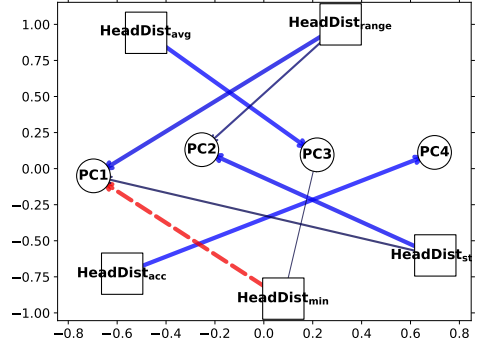
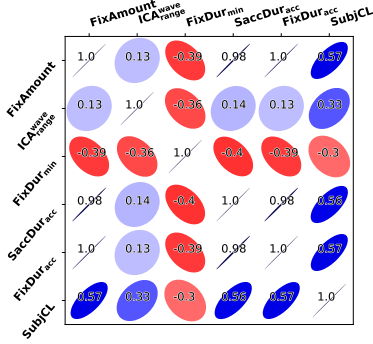


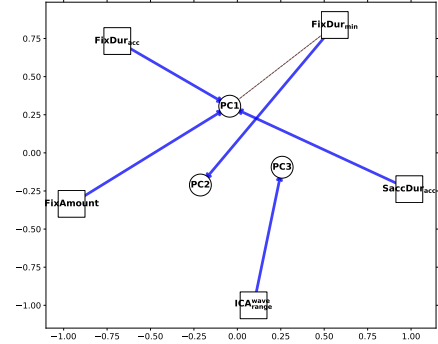
(a) Body Posture – Pearson



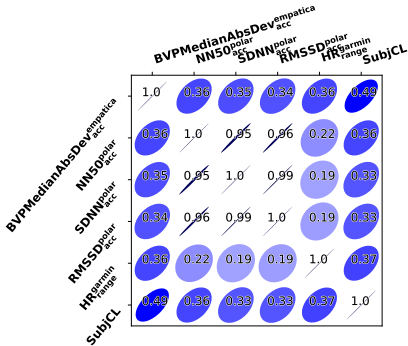
(b) Body Posture – PCA



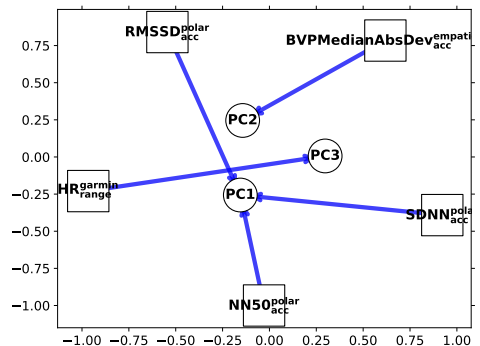
(c) Eyes – Pearson



(d) Eyes – PCA



(e) Heart – Pearson



(f) Heart – PCA

Figure 13: Correlations and PCA for body posture, eye, and heart modalities.

Table 1.3: Pairwise comparisons between the feature selected models *without LMEM/without participant and segment* (Table 1.1). * shows significance with $p < 0.05$ after Bonferroni-Holm correction. \tilde{t} is the test statistics for the modified paired t -test (Dietterich 1998).

Features	\tilde{t}
Time vs. combined	-4.06 *
Text vs. combined	-6.03 *
Keyboard vs. combined	-5.35 *
Body posture vs. combined	-6.32 *
Eyes vs. combined	-0.98
Heart vs. combined	-1.42
Skin vs. combined	-1.34

and color shows the strength of the loading; blue continuous lines represent positive loadings, while red dashed lines indicate negative loadings. For space reasons, we only summarize the most interesting results, which are all statistically significant.

For the *time features*, we see that PeTime and LNPeTime correlate very strongly and load on the same PC, but also that both show strong correlations to SubjCL.

For the *text features*, there expectedly are very strong correlations (-0.9) between TER and BLEU and between HTER and HBLEU, where each pair also loads on the same PC. Furthermore, strong correlations can be observed between TER and HTER, as well as between BLEU and HBLEU.

For the *keyboard features*, we see a very strong correlation between APR and PWR, however, both load on distinct PCs. PWR correlates more strongly to SubjCL than APR, indicating that PWR is by itself a better estimator of SubjCL than APR.

As expected, the most relevant *eye features* FixAmount, SaccDur_{acc}, and FixDur_{acc} correlate by almost 1, load on the same PC, and strongly relate to SubjCL.

For the *heart features*, the correlations between NN50_{acc}^{polar}, SDNN_{acc}^{polar}, and RMSSD_{acc}^{polar} are again very close to 1, and the PCA plot nicely visualizes that they cluster together. BVPMedAbsDev shows the strongest correlation to SubjCL.

Inspecting the most relevant *skin features*, we see very strong correlations between FreqFrameGSR_{avg}^{64,Empatica} and Leda_{avg}, as well as medium to strong correlations between the frequency frame and SkinTemp_{acc}^{Garmin} features.

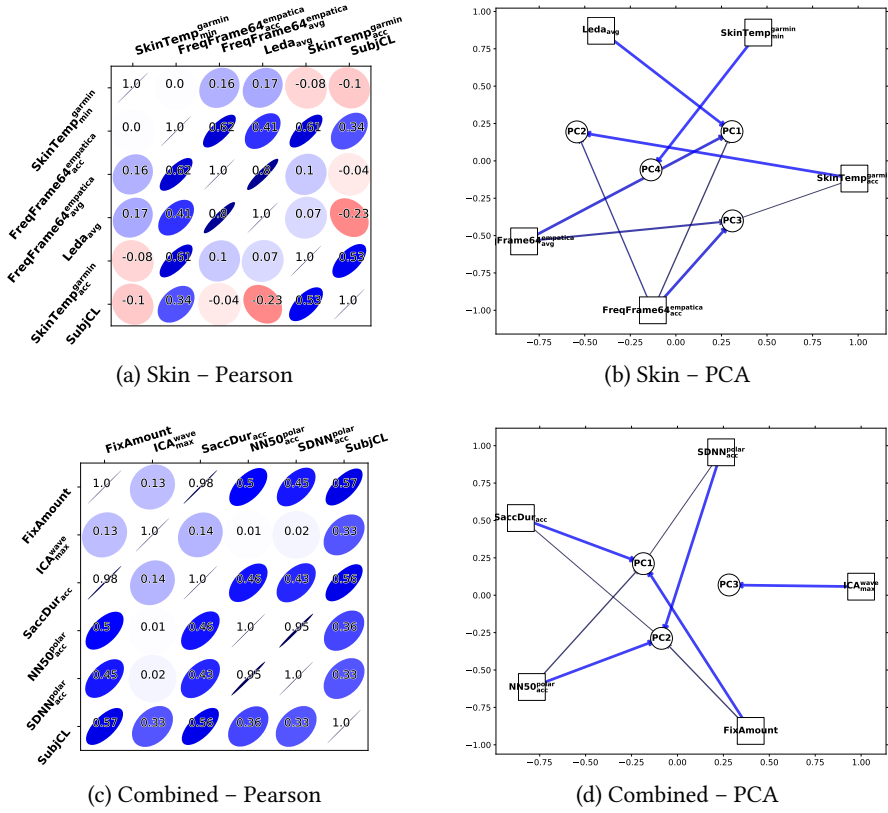


Figure 1.4: Correlations and PCA for the skin and combined modalities.

Most interestingly, for the *combined features* we can again see that $SDNN^{polar}_{acc}$ and $NN50^{polar}_{acc}$, as well as $FixAmount$ and $SaccDur_{acc}$, correlate with almost a value of 1. There also seems to be a strong link between the HRV measures and the eye measures $SaccDur_{acc}$ and $FixAmount$. The PCA further shows that there is one PC for the HRV measures, one for the ICA, and another one for the eye features $FixAmount$ and $SaccDur_{acc}$.

4.4 Discussion

Overall, very good regression results of up to 0.7 MSE on a 9-point scale were achieved by our regression models. This amount of error should be acceptable for most possible applications discussed in Herbig, Pal, van Genabith, et al. 2019. While the 5 by 2-fold CV results are often slightly worse, which might be be-

cause less training data was seen, the results of 1 by 10-fold and 5 by 2-fold are comparable, and the very small standard deviations indicate model robustness.

When comparing the regression results without adding participant and segment to Herbig, Pal, Vela, et al. (2019), whose approach is almost the same apart from having fewer sensors and features, we note a few similarities and differences: first of all, we found consistently better results across all modalities; however, already the baseline yields better results on our dataset. While the time features in Herbig, Pal, Vela, et al. (2019) were rather good, they are among the worst modalities here. A reason might be that we considered many more features, that helped the other modalities improve over the time as a feature. Furthermore, while in Herbig, Pal, Vela, et al. (2019) the eyes were by far the best among the three main categories eye, skin, and heart, all three show similar results here. This could be due to the numerous additional skin and heart features considered in our analysis. Whereas in both studies the combined approach leads to the best results, the performance gains when combining multiple modalities were much stronger in Herbig, Pal, Vela, et al. (2019), probably again because the three main categories are already very good by themselves.

So when we do not consider the individual participant and the segment they are post-editing (Table 1.1 or Herbig, Pal, Vela, et al. 2019), we can achieve the best results only with our main categories, eyes, heart, skin, or by combining features from several modalities. This is relevant for less controlled and more practical applications, e.g. adapting the user interface to perceived CL, where it is impossible to use participant and segment information, as ideally no two translators should post-edit the same sentence (which would otherwise be contained in TM).

In contrast, when we do consider participant and segment (Table 1.2), modalities of lesser quality, like time, text, keyboard, or body posture can also achieve good results. So considering *who is editing what* seems to yield enough information to learn from when combined with these features, while without considering participant and segment, the generalization is impeded. However, if the goal is to conduct a controlled experiment, e.g. to investigate the impact of different sentence features on subjectively felt CL, integrating participant and segment into the models allows to also achieve valuable estimates with these other modalities. The above experiment therefore also suggests that text quality, keyboard, and time measures, which are frequently used in the literature to estimate effort, only work well in controlled settings.

While we cannot compare all our correlation and PCA results to Vieira (2016), since we considered many more features, there is still some interesting overlap: The time features in both studies correlated strongly to SubjCL. Furthermore, the link between the PWR and SubjCL also seems comparable, while that between

APR and SubjCL appears weaker in our dataset. However, the correlation between these two keyboard features is similarly strong in both studies. The eye features FixAmount and FixDur also correlate to a similar extent with SubjCL in both studies. To summarize, we could both reproduce (except APR vs. SubjCL) and extend the findings by Vieira (2016), which strengthens our results.

The correlation and PCA especially revealed that many highly redundant features were selected by the feature selection approach (e.g. the HRV measures). The reason for this probably is their strong correlation to SubjCL; however, due to the redundancy, it is unclear whether incorporating multiple such features really helps. Therefore, we want to explore if handcrafting a set of features with fewer redundancies, or using a more sophisticated feature selection approach than RFECV, could boost the performance further. Since space constraints allowed us to analyze only very few features in terms of correlations and PCA, we also plan to investigate the link to the non-selected features, as well as a PCA including more features from all different modalities than the few reported here.

4.5 Limitations

The results presented in this study are subject to the following limitations: The data sample is relatively small, since only 10 subjects participated in our study. Next, while we performed CV and only report results on segments unseen during training, we did not completely leave out participants and then predict those participants' perceived CL from the data gathered by the other participants. Thus, to achieve these results in practice one may need to fine-tune and train for new users. Moreover, one should also note that our eye tracker only samples at 90 Hz, which could affect the peak velocity reconstruction and thereby saccades (Mack et al. 2017). Last, while our predictive approach yields interesting first insights, it is only an automatic "top-down" approach that might be improved by selecting an optimal set of features and tuning the hyper-parameters.

5 Conclusions and future work

In this paper, we have focused on perceived cognitive PE effort and argued for the need to robustly measure CL during PE. In contrast to most related work, we investigated whether and how multiple modalities to measure CL can be combined and used for the task of predicting the level of perceived CL during PE of MT. To the best of our knowledge, our analyzed feature set comprises the most

diverse set of features from a variety of modalities that has to date been investigated in the translation domain, considering even more factors than Herbig, Pal, Vela, et al. (2019).

Based on the data gathered from 10 professional translators, we report how well subjective CL can be predicted depending on the various features: When the models are unaware of which participant and segment the data belongs to, eye, skin, and heart features, or a combination of different modalities, performed best. In contrast, for regression models that can react differently depending on participant and segment, the less well performing categories time, text, keyboard, and body posture also achieved good results, probably due to overfitting on the participant. While this finding is very interesting for controlled experiments, it is less relevant for practical use, where no two participants should PE the same segment. Overall, the trained models can estimate CL during PE without interrupting the actual process through manual ratings with comparably low error of at best 0.7 MSE on a 9-point scale. However, further data analysis is needed to understand the required steps to achieve such results in practice.

We also report how strongly the different measures correlate and which features cluster together, where we reproduce almost all the findings of Vieira (2016) and extend them further by considering many more features.

In the future, we want to conduct more detailed investigations, e.g. in terms of a more complex feature selection approach or hand-crafting a subset of features based on the correlation and PCA findings, in combination with hyper-parameter tuning, to make better use of the available data than the chosen “top-down” regression approach. Furthermore, we want to use the captured continuous signals to already predict perceived CL while still editing the segment (i.e. based on a time window of the data), to allow for more real-time applications.

The long-term goal is to be able to decrease the perceived CL, and thereby stress and exhaustion, during PE. As discussed in Herbig, Pal, van Genabith, et al. (2019), this could be achieved by fine-tuning MT systems on the user’s CL measurements to produce less demanding outputs, or by automatically showing alternative translations or other forms of assistance. The measurement techniques explored within this paper form the basis for future research towards this goal.

Acknowledgments

This research was funded in part by the Deutsche Forschungsgemeinschaft (DFG) under grant number GE 2819/2-1/AOBJ: 636684. The responsibility lies with the authors.

References

- Arshad, Syed, Yang Wang & Fang Chen. 2013. Analysing mouse activity for cognitive load detection. In *Proceedings of the 25th Australian computer-human interaction conference: Augmentation, application, innovation, collaboration*, 115–118.
- Asteriadis, Stylianos, Paraskevi Tzouveli, Kostas Karpouzis & Stefanos Kollias. 2009. Estimation of behavioral user state based on eye gaze and head pose: Application in an e-learning environment. *Multimedia Tools and Applications* 41(3). 469–493.
- Benedek, Mathias & Christian Kaernbach. 2010. A continuous measure of phasic electrodermal activity. *Journal of Neuroscience Methods* 190(1). 80–91.
- Carl, Michael, Barbara Dragsted, Jakob Elming, Daniel Hardt & Arnt Lykke Jakobsen. 2011. The process of post-editing: A pilot study. *Copenhagen Studies in Language* 41. 131–142.
- Chen, Fang, Jianlong Zhou, Yang Wang, Kun Yu, Syed Z. Arshad, Ahmad Khawaji & Dan Conway. 2016. *Robust multimodal cognitive load measurement*. Cham: Springer International Publishing.
- Chen, Siyuan & Julien Epps. 2013. Automatic classification of eye activity for cognitive load measurement with emotion interference. *Computer Methods and Programs in Biomedicine* 110(2). 111–124.
- Daems, Joke. 2016. *A translation robot for each translator?: A comparative study of manual translation and post-editing of machine translations: Process, quality and translator attitude*. Ghent University. (Doctoral dissertation).
- Demberg, Vera & Asad Sayeed. 2016. The frequency of rapid pupil dilations as a measure of linguistic processing difficulty. *PloS One* 11(1). 1–29.
- Dietterich, Thomas G. 1998. Approximate statistical tests for comparing supervised classification learning algorithms. *Neural Computation* 10(7). 1895–1923.
- Doherty, Stephen, Sharon O'Brien & Michael Carl. 2010. Eye tracking as an MT evaluation technique. *Machine Translation* 24(1). 1–13.
- Gehring, Jonas, Michael Auli, David Grangier, Denis Yarats & Yann N. Dauphin. 2017. Convolutional sequence to sequence learning. In *Proceedings of the 34th international conference on machine learning*, vol. 70, 1243–1252.
- Goldberg, Joseph H. & Xerxes P. Kotval. 1999. Computer interface evaluation using eye movements: Methods and constructs. *International Journal of Industrial Ergonomics* 24(6). 631–645.
- Haapalainen, Eija, SeungJun Kim, Jodi F. Forlizzi & Anind K. Dey. 2010. Psychophysiological measures for assessing cognitive load. In *Proceedings of the 12th ACM international conference on ubiquitous computing*, 301–310.

- Herbig, Nico, Santanu Pal, Josef van Genabith & Antonio Krüger. 2019. Multi-modal approaches for post-editing machine translation. In *Proceedings of the 2019 CHI conference on human factors in computing systems*, 1–11.
- Herbig, Nico, Santanu Pal, Mihaela Vela, Antonio Krüger & Josef van Genabith. 2019. Multi-modal indicators for estimating perceived cognitive load in post-editing of machine translation. *Machine Translation* 33(1–2). 91–115.
- Hockey, Robert. 1997. Compensatory control in the regulation of human performance under stress and high workload: A cognitive-energetical framework. *Biological Psychology* 45(1). 73–93.
- Hossain, Gahangir & Mohammed Yeasin. 2014. Understanding effects of cognitive load from pupillary responses using Hilbert analytic phase. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 375–380.
- Iani, Cristina, Daniel Gopher & Peretz Lavie. 2004. Effects of task difficulty and invested mental effort on peripheral vasoconstriction. *Psychophysiology* 41(5). 789–798.
- Iqbal, Shamsi T., Xianjun Sam Zheng & Brian P. Bailey. 2004. Task-evoked pupillary response to mental workload in human-computer interaction. In *Extended abstracts on human factors in computing systems*, 1477–1480.
- Koglin, Arlene. 2015. An empirical investigation of cognitive effort required to post-edit machine translated metaphors compared to the translation of metaphors. *The International Journal for Translation & Interpreting* 7(1). 126–141.
- Koponen, Maarit. 2016. Is machine translation post-editing worth the effort? A survey of research into post-editing and effort. *The Journal of Specialised Translation* 25. 131–148.
- Kramer, Arthur F. 1991. Physiological metrics of mental workload: A review of recent progress. In Diane L. Damos (ed.), *Multiple-task performance*, chap. 11, 279–328. Oxfordshire: Taylor & Francis.
- Krings, Hans P. 2001. *Repairing texts: Empirical investigations of machine translation post-editing processes*. Vol. 5. Kent, OH: Kent State University Press.
- Lacruz, Isabel & Gregory M. Shreve. 2014. Pauses and cognitive effort in post-editing. In Sharon O'Brien, Laura W. Balling, Michael Carl, Michel Simard & Lucia Specia (eds.), *Post-editing of machine translation: Processes and applications*, chap. 11, 246–274. Newcastle upon Tyne: Cambridge Scholars Publishing.
- Lacruz, Isabel, Gregory M. Shreve & Erik Angelone. 2012. Average pause ratio as an indicator of cognitive effort in post-editing: A case study. In *AMTA workshop on post-editing technology and practice*, 21–30.

- Lavie, Alon & Abhaya Agarwal. 2007. Meteor: An automatic metric for MT evaluation with high levels of correlation with human judgments. In *Proceedings of the second workshop on statistical machine translation*.
- Lin, Chin-Yew & Franz J. Och. 2004. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In *Proceedings of the 42nd annual meeting of the Association for Computational Linguistics*. Barcelona, Spain.
- Macizo, Pedro & M. Teresa Bajo. 2006. Reading for repetition and reading for translation: Do they involve the same processes? *Cognition* 99(1). 1–34.
- Mack, David J., Sandro Belfanti & Urs Schwarz. 2017. The effect of sampling rate and lowpass filters on saccades—a modeling approach. *Behavior Research Methods* 49(6). 2146–2162.
- Mellinger, Christopher Davey. 2014. *Computer-assisted translation: An empirical investigation of cognitive effort*. Kent, OH: Kent State University.
- Moorkens, Joss, Sharon O’Brien, Igor A. L. Da Silva, Norma B. De Lima Fonseca & Fabio Alves. 2015. Correlations of perceived post-editing effort with measurements of actual effort. *Machine Translation* 29(3). 267–284. DOI: 10.1007/s10590-015-9175-2.
- Mulder, Lambertus J. M. 1992. Measurement and analysis methods of heart rate and respiration for use in applied environments. *Biological Psychology* 34(2). 205–236.
- O’Brien, Sharon. 2005. Methodologies for measuring the correlations between post-editing effort and machine translatability. *Machine Translation* 19(1). 37–58.
- O’Brien, Sharon. 2006a. Eye-tracking and translation memory matches. *Perspectives: Studies in translatology* 14(3). 185–205.
- O’Brien, Sharon. 2006b. Pauses as indicators of cognitive effort in post-editing machine translation output. *Across Languages and Cultures* 7(1). 1–21.
- Paas, Fred G.W.C., Juhani E. Tuovinen, Huib Tabbers & Pascal W.M. Van Gerven. 2003. Cognitive load measurement as a means to advance cognitive load theory. *Educational Psychologist* 38(1). 63–71.
- Paas, Fred G.W.C. & Jeroen J.G. van Merriënboer. 1994. Instructional control of cognitive load in the training of complex cognitive tasks. *Educational Psychology Review* 6(4). 351–371.
- Papineni, Kishore, Salim Roukos, Todd Ward & Wei-Jing Zhu. 2002. BLEU: A method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the ACL*, 311–318. Philadelphia, Pennsylvania.

- Rowe, Dennis W., John Sibert & Don Irwin. 1998. Heart rate variability: Indicator of user state as an aid to human-computer interaction. In *Proceedings of the conference on human factors in computing systems*, 480–487.
- Shaffer, Fred & Jay P. Ginsberg. 2017. An overview of heart rate variability metrics and norms. *Frontiers in Public Health* 5. 1–7.
- Shi, Yu, Natalie Ruiz, Ronnie Taib, Eric Choi & Fang Chen. 2007. Galvanic skin response (GSR) as an index of cognitive load. In *Extended abstracts on human factors in computing systems*, 2651–2656.
- Snover, Matthew, Bonnie Dorr, Richard Schwartz, Linnea Micciulla & John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the association for machine translation in the Americas*, 223–231.
- Snover, Matthew, Nitin Madnani, Bonnie Dorr & Richard Schwartz. 2009. Fluency, adequacy, or HTER? Exploring different human judgments with a tunable MT metric. In *Proceedings of the 4th workshop on statistical machine translation*, 259–268.
- Soukupova, Tereza & Jan Cech. 2016. Real-time eye blink detection using facial landmarks. In *21st computer vision winter workshop*, 1–8.
- Stuyven, Els, Koen Van der Goten, André Vandierendonck, Kristl Claeys & Luc Crevits. 2000. The effect of cognitive load on saccadic eye movements. *Acta Psychologica* 104(1). 69–85.
- Sweller, John, Jeroen J.G. van Merriënboer & Fred G.W.C. Paas. 1998. Cognitive architecture and instructional design. *Educational Psychology Review* 10(3). 251–296.
- van den Berg, Marten, Peter Rijnbeek, Maartje Niemeijer, Albert Hofman, Gerard van Herpen, Michiel Bots, Hans Hillege, Kees Swenne, Mark Eijgelsheim, Bruno Stricker & Jan Kors. 2018. Normal values of corrected heart-rate variability in 10-second electrocardiograms for all ages. *Frontiers in Physiology* 9. 1–9.
- Van Orden, Karl F., Wendy Limbert, Scott Makeig & Tzyy-Ping Jung. 2001. Eye activity correlates of workload during a visuospatial memory task. *Human Factors* 43(1). 111–121.
- Vieira, Lucas Nunes. 2014. Indices of cognitive effort in machine translation post-editing. *Machine Translation* 28(3–4). 187–216.
- Vieira, Lucas Nunes. 2016. How do measures of cognitive effort relate to each other? A multivariate analysis of post-editing process data. *Machine Translation* 30(1–2). 41–62.

- Villarejo, María Viqueira, Begoña García Zapirain & Amaia Méndez Zorrilla. 2012. A stress sensor based on galvanic skin response (GSR) controlled by ZigBee. *Sensors* 12(5). 6075–6101.
- Yamakoshi, Takehiro, Ken-ichi Yamakoshi, Shinobu Tanaka, Masamichi Nogawa, Sang-Bum Park, Mariko Shibata, Yoshito Sawada, Peter Rolfe & Yasuo Hirose. 2008. Feasibility study on driver's stress detection from differential skin temperature measurement. In *Engineering in medicine and biology society*, 1076–1079.

Chapter 2

Comparing NMT and PBSMT for post-editing in-domain formal texts: A case study

Sergi Álvarez^a, Toni Badia^a & Antoni Oliver^b

^aUniversitat Pompeu Fabra ^bUniversitat Oberta de Catalunya

This paper details a comparative analysis between phrase-based statistical machine translation (PBSMT) and neural machine translation (NMT) for English-Spanish in-domain medical documents using human rankings, fluency and adequacy, and post-editing (technical and temporal) effort, performed by professional translators. When MT output is ranked against translations performed by professional translators, results show a clear preference for human translations, with NMT in the second position. Regarding MT outputs, NMT is perceived as more fluent and conveying better the meaning of the source sentence. Despite this preference, post-editing temporal effort does not improve significantly in NMT compared to PBSMT, although technical effort is reduced.

1 Introduction

Over the last years, post-editing of machine translation (PEMT) has become common practice in the translation industry. It has been included as part of the translation workflow because it increases productivity and reduces costs (Guerberof 2009a). A recent survey showed that more than half of the language service providers (LSPs) offered PEMT as a service (Lommel & DePalma 2016). Post-editors “edit, modify and/or correct pre-translated text that has been processed by an MT system from a source language into (a) target language(s)” (Allen 2003).



Yet, many professional translators state that after post-editing a few MT segments, they delete the remaining segments and translate everything from scratch if they consider it will take them less time (Parra Escartín & Arcedillo 2015).

Effective PE, therefore, requires sufficient quality of the MT output. The issue, then, is how to detect that a machine translation output is good enough to serve as input to PE. Very often, the usual automatic metrics do not always correlate to PE effort (Koponen 2016). Even translators' perception does not always match PE effort (Koponen 2012; Moorkens 2018). Research in this field has mainly focused on measuring the PE effort related to MT output quality (Guerberof 2009a,b; Specia 2011; 2010), productivity (O'Brien 2011; Parra Escartín & Arcedillo 2015; Plitt & Masselot 2010; Sanchez-Torron & Koehn 2016), translator's usability (Castilho et al. 2014; Moorkens & O'Brien 2013) and perceived PE effort (Moorkens et al. 2015).

Statistical machine translation (SMT) has been well established as the dominant approach in machine translation for many years. However, in the last few years, research has become more interested in neural machine translation after the computational limitations have been solved (Bahdanau et al. 2018; Cho et al. 2014). The first results obtained have been very successful in terms of quality, for example in WMT 2016 (Bojar et al. 2016), WMT 2017 (Bojar et al. 2017), and WMT 2018 (Bojar et al. 2018). These promising results have driven a technological shift from (phrase-based) statistical machine translation (SMT) to neural machine translation (NMT) in many translation industry scenarios.

All of the current research on post-editing machine translation output uses the division established by Krings (2001) regarding PE effort: temporal effort (time spent PE), technical effort (number of edits, often measured using keystroke analysis), and cognitive effort (usually measured with eye-tracking or think-aloud protocols). Even though no current measure includes all three dimensions, cognitive effort correlates with technical and temporal PE effort (Moorkens et al. 2015). In our experiments, we use automatic measures of both temporal and technical effort.

As this new approach to MT becomes more popular among LSPs and translators, it is essential to test what NMT can offer for PE in terms of quality compared to the results of PBSMT. Recent studies (Bentivogli et al. 2016; Castilho, Moorkens, Gaspari, Sennrich, et al. 2017; Toral & Sánchez-Cartagena 2017) have stated an improved quality of NMT for PE. In this paper, we continue in this direction, but we focus on in-domain formal documents, which are the ones usually post-edited by professional translators.

Our objectives with these experiments are threefold:

- Determine which MT method (PBSMT or NMT) yields better results for PE in-domain formal texts.
- Analyze the relation between human and automatic metrics for PE.
- Study translators perception as a prospective measure of PE effort.

In Section 2, we review previous work comparing SMT and NMT approaches. In Section 3 we describe the MT systems and the training corpus used. In Section 4 we include the automatic evaluation of the MT systems used. We give details about the methodology used for our experiments in Section 5. We explain the results obtained in Section 6 and, finally, we state the main conclusions and our plans for future work in Section 7.

2 Previous work

One of the first complete papers studying the impact of SMT and NMT in PE was Bentivogli et al. (2016). In it, they carry out a small scale study on post-editing NMT and SMT outputs of English to German translated TED talks. They conclude that NMT generally decreases the PE effort, but degrades faster than SMT with sentence length. One of the main strengths of NMT is the reordering of the target sentence.

Wu et al. (2016) evaluate the quality of NMT and SMT, in this case using BLEU (Papineni et al. 2002) and human scores for machine-translated Wikipedia entries. Results show that NMT systems outperform and improve the quality of MT results. Other studies have confirmed this diagnostics (Junczys-Dowmunt et al. 2016; Isabelle et al. 2017), as have the results of the automatic PE tasks at the Conference on Machine Translation (Bojar et al. 2016; 2017).

Toral & Sánchez-Cartagena (2017) broaden the scope of Bentivogli et al. (2016) adding different language combinations and metrics, and they conclude that although NMT yields better quality results in general, it is negatively affected by sentence length, and the improvement of the results is not always perceivable in all language pairs.

Castilho et al. (2017) discuss three studies using automatic and human evaluation methods. One of them includes in-domain formal texts for chemical patent titles and abstracts. In addition to the automatic metrics, two reviewers assess 100 random segments to rank the translations and to identify translation errors. Automatic evaluation doesn't give clear results, but the SMT system is ranked higher than NMT in human evaluation.

Castilho et al. (2017) report on a comparative study of PBSMT and NMT, with four language pairs and different automatic metrics and human evaluation methods. It highlights some strengths and weaknesses of NMT, which in general yields better results. The study focuses especially on PE and uses the PET interface (Aziz et al. 2012) to compare educational domain output from both systems using different metrics. They conclude that NMT reduces word order errors and improves fluency for certain language pairs, so fewer segments require PE, especially because there is a reduction in the number of morphological errors. However, they don't detect a decrease in PE effort nor a clear improvement in omission and mistranslation errors.

Our experiments study the differences of post-editing NMT and SMT outputs for formal in-domain texts. We compare the usual automatic scores for MT with direct and indirect PE effort metrics. Mainly, we study translators' perception regarding quality, and fluency and accuracy, and analyze temporal and technical post-editing effort.

3 MT systems and training corpus

3.1 MT systems

In order to help contextualise the results in our experiments, we have decided to use two MT systems as references to compare their results with the ones of the systems we trained. As reference MT systems, we have chosen Apertium (Forcada et al. 2011), a shallow transfer MT system, and Google Translate, a neural MT system for the English-Spanish language pair, which is the one we use in our experiments.

For training the PBSMT and neural MT systems we have used ModernMT (Germann et al. 2016) version 2.4. This version allows to train both statistical and neural MT systems. We have used the default options for this version. One of the salient characteristics of ModernMT is the fact that it can take into account the context of the sentence to be translated. In the evaluation results, we show figures for both cases: with and without taking the context into account. In the experiments we take context to be the previous and the next segment (except for the first and last segment, where we have taken into account the next and the previous segment only, respectively). Short contexts are usually enough to calculate the context vector used by ModernMT.

3.2 Data: Medical corpus

To train the system, we have compiled all of the publicly available corpora in the English-Spanish pair known to us. We have also created several corpora from websites with medical content:

- The EMEA¹ (*European Medicines Agency*) corpus.
- The IBECS² (*Spanish Bibliographical Index in Health Sciences*) corpus.
- Medline Plus:³ we have compiled our own corpus from the web and we have combined this with the corpus compiled in MeSpEn⁴.
- MSDManuals⁵ English-Spanish corpus, compiled for this project under permission of the copyright holders.
- Portal Clínic⁶ English-Spanish corpus, compiled by us for this project.
- The PubMed⁷ corpus.
- The UFAL Medical Corpus⁸ v1.0.

We have also treated as a corpus glossaries and glossary-like databases containing a lot of useful terms and expressions in the medical domain. Namely, we have used the English-Spanish glossary from MeSpEn, the 10th revision of the international statistical classification of ICD and SnowMedCT.

With all the corpora and glossaries we have created an in-domain training corpus of 2,836,580 segments and entries. We have split the corpus in two parts: 99% of the segments for training, and the remaining 1% for testing.

We have also used other general corpora for training the MT systems, namely the Scielo corpus, the Europarl corpus⁹ (Koehn 2005), Global Voices corpus¹⁰ and

¹<http://opus.npl.eu/EMEA.php>

²<http://ibecs.isciii.es>

³<https://medlineplus.gov/>

⁴<http://temu.bsc.es/mespen/>

⁵<https://www.msdmanuals.com/>

⁶<https://portal.hospitalclinic.org>

⁷<https://www.ncbi.nlm.nih.gov/pubmed/>

⁸https://ufal.mff.cuni.cz/ufal_medical_corpus

⁹<http://www.statmt.org/europarl/>

¹⁰<https://globalvoices.org/>

News Commentary. The IBECS, Scielo, Pubmed and a part of the MedlinePlus corpus have been obtained from the MeSpEn corpus¹¹ (Villegas et al. 2018).

In Table 2.1 the size of all corpora and glossaries used for training the MT systems are shown. The figures are calculated after eliminating all the repeated source segment – target segment pairs in the corpora.

Table 2.1: Size of the corpora and glossaries used to create the corpus to train the MT systems.

Corpus	Segments/Entries	Tokens eng	Tokens spa
EMEA	366,769	5,327,963	6,008,543
IBECS	628,798	13,432,096	14,879,220
MedLine Plus	15,689	209,074	234,660
MSD Manuals	241,336	3,719,933	4,467,906
Portal Clinic	8,797	159,717	169,294
PubMed	320,475	2,752,139	3,035,737
UFAL	258,701	3,202,162	3,437,936
Glossary MeSpEn	125,645	286,257	348,415
ICD10-en-es	5,202	25,460	30,580
SnowMedCT Denom.	887,492	3,509,062	4,457,681
SnowMedCT Def.	4,268	177,861	184,574
In-domain	2,836,580	32,479,955	36,893,257
Scielo	741,407	17,464,256	19,305,165
Europarl	1,961,672	50,008,219	52,489,142
Global Voices	559,418	10,717,938	11,496,683
News Commentary	259,412	5,898,912	6,903,975
Out-of-domain	3,521,363	84,087,899	90,193,659

4 Automatic evaluation of the MT systems

In Table 2.2 we can observe the evaluation values of the trained systems using MTEval¹² along with Apertium and Google Translate. This software allows to calculate BLEU, NIST, RIBES and WER using only one reference. We have used all

¹¹<http://temu.bsc.es/mespen/>

¹²<https://github.com/odashi/mteval>

the test sets of the corpus. As shown in the table, the systems trained in the experiment obtain better results in all metrics than the reference systems used, except for the Google Translate system, which obtains a slightly better NIST result than the MMT Phrase-Based system without context and a better WER result than the two MMT Phrase-Based systems. The MMT Neural system performs consistently better than the MMT Phrase-Based system. In the MMT Neural system, we do not see any significant difference between the results obtained when trained with or without context.

Table 2.2: Results of the automatic evaluation using mteval.

MT system	BLEU	NIST	RIBES	WER
Apertium	0.192577	6.442539	0.713117	0.702716
Google T.	0.402497	9.632268	0.809469	0.530053
MMT P.B. no context	0.424183	9.536248	0.814425	0.637821
MMT P.B. context	0.444832	9.801466	0.819303	0.621032
MMT Neural no context	0.503935	11.106222	0.836954	0.485474
MMT Neural context	0.505778	11.141294	0.836313	0.481039

5 Experiments

We carried out three different experiments with English-Spanish medical texts to assess human perception and evaluation of both PBSMT and NMT systems.

5.1 Translation ranking

In the first part, participants had to answer some questions about their previous experience in the translation industry. The survey was open both to students and professional translators as we were mainly interested in the perception of quality. In the second part of the survey, participants had to rank the translation of 40 segments (human translation, NMT and PBSMT), which had no context and were randomized to avoid bias. They were selected so there were no repeated translations and all had a minimum length of 100 characters. Then we applied a script to ensure there was a minimum editing distance of 15% between the human-PBSMT, human-NMT and PBSMT-NMT solutions. This reduced the number of segments from 230 to 145. We hand-picked 40 segments without typos nor any other problem.

5.2 Fluency and adequacy

We presented a survey with the same English segments as in the previous experiment. In the first part, participants (both students and professional translators) had to answer some questions about their previous experience in the translation industry. Afterwards, they had to evaluate the fluency and adequacy of the proposed translation on a four-point Likert scale. The translation was either PBSMT or NMT chosen randomly without any knowledge of the participants. The goal was to assess fluency and adequacy for in-domain formal texts.

5.3 PE time and technical effort

Finally, in the third experiment, participants had to post-edit 41 segments from a 2018 medical paper. They had to carry out the task in PET (Aziz et al. 2012)¹³, a computer-assisted translation tool that supports PE. It was used with its default settings. It logged both PE time and edits (keystrokes, insertions and deletions, that is, technical effort). Four professional translators with more than two years of experience post-editing carried out the task: two of them post-edited the PBSMT output and the other two post-edited the NMT output.

6 Results

6.1 Translation ranking

29 people answered the survey. From those, 86.21% had previous experience as translators and 58.62% had worked on PE tasks. Confirming the initial hypothesis, most respondents preferred the human translation. However, this percentage was only of 60.52%. The second most preferred translation was NMT, with 25.17%, and PBSMT was only considered the best translation for 14.31% of the segments. We calculated inter annotator agreement using Fleiss' kappa (Fleiss 1971), which showed a fair agreement among the annotators ($\kappa = 0.36$). These results were statistically significant in a one-way ANOVA comparison ($p < 0.05$).

Although the survey was conducted on a fairly small number of sentences, it seems to point in two directions: NMT is far from achieving the quality of human translation for medical texts, and NMT yields better translations than PBSMT. We conducted a manual analysis of the sentences in which NMT or PBSMT were selected as the best translation. It was observed the main reason for the selection was terminology precision and fluency of the MT output.

¹³<http://wilkeraziz.github.io/dcs-site/pet/index.html>

Table 2.3: Results of the human-NMT-PBSMT ranking survey.

Evaluation	Human	NMT	PBSMT
EN-ES (40)	60.52%	25.17%	14.31%

6.2 Fluency and adequacy

In the second experiment, eleven people answered the survey. Seven of them were translators with more than two years of experience and only four of them were students. Both fluency and adequacy obtained a higher rate for NMT after calculating the mean for both MT systems. We calculated inter annotator agreement using Fleiss’ kappa (Fleiss 1971). For fluency, it showed poor agreement among the annotators ($\kappa = 0.01$). Results were statistically significant in a one-way ANOVA comparison, with an F -ratio value of 2.75586 and a p -value of 0.04856 (significance at $p < 0.05$). For adequacy, there was also poor agreement among annotators. These results weren’t statistically significant, with an F -ratio value of 0.96767 and a p -value of 0.412816 ($p < 0.05$).

If we take a closer look at the sentences that had to be assessed, PBSMT segments often contain morphological problems (e.g. concordance) that we cannot spot in NMT segments, as in example (1). This way the generally higher ratings for fluency and adequacy of the NMT system are confirmed.

- (1) Source: Craniopharyngioma had more hormone deficiencies
 Gloss: Craneofaringioma tenían más déficits hormonales
 PBSMT: ‘Craneofaringioma/had (plural)/more/deficits/hormonal’

Table 2.4: Results of the ranking survey.

System	Fluency	Adequacy
PBSMT	2.28	2.24
NMT	2.46	2.50

6.3 PE time and technical effort

Results for the PE task by professional translators have been grouped in temporal effort and technical effort (see Tables 2.5 and 2.6). In both cases, the mean for

PBSMT is higher, though only technical effort shows a statistically significant difference (in a t -test with a p -value of 0.002054). It is worth highlighting that there was a considerable difference in time and keylogging between the translators, especially for the two professionals who post-edited PBSMT (as indicated by the standard deviation in Tables 2.5 and 2.6).

Table 2.5: Temporal PE effort (secs/segment).

System	Mean	SD
PBSMT	88.75	44.59
NMT	79.25	33.43

Table 2.6: Technical effort (keystrokes/segment).

System	Mean	SD
PBSMT	130.68	39.63
NMT	54.99	16.90

7 Conclusions and future work

Although the number of segments analyzed is quite small, for this language combination and text type, there seems to be a clear preference for human translations, which are considered better in more than half of the cases. Regarding MT engines, NMT presents more fluency and adequacy. This corresponds with the higher results in all automatic metrics. However, the results for the perception and automatic assessments do not correlate with PE time, even though there is a reduction in technical effort when post-editing NMT outputs. Thus, even though NMT produces more fluent results, this improvement does not always entail a reduction of the PE effort for professional translators, probably due to the added difficulty of error spotting in more fluent outputs.

In future research, we intend to further analyze PE, increasing the number of segments and language combinations to assess the correlation between automatic metrics and PE (technical and temporal) effort.

Acknowledgements

We want to thank the copyright holders for granting permission for the MSD-Manuals website and for using these texts to create an English-Spanish parallel corpus. The training of the neural MT systems has been possible thanks to the NVIDIA GPU grant programme.

References

- Allen, Jeffrey H. 2003. Post-editing. In Harold Sommer (ed.), *Computers and translation: A translator's guide*, 297–317. Amsterdam: John Benjamins. DOI: 10.1075/btl.35.19all.
- Aziz, Wilker, Sheila C. M. De Sousa & Lucia Specia. 2012. PET: A tool for post-editing and assessing machine translation. In *Proceedings of the eight international conference on language resources and evaluation (LREC'12)*, 3982–3987.
- Bahdanau, Dzmitry, Felix Hill, Jan Leike, Edward Hughes, Pushmeet Kohli & Edward Grefenstette. 2018. Jointly learning “what” and “how” from instructions and goal-states. In *6th international conference on learning representations, ICLR 2018, Vancouver, BC, Canada, April 30–May 3, 2018, workshop track proceedings*. OpenReview.net. <https://openreview.net/forum?id=BkmZvdkPM>.
- Bentivogli, Luisa, Arianna Bisazza, Mauro Cettolo & Marcello Federico. 2016. Neural versus phrase-based machine translation quality: A case study. In *Proceedings of the 2016 conference on empirical methods in Natural Language Processing*, 257–267. Austin, Texas: Association for Computational Linguistics. DOI: 10.18653/v1/D16-1025.
- Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia & Marco Turchi. 2017. Findings of the 2017 conference on machine translation (WMT17). In *Second conference on machine translation*, 169–214. <http://www.aclweb.org/anthology/W17-4717>.
- Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurelie Neveol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor & Marcos Zampieri. 2016. Findings of the 2016 conference on machine translation. *Proceedings of the First Conference on Machine Translation 2*. 131–198. <http://www.aclweb.org/anthology/W16-2301>.
- Bojar, Ondřej, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn & Christof Monz. 2018. Findings of the 2018 conference on machine translation (WMT18). In *Proceedings of the third conference on machine translation*, 272–303. <http://aclweb.org/anthology/W18-6401.pdf>.
- Castilho, Sheila, Joss Moorkens, Federico Gaspari, Iacer Calixto, John Tinsley & Andy Way. 2017. Is neural machine translation the new state of the art? *The Prague Bulletin of Mathematical Linguistics* 108(1). 109–120. DOI: 10.1515/pralin-2017-0013.

- Castilho, Sheila, Joss Moorkens, Federico Gaspari, Rico Sennrich, Vilemini Sisoni, Yota Georgakopoulou, Pintu Lohar, Andy Way, Miceli Barone & Maria Gialama. 2017. A comparative quality evaluation of PBSMT and NMT using professional translators. In *Proceedings of MT summit XVI, vol.1: Research track*, 116–131.
- Castilho, Sheila, Sharon O’Brien, Fabio Alves & Morgan O’Brien. 2014. Does post-editing increase usability? A study with Brazilian Portuguese as target language. In *Proceedings of the 17th annual conference of the European association for machine translation*, 183–190. Dubrovnik, Croatia: European Association for Machine Translation. <https://www.aclweb.org/anthology/2014.eamt-1.40>.
- Cho, Kyunghyun, Bart van Merriënboer, Dzmitry Bahdanau & Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8: Eighth workshop on syntax, semantics and structure in statistical translation*, 103–111. Doha, Qatar: Association for Computational Linguistics. DOI: 10.3115/v1/W14-4012.
- Fleiss, Joseph L. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin* 76(5). 378–382.
- Forcada, Mikel L., Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O’Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez & Francis M. Tyers. 2011. Apertium: A free/open-source platform for rule-based machine translation. *Machine translation* 25(2). 127–144.
- Germann, Ulrich, Eduard Barbu, M. Bentivoglio, Nikolay Bogoychev, C. Buck, D. Caroselli, L. Carvalho, A. Cattelan, R. Cattoni, Mauro Cettolo, Marcello Federico, Barry Haddow, David Madl, L. Mastrostefano, Prashant Mathur, A. Ruopp, A. Samiotou, V. Sudharshan, M. Trombetti & Jan van der Meer. 2016. Modern MT: A new open-source machine translation platform for the translation industry. *Baltic Journal of Modern Computing* 4. 397–397.
- Guerberof, Ana. 2009a. Productivity and quality in MT post-editing. In *Proceedings of MT Summit XII: Beyond translation memories: New tools for translators MT*, 1–9. Ottawa, Canada: AMTA. <http://www.mt-archive.info/MTS-2009-Guerberof.pdf>.
- Guerberof, Ana. 2009b. Productivity and quality in the post-editing of outputs from translation memories and machine translation. *The International Journal of Localisation* 7(1). 11–21. <https://www.tdx.cat/bitstream/handle/10803/90247/GuerberofThesis%20Final.pdf?sequence=1&isAllowed=y>.
- Isabelle, Pierre, Colin Cherry & George F. Foster. 2017. A challenge set approach to evaluating machine translation. *Computing Research Repository* abs/1704.07431. <http://arxiv.org/abs/1704.07431>.

- Junczys-Dowmunt, Marcin, Tomasz Dwojak & Hieu Hoang. 2016. Is neural machine translation ready for deployment? A case study on 30 translation directions. *CoRR* abs/1610.01108. <http://arxiv.org/abs/1610.01108>.
- Koehn, Philipp. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of the Xth MT summit*, vol. 5, 79–86.
- Koponen, Maarit. 2012. Comparing human perceptions of post-editing effort with post-editing operations. In *Proceedings of the seventh workshop on statistical machine translation (WMT '12)*, 181–190. Montréal, Canada: Association for Computational Linguistics. <https://www.aclweb.org/anthology/W12-3123>.
- Koponen, Maarit. 2016. Is machine translation post-editing worth the effort? A survey of research into post-editing and effort. *The Journal of Specialised Translation* 25. 131–148.
- Krings, Hans P. 2001. *Repairing texts: Empirical investigations of machine translation post-editing processes*. Vol. 5. Kent, OH: Kent State University Press.
- Lommel, Arle & Donald A. DePalma. 2016. *Europe's leading role in machine translation: How Europe is driving the shift to MT*. Tech. rep. Boston. <http://cracker-project.eu>.
- Moorkens, Joss. 2018. Eye tracking as a measure of cognitive effort for post-editing of machine translation. In Walker Calum & Federico M. Federici (eds.), *Eye tracking and multidisciplinary studies on translation*, 55–70. Amsterdam. DOI: 10.1075/btl.143.04moo.
- Moorkens, Joss & Sharon O'Brien. 2013. User attitudes to the post-editing interface. In *Proceedings of machine translation summit XIV: Second workshop on post-editing technology and practice, Nice, France*, 19–25. <http://www.mt-archive.info/10/MTS-2013-W2-Moorkens.pdf>.
- Moorkens, Joss, Sharon O'Brien, Igor A. L. Da Silva, Norma B. De Lima Fonseca & Fabio Alves. 2015. Correlations of perceived post-editing effort with measurements of actual effort. *Machine Translation* 29(3). 267–284. DOI: 10.1007/s10590-015-9175-2.
- O'Brien, Sharon. 2011. Towards predicting post-editing productivity. *Machine Translation* 25(3). 197–215.
- Papineni, Kishore, Salim Roukos, Todd Ward & Wj Zhu. 2002. BLEU: A method for automatic evaluation of machine translation. In vol. July, 311–318. DOI: 10.3115/1073083.1073135.
- Parra Escartín, Carla & Manuel Arcedillo. 2015. A fuzzier approach to machine translation evaluation: A pilot study on post-editing productivity and automated metrics in commercial settings. In vol. 1, 40–45. <https://aclweb.org/anthology/W/W15/W15-4107.pdf>.

- Plitt, Mirko & François Masselot. 2010. A productivity test of statistical machine translation post-editing in a typical localisation context. *The Prague bulletin of mathematical linguistics* 93. 7–16. <https://ufal.mff.cuni.cz/pbml/93/art-plitt-masselot.pdf>.
- Sanchez-Torrón, Marina & Philipp Koehn. 2016. Machine translation quality and post-editor productivity. In *Proceedings of AMTA 2016*, 16–26. <https://researchspace.auckland.ac.nz/handle/2292/31486>.
- Specia, Lucia. 2010. Combining confidence estimation and reference-based metrics for segment-level MT evaluation. In *The ninth conference of the association for machine translation in the Americas*. <https://amta2010.amtaweb.org/AMTA/papers/2-03-BanerjeeDuEtal.pdf>.
- Specia, Lucia. 2011. Exploiting objective annotations for measuring translation post-editing effort. In *Proceedings of the 15th conference of the European association for machine translation*, 73–80. <http://www.mt-archive.info/EAMT-2011-Specia.pdf>.
- Toral, Antonio & Víctor M. Sánchez-Cartagena. 2017. A multifaceted evaluation of neural versus phrase-based machine translation for 9 language directions. In *Proceedings of the 15th conference of the European chapter of the Association for Computational Linguistics*, vol. Volume 1: Long papers, 1063–1073. Valencia, Spain: Association for Computational Linguistics.
- Villegas, Marta, Ander Intxaurre, Aitor Gonzalez-Agirre, Montserrat Marimon & Martin Krallinger. 2018. The MeSpEN resource for English-Spanish medical machine translation and terminologies: Census of parallel corpora, glossaries and term translations. *Language Resources and Evaluation* 52. 32–39.
- Wu, Yonghui, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes & Jeffrey Dean. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR* abs/1609.08144. <http://arxiv.org/abs/1609.08144>.