



# WSP-NMT: Code-Switching with Word Senses for Pretraining in Neural Machine Translation

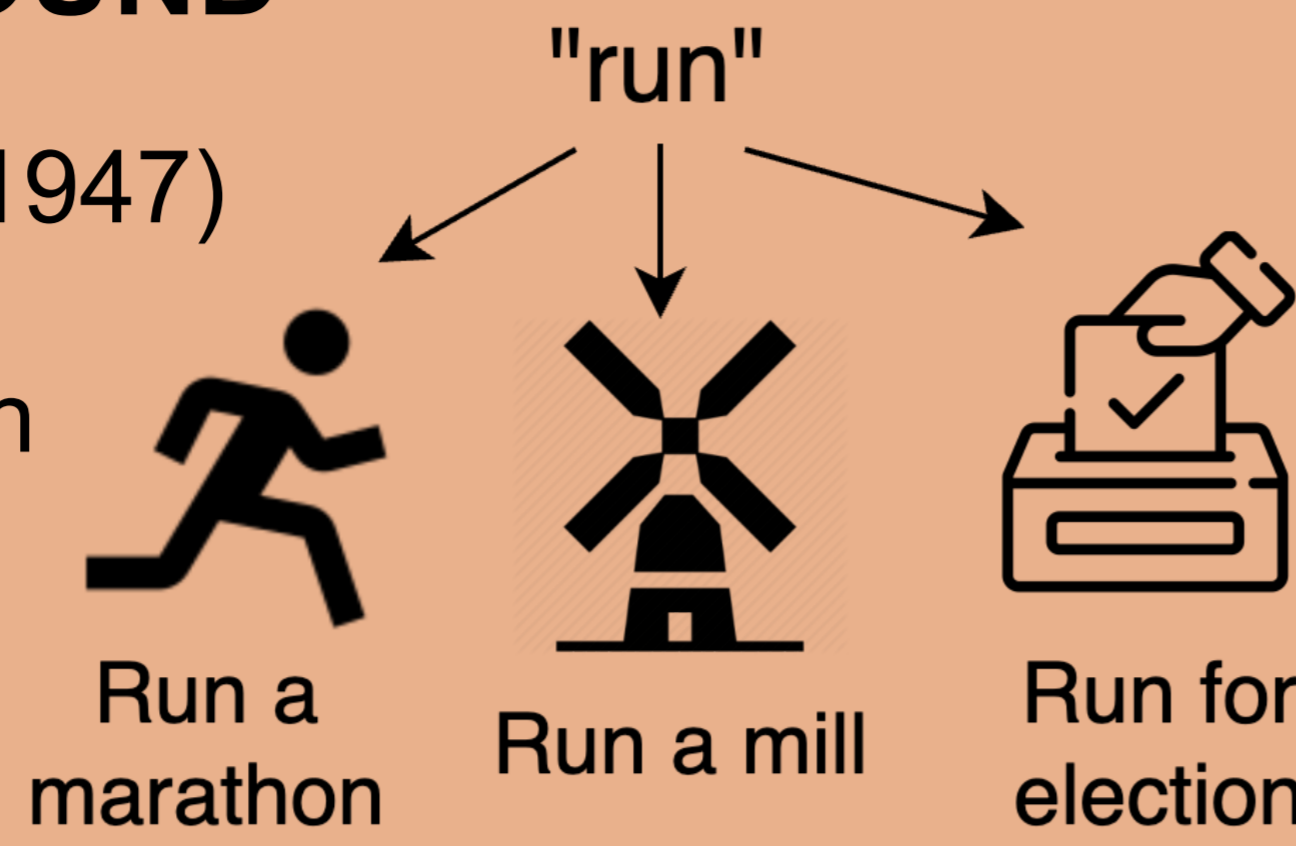
Vivek Iyer<sup>1</sup>, Edoardo Barba<sup>2</sup>, Alexandra Birch<sup>1</sup>, Jeff Pan<sup>1</sup>, Roberto Navigli<sup>2</sup>  
<sup>1</sup>The University of Edinburgh <sup>2</sup>Sapienza University of Rome  
vivek.iyer@ed.ac.uk



SAPIENZA  
UNIVERSITÀ DI ROMA

## 1. BACKGROUND

- Lexical ambiguity in MT (Weaver, 1947)
- Modern NMT systems struggle with WSD biases
- We re-examine NMT pretraining



## 2. CODE-SWITCHED PRETRAINING (CSP)

- Popular NMT pretraining approach, eg. AA (Pan et al., 2021)
- Synthetic Code-Switching: Words  $\Leftrightarrow$  Lexical Translations
- “Sense-agnostic” pretraining!!

Original sentence:	"If we don't win, there will be some inquiries of why we haven't," Graves told BBC Radio Leeds
AA-noised sentence:	" If noi annetada t ויטוריה , ٪ج ٪اسخ jet sometime αιτήσεις seine kuna bize haven't , " Graves erzählte BBC Radio Leeds.

Fig 1: Sourced from Figure 6, Pan et al., 2021

## 3. MOTIVATION

Idea: Disambiguate, then Code-Switch with word **sense** translations!

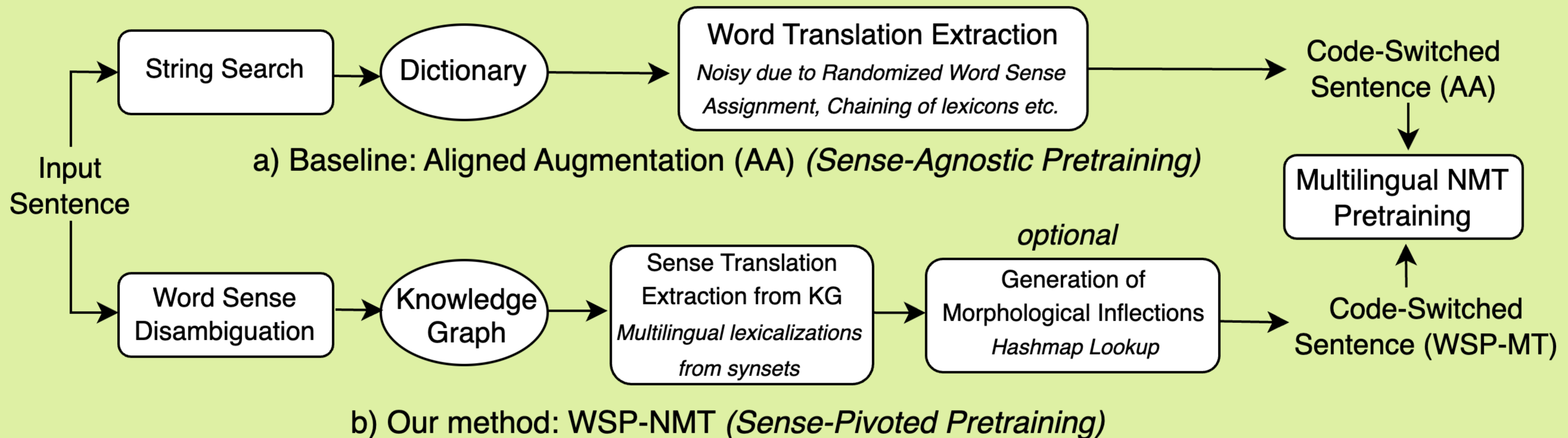
	<b>Source Sentence:</b>	He had an <b>edge</b> on the competition.
	<b>Baseline Translation (AA):</b>	Ha avuto un <b>margin</b> e alla concorrenza.
	<b>Our Translation (WSP-NMT):</b>	Aveva un <b>vantaggio</b> sulla concorrenza.

Fig 2: AA vs WSP-NMT. *Margine=edge, vantaggio=advantage*

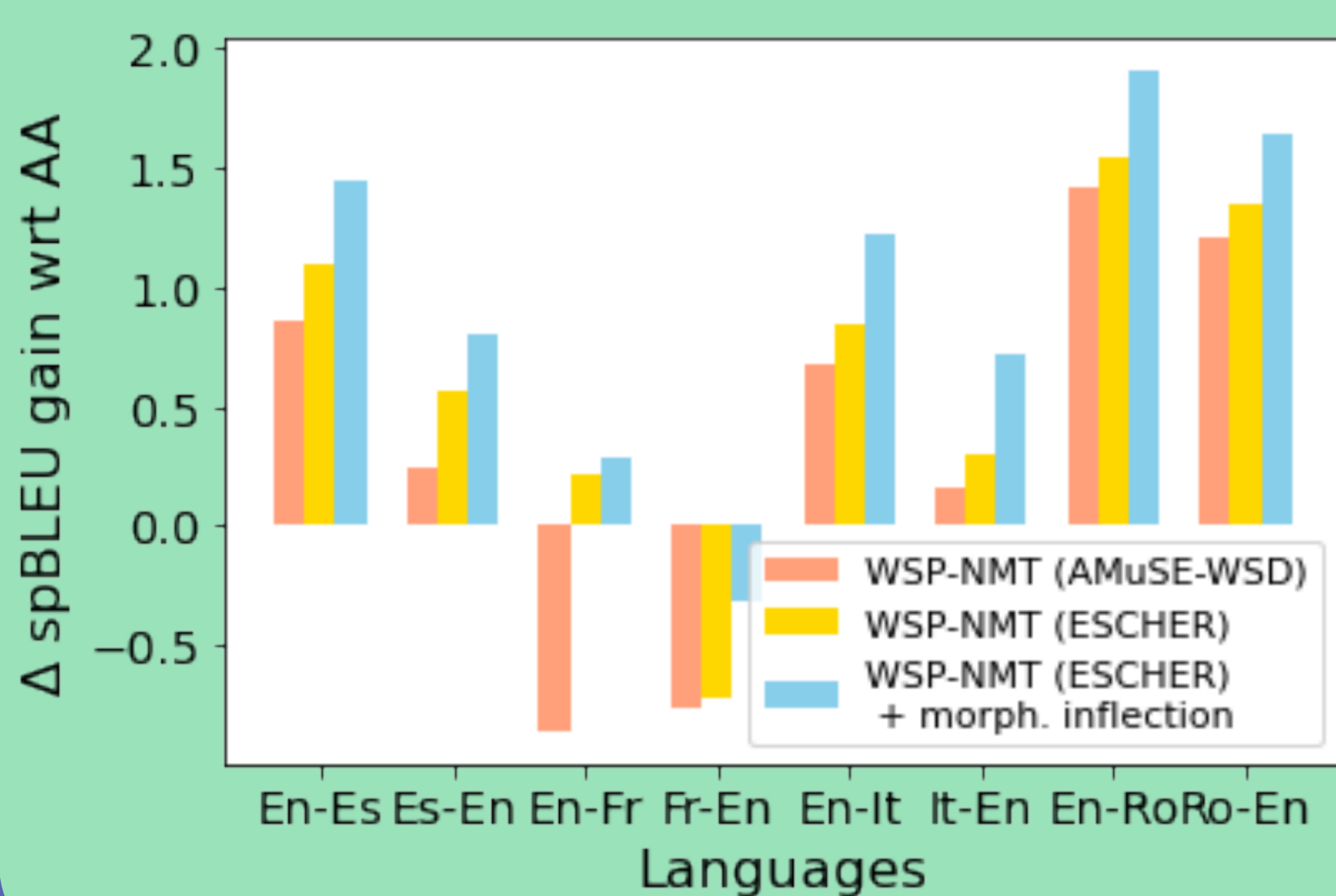
## 4. CONTRIBUTIONS

1. **Sense-pivoted pretraining** can improve overall MT quality and WSD performance
2. KGs + mNMT pretraining = better {reliability, accuracy}
3. Super effective in data-constrained scenarios!

## 5. APPROACH (WSP-NMT)



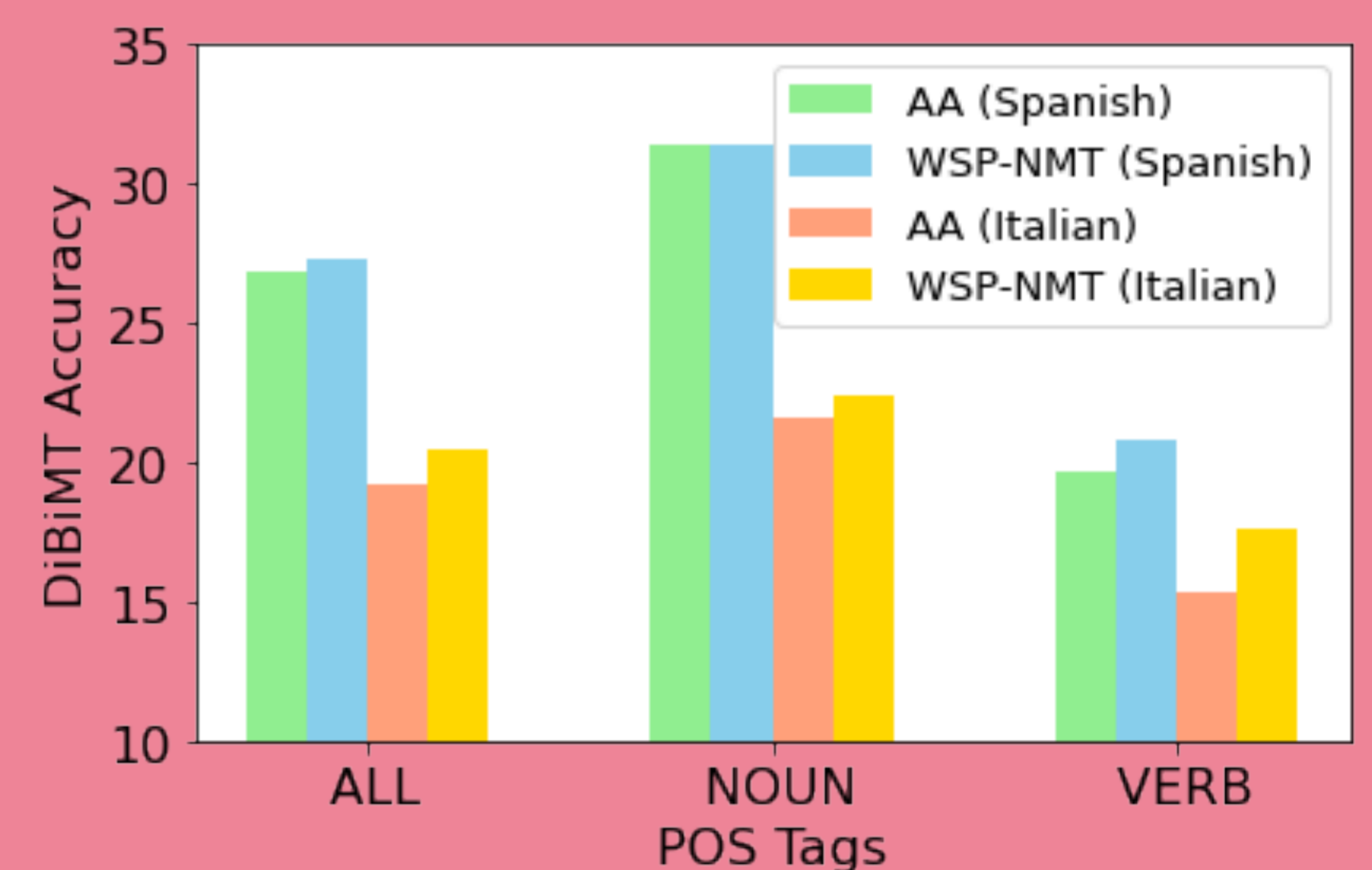
## 6. RESULTS (OVERALL MT)



- ✓ Consistent gains! Sense-pivoted pretraining helps :-)
- ✓ Better WSD (ESCHER) = better MT quality. But AMuSE-WSD is a good alternative too! (2.3x cheaper)
- ✓ Morph. Inflection Prediction w/ MUSE lexicons for {gender, tense} agreement
- ✓ Lower-resourced En-Ro (5x less data) gains the most!!

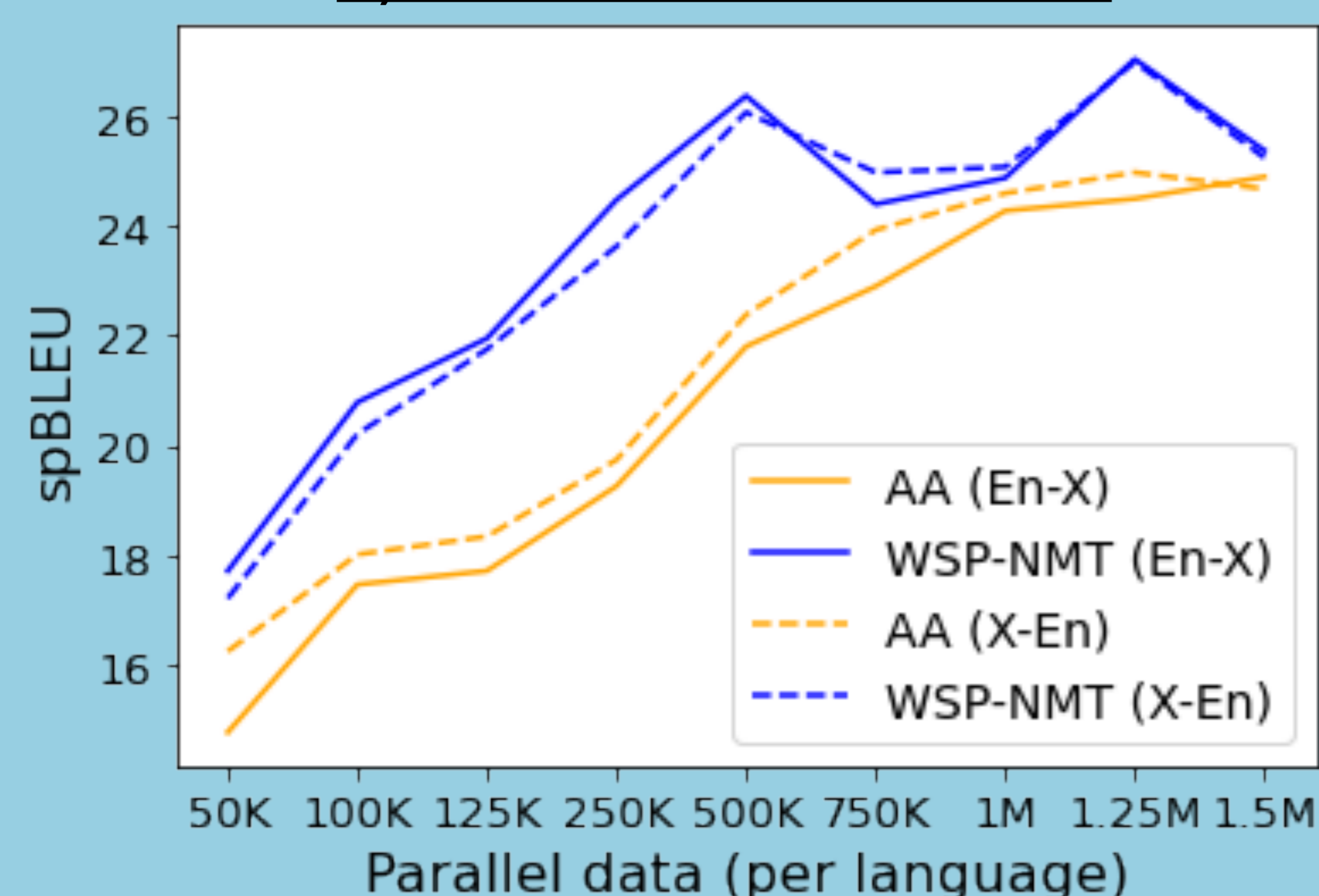
## 7. RESULTS (AMBIGUOUS MT)

Significant gains in verb disambiguation!



## 8. SCALING TO RESOURCE-CONSTRAINED SETTINGS

### A) Data size vs Performance



Highly effective in low & medium data setups!

### B) Zero-Shot Translation

Baseline	En-Pt	Pt-En
AA	2.92	6.88
WSP-NMT	3.60	8.52

Enhanced multilingual convergence!

### C) Zero-Shot WSD (Indo-Iranian)

Baseline	En-X	X-En
AA	22.79	20.49
WSP-NMT	22.71	20.23

Need disambiguation resources :(

## 9. APPLICATIONS

- ✓ Domain-specific MT
  - less data, well-resourced langs
- ✓ Information-centric domains
  - Healthcare, News etc.

## 10. CONCLUSION

### Advantages:

- ↑ Reliability, ↑ Quality, ↓ Errors
- Useful in low-data setups

### Disadvantages:

- Need WSD

