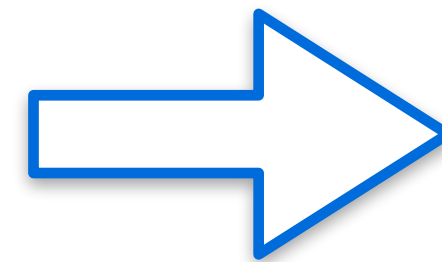
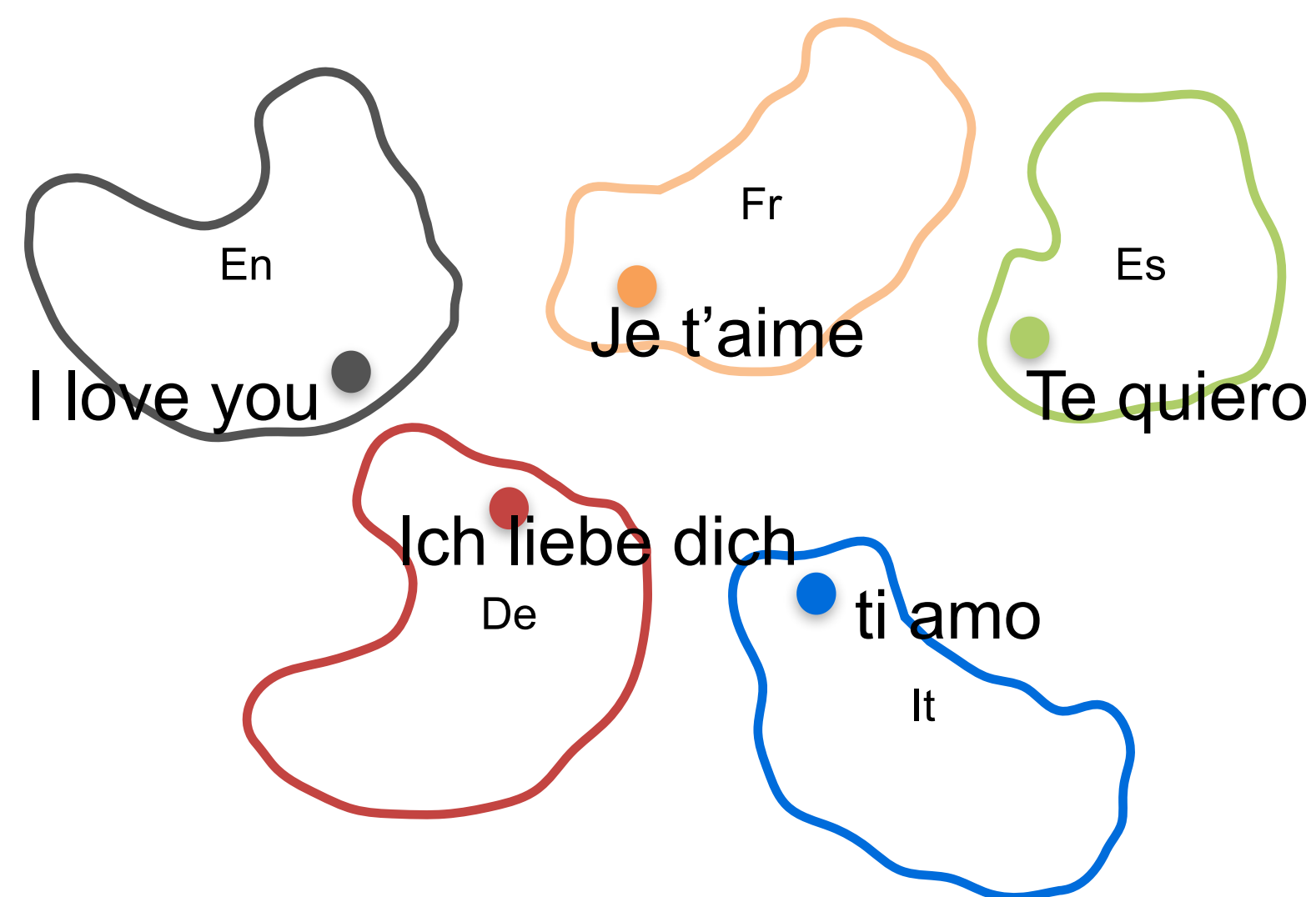
# Idea 1: Aligning Semantic Representations across Languages

- Key idea:
  - Words in different languages with the same meaning should have the same embedding, but the training objective does not necessarily encourage that!

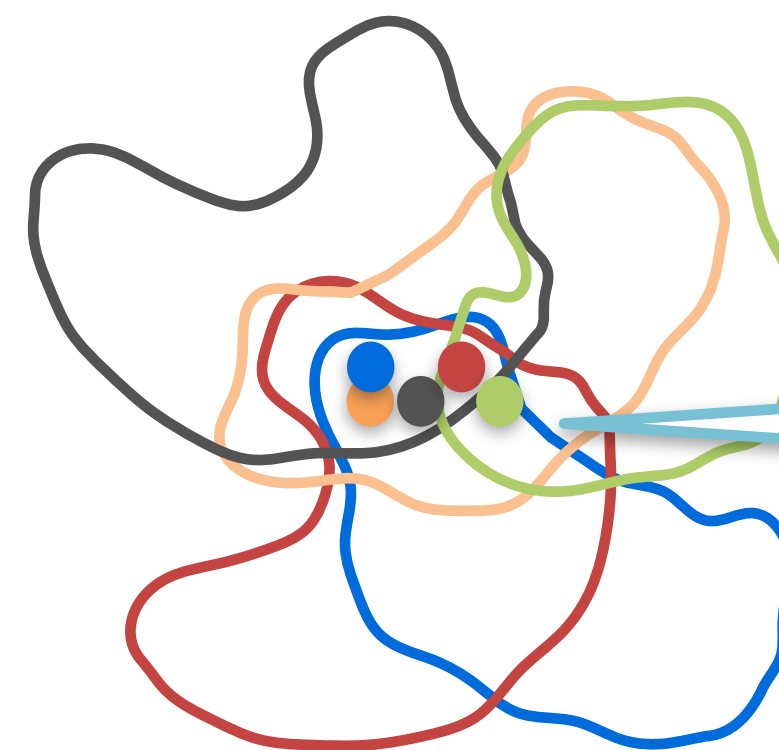


# Proposed mRASP: Aligning Semantic Representations across Languages

- Key idea:
  - Words in different languages with the same meaning should have the same embedding
  - Parallel sentences in different languages should have the same representation



ideally



<En> I **love** you.  
<Fr> Je t'**aime**.  
<De> Ich **liebe** dich.  
<Es> Te **quiero**.  
<It> ti **amo**.

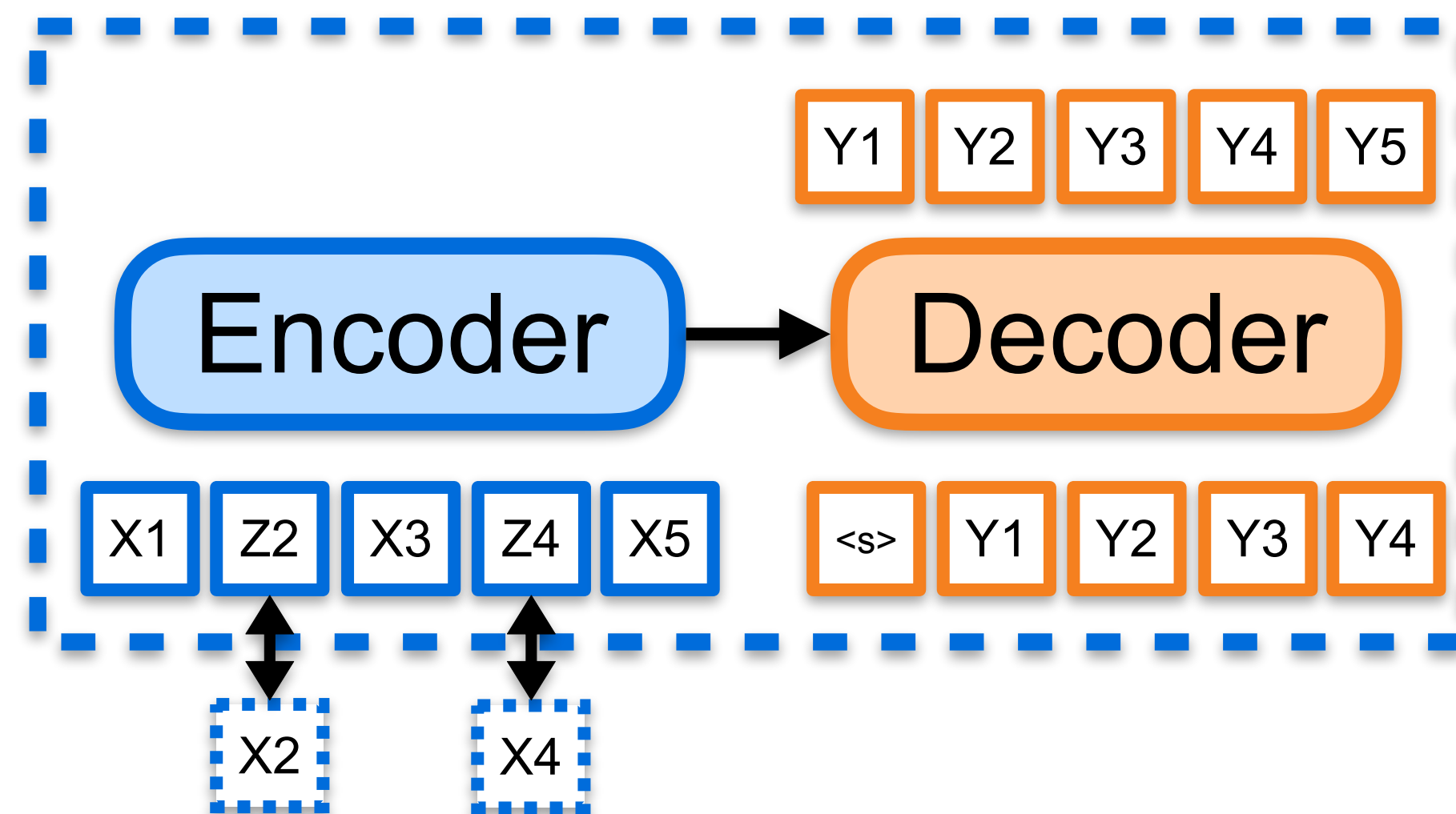
Pre-training Multilingual Neural Machine Translation by Leveraging Alignment Information

[Lin, Pan, Wang, Qiu, Feng, Zhou, **Lei Li**, EMNLP2020]

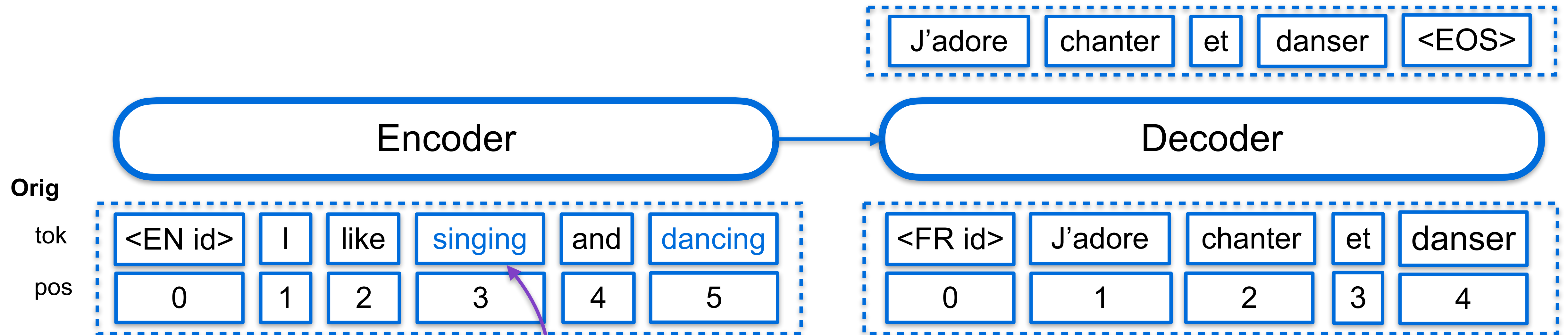
Contrastive Learning for Many-to-many Multilingual Neural Machine Translation [Pan, Wu, Wang, **Lei Li**, ACL 2021]<sup>27</sup>

# Aligning Semantic Representations across Languages

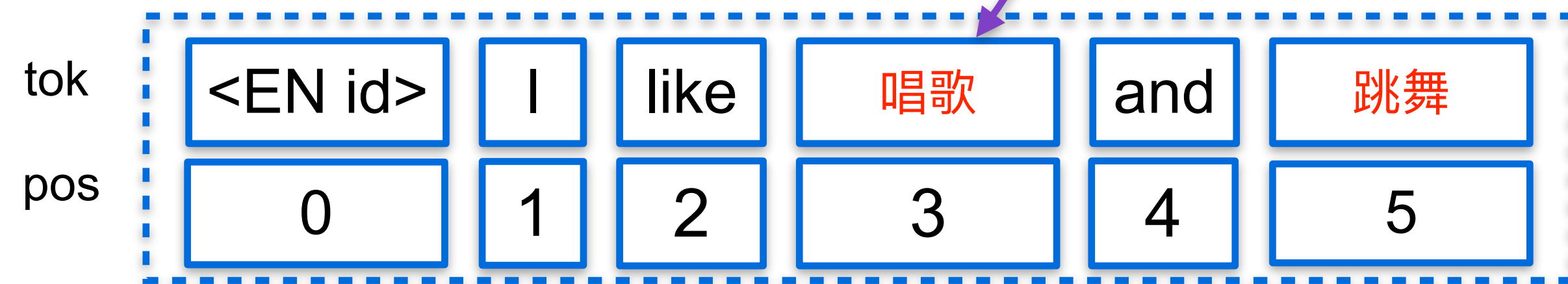
- **mRASP**: **m**ultilingual **R**andom **A**ligned **S**ubstitution **P**re-training
  - ▶ Multilingual Pre-training Approach
  - ▶ RAS: specially designed training method to align semantic embeddings



# mRASP: Random Aligned Substitution



## Random Aligned Substitution



$$\mathcal{L}_{RAS} = \sum_{i,j \in \mathcal{E}} \mathbb{E}_{(\mathbf{x}^i, \mathbf{x}^j) \sim \mathcal{D}_{i,j}} \left[ -\log P_{\theta}(\mathbf{x}^i | C(\mathbf{x}^j)) \right]$$

Randomly replace a source word to its synonym in different language.

# mRASP: Bringing Synonym Representations Closer

$$\mathcal{L}^{pre} = \sum_{i,j \in \mathcal{E}} \mathbb{E}_{(\mathbf{x}^i, \mathbf{x}^j) \sim \mathcal{D}_{i,j}} \left[ -\log P_{\theta} \left( \mathbf{x}^i \mid C(\mathbf{x}^j) \right) \right]$$

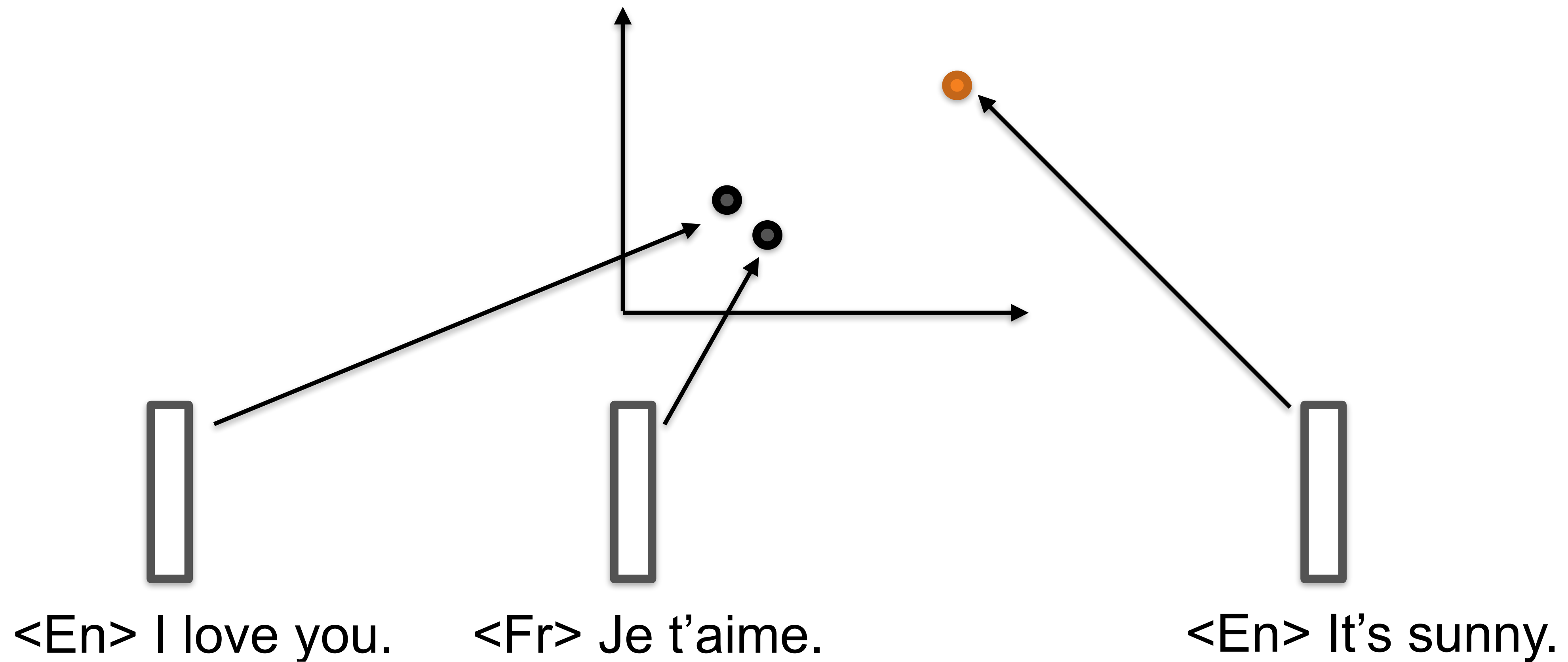
training with translation loss to bring closer



<EN id>	I	like	singing	and	dancing
0	1	2	3	4	5

<EN id>	I	like	chanter	and	danser
0	1	2	3	4	5

# Idea 2: Bring parallel sentence representations closer

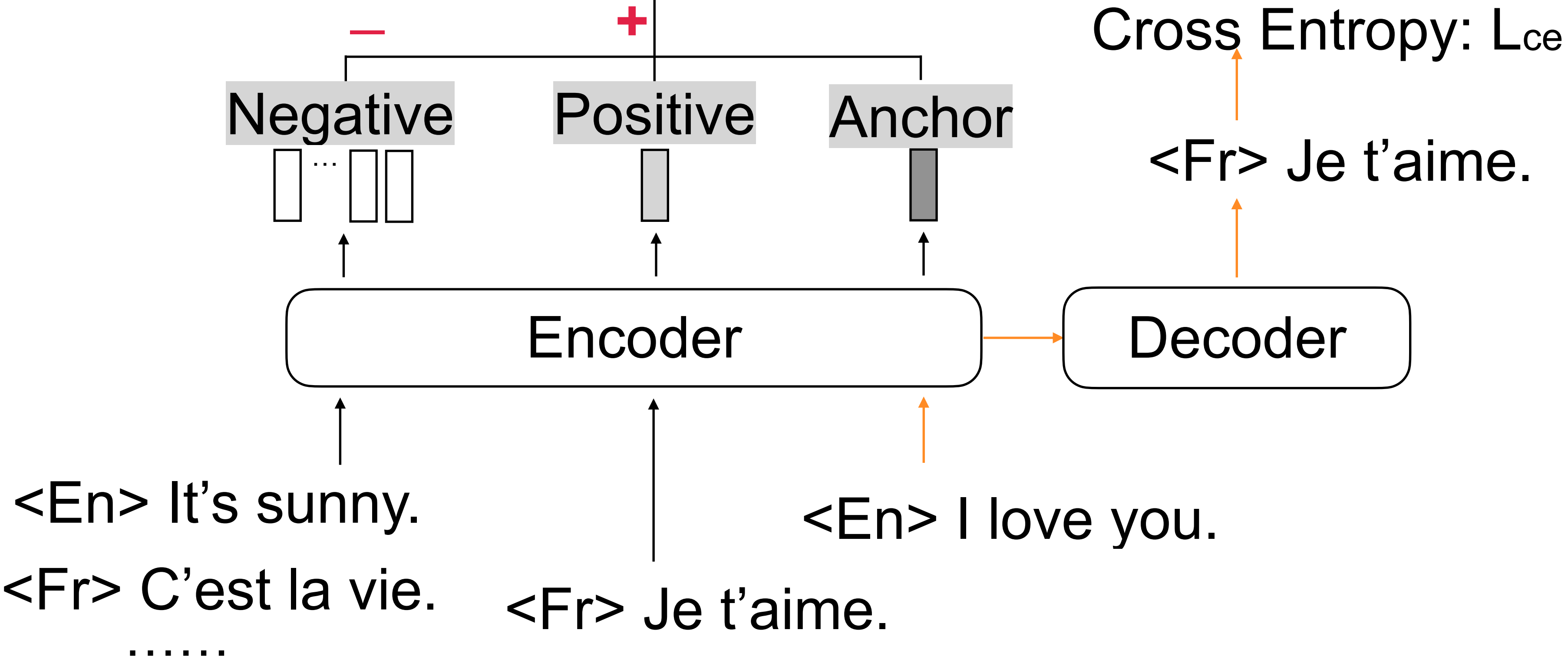


# mRASP2: Contrastive Learning to Bring Sentence Representations Closer

## Contrastive Loss: $\mathcal{L}_{ctr}$

$$\mathcal{L} = \mathcal{L}_{ce} + \lambda |s| \mathcal{L}_{ctr}$$

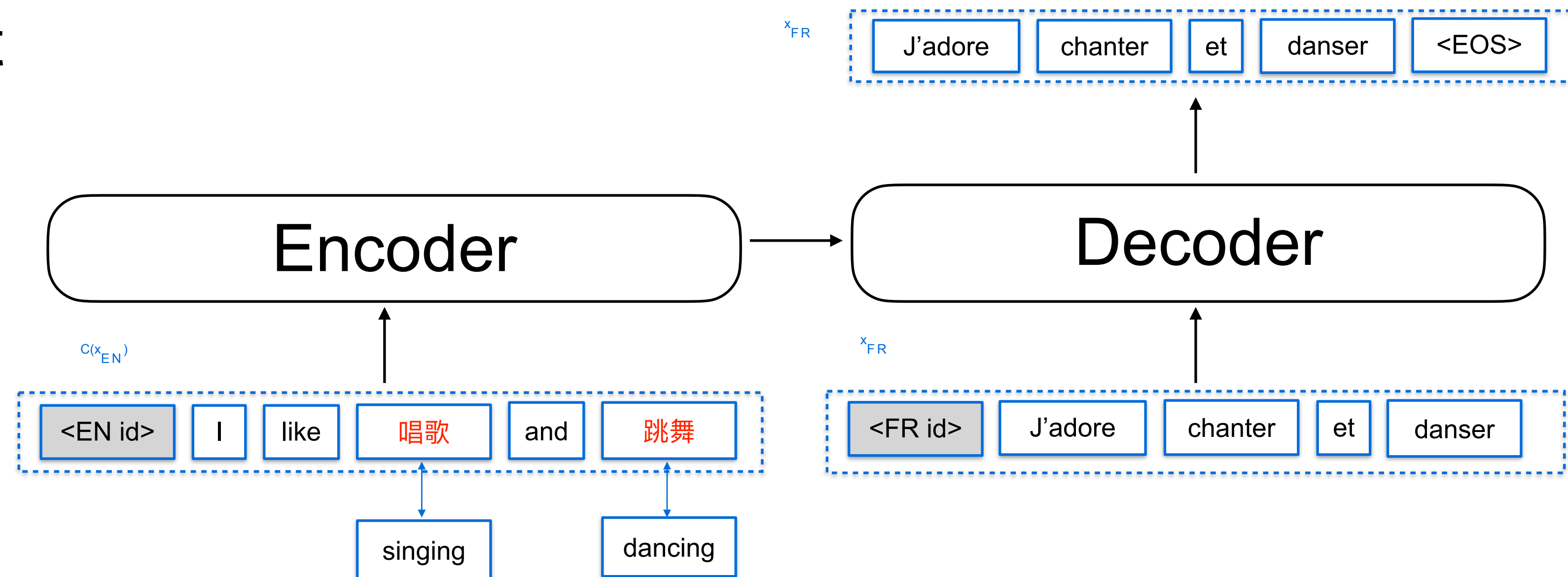
$$\mathcal{L}_{ctr} = - \sum_{\mathbf{x}^i, \mathbf{x}^j \in \mathcal{D}} \log \frac{e^{\text{sim}^+(\mathcal{R}(\mathbf{x}^i), \mathcal{R}(\mathbf{x}^j))/\tau}}{\sum_{\mathbf{y}^j} e^{\text{sim}^-(\mathcal{R}(\mathbf{x}^i), \mathcal{R}(\mathbf{y}^j))/\tau}}$$



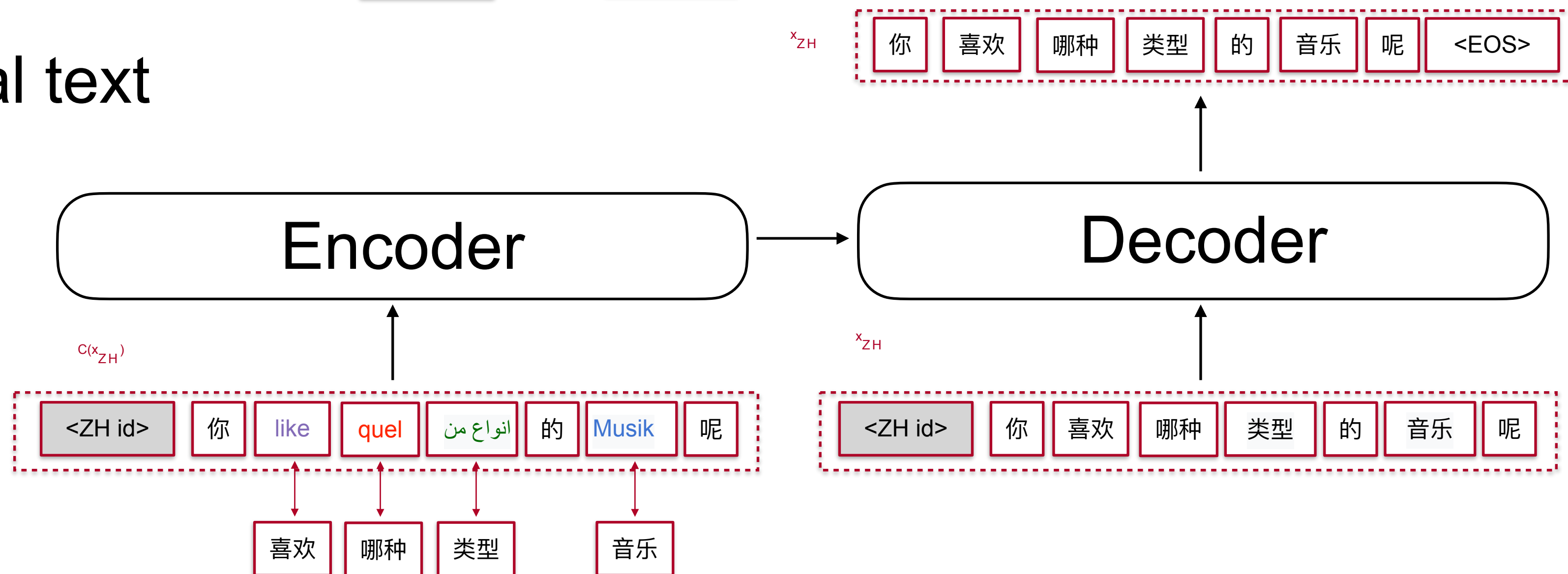


# mRASP2: Integrating Monolingual Data in Unified Training

- Parallel text

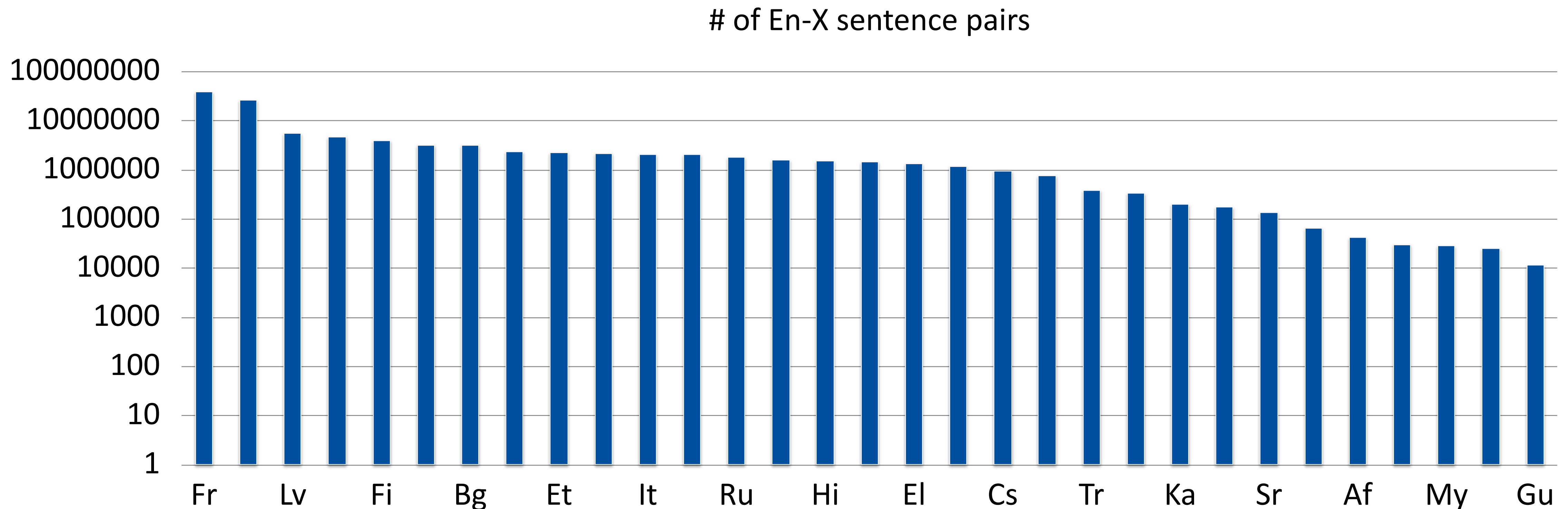


- Monolingual text



# Training Data for mRASP

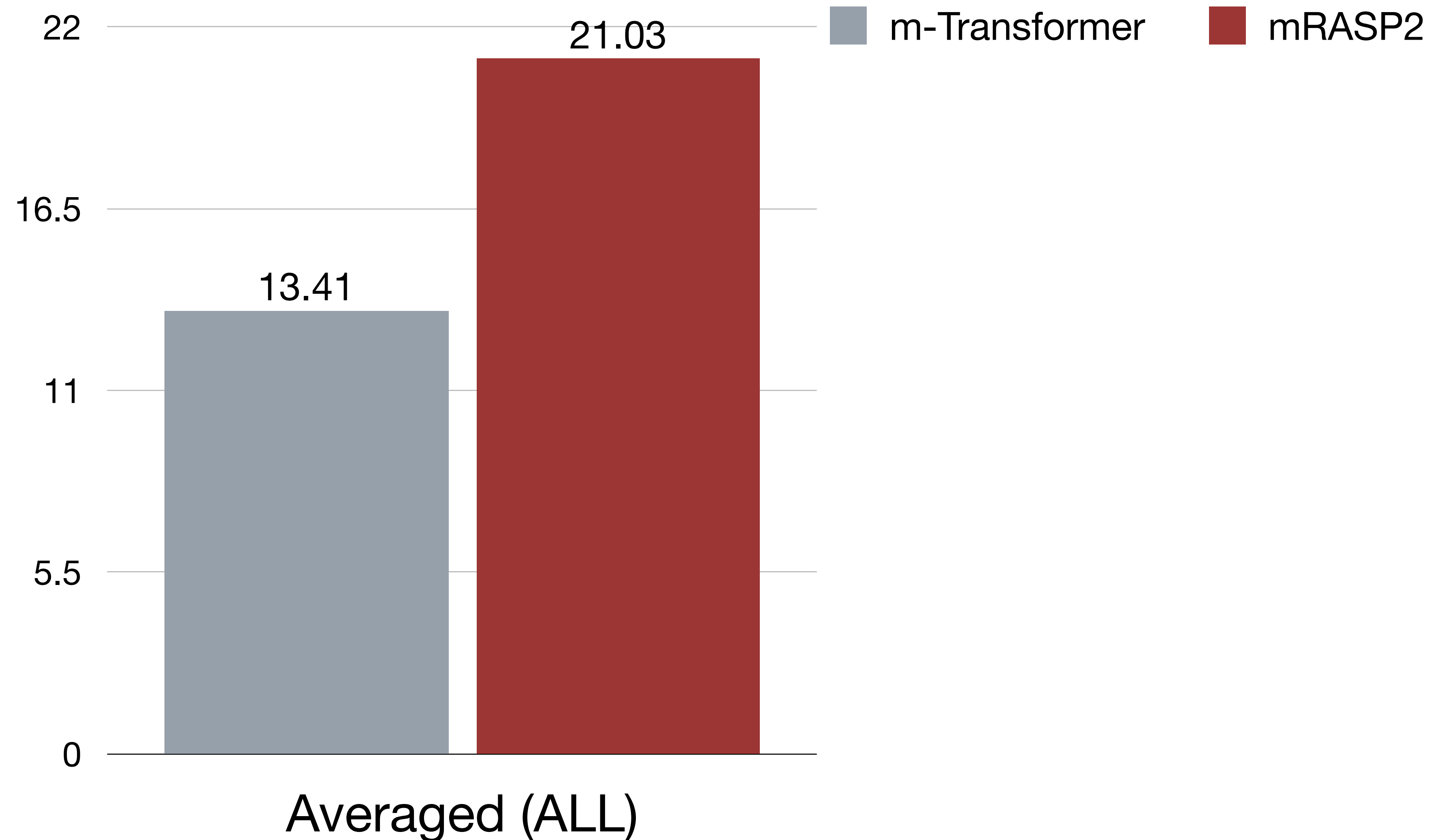
- Pre-training Dataset: PC32 (Parallel Corpus 32)
  - 32 English-centric language pairs, resulting in 64 directed translation pairs in total
  - Contains a total size of 110.4M public parallel sentence pairs



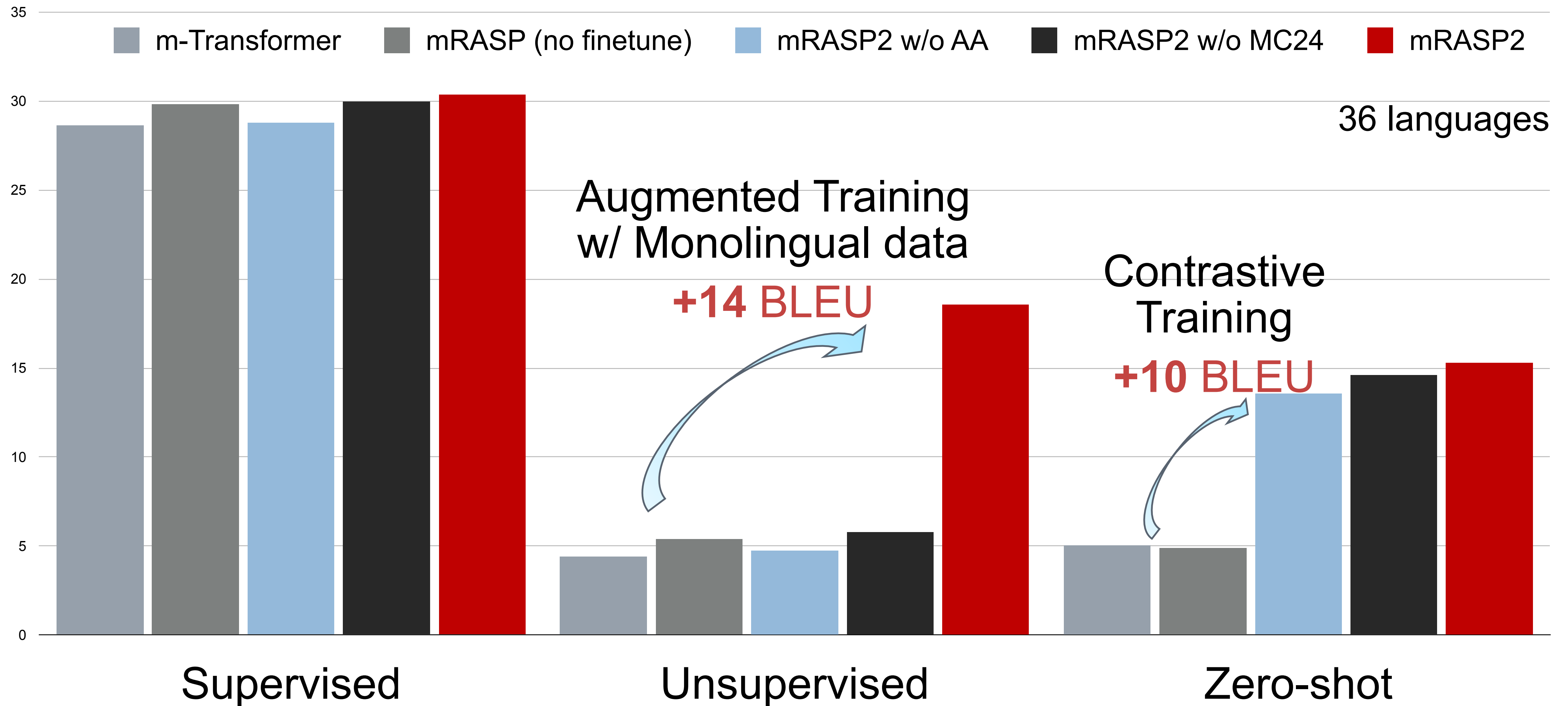
# mRASP2: a single MNMT model (no fine-tuning)

---

Overall Results in all scenarios: 56 directions



# mRASP significantly improves Zero-shot and Unsupervised Translation



Pre-training Multilingual Neural Machine Translation by Leveraging Alignment Information

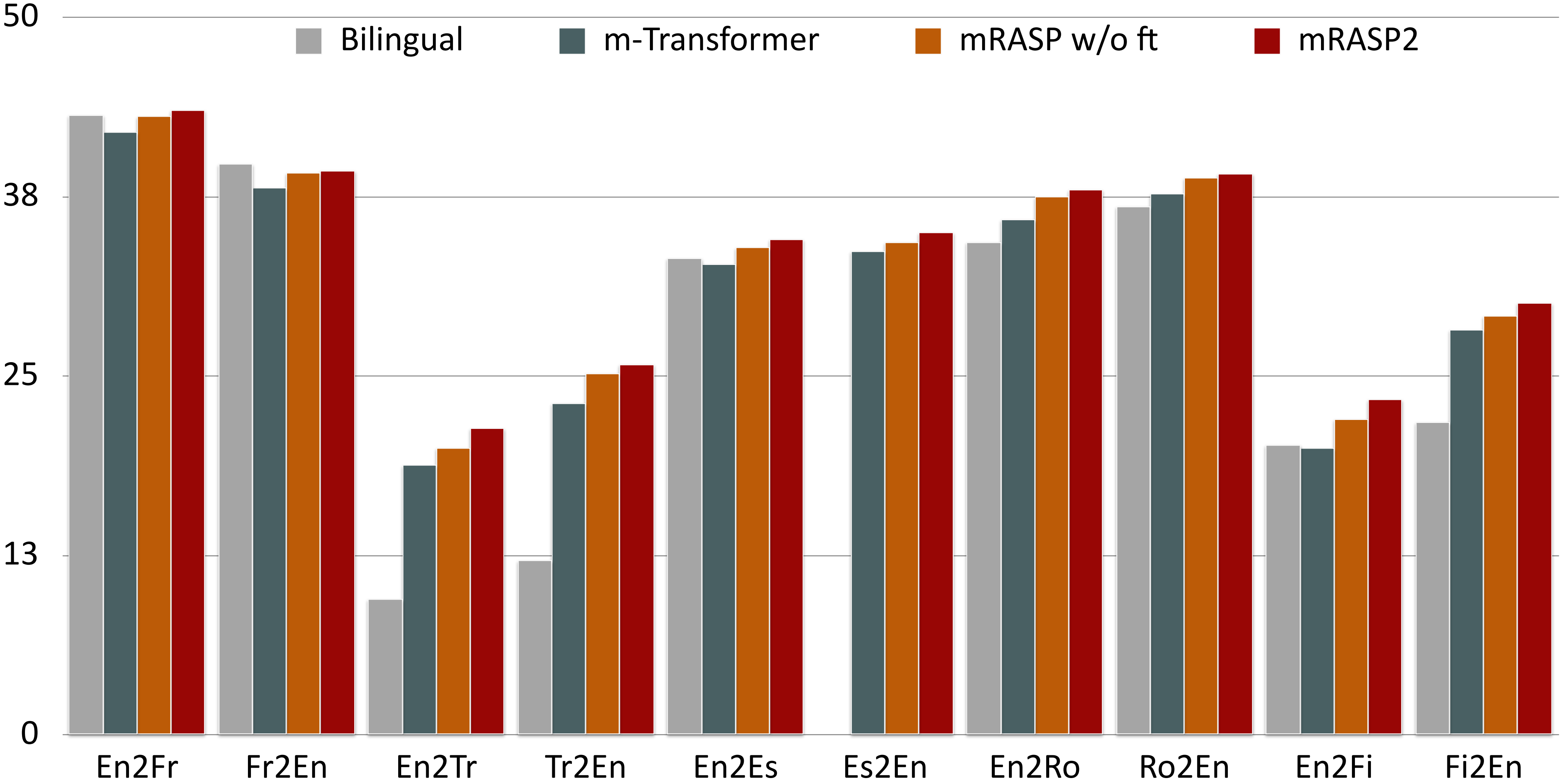
[Lin, Pan, Wang, Qiu, Feng, Zhou, **Lei Li**, EMNLP2020]

Contrastive Learning for Many-to-many Multilingual Neural Machine Translation [Pan, Wu, Wang, **Lei Li**, ACL 2021]<sup>36</sup>

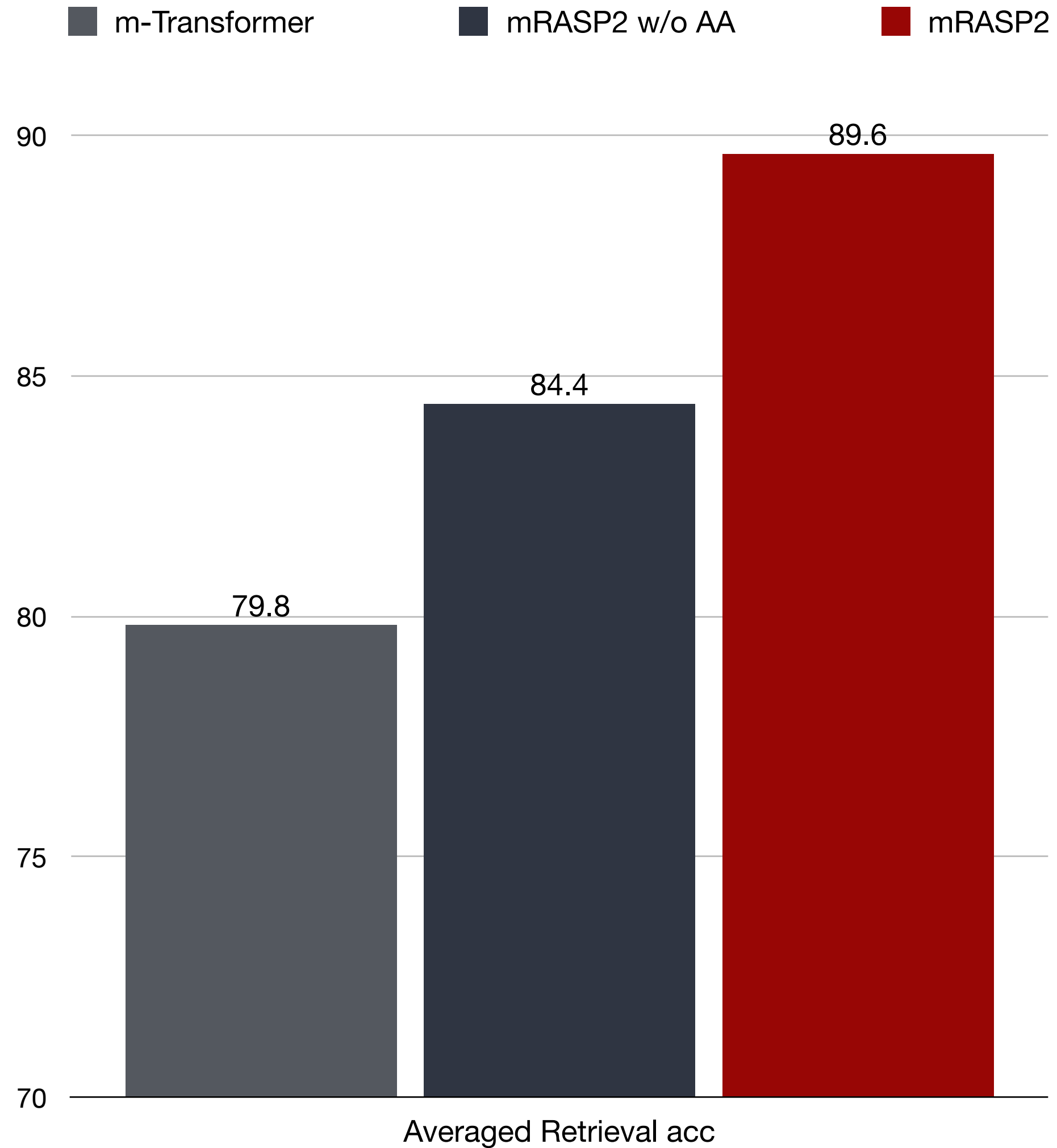


# mRASP2: Comparable or Better Performance on Supervised Directions

Tokenized BLEU on supervised directions



# Better Semantic Alignment: Sentence Retrieval

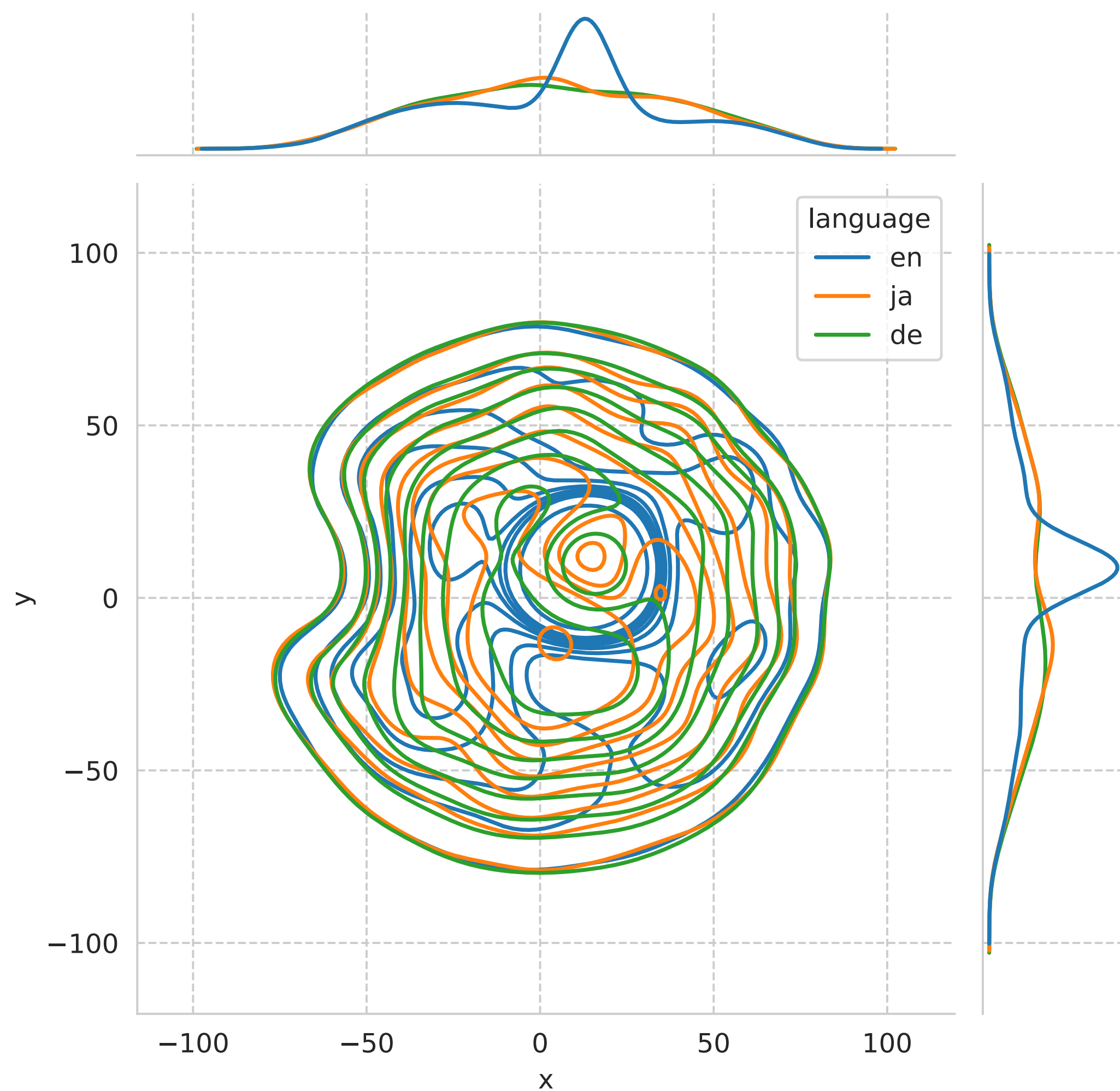


15-way parallel test set(Ted-M): 2284 samples

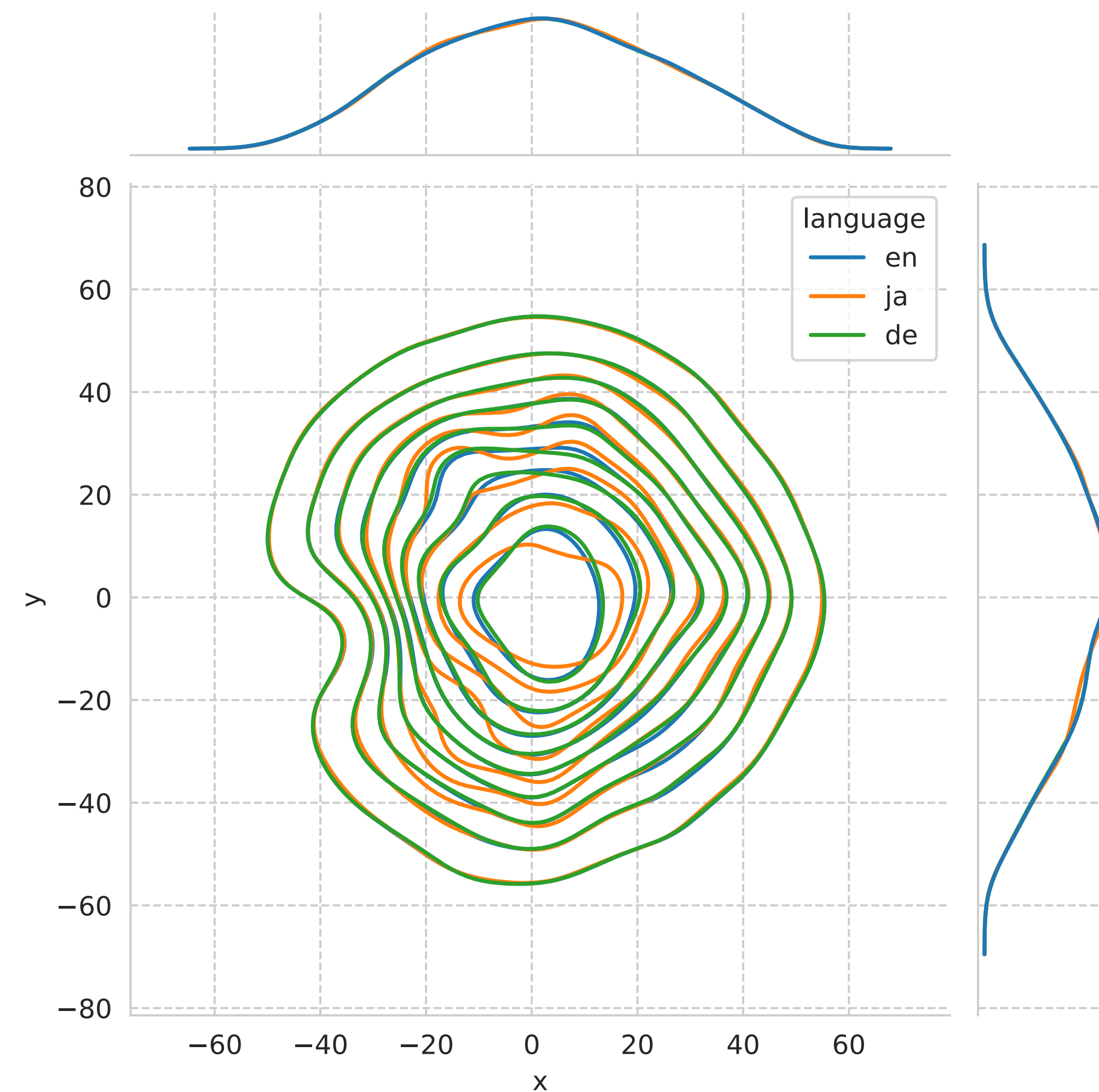
Contrastive Learning and Randomly Aligned Substitution both contribute to the improvement on sentence retrieval

# mRASP Produces Better Semantic Alignment

m-Transformer



mRASP



# mRASP Fine-tuning

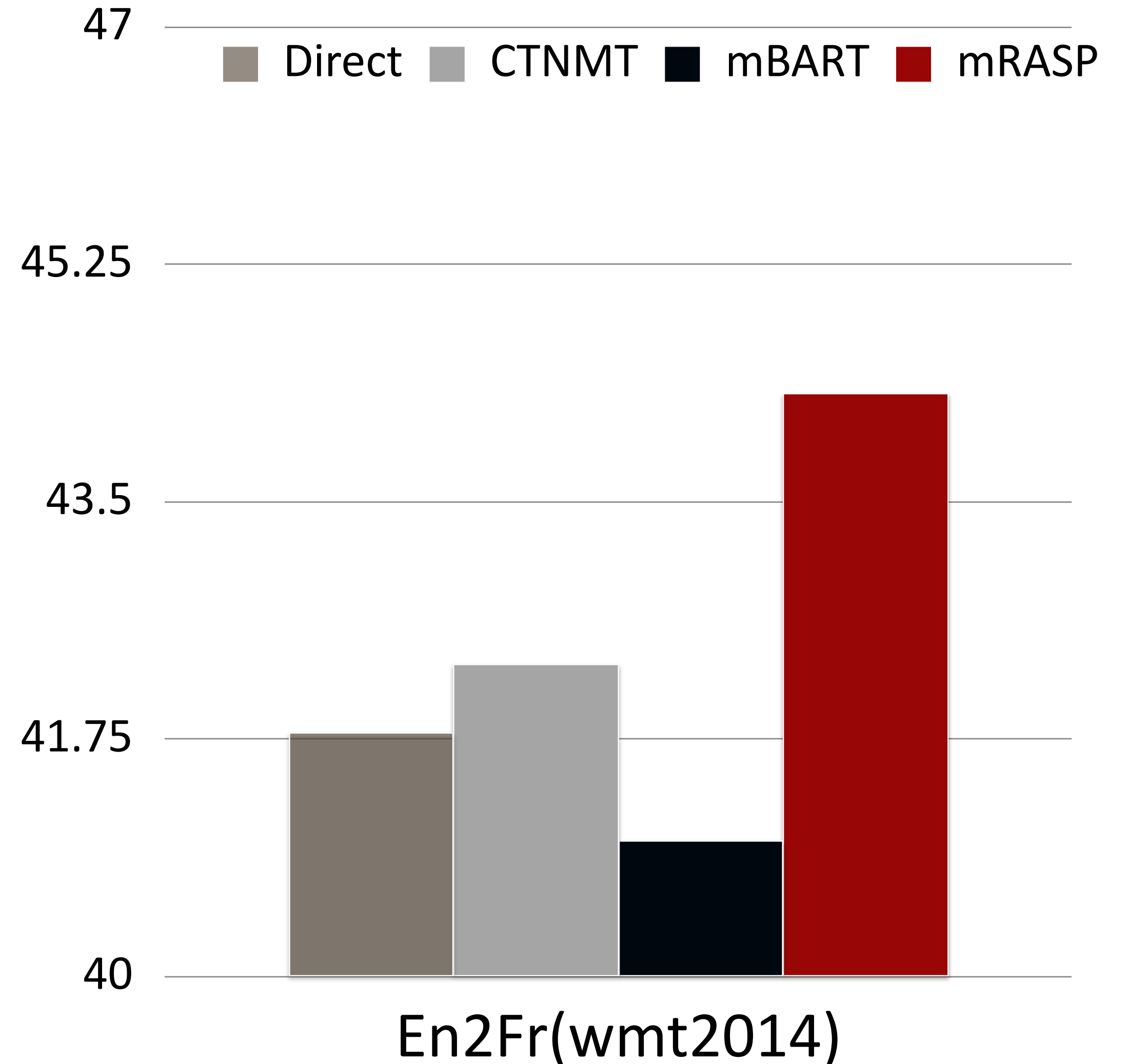
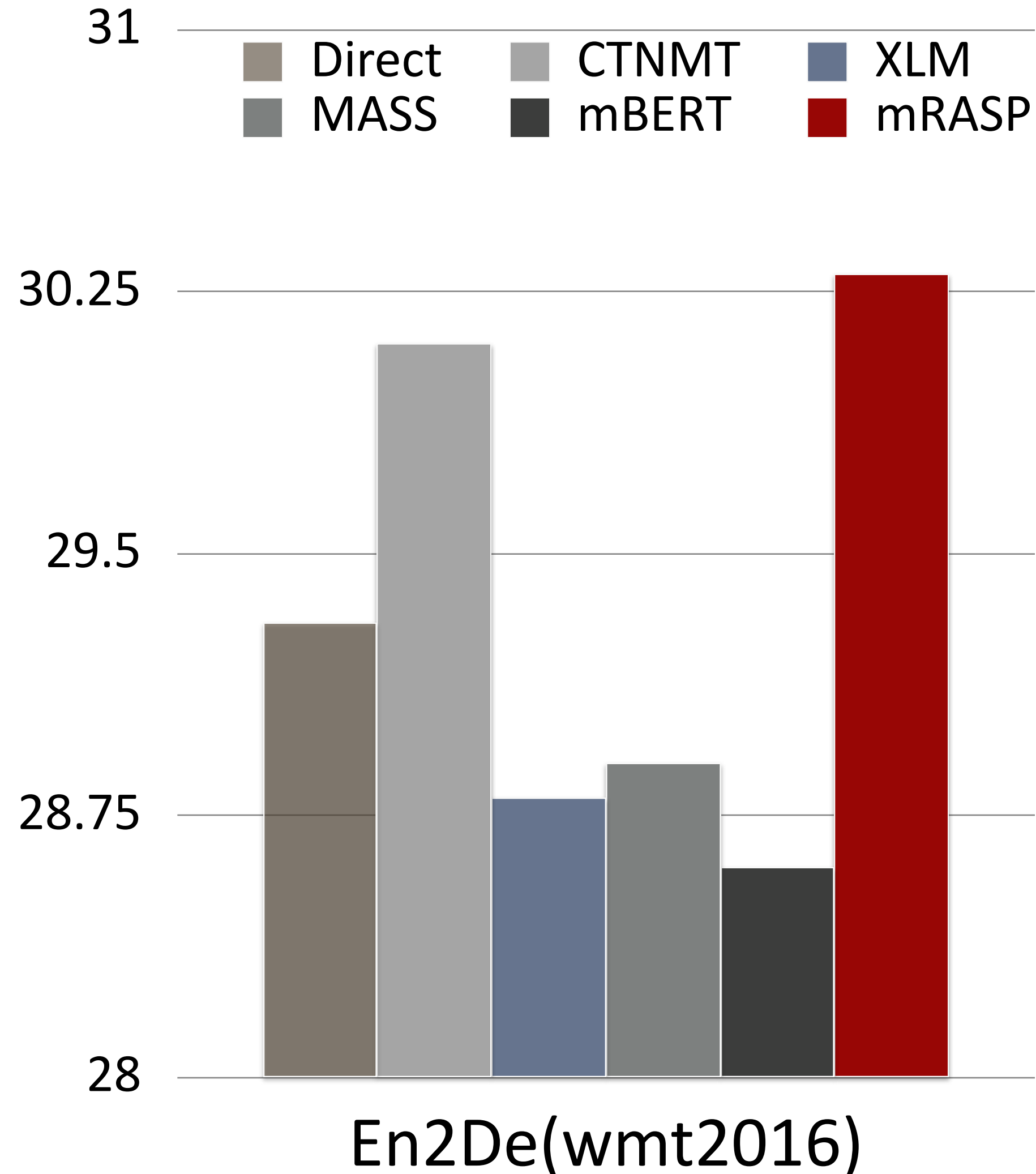
---

- Further fine-tuning based on mRASP model
- Fine-tuning Dataset
- Indigenous Corpus: included in pre-training phase
  - Extremely low resource (<100K) (Be, My, etc.)
  - Low resource(>100k and <1M) (He, Tr, etc.)
  - Medium resource (>1M and <10M) (De, Et, etc.)
  - Rich resource (>10M) (Zh, Fr, etc.)

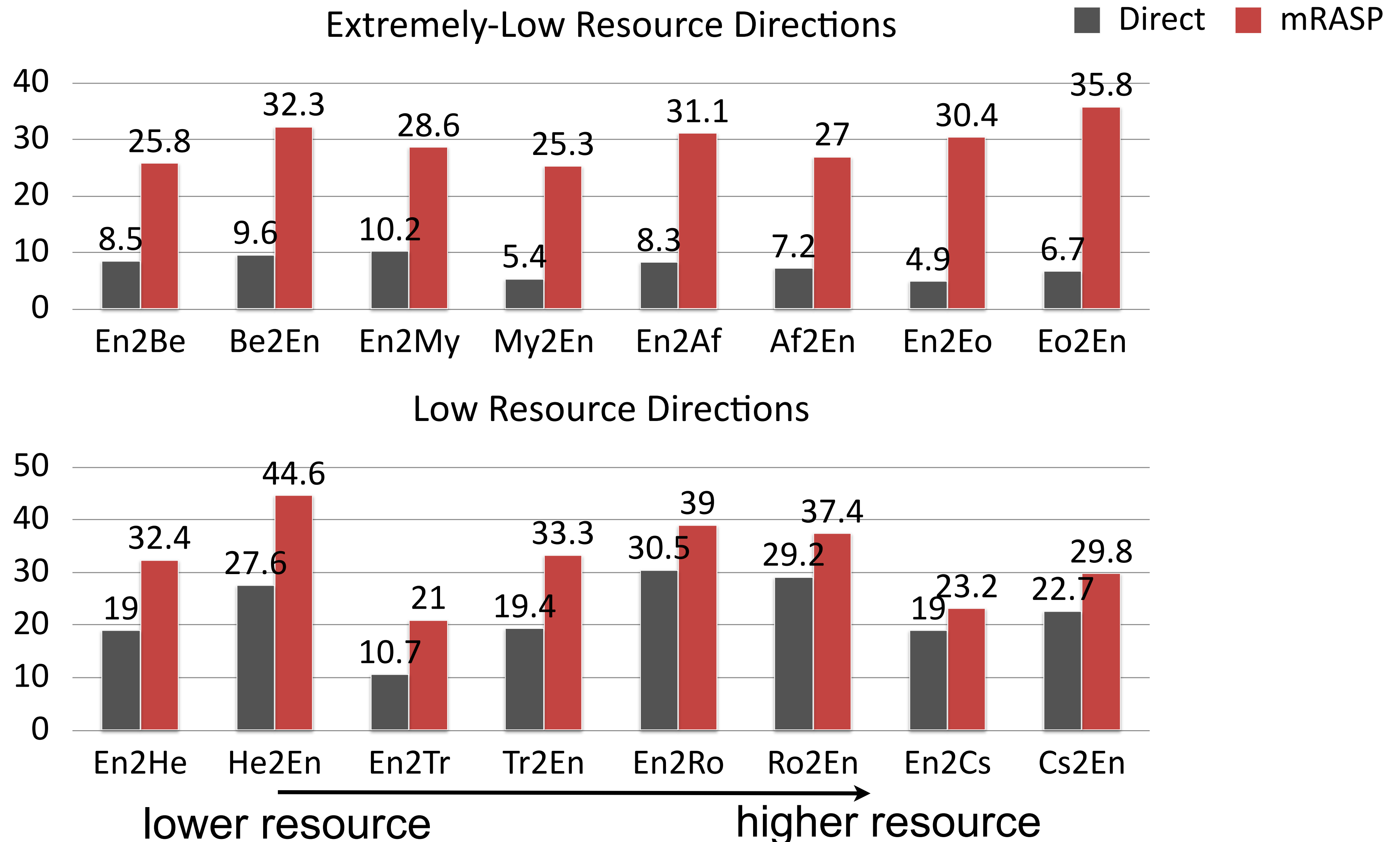


# mRASP Fine-tunes better: Rich resource works

- En->Fr +1.1BLEU.



# mRASP: Low resource works



# mRASP: Unseen languages

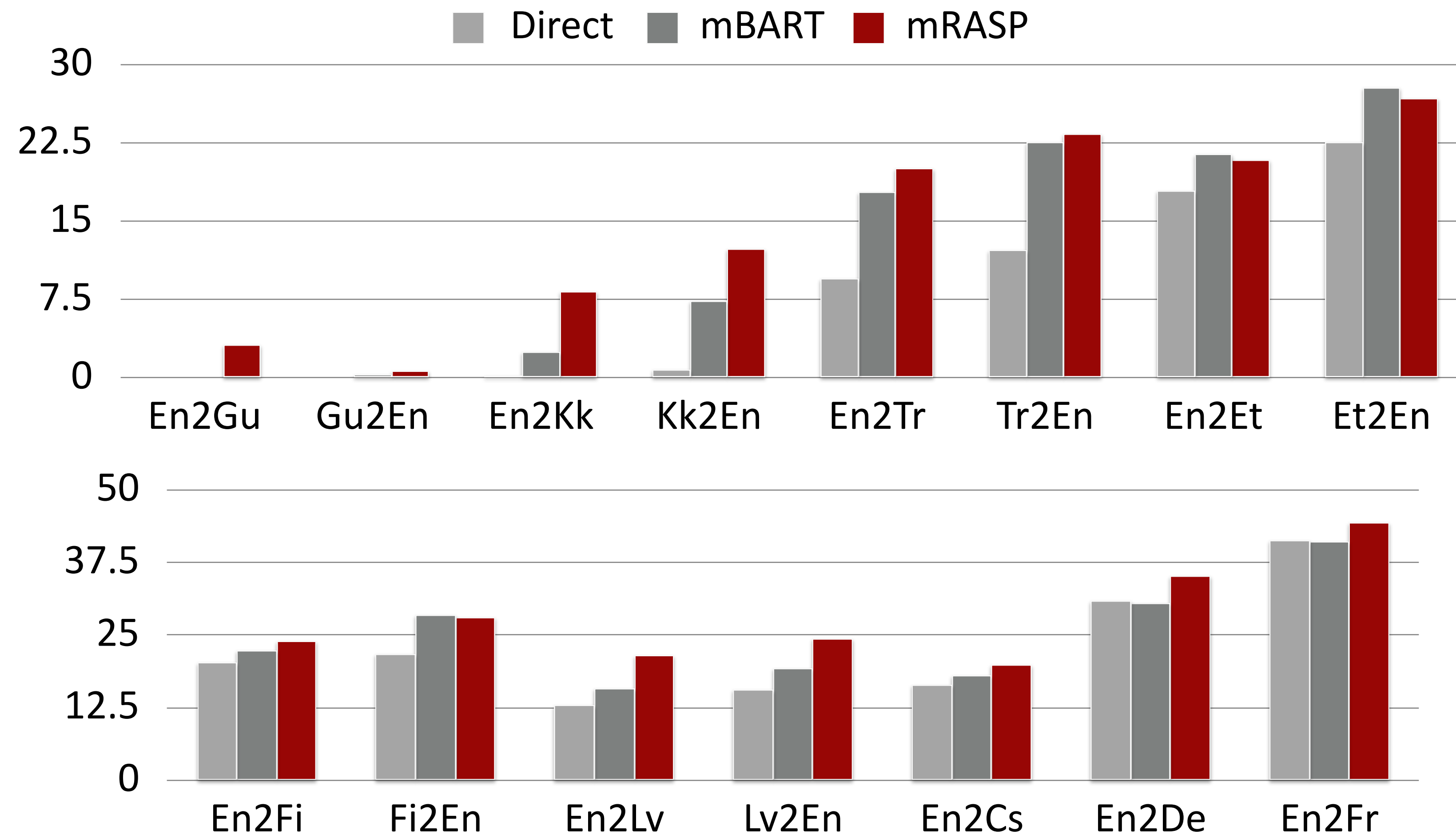
- mRASP generalizes on all exotic scenarios.

		Fr-Zh(20K)		De-Fr(9M)	
		→	←	→	←
Exotic Pair	Direct	0.7	3	23.5	21.2
	mRASP	25.8	26.7	29.9	23.4
		NI-Pt(12K)		Da-El(1.2M)	
		→	←	→	←
Exotic Full	Direct	0.0	0.0	14.1	16.9
	mRASP	14.1	13.2	17.6	19.9
		En-Mr(11k)		En-Gl(1.2M)	
		→	←	→	←
Exotic Source/ Target	Direct	6.4	6.8	8.9	12.8
	mRASP	22.7	22.9	32.1	38.1
		En-Eu(726k)		En-Sl(2M)	
		→	←	→	←
Exotic Source/ Target	Direct	7.1	10.9	24.2	28.2
	mRASP	19.1	28.4	27.6	29.5

12k: Direct not work **VS** mRASP achieves 10+ BLEU!!

# mRASP: Compare with other methods

- mRASP outperforms mBART for all but two language pairs.



# Summary

---

- Pre-training for NMT
  - sequence to sequence training objective
  - MASS: masked prediction using seq2seq
  - mBart: Recover original sentence from noised ones in multiple languages.
- Multilingual joint training
  - mRASP & mRASP2:
    - ▶ augmenting data with randomly substitute of words from bilingual lexicon
    - ▶ monolingual reconstruction
    - ▶ contrastive learning

# Discussion

---

- What strategies for training multilingual NMT