

A Post-Editing Dataset in the Legal Domain: Do we Underestimate Neural Machine Translation Quality?

**Julia Ive¹, Lucia Specia^{1,2}, Sara Szoc³, Tom Vanallemeersch³, Joachim Van den Bogaert³,
Eduardo Farah⁴, Christine Maroti⁴, Artur Ventura⁴, Maxim Khalilov⁴**

¹ Department of Computer Science, University of Sheffield, UK,

² Department of Computing, Imperial College London, UK,

³ CrossLang, Belgium

⁴ Unbabel, Portugal

j.ive@sheffield.ac.uk, l.specia@imperial.ac.uk

first.lastname@crosslang.com

eduardo.farah@unbabel.com, christine@unbabel.com, artur@unbabel.com, maxim@unbabel.com

Abstract

We introduce a machine translation dataset for three pairs of languages in the legal domain with post-edited high-quality neural machine translation and independent human references. The data was collected as part of the EU APE-QUEST project and comprises crawled content from EU websites with translation from English into three European languages: Dutch, French and Portuguese. Altogether, the data consists of around 31K tuples including a source sentence, the respective machine translation by a neural machine translation system, a post-edited version of such translation by a professional translator, and – where available – the original reference translation crawled from parallel language websites. We describe the data collection process, provide an analysis of the resulting post-edits and benchmark the data using state-of-the-art quality estimation and automatic post-editing models. One interesting by-product of our post-editing analysis suggests that neural systems built with publicly available general domain data can provide high-quality translations, even though comparison to human references suggests that this quality is quite low. This makes our dataset a suitable candidate to test evaluation metrics. The data is freely available as an ELRC-SHARE resource.

Keywords: Machine Translation, Post-editing Dataset, Legal Domain, Automatic Post-Editing, Quality Estimation

1. Introduction

Current state-of-the-art (SOTA) in Neural Machine Translation (NMT) has reached remarkable progress. This has also pushed the boundaries of manual and automatic evaluation procedures relying on human reference translations. For example, recent studies have shown that reference translations are often judged by humans as having lower translation quality than top NMT systems (Hassan et al., 2018), with follow-up studies showing that this is partly due to the limited quality of the human translations (Toral et al., 2018). Even though problems with independently collected (single) human reference translations for evaluation have been highlighted in the past (Fomicheva and Specia, 2016), this practice is more questionable with high quality NMT systems. In these conditions, metrics based on human post-editing of the machine translations become particularly important. This feedback can be used to assess MT quality directly, as well as to build and benchmark metrics for the automatic evaluation of Machine Translation (MT) output, and to build Quality Estimation (QE) and Automatic PE (APE) models.

We describe a machine translation dataset in the legal domain resulting from activities performed in the framework of the APE-QUEST project (<http://ape-quest.eu>). This project aims to integrate MT, QE and APE. The dataset focuses on the areas of online dispute resolution (ODR), procurement and justice. It continues the well-established tradition in MT of using the legal data resulting from the EU procedures that was started with the Europarl corpus (Koehn, 2005). The data consists of around 31K tuples

including an English source sentence, the respective machine translation by an NMT system (into Dutch, French and Portuguese), a post-edited version of such translation by professional translators, and an independently created reference translation (31,403 cases all together). Interestingly, our English-Dutch, English-French and English-Portuguese machine translations require very few post-edits – as expected. However, when compared to independently created human references using standard metrics like BLEU (Papineni et al., 2002), TER (Snover et al., 2009) and METEOR (Denkowski and Lavie, 2014) the resulting figures indicate rather low translation quality.

Existing PE datasets, e.g., the Autodesk dataset (Zhechev, 2012), are in their majority generated by editing the output of statistical MT systems and are no longer useful for NMT systems as the errors and nature of required corrections are different. The PE NMT datasets used in the annual WMT shared APE and QE tasks are created using the previous generation of NMT systems (mostly RNN-based, which exhibit inferior quality to the current systems) and/or cover only the IT or life sciences domains (Specia et al., 2017; Chatterjee et al., 2019). We propose a PE dataset built using translations from SOTA neural architectures and for a new domain. Finally, the language pairs we propose are also new (current NMT PE datasets propose English-German, English-Latvian, English-Italian and English-Russian translations).

In the remainder of this paper we first describe our data sources (Section 2) and the MT systems built (Section 3) to translate this data. We present the PE process and its results

and qualitative analysis in Section 4. Section 5 demonstrates two use cases of the dataset.

2. Post-Editing Data

The APE-QUEST project aims to align with the requirements of DSIs (Digital Service Infrastructures)¹ of the European Commission (EC) regarding the use of the translation system developed by the EC (CEF eTranslation system²). More specifically, the project’s use case deals with translating and post-editing data in the legal domain, which involves the ODR DSI, eProcurement DSI and eJustice DSI.

For the English-Dutch language pair, we first collected in-domain data using the publicly available eProcurement dataset³ and manually collected data from the ODR website.⁴ Using XenC (Rousseau, 2013), we trained a language model from this data, and carefully selected a subset from large generic-domain parallel corpora,⁵ i.e. the subset with the lowest perplexity according to the language model. Duplicates and very short sentences (under 5 words) were removed from those best-scored examples. A random subset of the resulting source and reference files was additionally manually cleaned, and sent for post-edition (see Section 4). For English-French and English-Portuguese, the datasets are composed of data which we scraped from the e-Justice website.⁶ We document-aligned, sentence-aligned and automatically cleaned the data, using the tools Malign,⁷ hunalign⁸ and Bicleaner.⁹ We then selected a random subset of the resulting data for the post-edition tasks.

3. Neural Machine Translation Systems

We trained a machine translation system for each language pair. For each of the pairs, we used NMT implementations from different toolkits: OpenNMT TensorFlow,¹⁰ tensor2tensor¹¹ and Marian.¹² This adds diversity to our dataset and provides an opportunity for comparative studies.

All three systems are built with publicly available data from the OPUS repository (Tiedemann, 2009). Table 1 summarises the statistics and the origin of the MT training data. For all the models we apply the BPE word segmentation

¹These infrastructures deliver networked cross-border services for citizens, businesses and public administrations.

²<https://ec.europa.eu/cefdigital/wiki/display/CEFDIGITAL/eTranslation>

³<https://github.com/ESPD/ESPD-Service/tree/master/espd-web/src/main/resources/i18n>

⁴<https://ec.europa.eu/consumers/odr>

⁵We took the English-Dutch parts of the Paracrawl Bicleaner v4.0 (<https://paracrawl.eu/releases.html>) and EUbookshop (<http://opus.nlpl.eu/EUbookshop.php>) corpora.

⁶<https://e-justice.europa.eu>

⁷<https://github.com/paracrawl/Malign>

⁸<http://mokk.bme.hu/en/resources/hunalign>

⁹<https://github.com/bitextor/bicleaner>

¹⁰<https://github.com/OpenNMT/OpenNMT-tf>

¹¹<https://github.com/tensorflow/tensor2tensor>

¹²<https://github.com/marian-nmt/marian>

approach (Sennrich et al., 2016) with 32K merge operations for EN-NL (32K subtokens joint vocabulary), with 90K merge operations for EN-FR (about 70K subtokens joint vocabulary after filtering the least frequent out) and with about 90K merge operations for EN-PT (about 90K subtokens joint vocabulary).

We train the EN-NL model for 25 iterations, using the `transformer_small` training parameters, a learning rate of 2.0 and 8K warmup steps. We average the last 8 checkpoints to obtain the final model. For EN-FR, we use the `transformer_big` parameters with a learning rate of 0.05 with 8K warmup steps. For EN-PT, we train an amun RNN model with default parameters for 4 epochs validating every 1,000 updates. The usage of the RNN architecture promotes the diversity of SOTA outputs in our dataset.

4. Post-Editing Process

Post-editing was done by professional translators, each being assigned a different portion of the data. They were asked to only correct actual errors and refrain from making stylistic improvements. The resulting data triplets comprise 10-11k sentences for each language pair. The statistics summarising the post-editing datasets are reported in Table 2. For the computation of statistics and automatic scores, the data was tokenised using the Moses toolkit scripts (Koehn et al., 2007).

Table 3 measures the edit distance (HTER), HBLEU and HMETEOR between MT and PE, and TER, BLEU and METEOR between MT and the independent reference (REF). HTER is defined as the minimum number of edits (substitution, insertion, deletion and shift) required to change an MT hypothesis so that it exactly matches a human post-edition of this hypothesis. HBLEU measures n -gram precision between MT hypotheses and post-edits, whereas HMETEOR – unigram precision and recall. The human-targeted HTER/HBLEU/HMETEOR (as compared to TER/BLEU/METEOR) variants are measured using actual post-edited MT rather than independent references. The data contains very few edits (11 HTER / 83 HBLEU / 89 HMETEOR on average).

Distributions over HTER bins confirm the overall high quality of the NMT outputs (Figure 1): around 67% of the sentences per language pair belong to the HTER bin between 0 and 10, and thus have minimum edits required. Moreover, the proportion of sentences requiring no correction at all (0 HTER) is rather high: 52% for EN-NL, 49% for EN-FR and 37% for EN-PT.

An interesting observation is the gap between the MT-PE HTER/HBLEU/HMETEOR and MT-REF TER/BLEU/METEOR scores (the difference (Δ) in (H)TER reaches 39 absolute points on average across languages), indicating extreme bias of popular automatic evaluation in this case. Note that these tendencies were not previously observed for the data translated by statistical MT (SMT) or by low-quality NMT: e.g., Specia et al. (2017) reports on average $\Delta 15$ (H)TER for SMT systems, $\Delta 11$ (H)TER for the low-quality NMT system. For the high-quality EN-DE language pair, Specia et al. (2017) reports a $\Delta 30$ (H)TER. We manually inspected the references and concluded that they are adequate and in

lang	# sent	corpora	model
EN-NL	40M	ECB, DGT, Europarl, EU, JRC, GlobalVoices, OpenSubtitles, NewsCommentary	transformer_small, OpenNMT Tensorflow
EN-FR	12M	DGT, ECB, Europarl, JRC-Acquis, subset of MultiUN	transformer_big, tensor2tensor
EN-PT	12M	ECB, DGT, Europarl, JRC, GlobalVoices, EMEA, EUbookshop, NewsCommentary	amun, Marian

Table 1: NMT training data statistics: Total number of sentences, corpora.

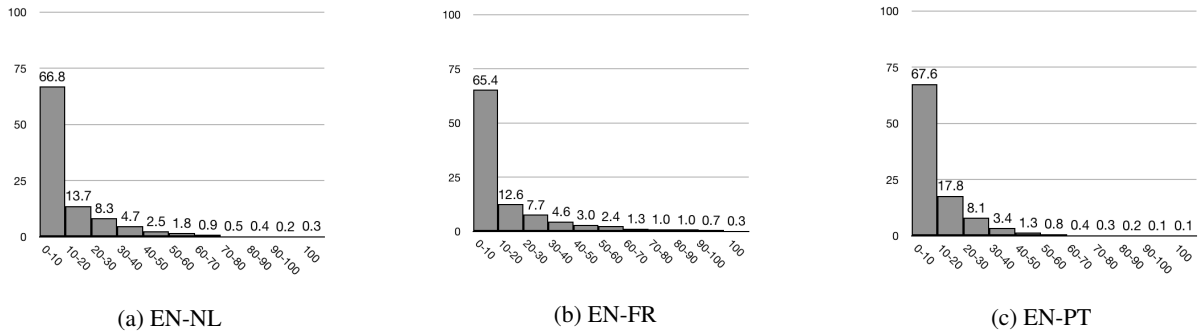


Figure 1: X-axis plots HTER bins. Y-axis – respective percentage of the dataset sentences.

lang	# sent	# tok., SRC	# tok., MT	# tok., PE
EN-NL	11,249	173,161	169,225	171,229
EN-FR	9,989	236,458	243,734	254,386
EN-PT	10,165	302,687	302,913	306,200

Table 2: Dataset statistics: Total number of sentences, average number of tokens in source, translation and post-edited sentences (using the Moses tokeniser).

	metric	EN – NL	EN – FR	EN – PT
MT-PE	HTER↓	9.8	13.2	10.1
	HBLEU↑	84.0	82.1	83.3
	HMETEOR↑	90.2	87.5	89.6
MT-REF	TER↓	48.4	52.9	49.7
	BLEU↑	31.4	32.1	33.0
	METEOR↑	48.3	51.6	50.4

Table 3: Edit distance (HTER↓), HBLEU↑ and HMETEOR↑ between PE and original MT, as well as TER↓, BLEU↑ and METEOR↑ between MT and independent references.

general have good quality. The issue is the fact that they differ significantly in phrasing from the NMT outputs and this type of variance significantly hinders reference-based automatic (or manual) evaluation.

Examples in Table 4 confirm that MT (requiring little or no edits) and REF convey the same meaning using different words.

5. Use Cases

In what follows we summarise the outcomes of benchmarking the presented datasets in SOTA Quality Estimation and Automatic Post-Editing tasks.

5.1. Quality Estimation

Quality Estimation (QE) (Blatz et al., 2004; Specia et al., 2018) predicts the quality of MT when automatic evaluation or human assessment is not possible (typically at system run-time). QE is mainly addressed as a supervised machine learning problem using quality-labelled data. Usually both source text and MT output serve as an input for a model that predicts a score for unseen MT units. Very often this score is HTER like in WMT QE shared tasks (Fonseca et al., 2019). Predictions for various types of units are possible: documents, paragraphs, sentences, words and phrases, with sentence-level predictions being the most common.

OpenKiwi QE Systems Our models utilise OpenKiwi (Kepler et al., 2019), an open-source framework for QE that implements a range of QE systems from the WMT 2015-18 shared tasks. We extend it to leverage recently proposed pre-trained models via transfer learning techniques. We follow OpenKiwi’s Predictor-Estimator implementation (Kim et al., 2017). Predictor-Estimator is a modular architecture that revolves around an encoder-decoder architecture (so-called Predictor), stacked with a bidirectional RNN (so called Estimator) that is trained to produce quality estimates. It predicts quality using the weights assigned by the Predictor to the words we seek to evaluate, which are concatenated with the representations of their left and right one-word contexts, and then used to feed the Estimator.

The training data for the Predictor is the same as the respective NMT systems, while the Estimator is trained on labelled PE data. To build the model, we randomly sample

FR	
MT-PE	Cette section les relie en leur fournissant des informations aux niveaux européen et national. 'This section links them by providing information at European and national levels.'
REF	Cette section offre un lien vers ces registres en fournissant des informations au niveau européen et national. 'This section provides a link to these registers by providing information at European and national level.'
NL	
MT-PE	Door de aanbestedende dienst toegekend registratienummer 'Allocated by the contracting authority registration number '
REF	Referentienummer van het dossier bij de aanbestedende dienst 'Reference number of the file at the contracting authority'
PT	
MT	Quando as regras em causa são favoráveis a uma parte, a parte vencida é condenada nas despesas judiciais por parte da parte vencida. 'Where the rules in question are favorable to a party, the unsuccessful party is to be ordered to pay the unsuccessful party's legal costs.'
PE	Quando as regras em causa são favoráveis a uma das partes, a parte vencida é condenada a pagar as despesas judiciais da parte vencedora. 'Where the rules in question are favorable to one of the parties, the unsuccessful party is ordered to pay the winning party's legal costs.'
REF	Quando o tribunal decide a favor de uma parte, ordena à parte vencida que pague à parte vencedora todas as custas judiciais suportadas pela parte em causa. 'When the court rules in favor of a party, it orders the losing party to pay the winning party all legal costs borne by the party concerned.'

Table 4: Examples of minor corrections performed to MT outputs. Independent reference expresses this meaning in different words.

from the human PE data 500 triplets for validation and 536 triples for test. The rest of the PE is used as the QE training data.

Results Table 5 reports results of our experiments. The best quality of the prediction ($r = 0.58$) is achieved for FR, what we tend to attribute to the fact that the quality of FR MT is slightly lower than for NL and PT, which makes the QE task less hard in this case. Following the shared task setup, Pearson’s r correlation coefficient is used as the primary evaluation metric for the scoring task (with Mean Absolute Error – MAE – as the secondary metric).

lang	r	MAE
EN–NL	0.38	0.14
EN–FR	0.58	0.14
EN–PT	0.38	0.08

Table 5: Pearson’s r correlation coefficient and MAE scores for the sentence-level QE task measured for the internal test set (536 sentences).

The datasets described in this paper open new avenues for research in QE, particularly in this challenging setting where a large proportion of the translations require no edits.

5.2. Automatic Post-Editing

Automatic Post-Editing (APE) (Simard and Foster, 2013; Chatterjee et al., 2017) seeks to reduce the burden of human post-editors and automatically corrects errors in MT out-

puts. APE is usually performed by monolingual translation models that “translate” from the raw MT to PE. Inputs to those systems are source sentences and raw MT that is expected to be corrected in the output. Given the high quality of current NMT outputs, the task has become particularly challenging. This increases the chance of APE systems to overfit or overcorrect new inputs at test time.

copycat APE Systems For our APE task we apply the recently introduced `copycat` networks (Ive et al., 2019), a Transformer-based pointer network framework. In the dual-source setting, the network can generate new words or copy words from either the source language or the original machine translation. The network has been shown to be rather conservative and make very few corrections to good quality raw MT – as a result of learning predominantly to copy.

We follow the procedure in (Ive et al., 2019) and train a dual-source with double attention `copycat` model. Again following Ive et al. (2019), we mimic MT data using 500K from general in-domain corpora available at OPUS¹³ and pre-train our models on this data. Zero HTER sentences are removed from the PE training data.

The data is tokenised and truecased using the Moses toolkit scripts (Koehn et al., 2007). We apply the BPE with 90K merge operations trained on legal corpora from OPUS for all the language pairs. For the experiments, we randomly select 1K lines for validation and 1K lines for test from each dataset. In the oracle setting, we also remove all 0 HTER sentences from the test set as if we had access to a perfect

¹³DGT for FR, Europarl for PT and NL

QE system. This results in test sets of 500 sentences.

Results Table 6 reports results of our experiments. We show HTER as the main metric and HBLEU as the secondary metric, as in the WMT APE shared task (Chatterjee et al., 2019). To highlight the difficulty of the task, the highest improvement of 0.5 HTER is obtained for EN-FR. Results for the test sets with 0 HTER sentences removed are more encouraging: with the 0.9 HTER reduction / 1.1 HBLEU increase on average.

lang	do nothing	Copypcat
full test set (1K)		
EN-FR	12.9 / 82.3	12.5 / 82.4
EN-NL	8.6 / 86.4	8.6 / 86.3
EN-PT	7.9 / 86.7	7.8 / 86.9
reduced test set (500, no 0 HTER)		
EN-FR	23.5 / 67.8	21.8 / 69.6
EN-NL	18.6 / 70.5	18.2 / 70.8
EN-PT	15.4 / 76.4	14.9 / 77.6

Table 6: HTER↓ / HBLEU↑ scores for the APE task. Bold-face values mark best results.

The datasets described in this paper address the lack of high-quality APE datasets in the domain. As indicated by the results of the most recent WMT APE challenge (Chatterjee et al., 2019), APE of such high-quality data remains a challenge since none of the submissions this year was able to beat such a high-quality “do nothing” baseline. By releasing our dataset we hope to stimulate research on the subject.

6. Conclusions

In this paper we described a new dataset with source, MT, PE triplets for three different language directions in the framework of the APE-QUEST project. The dataset is publicly available so it can be used for a variety of research and user-oriented purposes: ELRC-SHARE resource ID 2654 “Post-editing corpus English to Dutch/French/Portuguese, legal domain”.¹⁴ We applied the dataset to quality estimation and automatic post-editing. The results were more promising for English-French, the language pair for which the quality of the underlying neural machine translation system was the poorest.

7. Acknowledgements

APE-QUEST is funded by the EC’s CEF Telecom programme (Connecting Europe Facility, project 2017-EU-IA-0151).

8. Bibliographical References

Blatz, J., Fitzgerald, E., Foster, G., Gandrabur, S., Goutte, C., Kulesza, A., Sanchis, A., and Ueffing, N. (2004).

¹⁴<https://elrc-share.eu/repository/browse/post-editing-corpus-english-to-dutchfrench-portuguese-legal-domain/0e67e4f258a511ea913100155d02670675d8eb48f8eb4131baa5f9f4152ad2d2>

Confidence Estimation for Machine Translation. In *Proc. of the International Conference on Computational Linguistics*, page 315.

Chatterjee, R., Gebremelak, G., Negri, M., and Turchi, M. (2017). Online automatic post-editing for MT in a multi-domain translation environment. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 525–535.

Chatterjee, R., Federmann, C., Negri, M., and Turchi, M. (2019). Findings of the WMT 2019 shared task on automatic post-editing. In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 11–28, Florence, Italy, August. Association for Computational Linguistics.

Denkowski, M. and Lavie, A. (2014). Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the EACL 2014 Workshop on Statistical Machine Translation*.

Fomicheva, M. and Specia, L. (2016). Reference bias in monolingual machine translation evaluation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 77–82, Berlin, Germany, August. Association for Computational Linguistics.

Fonseca, E., Yankovskaya, L., Martins, A. F. T., Fishel, M., and Federmann, C. (2019). Findings of the WMT 2019 shared tasks on quality estimation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 1–10, Florence, Italy, August. Association for Computational Linguistics.

Hassan, H., Aue, A., Chen, C., Chowdhary, V., Clark, J., Federmann, C., Huang, X., Junczys-Dowmunt, M., Lewis, W., Li, M., et al. (2018). Achieving human parity on automatic Chinese to English news translation. *CoRR*, abs/1803.05567.

Ive, J., Madhyastha, P., and Specia, L. (2019). Deep copypcat networks for text-to-text generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3218–3227, Hong Kong, China, November. Association for Computational Linguistics.

Kepler, F., Trénous, J., Treviso, M., Vera, M., and Martins, A. F. T. (2019). OpenKiwi: An Open Source Framework for Quality Estimation. In *Proc. of the Annual Meeting of the Association for Computational Linguistics, System Demonstration*. Association for Computational Linguistics.

Kim, H., Lee, J.-H., and Na, S.-H. (2017). Predictor-estimator using multilevel task learning with stack propagation for neural quality estimation. In *Proceedings of the Second Conference on Machine Translation*, pages 562–568, Copenhagen, Denmark, September. Association for Computational Linguistics.

Koehn, P., Hoang, H., Birch, A., Callison-Burch, C., Federico, M., Bertoldi, N., Cowan, B., Shen, W., Moran, C., Zens, R., Dyer, C., Bojar, O., Constantin, A., and Herbst,

- E. (2007). Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180.
- Koehn, P. (2005). Europarl: A Parallel Corpus for Statistical Machine Translation. In *Proceedings of the X Machine Translation Summit*, pages 79–86.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. (2002). Bleu: a method for automatic evaluation of machine translation. In *Proceedings of 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318.
- Rousseau, A. (2013). XenC: An open-source tool for data selection in natural language processing. *The Prague Bulletin of Mathematical Linguistics*, (100):73–82.
- Sennrich, R., Haddow, B., and Birch, A. (2016). Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August. Association for Computational Linguistics.
- Simard, M. and Foster, G. (2013). PEPr: Post-edit propagation using phrase-based statistical machine translation. In *Proceedings of the XIV Machine Translation Summit*.
- Snover, M., Madnani, N., Dorr, B. J., and Schwartz, R. (2009). Fluency, adequacy, or HTER?: Exploring different human judgments with a tunable MT metric. In *Proceedings of the Fourth Workshop on Statistical Machine Translation*, pages 259–268.
- Specia, L., Harris, K., Blain, F., Burchardt, A., Mackentanz, V., Skadiņa, I., Negri, M., and Turchi, M. (2017). Translation quality and productivity: A study on rich morphology languages. In *Proceedings of the 16th Machine Translation Summit (MT Summit XVI)*, pages 282–299.
- Specia, L., Scarton, C., and Paetzold, G. H. (2018). Quality Estimation for Machine Translation. *Synthesis Lectures on Human Language Technologies*, 11(1):1–162.
- Tiedemann, J. (2009). News from OPUS - A collection of multilingual parallel corpora with tools and interfaces. In *Recent Advances in Natural Language Processing*, volume V, pages 237–248.
- Toral, A., Castilho, S., Hu, K., and Way, A. (2018). Attaining the unattainable? reassessing claims of human parity in neural machine translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 113–123, Brussels, Belgium, October. Association for Computational Linguistics.
- Zhechev, V. (2012). Machine translation infrastructure and post-editing performance at autodesk. In *AMTA 2012 Workshop on Post-Editing Technology and Practice (WPTP 2012)*.