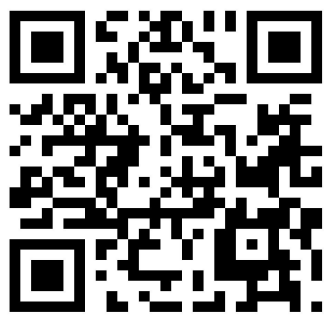# Guiding Large Language Models to Post-edit Machine Translation with Error Annotations

Dayeon Ki    Marine Carpuat

University of Maryland, College Park
dayeonki@umd.edu        marine@umd.edu
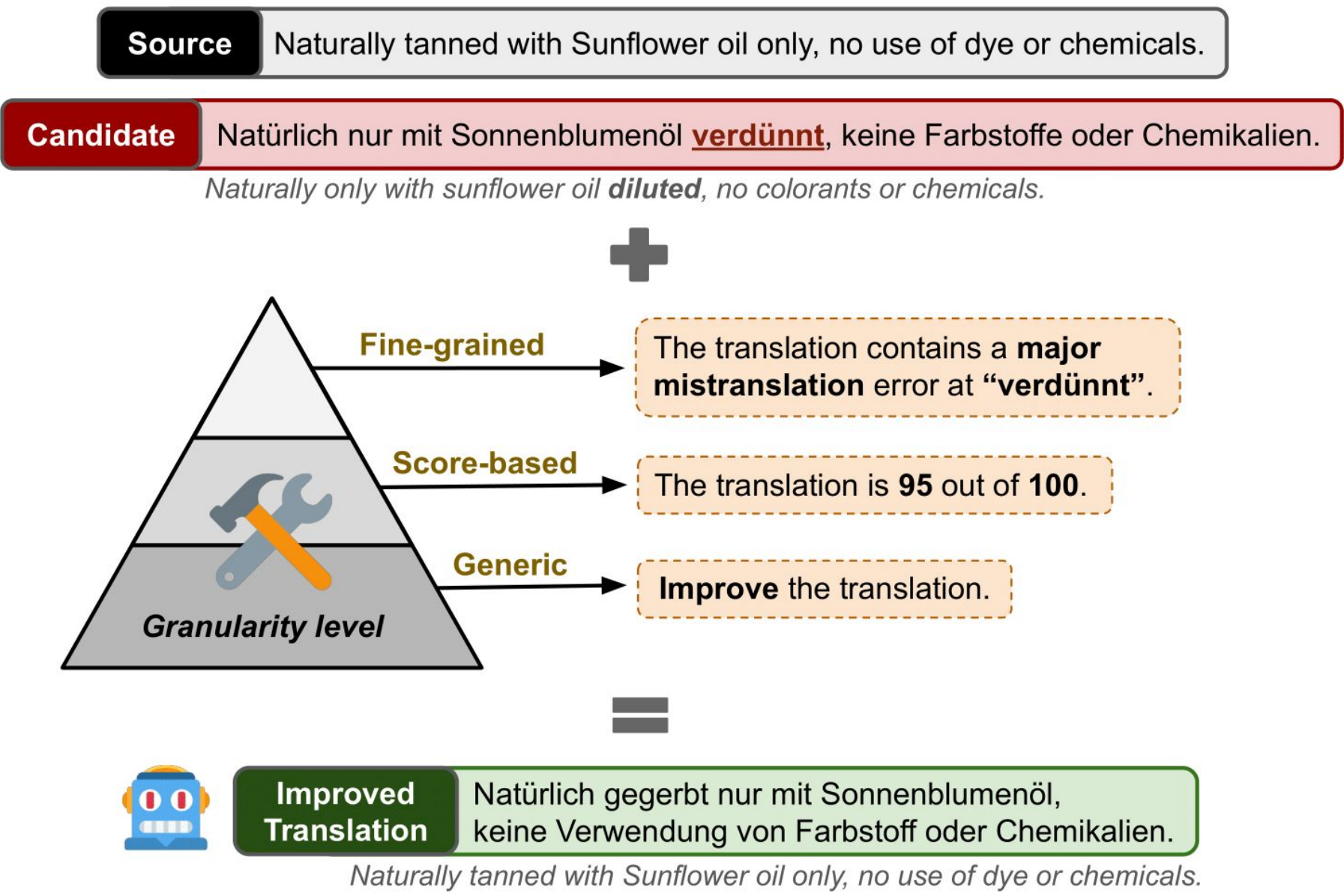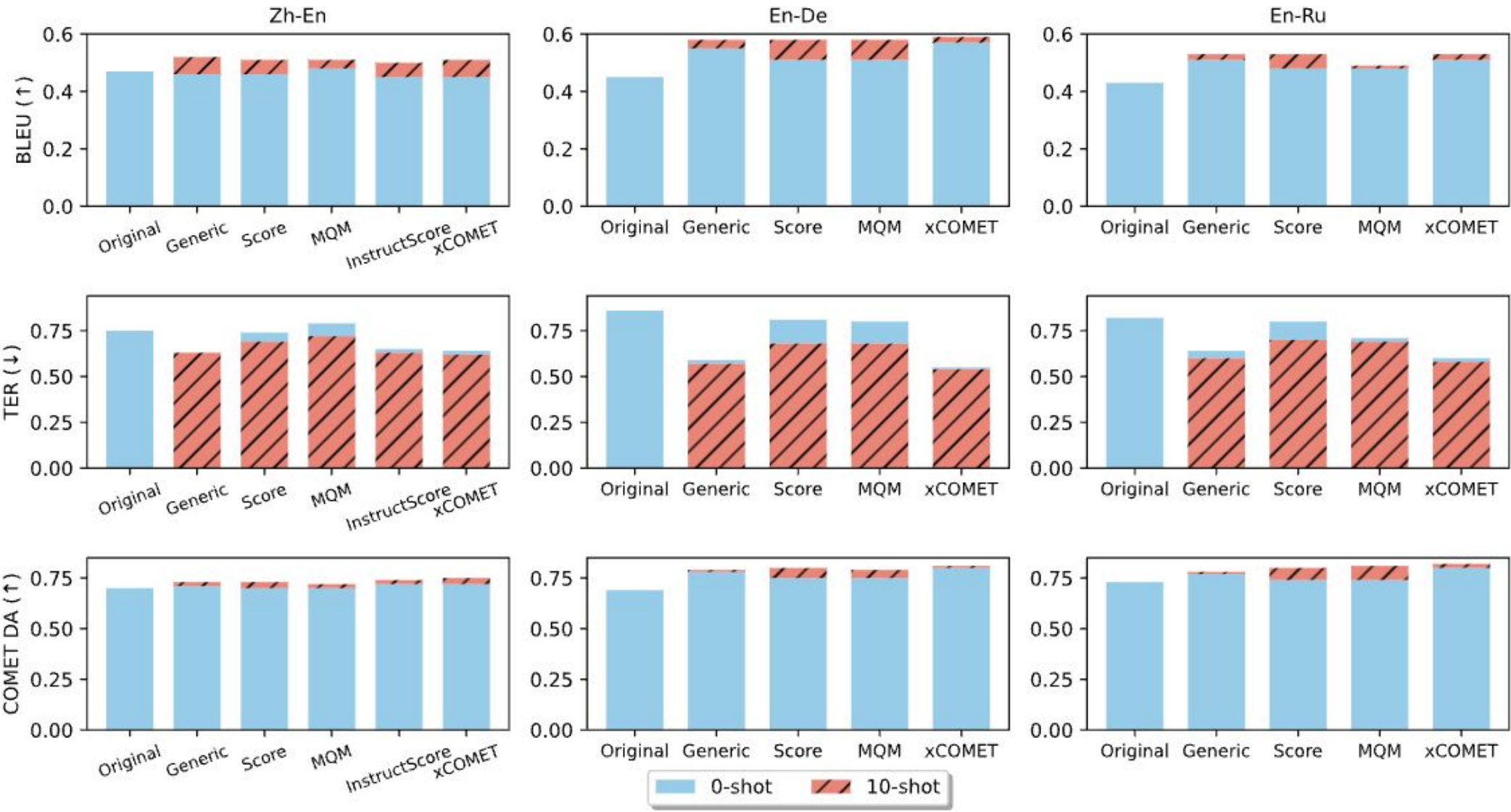
NAACL 2024

Paper

Twitter

## Motivation

- Performance of LLMs remains uneven with variation across models, languages, and translation directions.
- Exploit the complementary strengths of LLMs and supervised Machine Translation system.

**Can smaller LLMs use fine-grained feedback to refine their translations?**
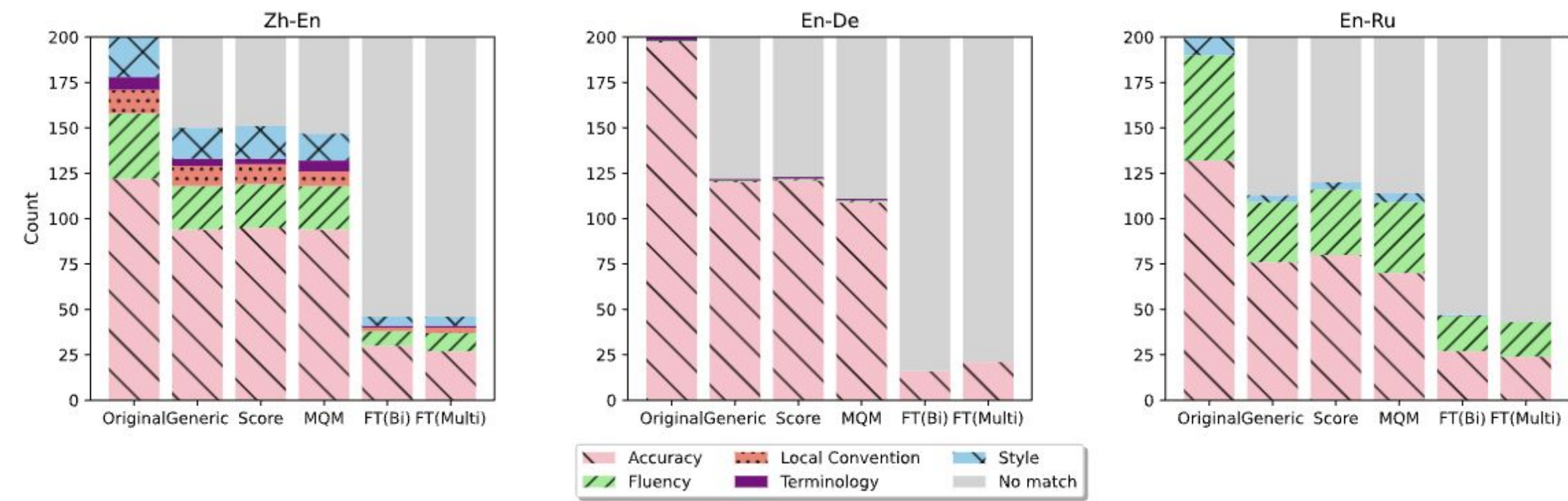


## Prompting Results

- Prompting LLaMA-2 to post-edit with any forms of feedback consistently improve MT performance.
- Most forms of our tested feedback, **regardless of granularity, converge to similar performance.**



## Error Analysis

- Fine-tuning with error annotations help **align LLM behavior** with the provided feedback.



## Types of Feedback

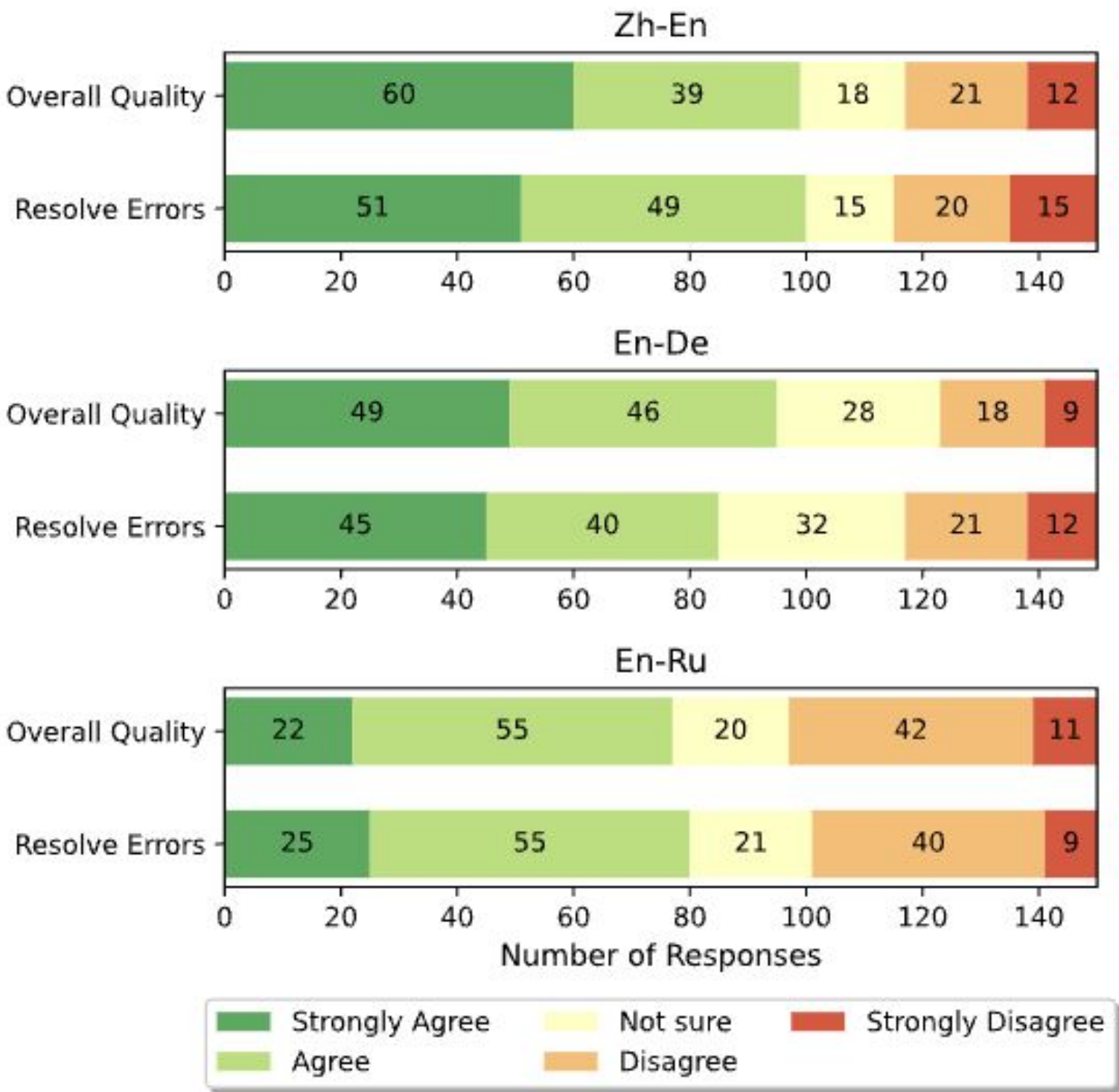| Category | Prompt |
|---|---|
| Generic | Improve the translation from English to German without any explanation. English: *The newer items are bagged only.* German: *Neue Gegenstände werden nur mit Gepäck versehen.* Improved German: |
| Score | Improve the translation from English to German without any explanation. This translation is scored 85 out of 100. English: *The newer items are bagged only.* German: *Neue Gegenstände werden nur mit Gepäck versehen.* Improved German: |
| Fine-grained | Improve the translation from English to German based on the identified errors without any explanation. (1) There is a major mistranslation error at "mit Gepäck versehen". English: *The newer items are bagged only.* German: *Neue Gegenstände werden nur mit Gepäck versehen.* Improved German: |

## Fine-tuning Results

### Automatic Evaluation

- Fine-tuning LLaMA-2 with error annotated translations gives an **extra boost** on top of prompting results.

| Language | Type | BLEU ($\uparrow$) | TER ($\downarrow$) | COMET$_{DA}$ ($\uparrow$) |
|---|---|---|---|---|
| Zh-En | Original | 0.47 | 0.75 | 0.70 |
| | prompt ($k$=0) | $0.48^\dagger$ | $0.72^\dagger$ | $0.70^\dagger$ |
| | prompt ($k$=10) | $0.51^\dagger$ | $0.65^\dagger$ | $0.72^\dagger$ |
| | FT (Bi) | $\mathbf{0.53}^\dagger$ | $0.63^\dagger$ | $\mathbf{0.76}^\dagger$ |
| | FT (Multi) | $\mathbf{0.53}^\dagger$ | $\mathbf{0.61}^\dagger$ | $\mathbf{0.76}^\dagger$ |
| En-De | Original | 0.45 | 0.86 | 0.69 |
| | prompt ($k$=0) | $0.51^\dagger$ | $0.68^\dagger$ | $0.75^\dagger$ |
| | prompt ($k$=10) | $0.58^\dagger$ | $0.56^\dagger$ | $0.79^\dagger$ |
| | FT (Bi) | $0.56^\dagger$ | $0.58^\dagger$ | $0.79^\dagger$ |
| | FT (Multi) | $\mathbf{0.59}^\dagger$ | $\mathbf{0.55}^\dagger$ | $0.79^\dagger$ |
| En-Ru | Original | 0.43 | 0.82 | 0.73 |
| | prompt ($k$=0) | $0.48^\dagger$ | $0.69^\dagger$ | $0.74^\dagger$ |
| | prompt ($k$=10) | $0.49^\dagger$ | $0.67^\dagger$ | $0.79^\dagger$ |
| | FT (Bi) | $0.51^\dagger$ | $0.65^\dagger$ | $\mathbf{0.80}^\dagger$ |
| | FT (Multi) | $\mathbf{0.52}^\dagger$ | $\mathbf{0.63}^\dagger$ | $\mathbf{0.80}^\dagger$ |

### Human Evaluation

- Fine-tuning **rewrites for naturalness** in the target language. (*better explains the context, flows better in target language*)



## Discussions

- Post-editing MT **does not require the largest proprietary LLMs** and can be done with **smaller open-source models.**
- Building on LLMs fine-tuned for many translation related tasks is a promising direction for encouraging transfer learning from limited amounts of annotation.