



EMNLP 2021

7th – 11th November | Online and in the Dominican Republic

Pushing the Right Buttons: Adversarial Evaluation of Quality Estimation

Diptesh Kanojia, Marina Fomicheva, Tharindu Ranasinghe,
Frédéric Blain, Constantin Orăsan, Lucia Specia



Roadmap

Motivation - “The Why?”

Key Contributions - “The What?”

Dataset

Probing Strategies - “The How?”

- Meaning-preserving Perturbations (MPPs)

- Meaning-altering Perturbations (MAPs)

- Quality Estimation Models

Results

Conclusions

Motivation

Quality Estimation (QE) - the task of predicting the quality of Machine Translation (MT) output in the absence of human reference translation.

Despite efforts towards NMT, important meaning errors in Machine Translation output still exist!

Can the QE models detect these meaning errors?

Key Findings

- SOTA QE models are robust to MPPs and are sensitive to MAPs.
- SOTA QE models fail to properly detect certain types of MAPs, such as negation omission.
- Our results on a set of QE models are consistent with their correlation with human judgements.

Dataset & Language Pairs

Dataset:

WMT 2020 Quality Estimation Shared Task 1

Language Pair (LP):

Russian (Ru) - English (En)

Romanian (Ro) - English (En)

Estonian (Et) - English (En)

Sinhala (Si) - English (En)

Nepali (Ne) - English (En)

| Language Pair | Ru-En | Ro-En | Et-En | Si-En | Ne-En |
|---------------|-------|-------|-------|-------|-------|
| #sentences | 1245 | 1035 | 766 | 404 | 100 |

Meaning-preserving Perturbations (MPPs)

Meaning-preserving Perturbation (MPP): a small change in the target-side translation that might affect the translation but does not affect the meaning of the sentence.

MPP1: Removal of Punctuations.

MPP2: Replacing Punctuations.

MPP3: Removal of Determiners.

MPP4: Replacing Determiners.

MPP5: Changing random words to UPPERCASE.

MPP6: Changing random words to lowercase.

Meaning-altering Perturbations (MAPs)

Meaning-altering Perturbation (MAP)

a change in the target-side translation which affects the overall meaning of the sentence.

MAP1: Removal of Negation Markers

MAP2: Removal of Random Content Words

MAP3: Duplication of Content Words

MAP4: Insertion of Content Words

MAP5: Replacing Content Words.

MAP6: BERT-based Sentence Replacement.

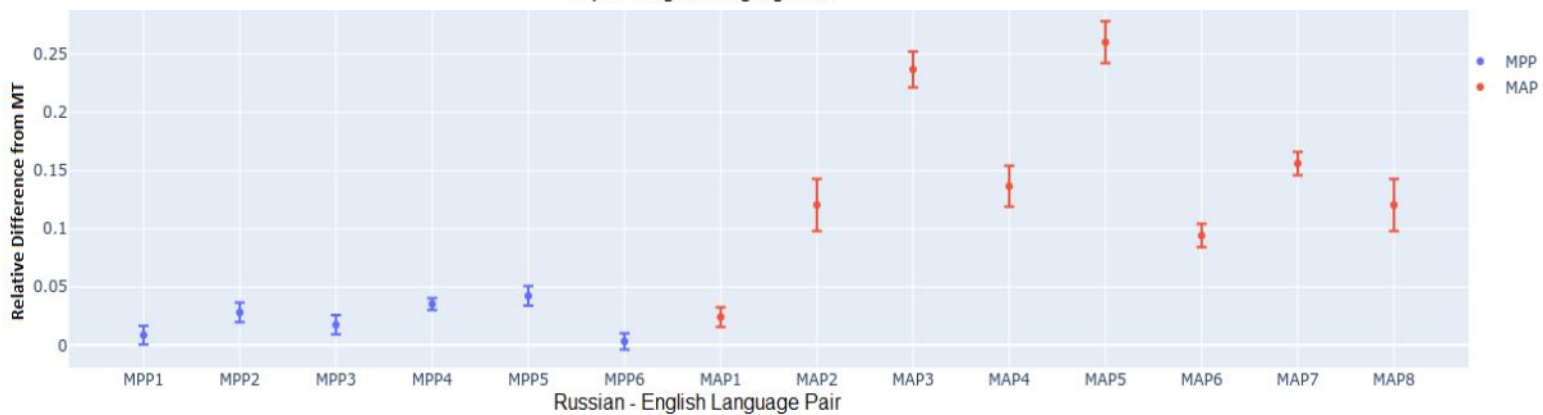
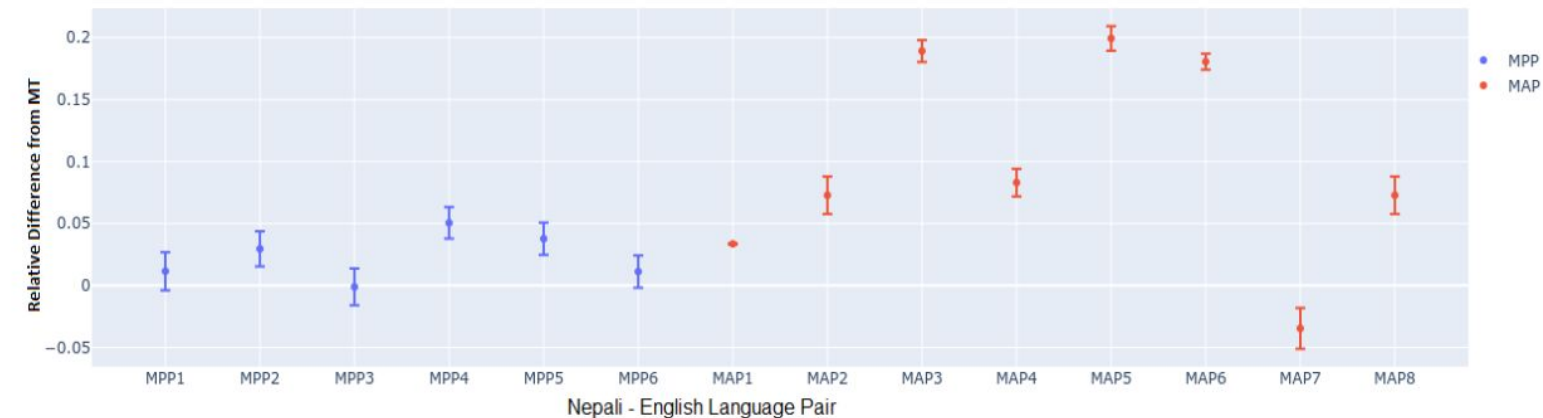
MAP7: Replacing word with Antonyms.

MAP8: Source-sentence as Target.

Quality Estimation Models

- MonoTransQuest (MonoTQ)
- SiameseTransQuest (SiameseTQ)
- MultiTransQuest (MultiTQ)
- Predictor-Estimator (OpenKiwi)
- SentSim (Unsupervised)

Do QE Models fail to detect MAPs?



Do perturbations affect SOTA QE Models?

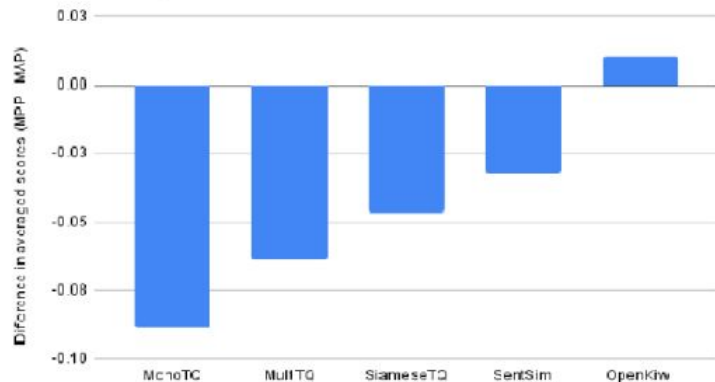
| | Ru-En | | | Ro-En | | | Et-En | | | Si-En | | | Ne-En | | |
|-----------|-------|-------------|-------------|-------|-------------|-------------|-------|-------------|-------------|-------|------|-------------|-------|-------------|-------------|
| | MT | MPP | MAP | MT | MPP | MAP | MT | MPP | MAP | MT | MPP | MAP | MT | MPP | MAP |
| MonoTQ | 0.81 | 0.78 | 0.66 | 0.82 | 0.80 | 0.74 | 0.81 | 0.79 | 0.73 | 0.71 | 0.65 | 0.64 | 0.75 | 0.74 | 0.68 |
| SiameseTQ | 0.86 | 0.85 | 0.86 | 0.58 | 0.57 | 0.52 | 0.92 | 0.91 | 0.91 | 0.58 | 0.57 | 0.52 | 0.68 | 0.68 | 0.65 |
| MultiTQ | 0.79 | 0.75 | 0.68 | 0.79 | 0.74 | 0.66 | 0.77 | 0.73 | 0.66 | 0.62 | 0.58 | 0.52 | 0.63 | 0.60 | 0.52 |
| OpenKiwi | 0.78 | 0.78 | 0.78 | 0.78 | 0.75 | 0.77 | 0.71 | 0.70 | 0.70 | 0.62 | 0.60 | 0.57 | 0.50 | 0.48 | 0.48 |
| SentSim | 0.54 | 0.57 | 0.57 | 0.78 | 0.76 | 0.72 | 0.50 | 0.53 | 0.52 | 0.41 | 0.43 | 0.41 | 0.47 | 0.52 | 0.50 |

Table 4 from the paper which shows average predicted scores by all the QE models on the test set for the unperturbed machine translation (MT), versus with meaning-preserving perturbations (MPP) and meaning-altering perturbations (MAP).

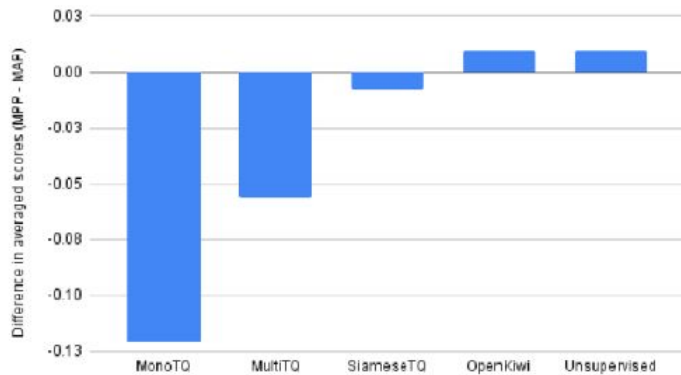
The lowest average scores (MPP/MAP) are boldfaced in each case, if lower than MT.

Can we use perturbations to rank QE models?

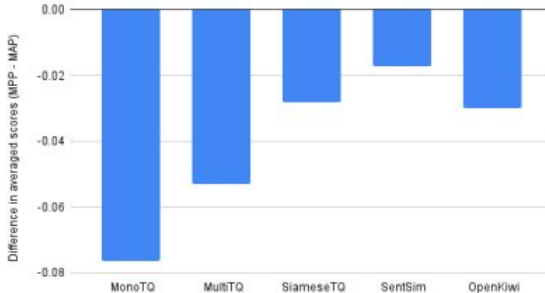
Romanian-English



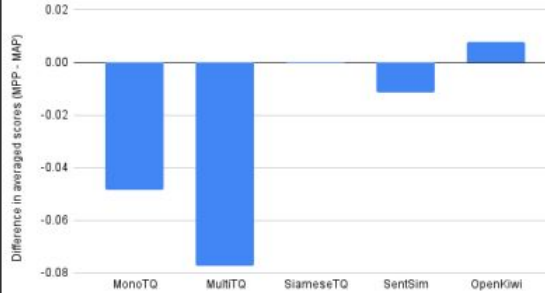
Russian-English



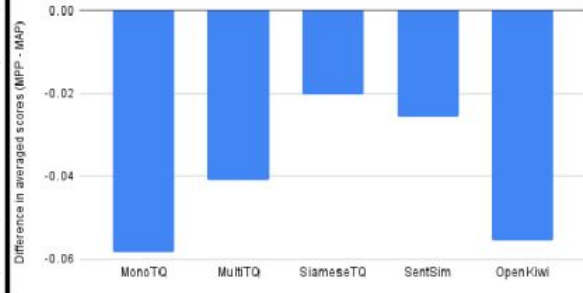
Nepali-English



Estonian-English



Sinhala-English



Conclusion and Future Work

- Probing the robustness of QE models.
- A perturbations-based method to detect failures of a QE model.
- Overall, predictive of the performance of a QE model.
- A method which does not rely on manual annotations.
- QE model ranking with this method.

Thank You!

Questions? :)

