

When the Gold Standard Isn't Necessarily Standard: Challenges of Evaluating the Translation of User-Generated Content



Lydia Nishimwe, Benoît Sagot, Rachel Bawden
Inria, France
{lydia.nishimwe, benoit.sagot, rachel.bawden}@inria.fr



PARIS Artificial Intelligence Research InstitutE

Which translation is correct?

his TOOOO funny!!

A Il est trop drôle !

B il est TROOOOP drôle !!

C il et TROOOOP drôle !!

Key Takeaways

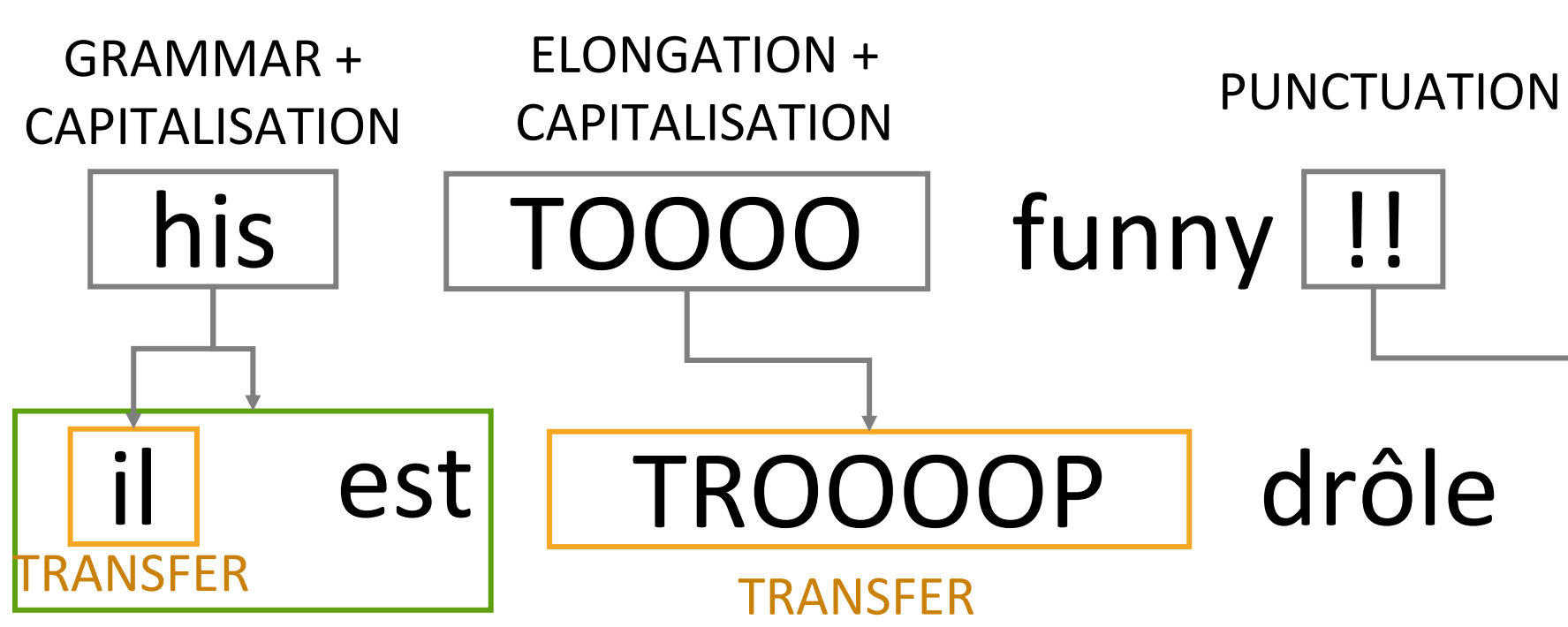
There is **no single gold standard** for UGC translation.

LLMs can be **steered toward different translation styles**.

Fair evaluation requires **guideline-aware metrics**.

RQ1 What is a "good" translation of user-generated content (UGC)?

Translation Guidelines



Defining a **taxonomy of phenomena and actions** from translation guidelines of UGC datasets

Different datasets expect **different levels of standardness in the gold translations**

Phenomenon

- Grammar
- Spelling
- Word elongation (e.g., *gooooaaaalllll*)
- Capitalisation (e.g., *NOPE, SoRry*)
- Informal abbreviations (e.g., *gonna, u*)
- Informal acronyms (e.g., *LOL, TBH*)
- Hashtags and subreddits
- URLs, user IDs, and retweet marks (*RT*)
- Emoticons and emojis
- Atypical punctuation
- Overt profanity (e.g., *fuck*)
- Self-censored profanity (e.g., *f*ck*)

	RoCS-MT	Footweets	MMTC	PFSMB
1	Normalise	Normalise	Normalise	Normalise
2	Normalise	Normalise	Normalise	Normalise
3	Normalise	Transfer	Transfer	Transfer
4	Normalise	Transfer	Transfer	Transfer
5	Normalise	Normalise	Normalise	Transfer
6	Normalise	Normalise	Transfer	Transfer
7	Copy	Copy	Transfer	Transfer
8	Copy	Copy	Copy	Copy
9	Copy	Copy	Copy	Copy
10	Normalise	Copy	Copy	Copy
11	Transfer	Transfer	Transfer	Transfer
12	Normalise	Normalise	Normalise	Transfer

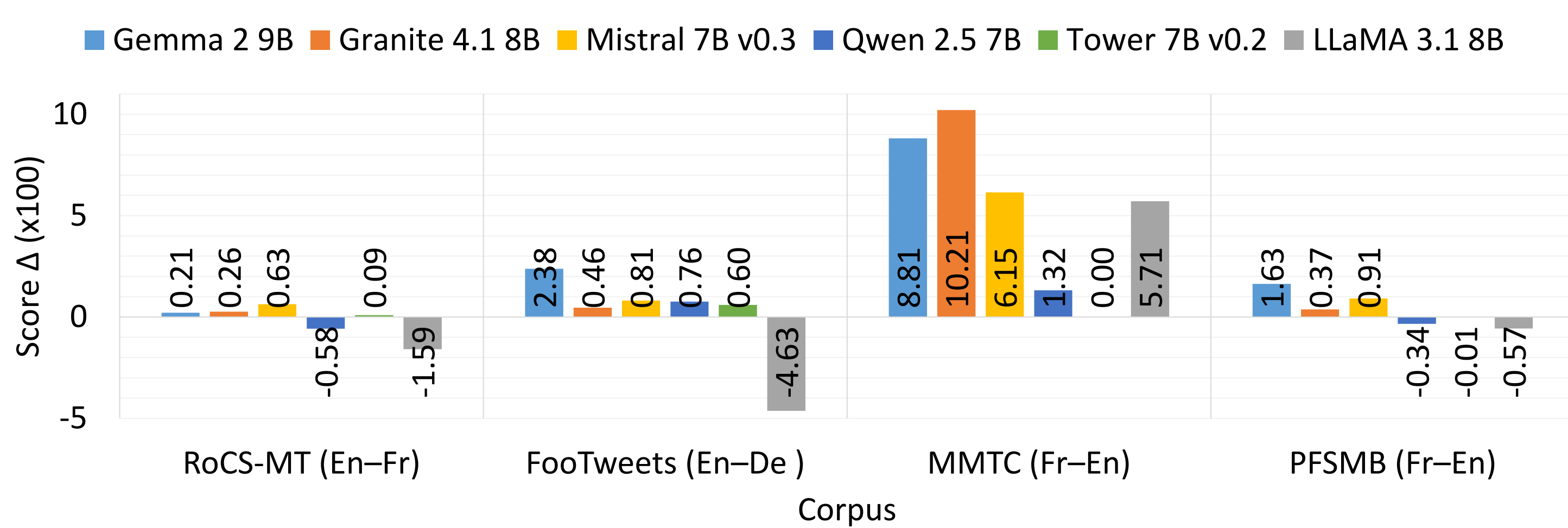
Highly standardising dataset ← → Minimally standardising dataset

RQ2 How can we ensure fair evaluation of UGC translation?

Prompting LLMs to translate with the **same guidelines** as the UGC datasets

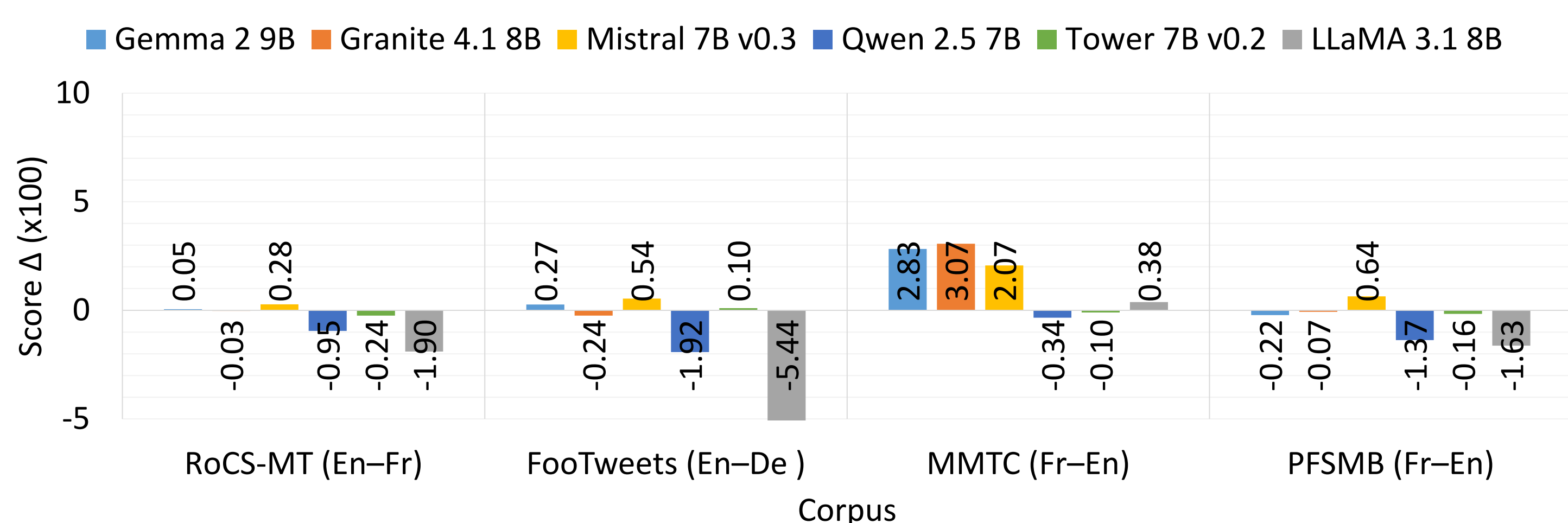
Reference-based Automatic Evaluation

COMET score variation (guided output vs. default)



Reference-less Automatic Evaluation

COMETKiwi score variation (guided output vs. default)



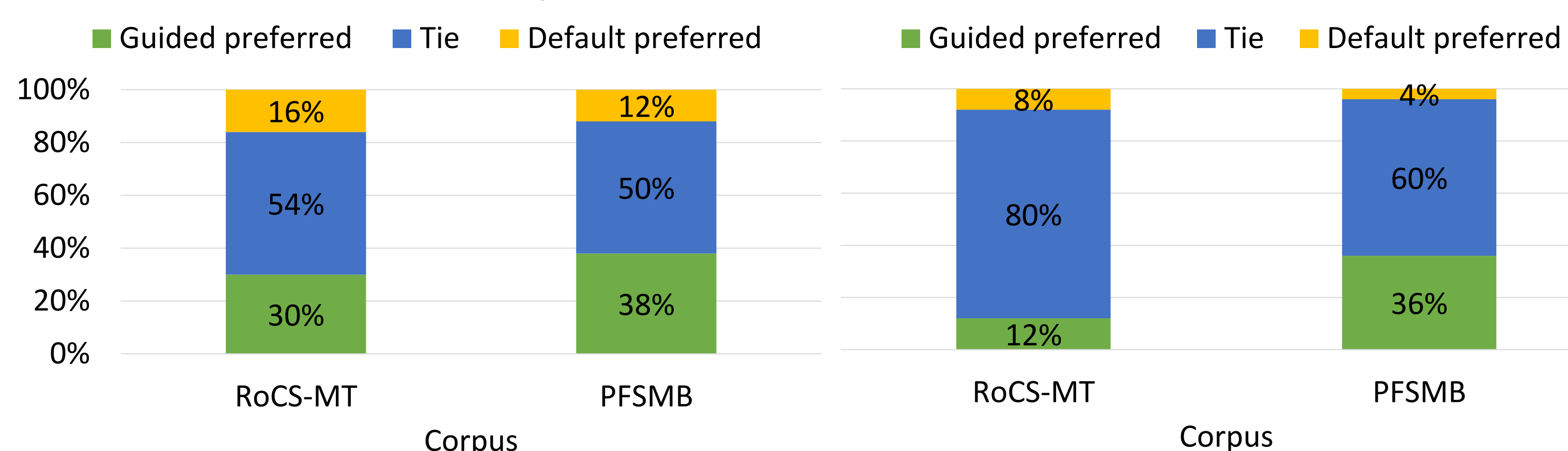
Qualitative Analysis (Example from MMTC)

Source: @JulieTom62 même pas besoins de regardé le match 🤔
 Normalised source: @JulieTom62 même pas besoin de regarder le match 🤔
 Reference: @JulieTom62 don't even need to watch the match 🤔
 Gemma's default output: No need to even watch the game 🤔
 + guidelines: @JulieTom62 no even need to watch the game 🤔

COMET
78.44
91.64 ↗

Human Evaluation

Overall Quality



Findings

- Most LLMs follow the UGC guidelines (Qwen follows them but hallucinates, Tower ignores them, LLaMA refuses to translate)
- Matching guidelines generally improve UGC translation but can introduce errors
- COMET aligns better with guideline adherence
- COMETKiwi aligns better with quality judgments

Future Work

- LLM-as-a-judge evaluation with guideline-specific rubrics
- More compositional guidelines and prompting strategies

References:

- RoCS-MT: Robustness Challenge Set for Machine Translation (Bawden & Sagot, WMT 2023)
- Footweets: A Bilingual Parallel Corpus of World Cup Tweets (Sluyter-Gäthje et al., LREC 2018)
- [MMTC] The Multilingual Microblog Translation Corpus: Improving and Evaluating Translation of User-Generated Text (McNamee & Duh, LREC 2022)
- [PFSMB] Comparison between NMT and PBSMT Performance for Translating Noisy User-Generated Content (Rosales Núñez et al., NoDaLiDa 2019)
- COMET-22: Unbabel-IST 2022 Submission for the Metrics Shared Task (Rei et al., WMT 2022)
- CometKiwi: IST-Unbabel 2022 Submission for the Quality Estimation Shared Task (Rei et al., WMT 2022)

Full paper here →

