

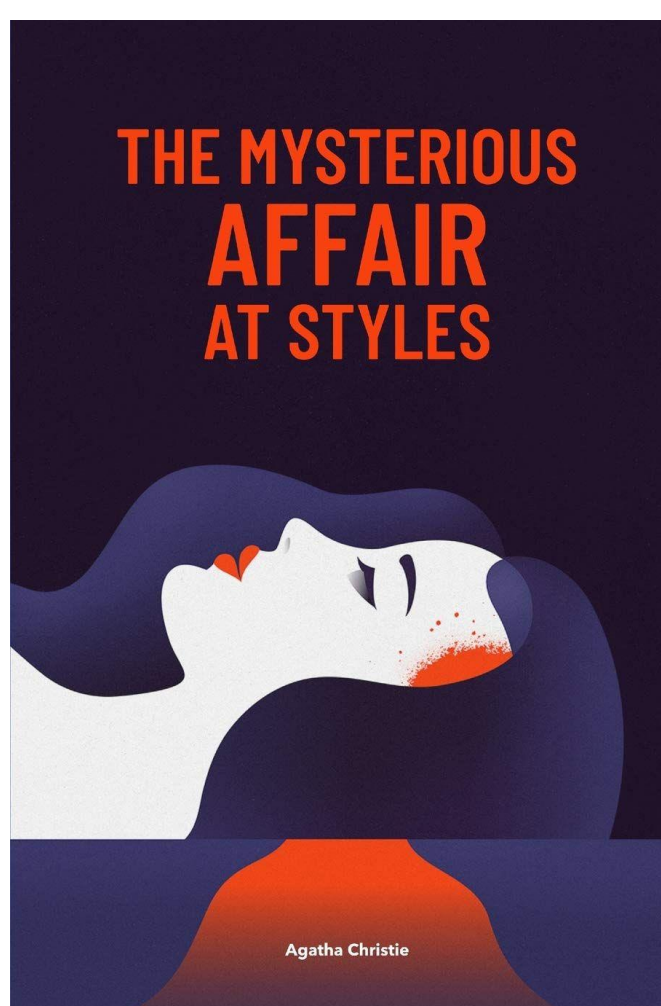


Introduction

Several studies have demonstrated that translation quality has improved enormously since the emergence of NMT^{1, 2, 3}

But: usually sentence-level evaluations of general text types

- What happens when we carry out the evaluation on machine translations of more creative text types, such as literature?
- What happens when we carry out the evaluation on document level?

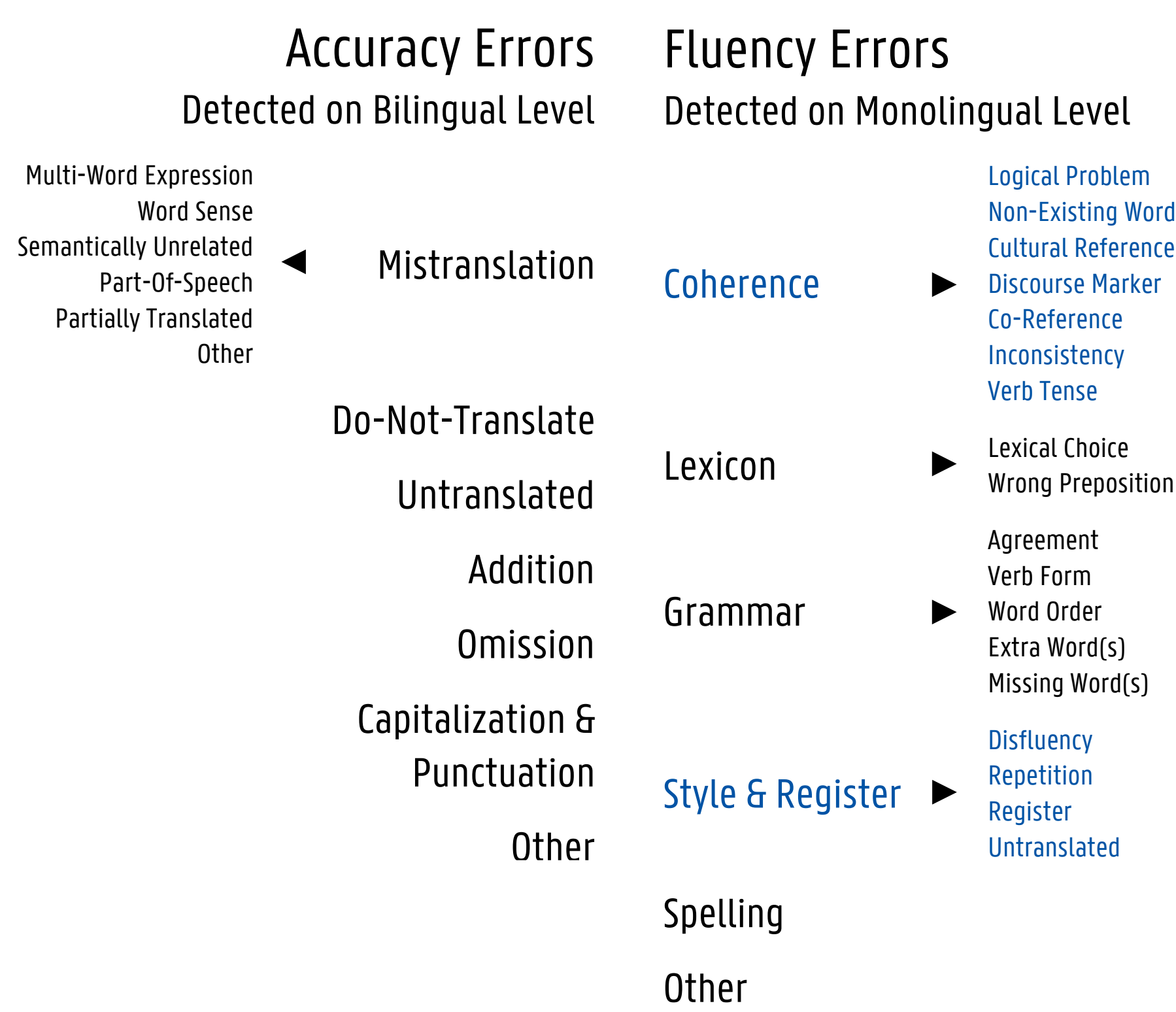


Text Selection

- Detective novel by Agatha Christie
- 58,110 words – 5,276 sentences
- Motivation: eye-tracking data for human translation available (future work: gather corresponding data for MT)
- Translated from English into Dutch by Google's NMT system

Fine-Grained Error Annotation

- Entire novel annotated by one annotator
- First chapter also independently annotated by a second annotator
- WebAnno annotation tool
- Errors classified within the hierarchical SCATE MT error taxonomy (3 levels), adapted to literary NMT on document level⁴
- Two-step annotation approach: first fluency, then accuracy
- Annotations spans are allowed to overlap



Extended SCATE MT error taxonomy with document-level features

64	Poirot was meteen ontuchttert	Grammar & Syntax Verb Form
65	"Kom, kom, mijn vriend," zei hij terwijl hij zijn armen door de mijne gleed.	Logical problem Coherence
66	"Ne vous fachez pas!"	Mistralation Other
67	"Come, come, my friend," he said, slipping his arms through mine.	Mistralation Other
68	"Kom, kom, mijn vriend," zei hij terwijl hij zijn armen door de mijne gleed.	Logical problem Coherence
69	"Ne vous fachez pas!"	Mistralation Other
70	"Ne vous fachez pas!"	Mistralation Other
71	"Ne vous fachez pas!"	Mistralation Other

Step 1: annotating all fluency errors without access to the source

Step 2: annotating & linking all accuracy errors in source & target

Inter-Annotator Agreement (IAA)

Do the annotators agree on their annotations on the first chapter?

(Dis)agreement on error detection

- Low IAA when calculated on word level (only 56% of the annotated words were annotated by both annotators)
- Decent IAA when calculated on annotation level (73% of the annotations overlap with an annotation by the other annotator)
- Explanation: tendency to disagree on the length of the error span

(Dis)agreement on error categorization

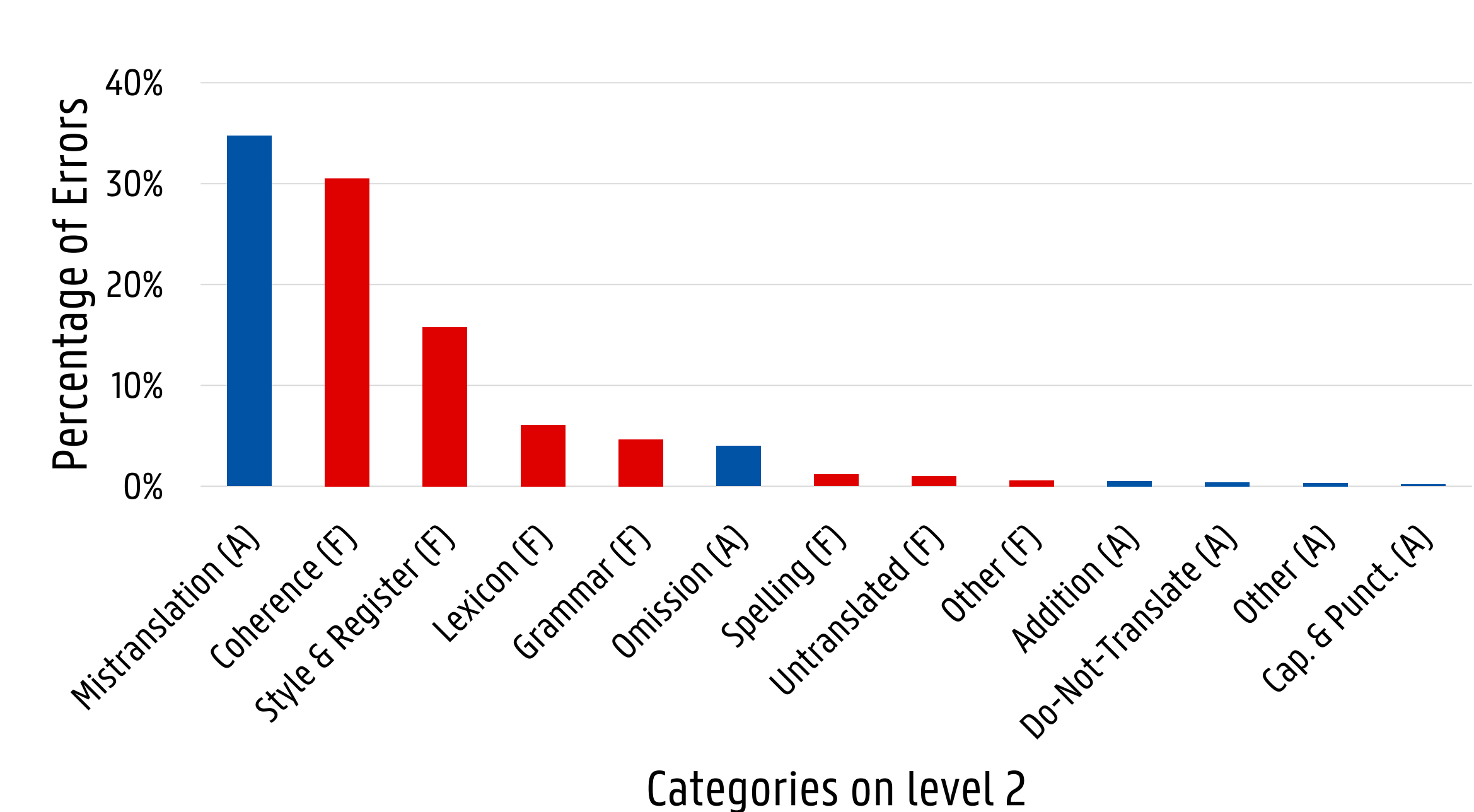
- 99% of the overlapping annotations agree on level 1 of the taxonomy, 90% on level 2, and 73% on level 3
- Most pronounced disagreement on level 2: 'coherence' vs. 'style & register' (5% of the annotation pairs agreeing on level 1)

Error Analysis

Overall Quality

- 44% of the sentences do not contain any errors
- On average, the erroneous source sentences are longer than the correct ones (17.11 vs. 9.56 words)
- Performance decreases from a source sentence length of 10 words onwards

Error Distribution

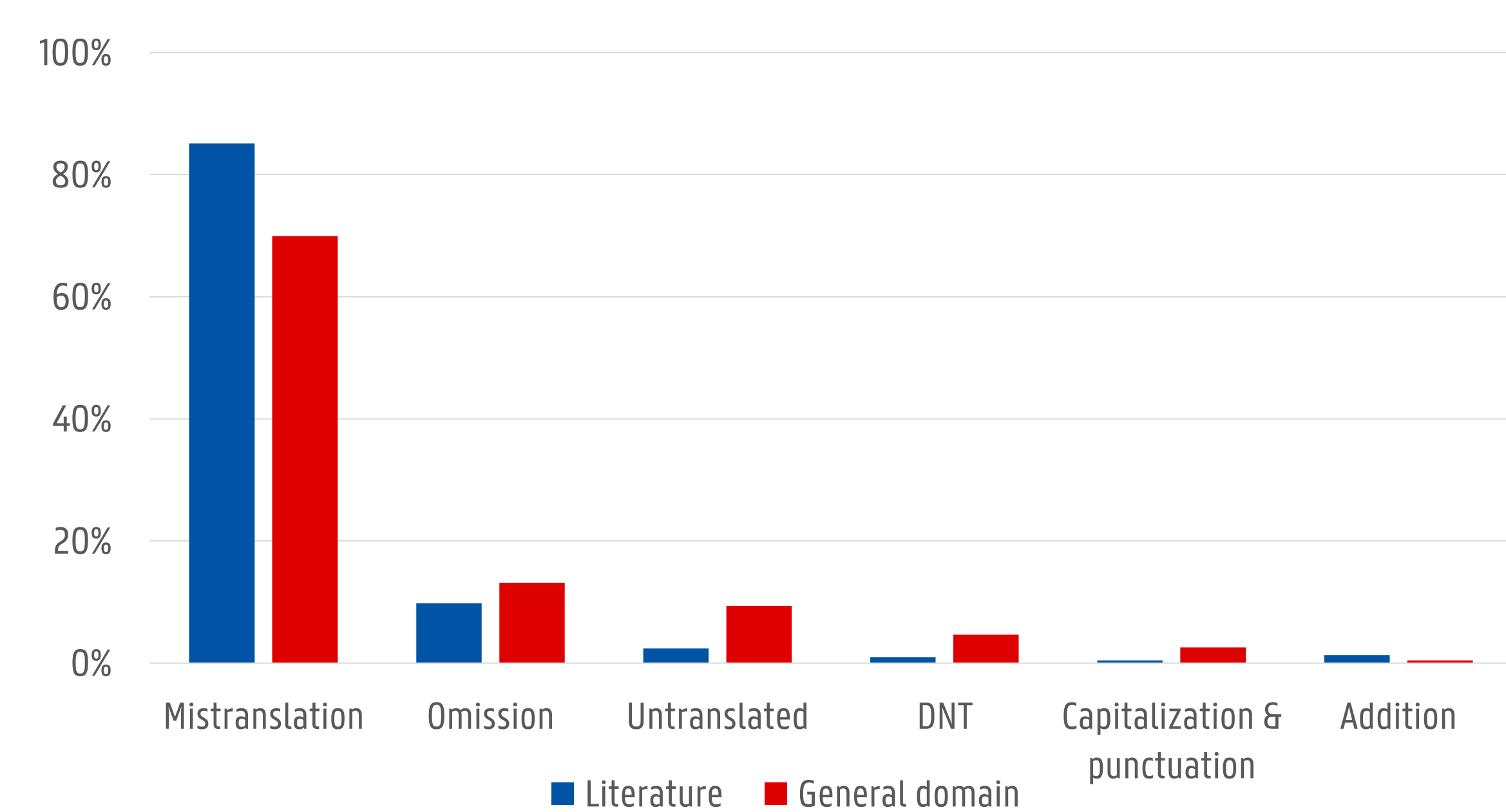


Co-Occurring Fluency & Accuracy Errors

Hypothesis: if accuracy errors tend to cause fluency errors, they will overlap regularly

- 27% of the fluency and 39% of the accuracy error spans fully overlap with the span of an error belonging to the other category
- 84% of those overlaps: 'coherence' x 'mistralation'

Accuracy Comparison with General Domain⁵



Conclusions

- Literary MT seems to be more promising for English-Dutch (44% correct sentences) than for English-Slovene (0%)⁶, English-Russian (17%)⁷ & German-English (33%)⁷
- Main problems for literary MT are 'mistralation', 'coherence' & 'style & register'
- Mistralations are the biggest accuracy issue for general-domain MT as well.
- Fluency & accuracy errors overlap regularly, especially 'coherence' & 'mistralation' errors