

It's Easier to Translate *out of* English than *into* it: Measuring Neural Translation Difficulty by Cross-Mutual Information

ACL 2020

Emanuele Bugliarello, Sabrina J. Mielke,
Antonios Anastasopoulos, Ryan Cotterell,
Naoaki Okazaki

UNIVERSITY OF COPENHAGEN





Evaluation Matrix

Translation quality of best system for test set

using metric

[Translations](#)
[Resources](#)
[Download](#)
[Info](#)
[Account](#)

		output language									
		Czech	German	English	Finnish	French	Gujarati	Kazakh	Lithuanian	Russian	Chinese
input language	Czech	19.3									
	German	20.1	42.8		37.3						
	English	29.9	44.9	27.4		28.2	11.1	20.1	36.3	44.6	
	Finnish		33.0								
	French		35.0								
	Gujarati			24.9							
	Kazakh			30.5							
	Lithuanian			36.3							
	Russian			40.2							
	Chinese			39.9							



Evaluation Matrix

Translation quality of best system for test set using metric

[Translations](#)

[Resources](#)

[Download](#)

[Info](#)

[Account](#)

		output language									
		Czech	German	English	Finnish	French	Gujarati	Kazakh	Lithuanian	Russian	Chinese
input language	Czech	19.3									
	German	20.1	42.8		37.3						
	English	29.9	44.9	27.4		29.2	11.1	29.1	26.2	44.6	
	Finnish			33.0							
	French		35.0								
	Gujarati			24.9							
	Kazakh			30.5							
	Lithuanian			36.3							
	Russian			40.2							
	Chinese			39.9							

Is **fi-en** easier than **en-fi**?



Evaluation Matrix

Translation quality of best system for test set using metric

[Translations](#)

[Resources](#)

[Download](#)

[Info](#)

[Account](#)

		output language									
		Czech	German	English	Finnish	French	Italian	Kazakh	Lithuanian	Russian	Chinese
input language	Czech	19.3									
	German	20.1	42.8		37.3						
	English	29.9	44.9	27.4		38.2	31.1	29.1	26.2	41.6	
	Finnish			33.0							
	French		35.0								
	Italian			24.9							
	Kazakh			30.5							
	Lithuanian			36.3							
	Russian			40.2							
	Chinese			39.9							

Is **fi-en** easier than **en-fi**?

We can't tell based on BLEU!

BLEU's shortcomings for cross-linguistic comparisons

BLEU's shortcomings for cross-linguistic comparisons

BLEU is a precision-based metric

BLEU's shortcomings for cross-linguistic comparisons

BLEU is a precision-based metric

1. BLEU depends on *tokenization* and the *notion of “word”*!

BLEU's shortcomings for cross-linguistic comparisons

BLEU is a precision-based metric

1. BLEU depends on *tokenization* and the *notion of "word"*!

Example:

“I will have been programming” English

“Programlayacağım” Turkish

BLEU's shortcomings for cross-linguistic comparisons

BLEU is a precision-based metric

1. BLEU depends on *tokenization* and the *notion of "word"*!

Example:

“I will have been programming” English

“Programlayacağım” Turkish

→ More partial credit for English!

BLEU's shortcomings for cross-linguistic comparisons

BLEU is a precision-based metric

1. BLEU depends on *tokenization* and the *notion of "word"*!

Example:

“I will have been programming” English

“Programlayacağım” Turkish

→ More partial credit for English!

Remedy: Look at the likelihood

BLEU's shortcomings for cross-linguistic comparisons

BLEU is a precision-based metric

1. BLEU depends on *tokenization* and the *notion of "word"*!

Example:

“I will have been programming” English

“Programlayacağım” Turkish

→ More partial credit for English!

Remedy: Look at the likelihood

2. We are still measuring: difficulty of **translation** *and* **generation**

Mutual Information expresses the act of translation

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$$H(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t | s))]$$

Mutual Information expresses the act of translation

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$$H(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t | s))]$$

$\underbrace{H(T)}$
uncertainty about T
a priori

Mutual Information expresses the act of translation

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$$H(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t | s))]$$

$H(T)$

uncertainty about T
a priori

$H(T | S)$

uncertainty about T
after knowing S

Mutual Information expresses the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$$H(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t | s))]$$

$$\underbrace{H(T)}_{\substack{\text{uncertainty about } T \\ \textit{a priori}}} - \underbrace{H(T | S)}_{\substack{\text{uncertainty about } T \\ \textit{after knowing } S}}$$

how much knowing S **reduced uncertainty** about T

Mutual Information expresses the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$H(T|S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t|s))]$

$$\underbrace{\text{MI}(S; T)}_{\substack{\text{mutual information} \\ \text{between } S \text{ and } T}} = \underbrace{H(T)}_{\substack{\text{uncertainty about } T \\ \text{a priori}}} - \underbrace{H(T|S)}_{\substack{\text{uncertainty about } T \\ \text{after knowing } S}}$$

how much knowing S **reduced uncertainty** about T

Mutual Information expresses the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$H(T|S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t|s))]$

$$\underbrace{\text{MI}(S; T)}_{\substack{\text{mutual information} \\ \text{between } S \text{ and } T}} = \underbrace{H(T)}_{\substack{\text{uncertainty about } T \\ \text{a priori}}} - \underbrace{H(T|S)}_{\substack{\text{uncertainty about } T \\ \text{after knowing } S}}$$

how much knowing S **reduced uncertainty** about T

symmetric!
assuming all entropies w.r.t. same joint $p(S, T)$

Mutual Information expresses the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$H(T|S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t|s))]$

$$\underbrace{\text{MI}(S; T)}_{\substack{\text{mutual information} \\ \text{between } S \text{ and } T}} = \underbrace{H(T)}_{\substack{\text{uncertainty about } T \\ \textit{a priori}}} - \underbrace{H(T|S)}_{\substack{\text{uncertainty about } T \\ \textit{after knowing } S}}$$

how much knowing S **reduced uncertainty** about T

symmetric!
assuming all
entropies w.r.t.
same joint $p(S, T)$

Example: en-zh

Mutual Information expresses the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$H(T|S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t|s))]$

$$\underbrace{MI(S; T)}_{\text{mutual information between } S \text{ and } T} = \underbrace{H(T)}_{\text{uncertainty about } T \text{ a priori}} - \underbrace{H(T|S)}_{\text{uncertainty about } T \text{ after knowing } S}$$

how much knowing S **reduced uncertainty** about T

symmetric!
assuming all entropies w.r.t. same joint $p(S, T)$

Example: en-zh

H(谢谢)

uncertainty about “谢谢”

Mutual Information expresses the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$H(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t | s))]$

$$\underbrace{MI(S; T)}_{\text{mutual information between } S \text{ and } T} = \underbrace{H(T)}_{\text{uncertainty about } T \text{ a priori}} - \underbrace{H(T | S)}_{\text{uncertainty about } T \text{ after knowing } S}$$

how much knowing S **reduced uncertainty** about T

symmetric!
assuming all entropies w.r.t. same joint $p(S, T)$

Example: en-zh

H(谢谢)

uncertainty about “谢谢”

H(谢谢 | Thanks)

uncertainty about “谢谢” after knowing its translation

Mutual Information expresses the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$H(T|S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t|s))]$

$$\underbrace{MI(S; T)}_{\text{mutual information between } S \text{ and } T} = \underbrace{H(T)}_{\text{uncertainty about } T \text{ a priori}} - \underbrace{H(T|S)}_{\text{uncertainty about } T \text{ after knowing } S}$$

how much knowing S **reduced uncertainty** about T

symmetric!
assuming all entropies w.r.t. same joint $p(S, T)$

Example: en-zh



uncertainty about “谢谢”

uncertainty about “谢谢” after knowing its translation

how much easier it has become to predict “谢谢”

Cross-Mutual Information measures models' performance on the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$$H(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t | s))]$$

$$\underbrace{MI(S; T)}_{\text{mutual information between } S \text{ and } T} = \underbrace{H(T)}_{\text{uncertainty about } T \text{ a priori}} - \underbrace{H(T | S)}_{\text{uncertainty about } T \text{ after knowing } S}$$

how much knowing S reduced uncertainty about T

Cross-Mutual Information measures models' performance on the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$$H(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t | s))]$$

$$\underbrace{MI(S; T)}_{\text{mutual information between } S \text{ and } T} = \underbrace{H(T)}_{\text{uncertainty about } T \text{ a priori}} - \underbrace{H(T | S)}_{\text{uncertainty about } T \text{ after knowing } S}$$

how much knowing S reduced uncertainty about T

Cross-Entropy: $H_q(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(q(t))]$ how surprised is model q in reality p ?

$$H_q(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(q(t | s))]$$

Cross-Mutual Information measures models' performance on the *act of translation*

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$$H(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t | s))]$$

$$\underbrace{MI(S; T)}_{\text{mutual information between } S \text{ and } T} = \underbrace{H(T)}_{\text{uncertainty about } T \text{ a priori}} - \underbrace{H(T | S)}_{\text{uncertainty about } T \text{ after knowing } S}$$

how much knowing S reduced uncertainty about T

Cross-Entropy: $H_q(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(q(t))]$ how surprised is model q in reality p ?

$$H_q(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(q(t | s))]$$

$$XMI(S \rightarrow T) := H_{q_{LM}}(T) - H_{q_{MT}}(T | S)$$

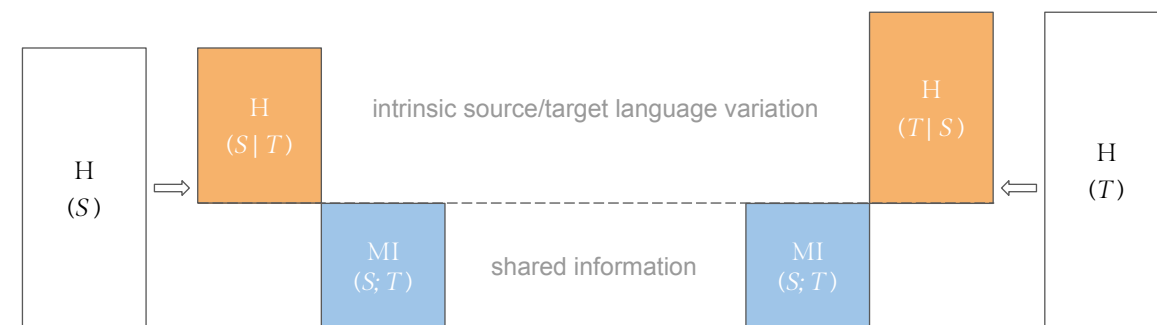
Cross-Mutual Information measures models' performance on the act of translation

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$$H(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t | s))]$$

$$\underbrace{\text{MI}(S; T)}_{\substack{\text{mutual information} \\ \text{between } S \text{ and } T}} = \underbrace{H(T)}_{\substack{\text{uncertainty about } T \\ \text{a priori}}} - \underbrace{H(T | S)}_{\substack{\text{uncertainty about } T \\ \text{after knowing } S}}$$

how much knowing S reduced uncertainty about T



Cross-Entropy: $H_q(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(q(t))]$ how surprised is model q in reality p ?

$$H_q(T | S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(q(t | s))]$$

$$\text{XMI}(S \rightarrow T) := H_{q_{LM}}(T) - H_{q_{MT}}(T | S)$$

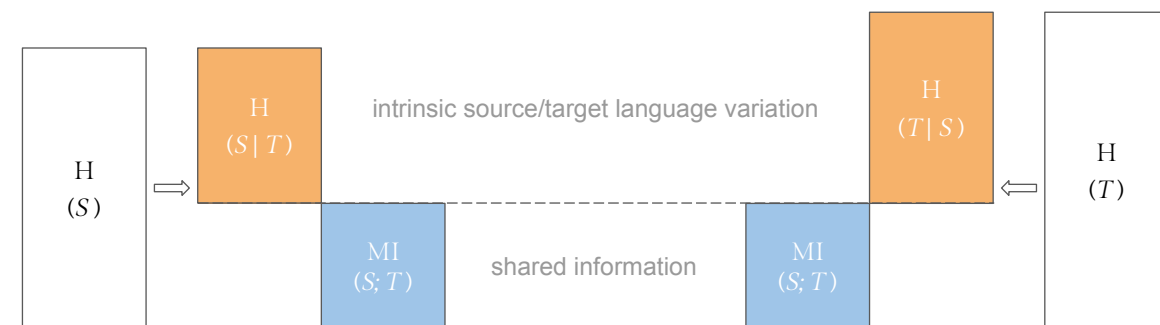
Cross-Mutual Information measures models' performance on the act of translation

Entropy: $H(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(p(t))]$ uncertainty

$$H(T|S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(p(t|s))]$$

$$\underbrace{\text{MI}(S; T)}_{\text{mutual information between } S \text{ and } T} = \underbrace{H(T)}_{\text{uncertainty about } T \text{ a priori}} - \underbrace{H(T|S)}_{\text{uncertainty about } T \text{ after knowing } S}$$

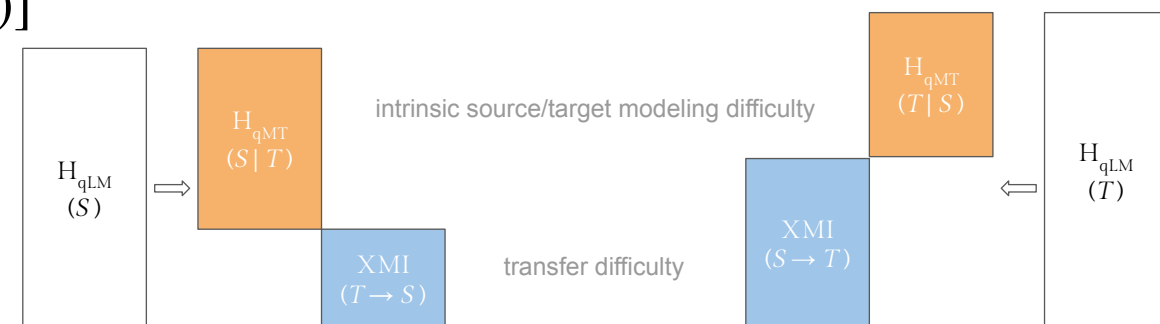
how much knowing S reduced uncertainty about T



Cross-Entropy: $H_q(T) = \mathbb{E}_{t \sim p(T)}[-\log_2(q(t))]$ how surprised is model q in reality p ?

$$H_q(T|S) = \mathbb{E}_{(s,t) \sim p(S,T)}[-\log_2(q(t|s))]$$

$$\text{XMI}(S \rightarrow T) := H_{q_{LM}}(T) - H_{q_{MT}}(T|S)$$



Experiments

Experiments

Setup

- *Data*: Fully 21-parallel subset of Europarl
- *Models*:
 - 20 [◦ → en] Transformers
 - 20 [en → ◦] Transformers

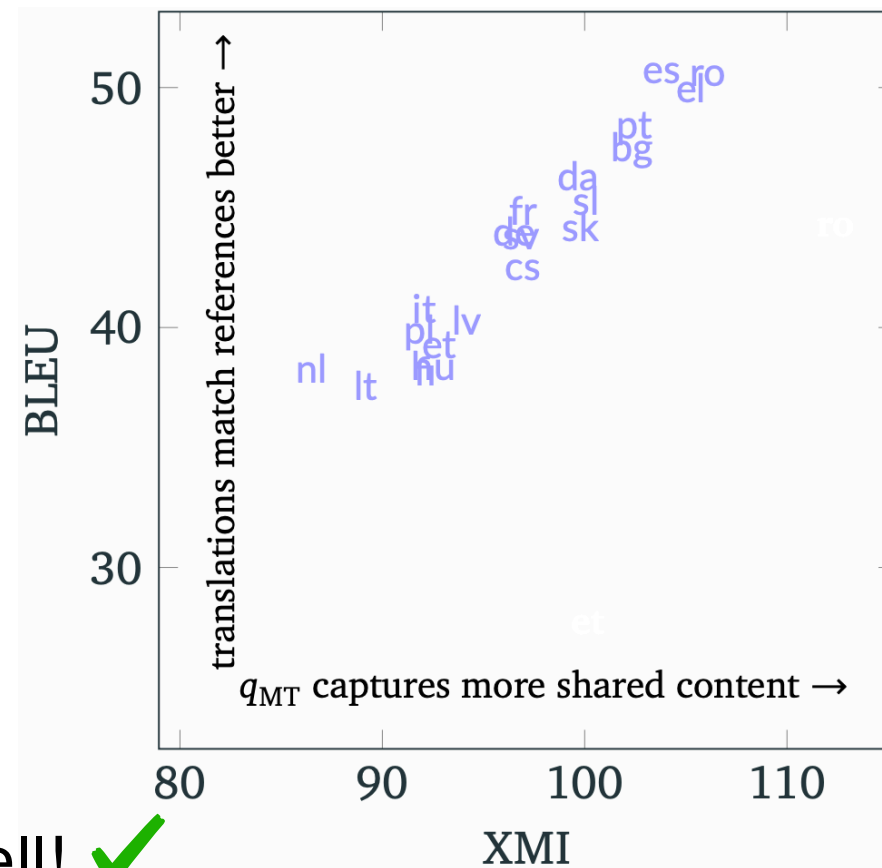
Experiments

Setup

- *Data*: Fully 21-parallel subset of Europarl
- *Models*:
 - 20 [$\circ \rightarrow \text{en}$] Transformers
 - 20 [$\text{en} \rightarrow \circ$] Transformers

Results

- For fixed target, BLEU and XMI correlate well! ✓



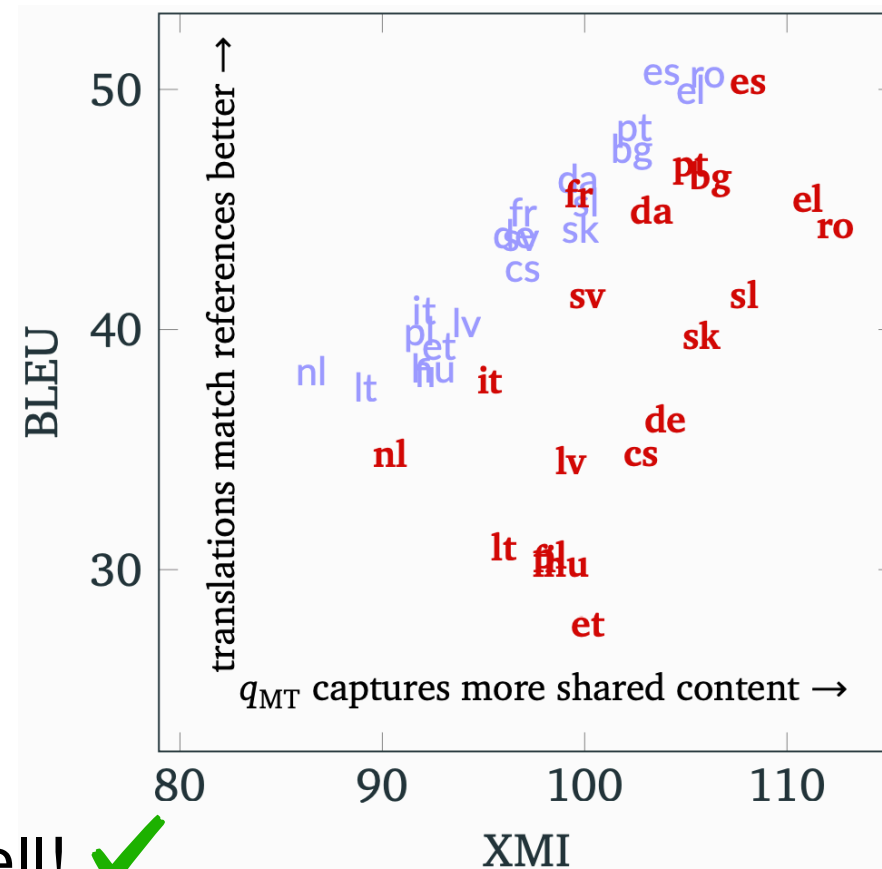
Experiments

Setup

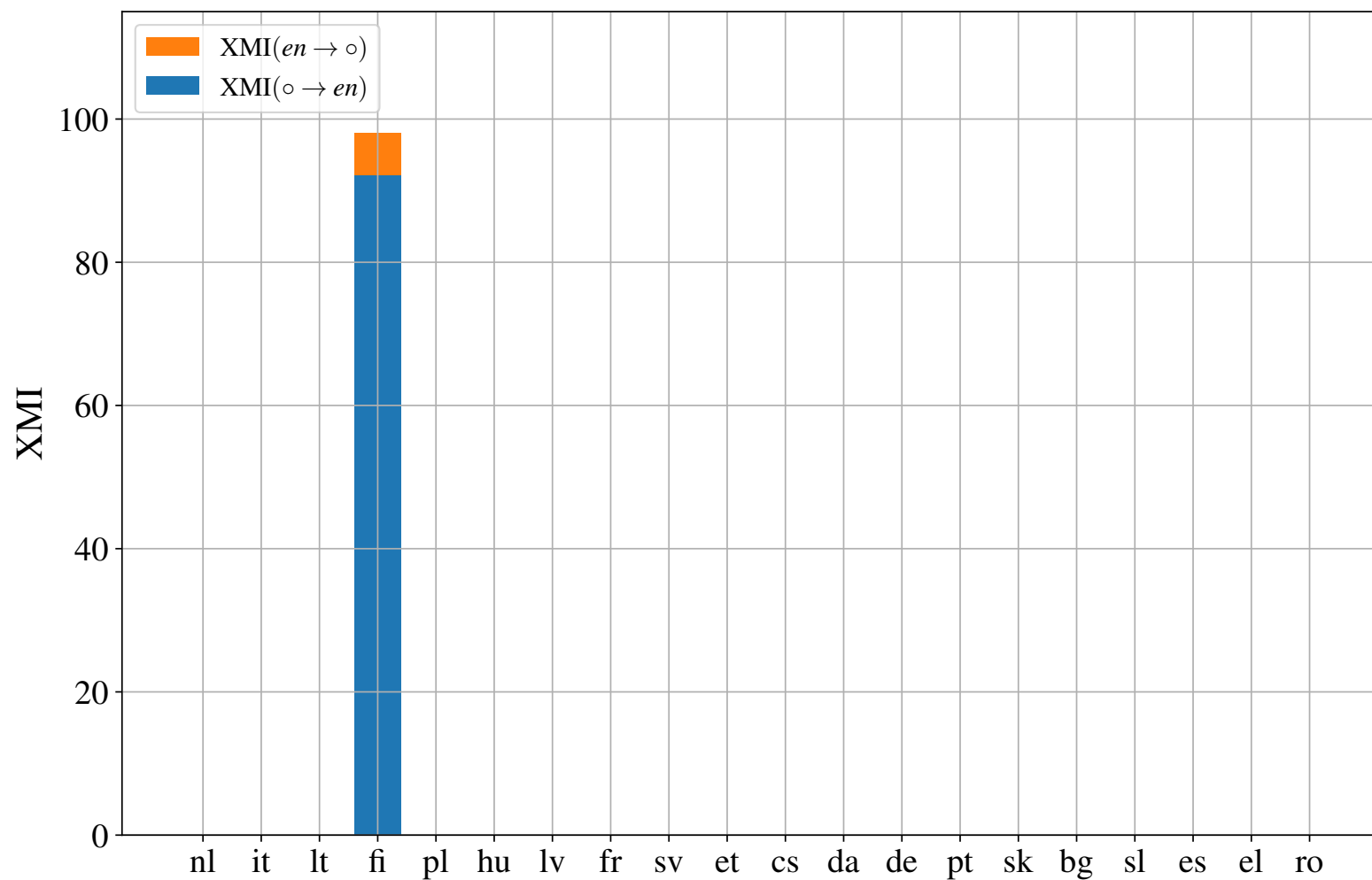
- *Data*: Fully 21-parallel subset of Europarl
- *Models*:
 - 20 [$\circ \rightarrow \text{en}$] Transformers
 - 20 [$\text{en} \rightarrow \circ$] Transformers

Results

- For fixed target, BLEU and XMI correlate well! ✓
- Check our paper for more correlations

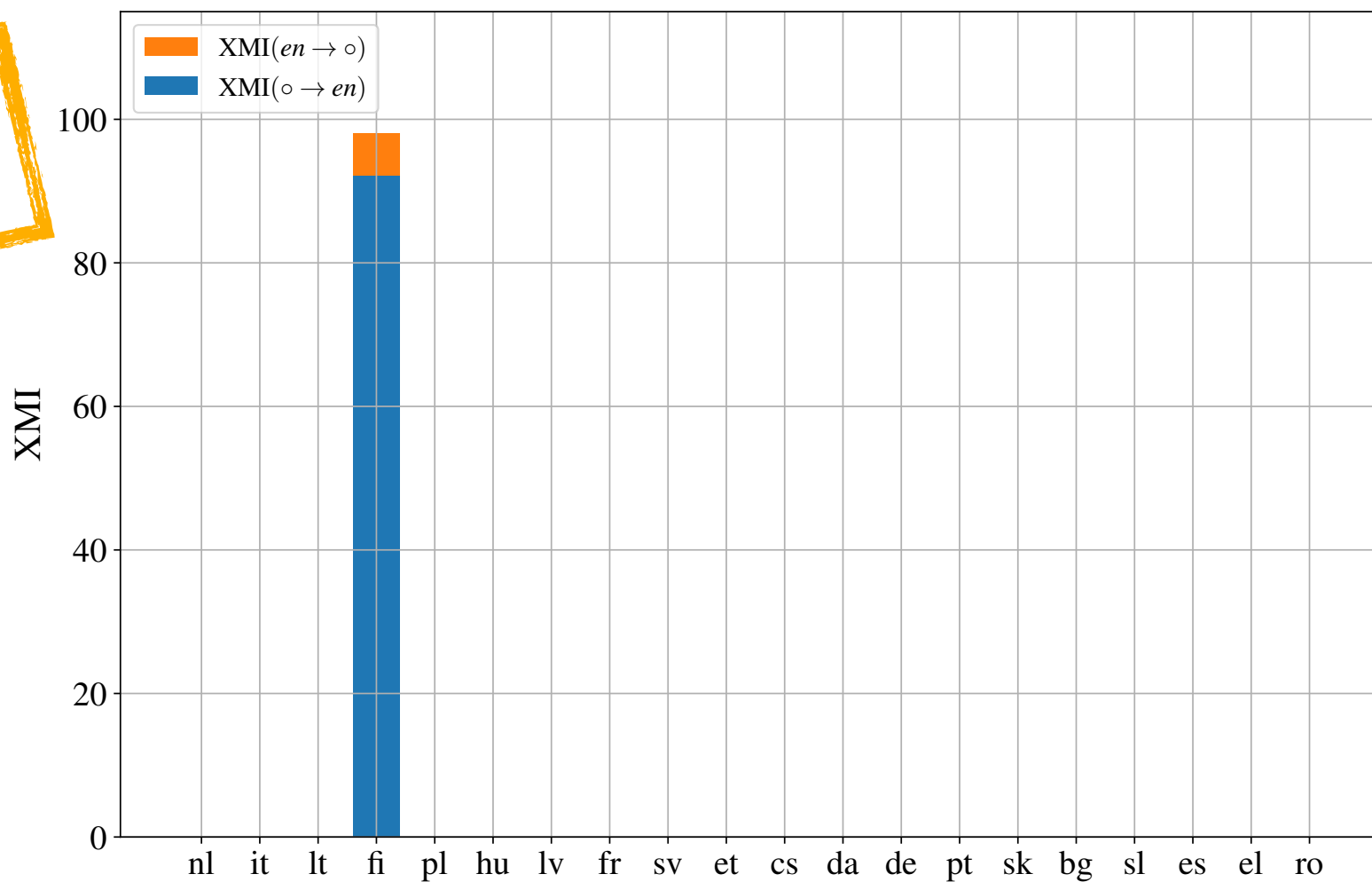


It's Easier to Translate *out of* English than *into* it!



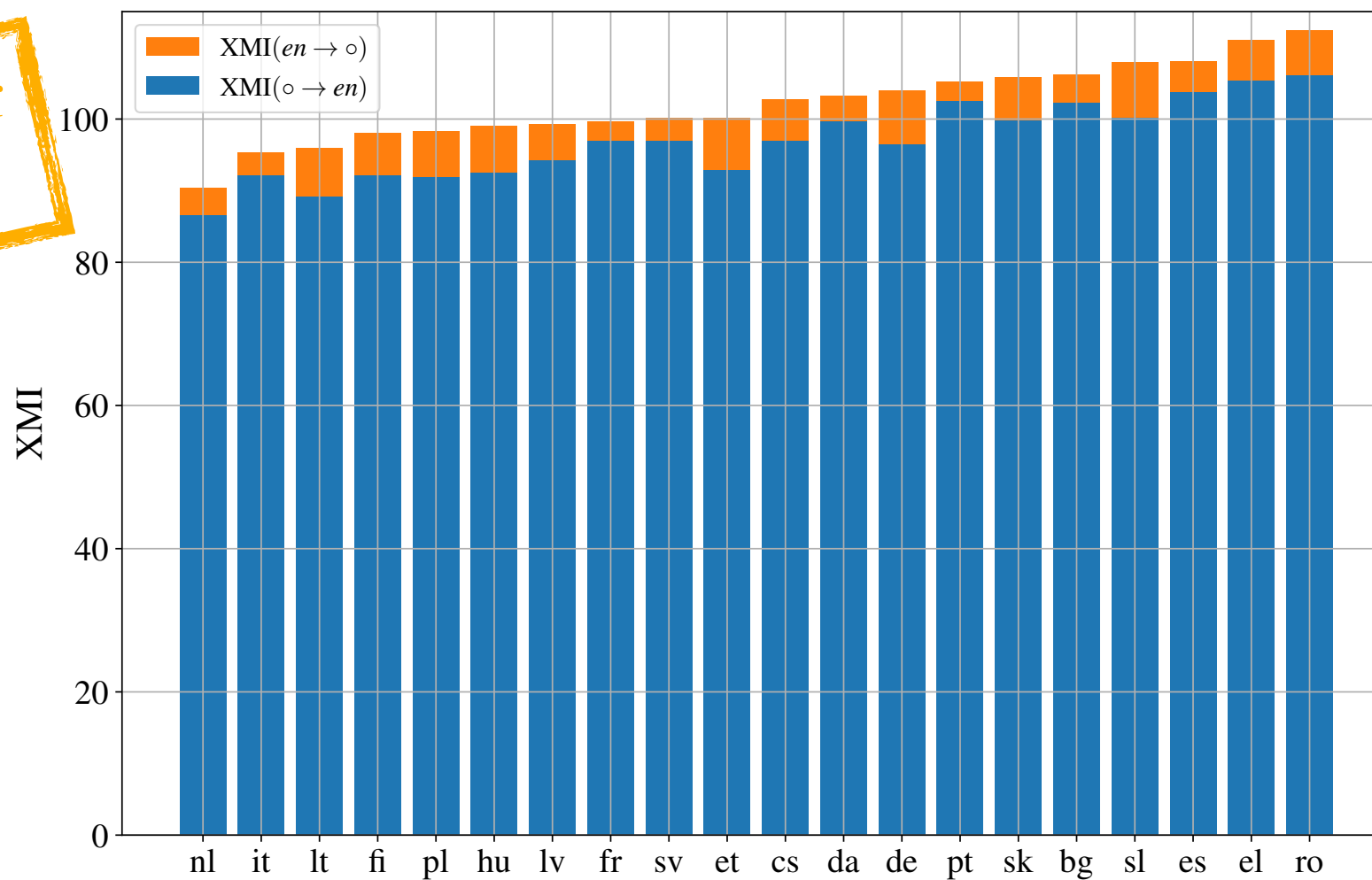
It's Easier to Translate *out of* English than *into* it!

en-fi is easier
than fi-en!



It's Easier to Translate *out of* English than *into* it!

en- \circ is easier
than \circ -en!



Correlations with XMI?

The usual: type-token ratio...
but on the source side!

Spearman's ρ	Metric	☉ → en	en → ☉	both
Mielke et al. (2019)	MCC_{src}	nope	nope	maybe?
	MCC_{tgt}	nope	nope	maybe?
	ADL_{src}	nope	nope	nope
	ADL_{tgt}	nope	nope	maybe?
	$HPE\text{-mean}_{src}$	nope	nope	maybe?
	$HPE\text{-mean}_{tgt}$	nope	nope	maybe?
Lin et al. (2019)	genetic	nope	nope	nope
	syntactic	nope	nope	nope
	featural	nope	nope	nope
	phonological	nope	nope	nope
	inventory	nope	nope	nope
	geographic	nope	nope	nope
Lin et al. (2019)	word number ratio	maybe?	nope	maybe?
	TTR_{src}	maybe?	–	-0.51
	TTR_{tgt}	–	nope	maybe?
	d_{TTR}	maybe?	nope	-0.47
	word overlap ratio	nope	nope	nope

Where to go from here?

Where to go from here?

- Cross-mutual information (XMI)

Where to go from here?

- Cross-mutual information (XMI)
 - A metric for translation difficulties between *any* two directions

Where to go from here?

- Cross-mutual information (XMI)
 - A metric for translation difficulties between *any* two directions
- No linguistic correlations, but TTR... again

Where to go from here?

- Cross-mutual information (XMI)
 - A metric for translation difficulties between *any* two directions
- No linguistic correlations, but TTR... again
 - Let's scale this up and evaluate more pairs!

Where to go from here?

- Cross-mutual information (XMI)
 - A metric for translation difficulties between *any* two directions
- No linguistic correlations, but TTR... again
 - Let's scale this up and evaluate more pairs!
 - Let's build better models!

Where to go from here?

Thanks!

- Cross-mutual information (XMI)
 - A metric for translation difficulties between *any* two directions
- No linguistic correlations, but TTR... again
 - Let's scale this up and evaluate more pairs!
 - Let's build better models!

Code available online at <https://github.com/e-bug/nmt-difficulty>