
Post-editing Effort of a Novel with Statistical and Neural Machine Translation

Antonio Toral,^{1,*} Martijn Wieling,¹ and Andy Way²

¹ Center for Language and Cognition, Faculty of Arts, University of Groningen, The Netherlands

² ADAPT Centre, School of Computing, Dublin City University, Ireland

Correspondence*:
Corresponding Author
a.toral.ruiz@rug.nl

2 ABSTRACT

3 We conduct the first experiment in the literature in which a novel is translated automatically and
4 then post-edited by professional literary translators. Our case study is *Warbreaker*, a popular
5 fantasy novel originally written in English, which we translate into Catalan. We translated one
6 chapter of the novel (over 3,700 words, 330 sentences) with two data-driven approaches to
7 Machine Translation (MT): phrase-based statistical MT (PBMT) and neural MT (NMT). Both
8 systems are tailored to novels; they are trained on over 100 million words of fiction. In the post-
9 editing experiment, six professional translators with previous experience in literary translation
10 translate subsets of this chapter under three alternating conditions: from scratch (the norm in the
11 novel translation industry), post-editing PBMT, and post-editing NMT. We record all the keystrokes,
12 the time taken to translate each sentence, as well as the number of pauses and their duration.
13 Based on these measurements, and using mixed-effects models, we study post-editing effort
14 across its three commonly studied dimensions: temporal, technical and cognitive. We observe
15 that both MT approaches result in increases in translation productivity: PBMT by 18%, and NMT
16 by 36%. Post-editing also leads to reductions in the number of keystrokes: by 9% with PBMT, and
17 by 23% with NMT. Finally, regarding cognitive effort, post-editing results in fewer (29% and 42%
18 less with PBMT and NMT respectively) but longer pauses (14% and 25%).

19 **Keywords:** literary translation, post-editing, neural machine translation, statistical machine translation, foreign literature, foreign
20 fiction

1 INTRODUCTION

21 Machine Translation (MT) is widely used in the translation industry today to assist professional human
22 translators, as using MT results in notable increases in translator productivity compared to translation
23 from scratch. This has been empirically shown in many use-cases over the last decade that rely on the
24 phrase- and rule-based paradigms to MT (PBMT and RBMT), for several text types, including technical
25 documents (Plitt and Masselot, 2010) and news (Martín and Serra, 2014), to mention just two.

26 The most common workflow employed is post-editing, a sequential pipeline in which the source document
27 is first translated with MT, and subsequently, a translator edits the MT translation (e.g. fixing errors) to
28 produce the final translation. In most of the use-cases explored in the literature the translation aim is

29 dissemination, and the translations obtained via post-editing have been found to be of equivalent or higher
30 quality (Plitt and Masselot, 2010; Green et al., 2013) to those produced from scratch.

31 Nonetheless, post-editing has been found to prime the translator, thus resulting in a final translation that is
32 similar to that initially suggested by the MT system (Green et al., 2013). Because the MT approaches most
33 widely used to date in post-editing translation workflows – RBMT and, above all, PBMT – are known to
34 lead to literal translations, post-edited translations are also perceived as being more literal than translations
35 from scratch (Martín and Serra, 2014).

36 While this is acceptable for text types such as technical documents, as the main objective of the translation
37 for these types of texts is to preserve the meaning of the original, it might not be the case for other text
38 types of a more creative nature, such as literary texts, because in this case the objective of the translation is
39 twofold: not only the meaning of the source text needs to be preserved but also its reading experience (Jones
40 and Irvine, 2013).

41 Recently, neural machine translation (NMT) has emerged as a new paradigm in MT, and has been
42 shown to considerably improve the translation quality achieved, regardless of the language pair (Toral and
43 Sánchez-Cartagena, 2017). In addition, the translations produced by NMT are much more fluent (Bentivogli
44 et al., 2016) than those derived by PBMT, until recently by far the most dominant paradigm in the field. In
45 addition, relevant to this work, it has been claimed that NMT does not lead to literal translations,¹ as is the
46 case with PBMT and RBMT.

47 At this point, because of (i) the maturity of post-editing in industry, and (ii) the rise of a new MT
48 paradigm (NMT) that results in more fluent and less literal translations than previous models (PBMT and
49 RBMT), it is timely to study the extent to which current MT technology can be useful in assisting with
50 professional translations of literary text. In this work we take the first steps in this direction by assessing
51 the effort involved in the post-editing of a novel, along the three dimensions commonly studied in the
52 literature (Krings and Koby, 2001), which constitute the research questions (RQs) underpinning this work:

- 53 • RQ1 (temporal effort). Does post-editing an MT output (using the NMT or PBMT paradigm) result in
54 shorter translation time compared to post-editing of outputs from the other type of MT system and/or
55 to translation from scratch?
- 56 • RQ2 (technical effort). Does post-editing on one of the two MT paradigms result in a lower number of
57 keystrokes than the other MT paradigm and/or than translation from scratch?
- 58 • RQ3 (cognitive effort). Does post-editing on one of the MT paradigms result in changes in cognitive
59 effort?

60 In this work we translate a fragment of a novel with NMT and PBMT. Subsequently, six professional
61 translators with previous experience in literary translation translate subsets thereof under three different
62 conditions: from scratch (the norm in the novel translation industry), post-editing the translation produced
63 by the PBMT system and post-editing that generated by the NMT system. For each sentence translated, we
64 record (i) the time spent to translate it, (ii) the number of keystrokes used, and (iii) the number of pauses
65 and time devoted to them. We then use these three measurements to attempt to provide answers to questions
66 RQ1, RQ2 and RQ3, respectively.

¹ “Neural network-based MT can, rather than do a literal translation, find the cultural equivalent in another language”, according to Alan Packer, Engineering Director at Facebook, in 2016, cf. <https://slator.com/technology/facebook-says-statistical-machine-translation-has-reached-end-of-life>

67 The rest of the paper is organised as follows. Section 2 outlines the state-of-the-art in MT of novels.
68 Next, Section 3 presents the MT systems (Section 3.1) and the novel (Section 3.2) used in our experiment,
69 followed by the experimental set-up (Section 3.3). Section 4 presents and discusses the results. Finally, in
70 Section 5, we draw our conclusions and propose lines of future work.

2 STATE-OF-THE-ART IN LITERARY TRANSLATION USING MT

71 Voigt and Jurafsky (2012) studied how referential cohesion is expressed in literary (short stories) and
72 non-literary (news stories) texts and how this cohesion affects translation. They found that literary texts use
73 more dense reference chains to express greater referential cohesion than news. They then compared the
74 referential cohesion of human versus machine translations of short stories from Chinese-to-English. MT
75 systems had difficulty in conveying the cohesion, which is attributed to the fact that they translate each
76 sentence in isolation while human translators can rely on information beyond the sentence level.

77 Jones and Irvine (2013) used generic PBMT systems to translate samples of French literature (prose and
78 poetry) including a fragment of Camus' *L'Étranger* into English. They analysed the translations from a
79 qualitative perspective to address what makes literary translation hard and to discover what the potential
80 role of MT could be.

81 Besacier and Schwartz (2015) presented a pilot study where a generic PBMT system followed by post-
82 editing was applied to translate a short story from English into French. Post-editing was performed by
83 non-professional translators, and the authors concluded that such a workflow can be a useful low-cost
84 alternative for translating literary works, albeit at the expense of lower translation quality.

85 Simultaneously to the previous work, Toral and Way (2015) built a PBMT system tailored to a
86 contemporary best-selling author (Ruiz Zafón) and then applied it to translate one of his novels, *El*
87 *prisionero del cielo*, between two closely-related languages (Spanish-to-Catalan). For 20% of the sentences,
88 the translations produced by the MT system and the professional translator (i.e. taken from the published
89 novel in the target language) were exactly the same. In addition, a human evaluation revealed that for
90 over 60% of the sentences, Catalan native speakers judged the translations produced by MT and by the
91 professional translator to be of the same quality.

92 Murchú (2017) machine translated the sci-fi novel *Air Cuan Dubh Drilseach* from Scottish Gaelic to Irish,
93 a pair of closely related languages, using the hybrid MT system Intergaelic and subsequently post-edited
94 the resulting MT output. Post-editing was 31% faster than translating from scratch and less than 50% of
95 the tokens in the MT output were corrected by the translator.

96 Toral and Way (2018) built PBMT and NMT systems tailored to novels for the English–Catalan language
97 pair. These were evaluated on a set of twelve widely known novels from the 20th and 21st centuries by
98 authors such as Joyce, Orwell, Rowling and Salinger, to name but a few. Overall, NMT resulted in an 11%
99 relative improvement (3 points absolute) over PBMT according to the BLEU evaluation metric (Papineni
100 et al., 2002). In a human evaluation conducted on the books by Orwell, Rowling and Salinger, the
101 translations generated by the NMT system were perceived by Catalan native speakers to be of equivalent
102 quality to the professional human translations for 14%, 29% and 32% of the sentences, respectively,
103 compared to 5%, 14% and 20% respectively, with the PBMT system. These findings have encouraged us to
104 expand this study to the post-editing of a novel, which we detail below.

3 MATERIAL AND METHODS

105 3.1 MT Systems

106 We trained two MT systems belonging to two different paradigms: PBMT and NMT. Both are tailored to
107 novels and a brief description of them follows. For a more detailed account, the reader is referred to Toral
108 and Way (2018).

109 The PBMT system is trained on a linear interpolation of in-domain (133 parallel novels from different
110 genres amounting to over 1 million sentence pairs) and out-of-domain (around 400,000 sentence pairs
111 of subtitles)² parallel data, with version 3 of the Moses toolkit (Koehn et al., 2007). The n -gram-based
112 language model, in addition to the training parallel data, uses monolingual in-domain (around 1,000 books
113 written in Catalan amounting to over 5 million sentences) and out-of-domain (around 16 million Catalan
114 sentences crawled from the web (Ljubešić and Toral, ???)) data. The system uses 3 reordering models
115 (lexical- and phrase-based, and hierarchical), an operation sequence model (Durrani et al., 2011) and an
116 additional language model based on continuous space n -grams (Vaswani et al., 2013). The last two models
117 are trained on the in-domain parallel data.

118 The NMT system follows the encoder-decoder approach and is built with Nematus (Sennrich et al.,
119 2017).³ This system is trained on the concatenation of the parallel in-domain training data (133 parallel
120 novels) and a synthetic corpus obtained by machine-translating the Catalan in-domain monolingual training
121 data (1,000 books) into English. The system uses sub-words as the basic translation unit; we segmented
122 the training data into characters and performed 90,000 operations jointly on both the source and target
123 languages (Sennrich et al., 2016). Finally, we generate an n -best list with the NMT system and rerank it
124 with a left-to-right NMT system.⁴

125 3.2 Novel

126 The novel used in this experiment is Sanderson's *Warbreaker*.⁵ This book fulfills our two requirements,
127 namely (i) literary quality, to make sure that the task is indeed challenging, and (ii) being freely
128 redistributable, to guarantee the reproducibility of our experiment. The first criterion is attested by its
129 reviews by critics, while the second is met as the book was published under a Creative Commons License
130 (CC-NC-ND specifically).

131 *Warbreaker* is pre-processed in the same way as the training data, namely it is sentence-split with
132 NLTK (Bird, 2006) and subsequently tokenised, truecased and normalised (in terms of punctuation) with
133 the corresponding Moses scripts.

134 In order to have an estimate of the difficulty posed by the translation of *Warbreaker*, we use two automatic
135 metrics. The first, type-token ratio (TTR), provides an indication of the richness of the vocabulary used in
136 the book. The second, n -gram overlap, corresponds to the percentage of n -grams in the novel that are also
137 found in the training data used to build the MT system. This measure thus provides an indication of the
138 degree of lexical divergence (or 'novelty') of the book that is to be translated with respect to the training
139 data.

² <http://opus.lingfil.uu.se/OpenSubtitles.php>

³ <https://github.com/rsennrich/nematus>

⁴ This system has the same settings as the regular NMT system, the only difference being that the target sentences of the training data are reversed at the word level.

⁵ <https://brandonsanderson.com/books/warbreaker/warbreaker/>

140 Table 1 shows the TTR and n -gram overlap (for $n = \{2, 3, 4\}$) of Warbreaker (both for the whole
141 book and for some individual chapters)⁶ as well as for the 12 books previously translated with our MT
142 systems (Toral and Way, 2018). For the latter we show the mean value for the 12 books as well as the 95%
143 confidence interval. In addition, we calculate the average sentence length (average number of words per
144 sentence) as previous research has shown that the performance of current NMT systems degrades with
145 increasing sentence length (Toral and Sánchez-Cartagena, 2017).

146 Comparing the scores of *Warbreaker* to those of the twelve well-known books we have previously
147 translated allows us to have an approximation as to how challenging translating *Warbreaker* is going to be.
148 The scores for *Warbreaker* (full book) fall inside the confidence intervals obtained for the twelve books for
149 two measures (2-gram overlap and TTR), they are slightly higher for another two (3- and 4-gram overlap)
150 and slightly lower for the remaining one (sentence length). According to these results we expect the novel
151 chosen to be slightly easier to translate than the average of the twelve novels we translated previously.

152 As for *Warbreaker*'s individual chapters, we select one for our experiment that has similar values to the
153 whole book, as that would make it (to some extent) representative of the book as a whole. We show the
154 values for the first three (prologue and Chapters 1 and 2) in Table 1 and pick Chapter 1 as it is the one
155 whose results are closest to the whole book for all the metrics considered (except sentence length, whose
156 value is longer than the average).

157 3.3 Experimental Setup

158 The professional translators performed the translation using PET v2.0 (Aziz et al., 2012),⁷ a computer-
159 assisted translation tool that supports both translating from scratch and post-editing. PET is used with its
160 default settings. A snapshot of the tool, as used in our experiment, is shown in Figure 1.

161 The source text translated in the experiment (*Warbreaker*'s Chapter 1) is made up of 3,743 words
162 distributed in 330 sentences. We divided it into 33 translation jobs, each of which is made of 10 consecutive
163 sentences (translation segments). There are three types of translation jobs (translation conditions):
164 translation from scratch (HT), and post-editing the translation produced by the PBMT (MT1) and NMT
165 systems (MT2).⁸

166 Six translators (henceforth T1 to T6) took part in the study. They saw all factors but not all combinations,
167 since they translated each job in one translation condition. The type of translation to be carried out for each
168 job by each translator is chosen randomly, with the following three constraints:

- 169 1. The first job is set to translation condition HT for translators T1 and T2, to MT1 for T3 and T4 and to
170 MT2 for T5 and T6.
- 171 2. Two consecutive jobs by a translator cannot follow the same translation condition.
- 172 3. For each translator the number of jobs under each translation condition is equal, i.e. each translator
173 translates 11 jobs under translation condition HT, 11 under MT1 and 11 under MT2.

174 We provided the translators with comprehensive translation guidelines,⁹ where it is stated that the aim is to
175 achieve publishable professional quality translations, both for translations from scratch and for post-editing.
176 With respect to post-editing, the guidelines encourage the translator to try to fix the translation provided

⁶ TTR scores are not shown for chapters as it is computed on 20,000 words, a much bigger amount of text than what makes up a chapter.

⁷ <http://rgcl.wlv.ac.uk/projects/PET/resources/PET-v2.0.tgz>

⁸ We referred to the two MT systems as MT1 and MT2 throughout the experiments so that the translators could not know anything about the MT paradigm into which they fell.

⁹ The manual is available at [TODO](#)

177 by the MT system. Only if this is deemed too time-consuming to fix (e.g. because the quality of the MT
178 output is too low) were the translators instructed to delete it and carry out the translation from scratch.

179 As in other computer-assisted tools, translations in PET are related to source sentences on a one-to-one
180 basis. In other words, each source sentence corresponds to one target sentence (see Figure 3.3). However, in
181 the translation of novels it is not that uncommon to have some cases of many-to-one (more than one source
182 sentence translated as one target sentence) or one-to-many (one source sentence translated as more than
183 one target sentence) translations. Due to this characteristic of literary translation, translators were told that
184 they could, in addition to one-to-one translations, perform one-to-many and/or many-to-one translations.
185 Details on how they could go about this are provided in the translator's manual.

186 For each research question (temporal, technical and cognitive effort), we first report the (descriptive)
187 results for the samples. For example, for temporal effort, the relative change in translation productivity
188 with post-editing versus translating from scratch is provided. Subsequently, we aim to generalise from the
189 samples (the translators that participate in the study and the sentences they translate) to populations (any
190 translator and any similar text) by using mixed-effects regression models (Baayen, 2008).¹⁰ Mixed-effects
191 regression models distinguish both fixed effects (i.e. the effects we are usually interested in) from random
192 effects (i.e. the factors we would like to generalise over). With respect to random effects, a distinction can
193 be made between random intercepts (i.e. the value of the dependent variable varies on the basis of the level
194 of the random-effect factor), and random slopes (i.e. the strength of the effect of a predictor varies on the
195 basis of the level of the random-effect factor). Specifically, we will build models where we contrast the
196 three translation conditions by including them as fixed effects, while including the translators (6 levels) and
197 translation segments, or sentences (330 levels) as random effects.

198 Previous studies in post-editing have shown that results vary considerably between translators and
199 segments. By taking a mixed-effects regression approach, we are able to include the by-translator and
200 by-segment random intercepts and slopes to model the variability associated with translator and segment.
201 For example, one individual translator may tend to take longer, or rewrite a larger part of a sentence than
202 another, which is modelled by a by-translator random intercept. Similarly, one sentence (due to its structure)
203 may be more likely to be rewritten than another, or may take more cognitive effort to translate, which is
204 modelled by a by-segment random intercept. In addition, one translator might show a greater difference
205 between the three conditions than another, which is modeled by by-translator random slope for translation
206 condition. Similarly, a by-segment random slope for translation condition is able to model that translation
207 condition may show a greater difference for one sentence than for another.

208 We conduct exploratory analyses, in which we first include random intercepts for translator and segment,
209 and subsequently add fixed-effect predictors one by one. For each predictor, we check whether its addition
210 results in a significantly better statistical model by comparing the model that adds that predictor to a
211 simpler model without that predictor. Any pair of models is subsequently compared in terms of Akaike's
212 Information Criterion (AIC; Akaike, 1973). If the AIC of the model that includes the predictor is at least 2
213 points lower than the model without the predictor then we consider the first model significantly better. The
214 evidence ratio can be calculated on the basis of the AIC difference¹¹ and represents the relative probability
215 that the model with the lowest AIC is more likely to provide a more precise model of the data. By using a
216 threshold of 2 (see also Groenewold et al. (2014)), we only select a more complex model if it is 2.7 times
217 more likely than the simpler model. After including the fixed effect predictors separately, we evaluate

¹⁰ For our analysis we use the `lme4` R package, for a normal linear regression model or a Poisson generalised linear regression model, but for a ratio as the dependent variable, we use the package `mgcv`, as beta regression is not implemented in the `lme4` package.

¹¹ Evidence ratio: $e^{\frac{\delta AIC}{2}}$

218 (using AIC comparisons) if interactions between the fixed-effect predictors are necessary. After obtaining
219 the best fixed-effects structure, we evaluate the optimal random-effects structure (i.e. by including random
220 effects, and evaluating their inclusion again using AIC comparison) and retain all fixed-effect factors which
221 are significant when the appropriate random effects structure is included. This approach is similar to that
222 used by Wieling et al. (2011).

223 An ethics approval for this study was obtained from The Research Ethics Committee of the Faculty of
224 Arts, University of Groningen. The professional translators involved in the study gave written informed
225 consent in accordance with the Declaration of Helsinki.¹²

4 RESULTS AND DISCUSSION

226 As was previously mentioned in Section 1, this work has three research questions, concerning temporal
227 (RQ1), technical (RQ2) and cognitive effort (RQ3). Next we detail the pre-processing of the data. The
228 subsequent three subsections attempt to provide answers to these three questions, based on the experimental
229 data collected.

230 4.1 Pre-processing

231 For each translated sentence by each translator, we extract the following elements from the PET logs:
232 length of the source and target text (in words and characters), translation condition (HT, MT1 or MT2),
233 translation time, number of keystrokes (total as well as belonging to different categories: letters, digits,
234 whitespace, symbols, navigation, deletion, copy, cut and paste), and number of pauses and their duration.
235 Following the findings by Lacruz et al. (2014), we include only pauses longer than 300 milliseconds.

236 We also pre-processed those translations without a 1-to-1 sentence equivalence. None of the translators
237 produced any 1-to-many translations, and only three out of the six translators generated many-to-1
238 translations. Moreover, these translators performed such translations in very few cases: from 6 to 10
239 sentences, i.e. from 1.8% to 3% of the translation units. The reason given by the translators as to why some
240 of them produced many-to-1 translations but no 1-to-many was due to the fact that sentences in novels in
241 Catalan tend to be longer than in English. Accordingly, conflating more than one English sentence into a
242 single Catalan translation equivalent made sense, albeit on rare occasions. The fact that the vast majority of
243 translations were 1-to-1 could be attributed to either of the two following reasons (or a combination of
244 both):

- 245 • While the instructions allowed for translations beyond the 1-to-1 sentence equivalence, the computer-
246 assisted tool used expects 1-to-1 sentence equivalence, so translators may feel discouraged to do
247 otherwise;
- 248 • While there may be the perception that in original novels and their translations, sentences do not tend
249 to correspond 1-to-1, this is actually the most frequent case, at least for the language pair we cover
250 in this study. In Toral and Way (2018), we sentence-aligned over 100 novels in English and their
251 translations in Catalan. Overall, 77% of the sentences were successfully aligned 1-to-1. The remaining
252 23% is made up of alignments that are not 1-to-1 but also of 1-to-1 alignments that the alignment tool
253 could not align confidently.

¹² <https://web.archive.org/web/20091015082020/http://www.wma.net/en/30publications/10policies/b3/index.html>

254 **4.2 Temporal Effort**

255 First, we report on translation productivity (measured as words per hour) per translation condition, as this
256 is a metric commonly used in related work, e.g. Plitt and Masselot (2010). Overall, translators produce 503
257 words per hour when translating from scratch (condition HT). Compared to this, post-editing the translation
258 produced by the PBMT results in 594 words per hour, an 18% increase in productivity, while post-editing
259 the NMT output leads to double that figure: 36% (685 words per hour). This is clearly indicative of the fact
260 that NMT outputs were superior to those from PBMT.

261 We now zoom in and look at each translator individually. Results are shown in Figure 2. We can observe a
262 large variability in translation speed, from the lowest value of 402 words/hour (translator T3, condition HT)
263 to the highest of 1,140 (T2, MT2). Despite this variability, we can observe clear trends when comparing
264 translation conditions: all translators are faster in condition MT1 compared to HT (relative increases range
265 from 1% for T6 to 46% for T2), and all are faster with MT2 than with MT1 (increases range from to
266 0.001% for T1 to 37% for T3).

267 Next, in order to generalise from samples to populations and to find out whether differences are statistically
268 significant, we build a linear mixed-effects regression model in which we predict translation time¹³ given
269 two (fixed-effect) numerical predictors (length of the segment in characters and trial number), one fixed-
270 effect factorial predictor (translation condition) and two random-effect factors (translators and segments).
271 Numerical predictors are centred and scaled. After fitting the final model, we conduct model criticism by
272 excluding data points which have an observed value deviating more than 2.5 standard deviations from the
273 predicted value by the model¹⁴ and refit the model. In this way, we prevent that potentially significant
274 effects are “carried” by these outliers (which are not well represented by the model; Baayen, 2008). We
275 assessed that the residuals of our final model approximately followed a normal distribution and were
276 homoscedastic.

277 In the best model, the two numerical fixed predictors are significant: translators take longer time the
278 longer the input text and shorter time as they advance through the experiment (trial number). The effect of
279 translation condition is also significant: compared to HT, translation time in condition MT1 is significantly
280 reduced, and so this is also the case for MT2 compared to MT1.

281 We find a significant interaction between input length and translation condition. Figure 3 shows that
282 the longer the input sentence the lower the advantage of MT2 over HT. There is no such effect for MT1
283 though. The fact that post-editing NMT is not advantageous over translating from scratch for long sentences
284 corroborates the finding that the translation quality provided by NMT degrades with sentence length (Toral
285 and Sánchez-Cartagena, 2017). Table 3 shows the significance level for each predictor and interaction
286 between predictors, not only for the model built for temporal effort but also for those used for technical
287 and cognitive effort (see Sections 4.3 and 4.4, respectively).

288 In terms of the random-effects structure, the final model included both random intercepts (by segment
289 and translator), and a by-item (segment) random slope for translation condition. The random slope reflects
290 that the difference in temporal effort between the three conditions varies per segment.

¹³ We transform it logarithmically, since its distribution is heavily skewed to the right.

¹⁴ A total of 56 out of 1,980 data points (2.8%) were removed.

291 **4.3 Technical effort**

292 We measure the technical effort by means of the number of keystrokes used to produce the final translation.
293 Similarly to what was done for temporal effort (cf. Section 4.2), we calculate the number of keystrokes
294 per character in the source sentence and per translation condition (HT, MT1 and MT2), i.e. the number
295 of keystrokes that it takes to translate one character with each translation method. Overall, it takes 1.94
296 keystrokes to translate each character when translating from scratch (condition HT). Compared to this,
297 post-editing PBMT (condition MT1) results in a 9% reduction (1.76 keystrokes per character), while NMT
298 leads to more than double that reduction, 23% (1.49 keystrokes per character).

299 We now zoom in and look at each translator individually. Results are shown in Figure 4. As with temporal
300 effort, there is large variability across translators and conditions, the lowest value being 0.8 keystrokes
301 per second (translator T2, condition MT2) and the highest 2.9 (translator T5, condition MT1). Some
302 trends arise but they are not as clear as was the case with temporal effort. Compared to HT, the number
303 of keystrokes is reduced with MT1 for three translators (maximum reduction: 45%, T2) and is increased
304 with the other three (maximum increase: 13%, T5). Compared to HT, MT2 results in a reduced number of
305 keystrokes for all translators except T6, for whom it increases slightly (2%). The maximum reduction is, as
306 in the case of MT1, for T2 (59%).

307 Next, as for temporal effort, we build a statistical mixed model to predict technical effort, for which
308 we consider the same set of predictors. Our dependent variable is the total number of keystrokes. As this
309 dependent variable reflects count data, we use Poisson generalised linear mixed-effects regression. As
310 in temporal effort, all the fixed predictors are significant. The longer the input, the more keystrokes are
311 used and the further a translator advances in the experiment, the fewer keystrokes s/he uses. The effect of
312 post-editing is significant, fewer keystrokes are required with MT1 compared to HT, and the same occurs
313 when we compare MT2 to MT1.

314 The interaction between input length and translation condition, which was significant for temporal effort,
315 is significant here too, but again only shows a difference between HT and MT2. The interaction is shown
316 in Figure 5. The longer the input sentence, the smaller the difference becomes between the number of
317 keystrokes used in conditions HT and MT2.

318 The optimal random-effects structure, in this case, consists of both a by-translator and a by-segment
319 random slope for translation condition, and a by-translator random slope for trial (reflecting that the trial,
320 i.e. learning effects, are different per translator).

321 In the experiment we not only logged the number of keystrokes used but also their type. We now delve
322 deeper into the keystroke results by differentiating the keystrokes into three groups: content (digits, letters,
323 white space and symbols), navigation keys and erase keys.¹⁵ Figure 4 showed the average number of keys
324 per source character for each translator under each of the translation conditions. Now, we show a different
325 perspective in Table 2, where we break up the average number of total keys into three groups of keys and
326 we aggregate the data for all the translators.

327 It has been previously shown that post-editing leads to a very different usage of the keyboard compared
328 to translation from scratch (Carl et al., 2011). Our results corroborate this: while post-editing reduces
329 considerably the number of content keywords used (-55% with PBMT and -63% with NMT), that translation

¹⁵ Other types of keystrokes were logged too, for the operations copy, cut, paste and undo. However, their usage was negligible in the experiment; they account for just 0.1% of the total number of keystrokes used, so have not been included in our analysis.

330 pipeline results in a massive increase in the use of navigation keys (228% with PBMT and 195% with
331 NMT) and, to a lesser extent, erase keys (105% for PBMT and 72% with NMT).

332 Figure 7 shows a complementary view of this data. For each translation condition, we depict the
333 proportion of keys that belong to each of the three groups considered (content, navigation and erase).
334 In translation from scratch content keystrokes make up 79% of the total, navigation 9% and erase the
335 remaining 12%. Post-editing leads to roughly equal percentages for each keystroke category: 39% content,
336 34% for navigation and 27% for erase with PBMT and 38%, 36% and 27%, respectively with NMT.

337 Finally, we show the complete picture with three variables at once (translators, translation condition
338 and keystroke type) in Figure ???. The trend is similar across translators for content keys; all of them
339 use substantially more keystrokes when translating from scratch than when post-editing. Navigation is
340 the type of keystrokes for which we observe the highest variation across translators; on one extreme two
341 translators (T2 and T6) use very few navigation keys, regardless of the translation condition. On the other,
342 one translator (T5) uses more than double the number of navigation keys than the second translator in
343 number of navigation keys (T3). For erase keys, we see similar trends across translators; all of them except
344 T2 use more erase keys when post-editing than when translating from scratch.

345 4.4 Cognitive Effort

346 We use pauses as a proxy to measure cognitive effort (Schilperoord, 1996; O’Brien, 2006). We consider
347 three different ways of expressing the dependent variable (Green et al., 2013):

- 348 • Count: the number of pauses.
- 349 • Mean duration: how long pauses take on average.
- 350 • Ratio: the amount of time devoted to pauses divided by the total translation time.

351 The number of pauses correlates strongly with the number of keystrokes ($R = 0.87$). Due to this and
352 because number of pauses is a count-dependent variable, we fit number of pauses as the dependent variable
353 with the Poisson regression model previously built for technical effort (see Section 4.3). According to the
354 model, there are 15.3 pauses per sentence when translating from scratch. Condition MT1 significantly
355 reduces this by 29% (10.9) and MT2 by 42% (8.8).

356 The mean duration of pauses correlates weakly with translation time ($R = 0.25$) and has no correlation
357 with number of keystrokes ($R = -0.02$). We fit the mean duration of pauses¹⁶ with the model previously
358 built to predict translation time (see Section 4.2). Pauses have a mean duration of 2,243 milliseconds
359 in the translating-from-scratch condition. In condition MT1 this significantly increases by 14% (2,559
360 milliseconds), while in MT2 this increases further, by 25% (2,810 milliseconds).

361 The ratio of pauses is a proportion, and thus we use beta regression. Pause ratio correlates with translation
362 time ($R = 0.57$) and hence we will use the same predictors, interactions and slopes as in the model
363 previously built to predict time (see Section 4.2). According to the model, pauses take 63% of the
364 translation time in condition HT. Post-editing, be it with MT1 or MT2, leads to significant increments of
365 around 2.5 percentage points (65.6% and 65.3%, respectively) of the time devoted to pauses. The difference
366 between MT1 and MT2 is not significant.

¹⁶ We transform it logarithmically, since its distribution is heavily skewed to the right.

5 CONCLUSIONS AND FUTURE WORK

367 We have conducted the first experiment in the literature in which a fragment of a novel is translated
368 automatically and then post-edited by professional translators. Specifically, we have translated one chapter
369 of *Warbreaker* (over 3,700 words) from English into Catalan with domain-specific PBMT and NMT
370 systems.

371 The experiment has been conducted by six professional translators, who translated consecutive fragments
372 of 10 sentences each in three alternating conditions: from scratch, post-editing PBMT, and post-editing
373 NMT. The time taken for each segment as well as the keystrokes used, the number of pauses and the
374 duration of pauses were recorded, which has allowed us to analyse the translation logs and study how
375 post-editing with PBMT and NMT affects temporal, technical and cognitive effort.

376 Regarding temporal effort, compared to translation from scratch, both PBMT and NMT lead to substantial
377 increases in translation productivity (measured as word per hour), of 18% and 36%, respectively. This
378 demonstrates convincingly that post-editing MT output – whatever the system – makes translators faster
379 than when they translate from scratch. Furthermore, it indicates that translations output by NMT engines
380 were better than those from the corresponding PBMT systems. In addition, we found that the gain with
381 PBMT remains constant regardless of the length of the input sentence, while the gain with NMT decreases
382 with long sentences.

383 With respect to the number of keystrokes used (the measure used for technical effort), NMT again resulted
384 in a more substantial reduction (23%) than PBMT (9%). As with temporal effort, the reduction in the
385 number of keystrokes for PBMT remains constant across input sentences of different length, while the
386 reduction with NMT decreases for long sentences. Finally, we have observed that the distribution of types
387 of keystrokes is very different in post-editing compared to translation from scratch. While the first results in
388 considerably fewer content keywords, it notably increases the number of navigation and erase keystrokes.

389 As for cognitive effort, which we measured using pauses as proxies, we found that NMT – and to a
390 lesser extent PBMT – significantly reduce the number of pauses (42% and 29%, respectively). Pauses
391 are considerably longer when post-editing (14% with PBMT and 25% with NMT) than when translating
392 from scratch. Finally, we observed that pauses take a longer fraction of the total translation time when
393 post-editing, and that the difference between PBMT and NMT is not significant.

394 In this study we have looked at post-editing effort, covering its three dimensions: temporal, technical
395 and cognitive. In the next phase of this work, we will explore translators' perceptions, which we recorded
396 during the experiments by means of pre- and post-experiment questionnaires and a debriefing session, and
397 compare these perceptions to the results and conclusions from the current study.

398 Finally, we will assess the quality of the resulting post-edited translations. In previous post-editing studies
399 this is commonly measured by assessing the translations in terms of adequacy and fluency. For literary
400 texts, however, there is an additional requirement, namely that the translation should preserve the reading
401 experience of the source text. Accordingly, we aim to measure this in our future work.

CONFLICT OF INTEREST STATEMENT

402 The authors declare that the research was conducted in the absence of any commercial or financial
403 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

404 AT conceptualised the research, co-designed and conducted the experiments and wrote the manuscript.
405 MW directed the statistical analyses and reviewed/edited the manuscript. AW co-designed the experiments
406 and reviewed/edited the manuscript. All authors listed have made a substantial, direct, and intellectual
407 contribution to the work and approved it for publication.

FUNDING

408 The research leading to these results has received funding from the European Association for Machine
409 Translation through its 2015 sponsorship of activities programme, proposal named "Pilot on Post-editing
410 Novels (PiPeNovel)". The ADAPT Centre for Digital Content Technology at Dublin City University is
411 funded under the Science Foundation Ireland Research Centres Programme (Grant 13/RC/2106) and is
412 cofunded under the European Regional Development Fund.

ACKNOWLEDGMENTS

413 We would like to thank the six professional translators that took part in this study, in alphabetical order:
414 X, Y, Z, and N translators that preferred to remain anonymous. In addition, we would like to thank Sheila
415 Castilho and Joss Moorkens for their feedback on the experiment set up and the translation guidelines and
416 Wilker Aziz for his help on processing the PET log files.

SUPPLEMENTAL DATA

417 The translation guidelines provided to translators, the raw logs from PET and an R notebook (source
418 code and HTML report) with all the statistical analyses conducted are provided with this manuscript as
419 supplementary data.

REFERENCES

- 420 Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *Second*
421 *International Symposium on Information Theory* (Budapest, Hungary), 267–281
- 422 Aziz, W., Castilho, S., and Specia, L. (2012). PET: a Tool for Post-editing and Assessing Machine
423 Translation. In *Proceedings of the 8th International Conference on Language Resources and Evaluation*
424 (Istanbul, Turkey), 3982–3987
- 425 Baayen, R. H. (2008). *Analyzing linguistic data: A practical introduction to statistics using R* (Cambridge,
426 UK: Cambridge University Press)
- 427 Bentivogli, L., Bisazza, A., Cettolo, M., and Federico, M. (2016). Neural versus phrase-based machine
428 translation quality: a case study. In *Proceedings of the 2016 Conference on Empirical Methods in*
429 *Natural Language Processing* (Austin, Texas), 257–267
- 430 Besacier, L. and Schwartz, L. (2015). Automated translation of a literary work : a pilot study. In
431 *Proceedings of the Fourth Workshop on Computational Linguistics for Literature* (Denver, Colorado,
432 USA), 114–122
- 433 Bird, S. (2006). NLTK: the natural language toolkit. In *Proceedings of the COLING/ACL on Interactive*
434 *presentation sessions* (Sydney, Australia), 69–72

- 435 Carl, M., Dragsted, B., Elming, J., Hardt, D., and Jakobsen, A. L. (2011). The Process of Post-Editing : a
436 Pilot Study. In *Proceedings of the 8th international NLPSC workshop. Special theme: Human-machine
437 interaction in translation* (Copenhagen, Denmark), 131–142
- 438 Durrani, N., Schmid, H., and Fraser, A. (2011). A joint sequence translation model with integrated
439 reordering. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics:
440 Human Language Technologies-Volume 1* (Portland, Oregon, USA), 1045–1054
- 441 Green, S., Heer, J., and Manning, C. D. (2013). The efficacy of human post-editing for language translation.
442 In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Paris, France),
443 439–448
- 444 Groenewold, R., Bastiaanse, R., Nickels, L., Wieling, M., and Huiskes, M. (2014). The effects of direct and
445 indirect speech on discourse comprehension in dutch listeners with and without aphasia. *Aphasiology*
446 28, 862–884
- 447 Jones, R. and Irvine, A. (2013). The (Un)faithful Machine Translator. In *Proceedings of the 7th Workshop
448 on Language Technology for Cultural Heritage, Social Sciences, and Humanities* (Sofia, Bulgaria),
449 96–101
- 450 Koehn, P., Hoang, H., Birch, A., Callison-Burch, C., Federico, M., Bertoldi, N., et al. (2007). Moses: open
451 source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the ACL
452 on Interactive Poster and Demonstration Sessions* (Prague, Czech Republic), 177–180
- 453 Krings, H. and Koby, G. (2001). *Repairing Texts: Empirical Investigations of Machine Translation
454 Post-editing Processes*. Translation studies (Kent State University Press)
- 455 Lacruz, I., Denkowski, M., and Lavie, A. (2014). Cognitive demand and cognitive effort in post-editing. In
456 *AMTA 2014: proceedings of the eleventh conference of the Association for Machine Translation in the
457 Americas, Workshop on Post-editing Technology and Practice (WPTP-3)* (Vancouver, BC), 73–84
- 458 Ljubešić, N. and Toral, A. (????). In *Proceedings of the Ninth International Conference on Language
459 Resources and Evaluation (LREC'14), title = caWaC - a Web Corpus of Catalan and its Application to
460 Language Modeling and Machine Translation, year = 2014, pages = 1728–1732* (Reykjavik, Iceland)
- 461 Martín, J. A. A. and Serra, A. C. (2014). Integration of a machine translation system into the editorial
462 process flow of a daily newspaper. *Procesamiento del Lenguaje Natural* 53, 193–196
- 463 Murchú, E. Ó. (2017). Bearná i litríocht na gaeilge a líonadh: Réiteach úr? (filling gaps in irish-language
464 literature: A novel approach). Presented at AR AN IMEALL I LÁR AN DOMHAIN: An tairseachúlacht
465 i litríocht agus i gcultúr na hÉireann agus na hEorpa, Prague, 14 Sept 2017
- 466 O'Brien, S. (2006). Pauses as indicators of cognitive effort in post-editing machine translation output.
467 *Across Languages and Cultures* 7, 1–21
- 468 Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. (2002). BLEU: A Method for Automatic Evaluation of
469 Machine Translation. In *Proceedings of the 40th Annual Meeting on Association for Computational
470 Linguistics* (Philadelphia, PA, USA), 311–318
- 471 Plitt, M. and Masselot, F. (2010). A productivity test of statistical machine translation post-editing in a
472 typical localisation context. *Prague Bull. Math. Linguistics* 93, 7–16
- 473 Schilperoord, J. (1996). *It's about Time: Temporal Aspects of Cognitive Processes in Text Production*.
474 Utrecht studies in language and communication (Rodopi)
- 475 Sennrich, R., Firat, O., Cho, K., Birch, A., Haddow, B., Hitschler, J., et al. (2017). Nematus: a toolkit for
476 neural machine translation. In *Proceedings of the Software Demonstrations of the 15th Conference of
477 the European Chapter of the Association for Computational Linguistics* (Valencia, Spain), 65–68

- 478 Sennrich, R., Haddow, B., and Birch, A. (2016). Neural Machine Translation of Rare Words with Subword
 479 Units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*
 480 (*Volume 1: Long Papers*) (Berlin, Germany), 1715–1725
- 481 Toral, A. and Sánchez-Cartagena, V. M. (2017). A multifaceted evaluation of neural versus phrase-based
 482 machine translation for 9 language directions. In *Proceedings of the 15th Conference of the European*
 483 *Chapter of the Association for Computational Linguistics: Volume 1, Long Papers* (Valencia, Spain),
 484 1063–1073
- 485 Toral, A. and Way, A. (2015). Translating Literary Text between Related Languages using SMT. In
 486 *Proceedings of the Fourth Workshop on Computational Linguistics for Literature* (Denver, Colorado,
 487 USA), 123–132
- 488 Toral, A. and Way, A. (2018). What level of quality can neural machine translation attain on literary
 489 text? In *Translation Quality Assessment: From Principles to Practice*, eds. J. Moorkens, S. Castilho,
 490 F. Gaspari, and S. Doherty (Berlin/Heidelberg: Springer)
- 491 Vaswani, A., Zhao, Y., Fossum, V., and Chiang, D. (2013). Decoding with large-scale neural language
 492 models improves translation. In *Proceedings of the 2013 Conference on Empirical Methods in Natural*
 493 *Language Processing* (Seattle, Washington, USA), 1387–1392
- 494 Voigt, R. and Jurafsky, D. (2012). Towards a Literary Machine Translation: The Role of Referential
 495 Cohesion. In *NAACL-HLT Workshop on Computational Linguistics for Literature* (Montréal, Canada),
 496 18–25
- 497 Wieling, M., Nerbonne, J., and Baayen, R. H. (2011). Quantitative social dialectology: Explaining
 498 linguistic variation geographically and socially. *PloS one* 6(9). 14pp.

TABLES

Document	N-gram overlap			TTR	Sentence length
	2	3	4		
Warbreaker	0.86	0.67	0.41	0.15	12.54
Prologue	0.86	0.63	0.38		13.14
Chapter 1	0.87	0.66	0.41		13.81
Chapter 2	0.89	0.67	0.42		13.08
12 books	0.86±0.03	0.63±0.03	0.37±0.03	0.17±0.03	16.78±3.03

Table 1. N-gram overlap ($n = \{2, 3, 4\}$), TTR and sentence length for *Warbreaker* and the means and 95% confidence intervals of those measures for the 12 books previously translated by Toral and Way (2018).

Keystroke type	Task Type				
	ht	mt1	$\Delta\%$	mt2	$\Delta\%$
Total	1.94	1.76	-9%	1.49	-23%
Content	1.52	0.69	-55%	0.56	-63%
Navigation	0.18	0.59	228%	0.53	195%
Erase	0.23	0.47	105%	0.40	72%

Table 2. Average number of different types of keystrokes used to translate each source character in each translation condition. For conditions MT1 and MT2, the relative changes with respect to translation from scratch (HT) are shown alongside the absolute values.

FIGURE CAPTIONS

Predictor	Temporal (time)	Technical (keystrokes)	Cognitive (pauses)		
			number	mean duration	ratio
Source length	↑***	↑***	↑***	↑***	↑*
Trial	↓***	↓*	↓.	↓**	-
Condition (MT1 vs HT)	↓***	↓**	↓***	↑***	↑***
Condition (MT2 vs HT)	↓***	↓**	↓***	↑***	↑*
Condition (MT2 vs MT1)	↓*	↓**	↓***	↑**	-
Length:MT1	-	-	-	-	-
Length:MT2	↑**	↑***	-	-	-

Table 3. Significance of predictors in the mixed models built for each effort dimension. Significance levels: - ($p > 0.1$), . ($p \leq 0.1$), * ($p \leq 0.05$), ** ($p \leq 0.01$), *** ($p \leq 0.001$). Direction: ↑ (significantly higher), ↓ (significantly lower). Two comparisons are carried out for level MT2 of the predictor condition (i.e. against levels HT and MT1), hence we correct these p -values with Holm-Bonferroni.

ready to edit! revisions: 0 total: 0s

Mentre Bevalis era tècnicament la capital d'Idris, no era tan gran, i tothom la coneixia de vista.

0/10

0 saved

To another, it might have been offensive.	A un altre, podria haver estat ofensiu.	
To Siri it was a blessing.	A Siri va ser una benedicció.	
She smiled, walking into the city proper.	Ella va somriure, caminant cap a la ciutat com Déu mana.	
She drew the inevitable stares.	Va dibuixar les mirades inevitables.	
<i>While Bevalis was technically the capital of Idris, it wasn't that big, and everyone knew her by sight.</i>	<i>Mentre Bevalis era tècnicament la capital d'Idris, no era tan gran, i tothom la coneixia de vista.</i>	
Judging by the stories Siri had heard from passing ramblemen, her home was hardly even a village compared with the massive metropolises in other nations.	Si s'havia de jutjar per les històries que Siri havia sentit des de divagacions, la seva llar amb prou feines era ni tan sols un poble comparat amb la immensa metròpolis en altres nacions.	
She liked it the way it was, even with the muddy streets, the thatched cottages, and the boring - yet sturdy - stone walls.	Li agradava com era, fins i tot amb els carrers enfangats, les casetes de palla, i els avorrits - però robustos - murs de pedra.	
Women chasing runaway geese, men pulling donkeys laden with spring seed, and children leading sheep on their way to pasture.	Dones que perseguien oques desbocades, homes que arrossegaven ases carregats amb llavor de primavera, i nens que duïen ovelles de camí cap a pastura.	
A grand city in Xaka, Hudres, or even terrible Hallandren might have exotic sights, but it would be crowded with faceless, shouting, jostling crowds, and haughty noblemen.	Una gran ciutat de Xaka, Hudres, o fins i tot terrible Hallandren podia tenir vistes exòtiques, però seria ple de gent sense cara, crits, empentes i nobles nobles.	

Figure 1. Snapshot from the translation environment showing the third task, in which the translator is to post-edit the translations produced by the NMT system for sentences 21–30 of the chapter.

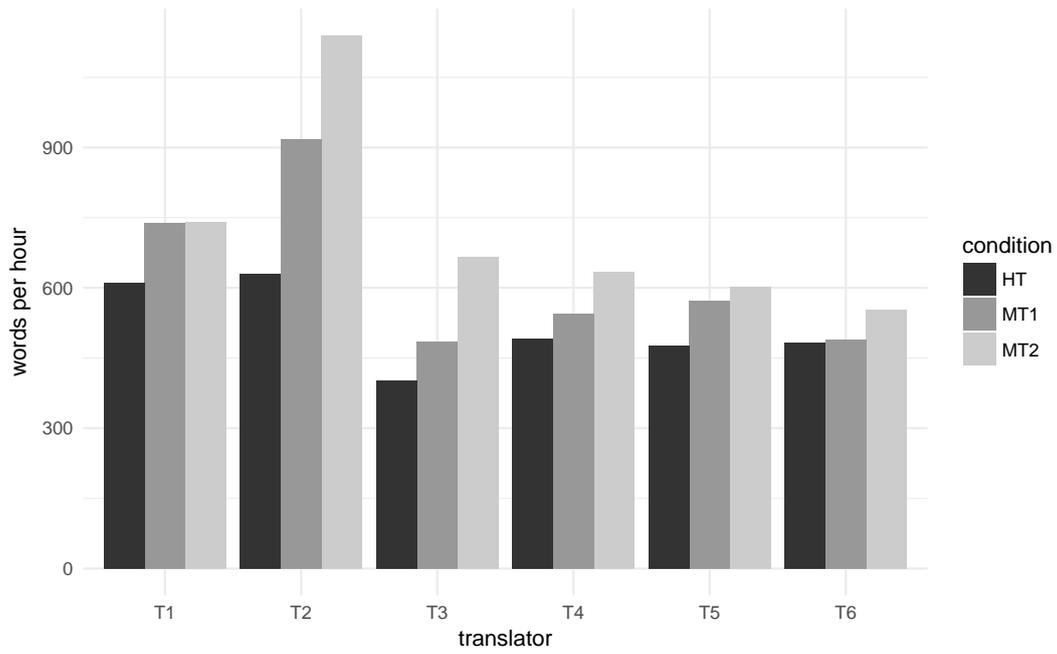


Figure 2. Translation productivity measured as words per hour for each of the translators in each of the translation conditions.

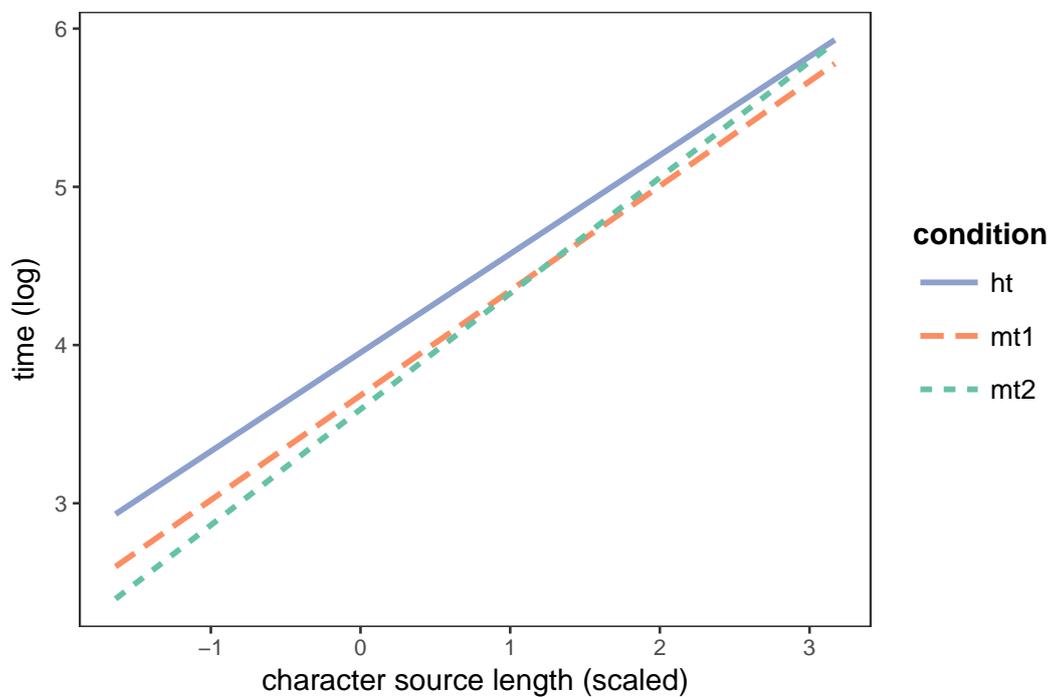


Figure 3. Interaction between the length of the source sentence and the translation condition on temporal effort.

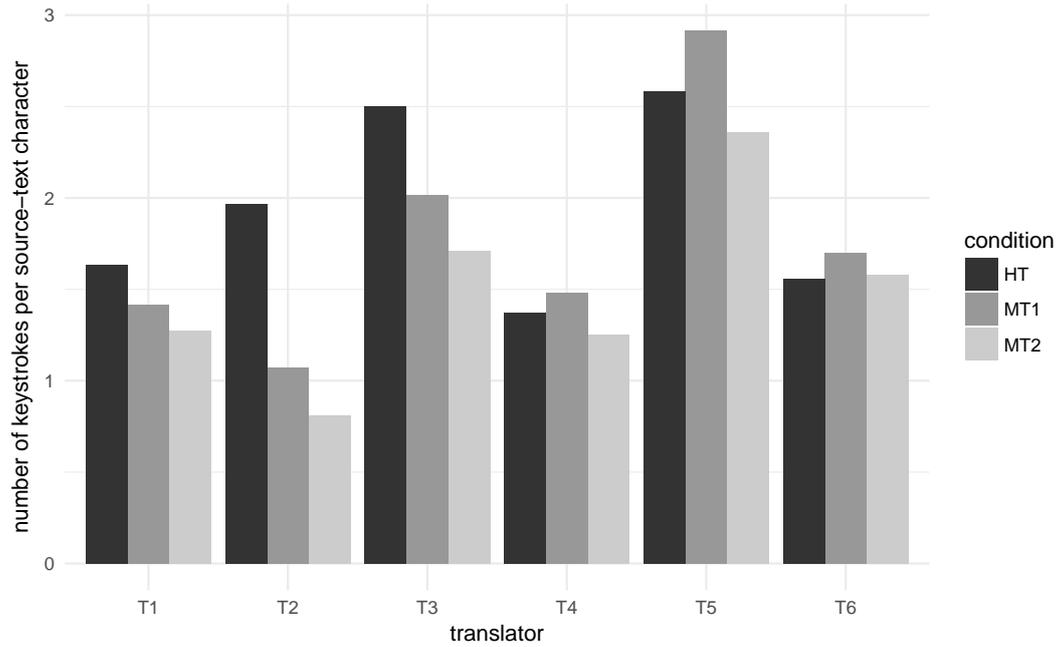


Figure 4. Technical effort measured as number of keystrokes per source character for each of the translators under each of the translation conditions.

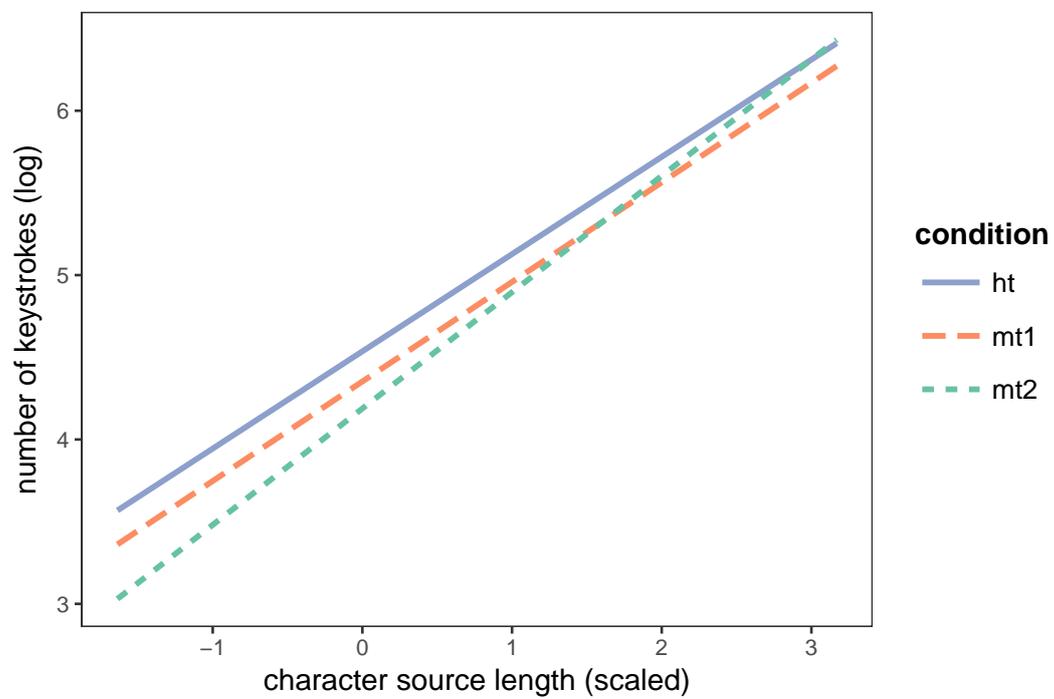


Figure 5. Interaction between the length of the source sentence and the translation condition on technical effort.

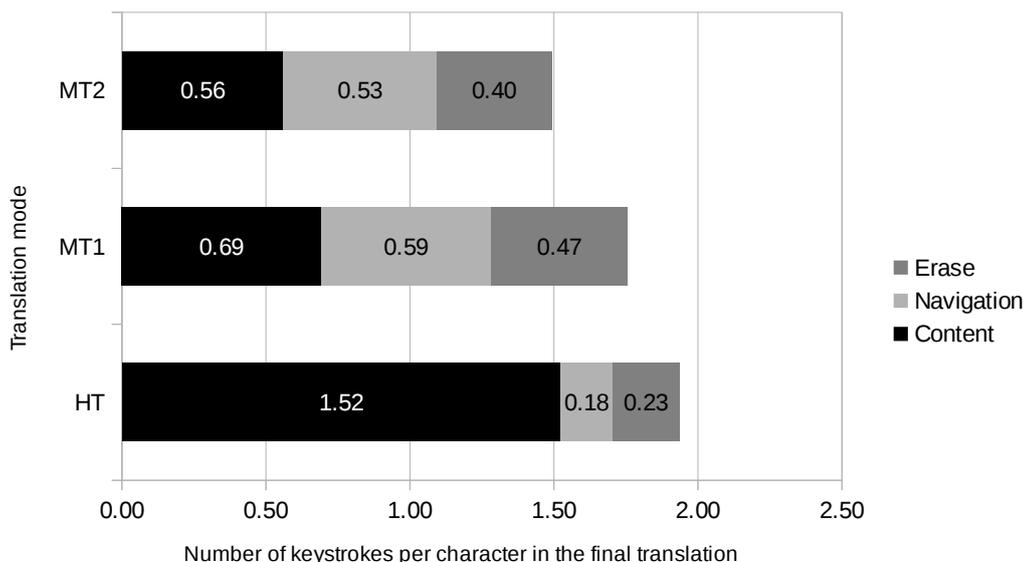


Figure 6. Proportion of each keystroke type (content, navigation and erase) in each translation condition (HT, MT1 and MT2) aggregating all the translators.

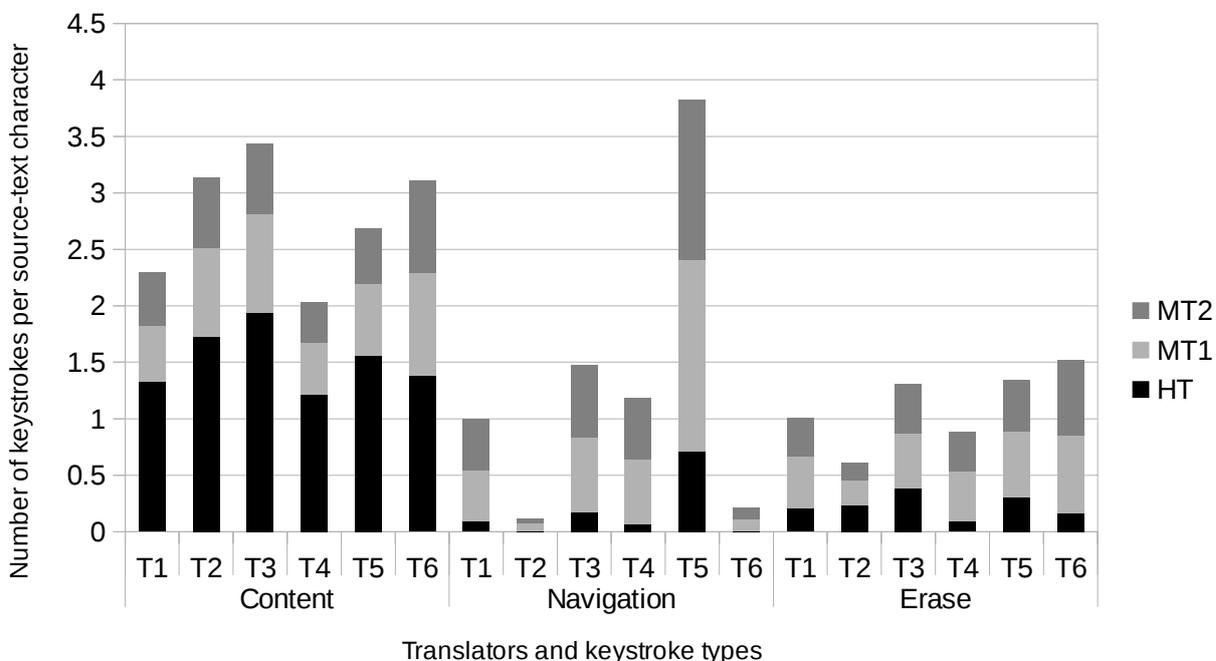


Figure 7. Proportion of each keystroke type (content, navigation and erase) in each translation condition (HT, MT1 and MT2) and for each translator.