



SCAN ME

## Introduction:

- Natural language types are imbalanced
  - A few types are frequent  $\Rightarrow$  less *information content*
  - Most types are rare  $\Rightarrow$  more *information content*
- Current eval metrics do not address imbalance
- MacroF measure is used for classifier evaluation on imbalanced test sets; we apply MacroF for machine translation evaluation
- MacroF shows differences in quality of supervised and unsupervised NMT in a way BLEU cannot

## Methods

$$\text{Match}(c) = \sum_{i=1}^m \min\{\mathcal{C}(c, h^{(i)}), \mathcal{C}(c, y^{(i)})\}$$

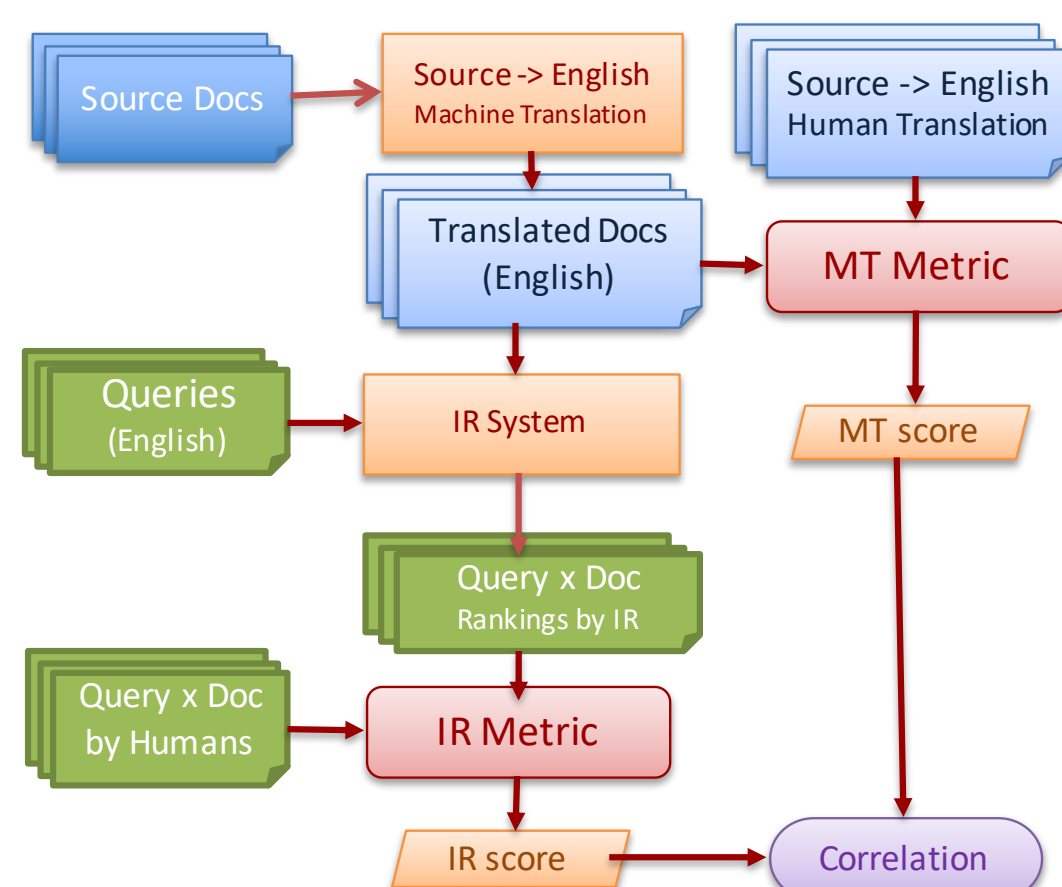
$$P_c = \frac{\text{Match}(c)}{\text{Preds}(c)} \quad R_c = \frac{\text{Match}(c)}{\text{Refs}(c)}$$

$$F_{\beta;c} = (1 + \beta)^2 \frac{P_c \times R_c}{\beta^2 \times P_c + R_c}$$

$$\text{Macro}F_{\beta} = \frac{\sum_{c \in V} F_{\beta;c}}{|V|}$$

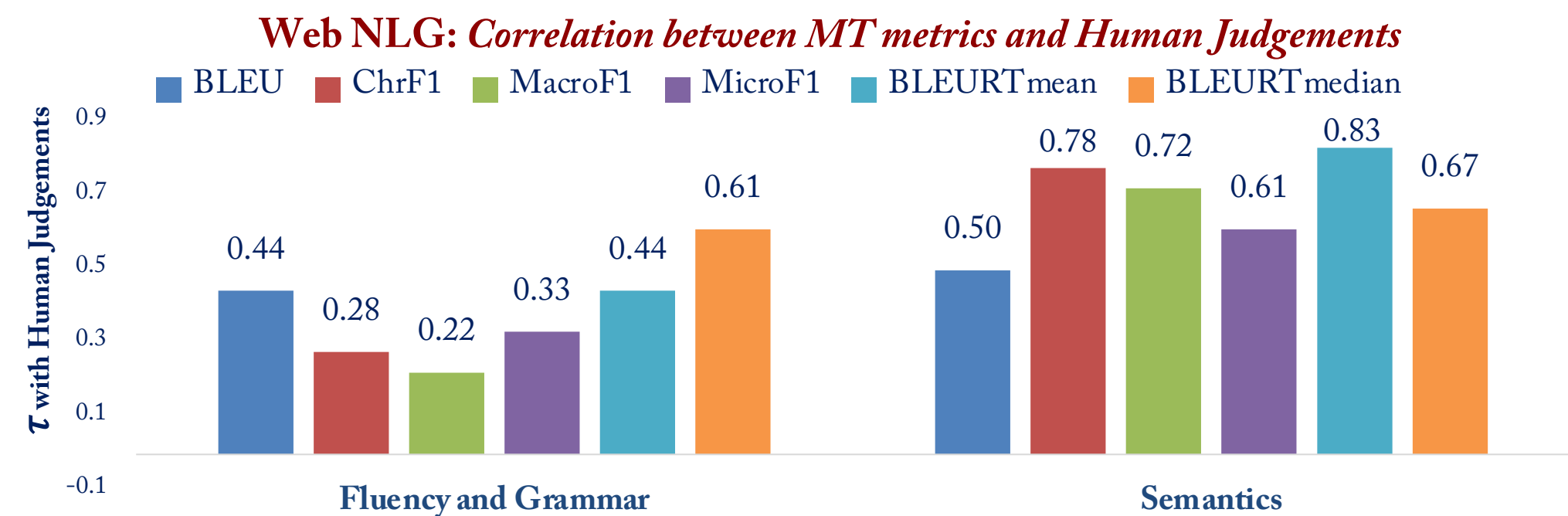
$$\text{Micro}F_{\beta} = \frac{\sum_{c \in V} f(c) \times F_{\beta;c}}{\sum_{c' \in V} f(c')} \text{ where } f(c) = \text{Refs}(c) + k; k \geq 1$$

## CLIR Pipeline

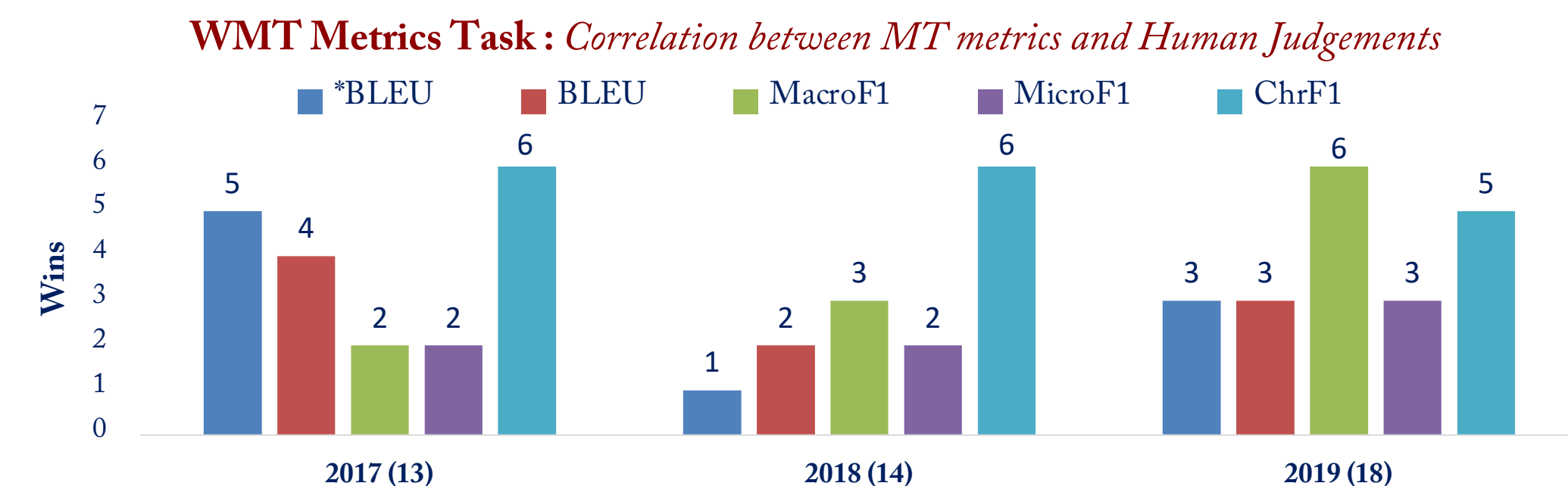


## Justification

- Direct assessment : WebNLG , WMT Metrics Task
- Downstream task performance indication



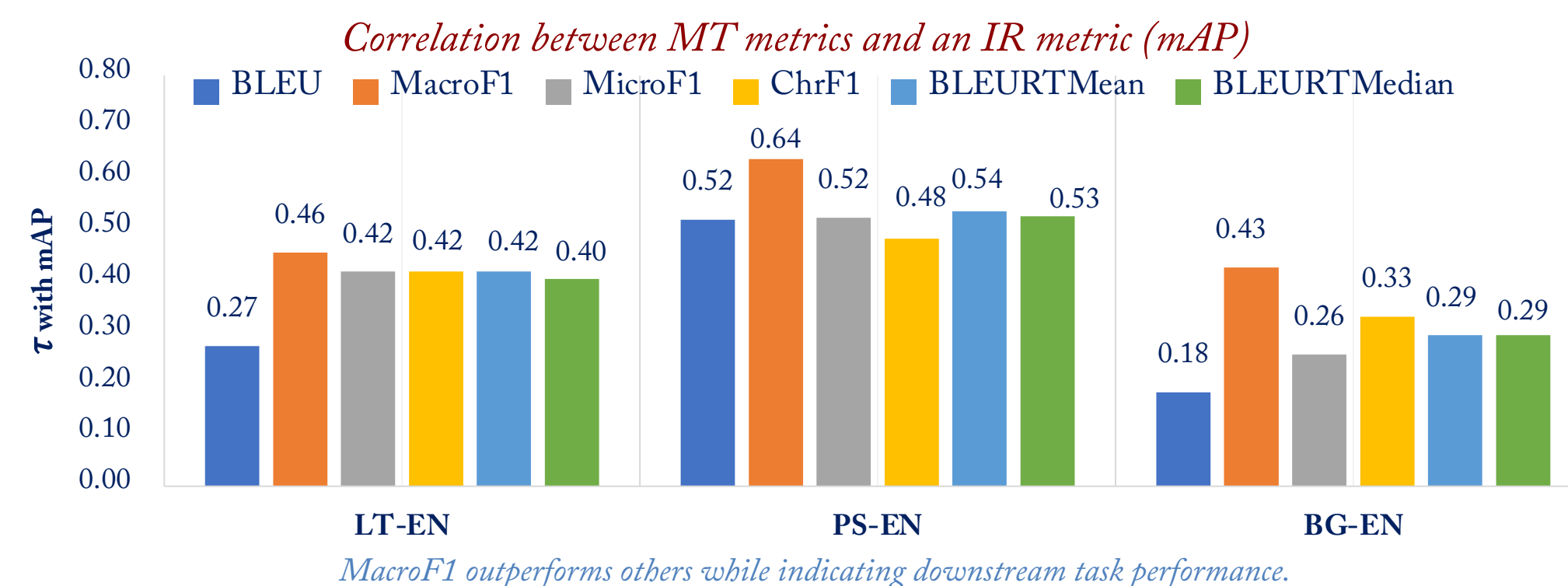
MacroF1 is a poor indicator of fluency and grammar; but a strong indicator of semantics; These results are based on English only.



MacroF1 has highest number of wins in the recent year when most systems are fluent, and adequacy is the key discriminating factor.

\*BLEU is precomputed values in metrics package. MacroF1, MicroF1 share the same tokenizer as BLEU, obtained using SacreBLEU

## Cross-lingual Information Retrieval Task: CLSSTS 2020



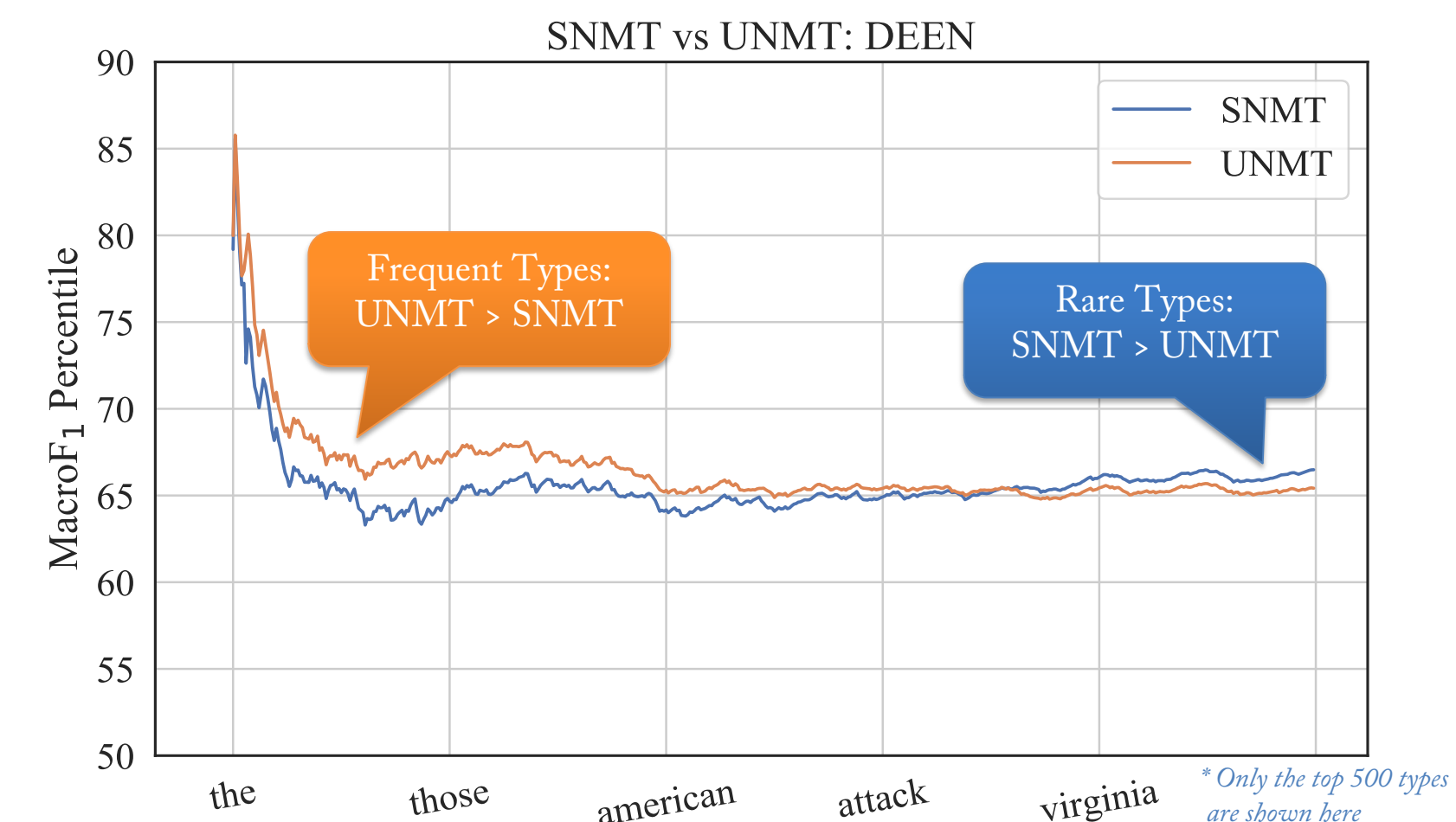
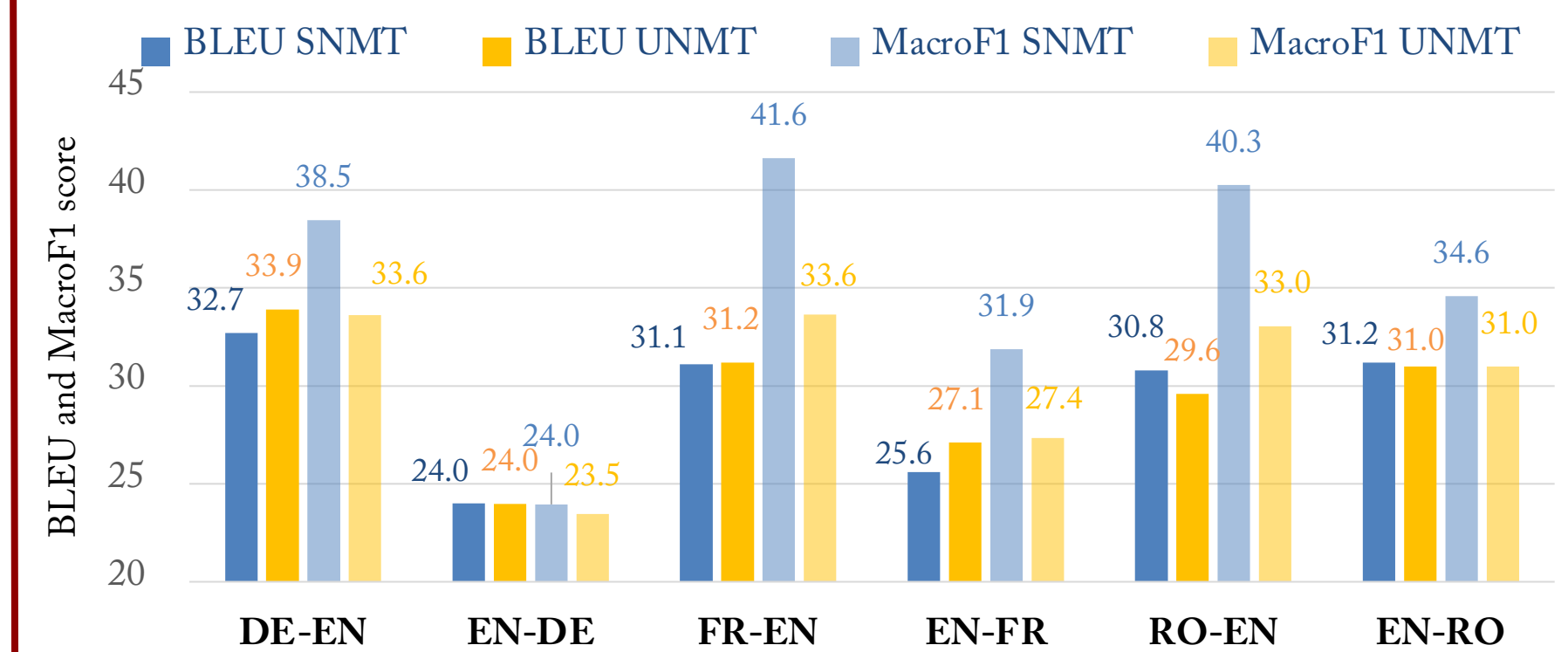
MacroF1 outperforms others while indicating downstream task performance.

## Usage

```

$ pip install git+https://github.com:isi-nlp/sacrebleu.git
$ sacrebleu $REF -m macrof microf < $HYP.detok
$ sacrebleu $REF -m macrof --report report.txt < $HYP.detok
    
```

## SNMT vs UNMT Quality Diff



## Manual Analysis:

Untranslation and truncation are heavily penalized by MacroF1 than BLEU

## Conclusion

- MacroF1 is a strong indicator of semantics; competitive on direct assessments and outperforms others on a downstream CLIR task
- MacroF1 is easily computable and interpretable, does not appear to have uncontrollable biases resulting from data (unlike model-based metrics)
- Macro-averaged evaluation is a useful technique for addressing the importance of long tail of language