

Interactive Translation Prediction vs. Conventional Post-editing in Practice: A Study with the CASMACAT Workbench

Germán Sanchis-Trilles · Vicent Alabau · Christian Buck ·
Michael Carl · Francisco Casacuberta · Mercedes García-
Martínez · Ulrich Germann · Jesús González-Rubio ·
Robin L. Hill · Philipp Koehn · Luis A. Leiva · Bartolomé Mesa-
Lao · Daniel Ortiz-Martínez · Herve Saint-Amand ·
Chara Tsoukala · Enrique Vidal

Received: date / Accepted: date

Abstract We conducted a field trial in computer-assisted professional translation to compare Interactive Translation Prediction (ITP) against conventional post-editing (PE) of machine translation (MT) output. In contrast to the conventional PE set-up, where an MT system first produces a static translation hypothesis that is then edited by a professional (hence “post-editing”), ITP constantly updates the translation hypothesis in real time in response to user edits. Our study involved nine professional translators and four reviewers working with the web-based CASMACAT workbench. Various new interactive features aiming to assist the post-editor/translator were also tested in this trial. Our results show that even with little training, ITP can be as productive as conventional PE in terms of the total time required to produce the final translation. Moreover, translation editors working with ITP require fewer key strokes to arrive at the final version of their translation.

Keywords CAT, SMT, interactive translation prediction, post-editing, field trial, user studies

G. Sanchis-Trilles · V. Alabau · F. Casacuberta · J. González-Rubio
L. A. Leiva · D. Ortiz-Martínez · E. Vidal
Pattern Recognition and Human Language Technologies Center
Universitat Politècnica de València
Valencia, Spain
Tel.: +34 963 878 172
E-mail: gsanchis@dsic.upv.es

C. Buck · U. Germann · R. Hill · P. Koehn · H. Saint-Amand · C. Tsoukala
School of Informatics
University of Edinburgh
Edinburgh, United Kingdom

M. Carl · M. García-Martínez · B. Mesa-Lao
Department of International Business Communication
Copenhaguen Business School
Copenhaguen, Denmark

1 Introduction

Contemporary professional translators rarely produce translations entirely from scratch. Instead, they increasingly rely on Translation Memories (TM), that is, data bases of texts that have already been translated, and their translations. At translation time, translations of text fragments similar to the actual source text are retrieved from the data base and edited by the translator to bridge the mismatch between retrieved text fragments and an actual correct translation of the current source text. As the quality of the raw output of fully automatic machine translation (MT) systems is on the rise, so is the commercial interest in integrating MT as an alternative or supplement to traditional TMs into the professional translation workflow. Recent studies (Koehn 2009a; Flourney and Duran 2009; Plitt and Masselot 2010; Federico et al. 2012; Green et al. 2013) have concluded that post-editing is, on average, more efficient than translating from scratch. However, the optimal form of human-computer interaction in the context of translation is still an open research question.

The open-source project CASMACAT addresses two needs in this area: first, it provides a new post-editing workbench for professional translators that is unobtrusive, yet provides support to the translator when it is relevant to do so; and second, it is able to log user activity in detail and thus record research data that can shed light on the mental processes underlying human translation in a computer-assisted translation (CAT) setting.

CASMACAT builds on the open-source, web-based post-editing tool MATECAT¹ and adds several major capabilities to the framework:

1. It offers *interactive translation prediction* (ITP; Barrachina et al. 2009) as an alternative to classical post-editing. The ITP functionality used in this study has been implemented by means of the Thot toolkit for statistical MT (Ortiz-Martínez and Casacuberta 2014). Various auxiliary features and customisations have been implemented to help tailor the MATECAT tool to the individual translator's preferences. They are described in Section 2.
2. CASMACAT can log user activity in detail and with precise timing information: key strokes, mouse activity, and translator's gaze (if used in combination with an eye tracker). Without eye tracking, the tool can be easily deployed in a web browser, eliminating the need for specialised hardware or software to run experiments. The logs from the user study discussed in this paper are available online for further analysis at http://bridge.cbs.dk/platform/?q=CRITT_TPR-db.
3. CASMACAT can be used with an e-pen as an alternative input device (Alabau et al. 2014). Such an interface is comfortable and effective in a number of situations. First, it is suited for post-editing sentences with only few errors, as it is often the case for sentences with strong fuzzy matches in translation memories, or during revision of human post-edited sentences. Second, it allows to perform such tasks while commuting, travelling or away from the desk for other reasons. The e-pen interface is also able to recognise gestures for interactive text editing, using a highly accurate, high-performance gesture recogniser (Leiva et al. 2013).

¹ <http://www.matecat.com>

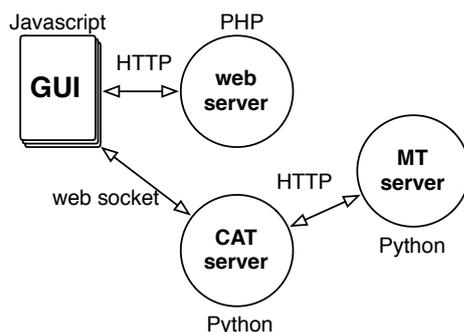


Figure 1 Components of the workbench

In the following, we present the results of a focused controlled user study of the CASMACAT workbench with professional translators that addressed the following questions:

- Does ITP boost or hinder overall translation productivity, especially when compared to conventional post-editing?
- What effect do different ITP visualisation options have on the interactive translation process?
- How satisfied are users with regard to the translations produced?

2 The CASMACAT workbench

The CASMACAT workbench consists of several components (Figure 1).

1. a graphical user interface (GUI) implemented as a web browser plugin in JavaScript;
2. a web server backend implemented in PHP that retrieves translation jobs from a MySQL database;
3. a CAT server that manages interactive translation prediction and event logging during an edit session; and
4. an MT server that provides raw translations as well as the underlying search graphs (compact representations of all translation options considered) to the CAT server.

The latter two components are implemented in Python but interface and interact with additional third-party components written in a variety of programming languages.

The GUI and the CAT server communicate via web sockets for speed; the other communication pathways are handled over HTTP for maximum compatibility with other software components. For example, the communication between CAT server and MT server relies on an extension of the Google Translate API, so that other MT engines compliant with the Google Translate API can easily be swapped in if desired. The web back-end accepts translation job uploads and offers file downloads in standard XML Localisation Interchange File Format (XLIFF).

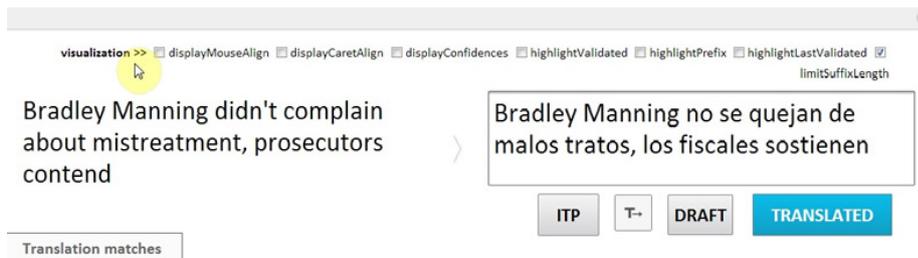


Figure 2 Screenshot of CASMACAT with optional visualisation features disabled

The CASMACAT workbench offers numerous user customisation options. In its most basic form (Figure 2), the tool is reminiscent of standard CAT tools. The source text is partitioned into a series of translation segments (typically individual sentences), with the source text shown on the left and an edit window on the right that allows editing of the translation of the “current” segment.

This basic interface is augmented by additional functionalities and display customisation options:

- **Intelligent autocompletion:** This is the fundamental interactive prediction feature of the CASMACAT workbench. Every time a keystroke is detected, the system produces a translation prediction for the entire sentence in accordance with the text that the user is writing or editing. The text to the left of the cursor is assumed to be approved by the human translator and serves as a prefix to identify the highest-scoring automatic translation that overlaps in this prefix. The remainder of the current translation prediction (to the right of the cursor) is then replaced with the updated prediction. Post-editing and ITP are mutually exclusive processing modes in the tool; the user can easily toggle between them by pressing a button below the text edit box (currently labeled ‘ITP’ in Figure 2). ITP — which could also be described as concurrent editing — and PE are quite different modes of operation cognitively as well. To facilitate the discussion in the remainder of this paper, we will henceforth subsume both forms of editing raw MT output under the term *editing*² whenever this distinction is of no relevance to the point of discussion.
- **Prediction rejection:** The current CASMACAT prototype also allows the human editor to use the mouse wheel to scroll through translation options (Sanchis-Trilles et al. 2008). When the mouse wheel is turned over a word, the system invalidates the current prediction and provides the user with an alternate translation option in which the first new word is different from the one at the current mouse position. This option is one of the advanced ITP features.
- **Search and replace** (even in future predictions): the workbench extends standard search-and-replace functionality to future translation predictions. Whenever a new replacement rule is created, it is automatically propagated to the forthcoming predictions made by the system, so that the user only needs to

² In anticipation of the discussion of the CASMACAT field trial in Section 4, we should also point out that editing is different from *reviewing*, in which a third party reviews and revises the text produced by the initial editor for quality assurance.

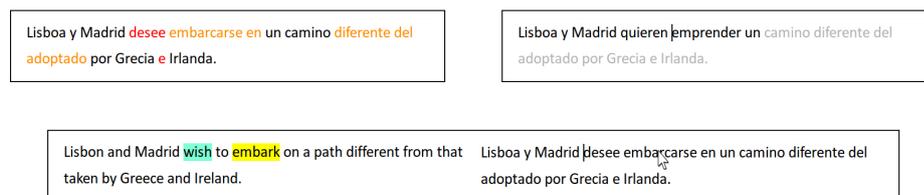


Figure 3 Advanced display options in CASMACAT: colour-coding of confidence estimates (top left), limited prediction horizon (top right), and word alignment visualisation (bottom).

specify them once. This specific function was implemented in response to user feedback in the first field trial of the tool. Note that this option implements a collection of replace rules, but does not resort to a fully fledged SMT system for doing so as in (Simard and Foster 2013).

The user can also choose a number of advanced visualisation options (Figure 3):

- **Visualisation of MT system confidence.** Automatic estimation of the reliability of the MT system output, also known as *confidence estimation*, is currently an active area of research. The CASMACAT workbench is able to colour-code such confidence estimates in the prediction. MT output identified as probably incorrect is marked in red; MT output of questionable reliability is coloured orange.
- **Limited prediction horizon.** Providing the human editor with a new prediction whenever a key is pressed has been shown to be cognitively demanding (Alabau et al. 2012). In the current prototype, when this option is active, predictions are shown only up to the first word of low confidence according to the confidence estimates associated with the prediction. Pressing the TAB key allows the user to ask the system for the next set of predicted words, displaying the remaining words in the suggested translation in grey.
- **Word alignment information.** Alignment of source and target information is an important part of the translation process (Brown et al. 1993). In order to display the correspondences between both the source and target words, this feature was implemented so that every time the user places the mouse (yellow) or the text cursor (cyan) on a word, the alignments made by the system are highlighted. The user can enable this visualisation option by activating *displayCaretAlign* for the alignments with the cursor and *displayMouseAlign* for the alignments with the mouse.
- **Visualisation of user edits** (not shown in Figure 3). This visualisation option comes in three variants, all of them implemented with the purpose of helping the user locate which changes were introduced by him, or what was produced by the system without interaction.
 - *changed words only*: the system highlights in green the words that the user has modified.
 - *entire prefix*: the system highlights the prefix, i.e. the first part of the segment that the user has validated.
 - *last edit only*: the system highlights the last word that the user has modified.

Table 1 Translation / post-editing process data logged and stored by CASMACAT. For details about these features, see Carl and Kay (2011); Carl (2012b; 2014).

<p>Keystrokes (KD): basic text modification operations (insertions or deletions), together with time of stroke, and the word in the final text to which the keystroke contributes.</p> <p>Fixations (FD): basic gaze data of text fixations on the source or target text, defined by the starting time, end time and duration of fixation, as well as the offset of the fixated character and word in the source or target window.</p> <p>Production units (PU): coherent sequence of typing, defined by starting time, end time and duration, percentage of parallel reading activity during unit production, duration of production pause before typing onset, as well as number of insertions and deletions.</p> <p>Fixation units (FU): coherent sequences of reading activity, including two or more subsequent fixations, characterised by starting time, end time and duration, as well as scan path indexes to the fixated words.</p> <p>Activity Units (CU): exhaustive segmentation of the session recordings into activities of typing, reading of the source or reading of the target text.</p> <p>Source tokens (ST): as produced by a tokeniser, together with TT correspondence, number, and time of keystrokes (insertions and deletions) to produce the translation and micro unit information (see below).</p> <p>Target tokens (TT): as produced by a tokeniser, together with ST correspondence, number, and time of keystrokes (insertions and deletions) to produce the token, micro unit information, amount of parallel reading activity during.</p> <p>Alignment units (AU): transitive closure of ST-TT token correspondences, together with the number of keystrokes (insertions and deletions) needed to produce the translation, micro unit information, amount of parallel reading activity during AU production, etc.</p> <p>Segments (SG): aligned sequences of source and target text segments, including duration of segment production, number of insertions and deletions, number and duration of fixations, etc.</p> <p>Session (SS): global properties of the session, such as source and target languages, total duration of the session, beginning and end of drafting, etc.</p>
--

3 Translation and post-editing process data

Another important feature of the CASMACAT workbench is its ability to record user activity in fine detail for analysing human and computer-assisted translation processes scientifically. That is, the tool not only stores translation product information (the source, raw MT output and final translation), but can also provide detailed editing process data with precise timing information, including eye tracking data if used in combination with an eye tracker.³ A gaze-to-word mapping algorithm runs in real time, and maps gaze samples and fixation points to the nearest letter on the screen; the character offset is then logged together with the gaze data. The tool also keeps a record of the different translation options that were presented to the user at the time. At storage time, CASMACAT aggregates and stores information about phases of coherent writing (production units; PU) and reading (fixation units; FU) from the raw user activity data (UAD). Table 1 summarises the information stored during interactive translation and post-editing sessions. During analysis, we derived further aggregate information from the stored UAD. These derived measures are described in Section 5.

³ In our experiments, we use an EyeLink1000 eye tracker.

4 Field trial

In June 2013, we conducted the second CASMACAT field trial in cooperation with Celer Soluciones SL, a language service provider based in Madrid. The trial involved nine freelance translators and four reviewers, all native speakers of Spanish who provide translation and post-editing services on a regular basis for this company. Detailed information about participants' age (46 years old on average), years of experience in translation (23 years on average), education (Translation, Philology or another degree), etc., can be found in the CRITT Translation Process Research database under the metadata folder.⁴

The trial texts were chosen from the WMT-2012 *news commentary* corpus (Callison-Burch et al. 2012), a collection of general news stories and opinion pieces. The texts consisted of 30 to 63 segments each, totalling approximately 1,000 words per document, as shown in Table 2. Each English source text was automatically translated into Spanish by a statistical MT system and then automatically loaded into the CASMACAT workbench for the participants to post-edit or translate interactively.

In an attempt to unify post-editing and translation criteria among participants, all of them were instructed to follow the same guidelines aiming at a final high-quality target text (publishable quality). The guidelines distributed in hard copy were: i) Retain as much raw MT as possible; ii) Do not introduce stylistic changes; iii) Make corrections only where absolutely necessary, i.e. correct words and phrases that are clearly wrong, inadequate or ambiguous according to Spanish grammar; iv) Make sure there are no mistranslations with regard to the English source text; v) Publishable quality is expected. The task of the four reviewers was to proofread and correct the texts produced by the initial translation editors. Below, we will refer to translations as produced by PE or (A)ITP as *(edited) draft translations* (as opposed to raw MT output), and to their final versions after review as *final translations*.

4.1 Experimental design

Three system setups were evaluated in the trial: conventional post-editing (PE), basic interactive translation prediction (ITP), and interactive translation prediction with advanced features (AITP). The test set consisted of three sets of three texts each. Each text was processed by three editors in each of the three conditions. Each editor processed each document exactly once under one of the three conditions. Table 2 shows the task assignments. Dataset 1 was processed under laboratory conditions with additional eye tracking at the office of Celer Soluciones SL. Datasets 2 and 3 were delivered over the Internet and processed at home. In each instance, keyboard and mouse activity was logged. For the PE setup, the highest-scoring translation hypothesis was used; ITP and AITP relied on a translation search graph delivered by the MT system. In the AITP condition, study participants were given access to all the advanced ITP features described in Section 2 and could freely choose which ones to enable and use.

The edited draft translations of Dataset 1 were subsequently proofread at Celer Soluciones SL, where each of the reviewers was assigned to review the work done

⁴ http://bridge.cbs.dk/platform/?q=CRITT_TPR-db

Table 2 Task assignments in the field trial

	text								
	dataset 1			dataset 2			dataset 3		
	T1.1	T1.2	T1.3	T2.1	T2.2	T2.3	T3.1	T3.2	T3.3
segments	49	30	45	63	55	51	59	61	47
source words	952	861	1121	1182	1216	1056	1396	1427	1258
Editor 1	PE	ITP	AITP	ITP	AITP	PE	AITP	PE	ITP
Editor 2	AITP	PE	ITP	PE	ITP	AITP	ITP	AITP	AITP
Editor 3	ITP	AITP	PE	AITP	PE	ITP	PE	ITP	PE
Editor 4	ITP	AITP	PE	AITP	PE	ITP	PE	ITP	AITP
Editor 5	PE	ITP	AITP	ITP	AITP	PE	AITP	PE	ITP
Editor 6	AITP	PE	ITP	PE	ITP	AITP	ITP	AITP	PE
Editor 7	AITP	PE	ITP	PE	ITP	AITP	ITP	AITP	PE
Editor 8	ITP	AITP	PE	AITP	PE	ITP	PE	ITP	AITP
Editor 9	PE	ITP	AITP	ITP	AITP	PE	AITP	PE	ITP

by a maximum of three editors. Gaze and keyboard activity for reviewers was also logged.

Before starting their tasks, participants were introduced to the CASMACAT workbench and the three different conditions under consideration during the trial. They were given time to familiarise themselves with the tool and try out the different visualisation options, and to decide which options they would enable when post-editing using AITP. After each session, participants were asked to complete an online questionnaire (see Section 5.3). When all sessions at Celer Soluciones SL were completed, an additional in-depth interview was conducted with each of the participants. Table 3 summarises the data collected during the trial.

5 System evaluation and results

User performance and evaluation is a central part of the CASMACAT project, and a rich dataset for analysis was collected during the field trial. This section provides several kinds of evaluation:

- Section 5.1 looks at the collected activity data, i.e. keystrokes and gaze data. In Section 5.1.1 we look at the amount of coherent typing activity needed to perform the task. Section 5.1.2 analyses the effort made by the study participants in terms of the number of insertions and deletions, and Section 5.1.3 the gaze behavior.
- Section 5.2 describes several paths to assess the linguistic quality of the final product. Section 5.2.1 computes the edit distance between edited draft translations and their final counterparts after review, and section 5.2.2 correlates processing time, number of text modifications, and edit distance between edited translation drafts and final translations.
- Section 5.3 presents the feedback provided by the study participants in the form of questionnaires after completing each task.

Table 3 Data collected during the field trial. 460 distinct source segments were post-edited or interactively translated by 9 translators. The data of 54 segments was lost due to technical failure.

# of processing logs collected					
condition	total	w/ gaze data		English	Spanish
PE	1,345	372			
ITP	1,368	372			
AITP	1,373	372			
lost data	54	—			
total	4,086	1,116			
			total tokens	94,865	101,671
			mean segment length	23	25

5.1 Evaluation of activity data

Table 3 summarises the user activity data that were collected during the field trial. For Dataset 1, gaze data were collected from all post-editors/translators and reviewers. We analyzed the processing logs with respect to overall post-editing times, user effort in terms of edit operations, and gaze behavior.

5.1.1 Overall processing time

In principle, the total processing time for a segment is the time lapsed between the moment the post-editor or translator enters the edit box for a segment and the time he or she proceeds to the next one. However, in some of the logs from the sessions conducted from home, we observed very long pauses (up to several hours) suggesting that the respective participant interrupted these sessions and then returned to them later. By analysing the intervals between recorded edit events (recall that gaze data was not recorded for Datasets 2 and 3), we can make inferences about the underlying translation activity. Since the events themselves do not consume a significant amount of time (they are considered to have no duration), the total edit time is equivalent to the total of the time lapsed *between* events, i.e., during pauses (including intervals between keystrokes during continuous typing).

In our data, the vast majority of the pauses had a duration of a few seconds or less. Figure 4 shows the pause duration by means of a box plot. Box plots visualise data by means of a box that includes the first and third quartiles of the distribution as well as two arms or whiskers containing the extreme values. Box plots can also represent outliers⁵ as isolated points at the left or at the right of the whiskers. Our data contain outliers so extreme that they could not be represented in the box plot without negatively affecting its legibility. Because of this, they have not been included in the diagram. Excluding outliers, all inter-keystroke intervals had a duration of 0.8 seconds or less. However, this does not mean that all of the outliers corresponded to noisy observations. Therefore, it is necessary to analyse the pauses more carefully. Here, we present two techniques to filter pause data in a meaningful way.

The first technique assumes that processing consists of alternating periods of typing and cognitive processing activity. Based on cognitive language processing and production theory (Alves and Vale 2009; Lacruz et al. 2012; Carl 2012a),

⁵ Outliers are defined here as those points that exceed $Q3+1.5$ times the inter-quartile range, cf. Montgomery (2004).

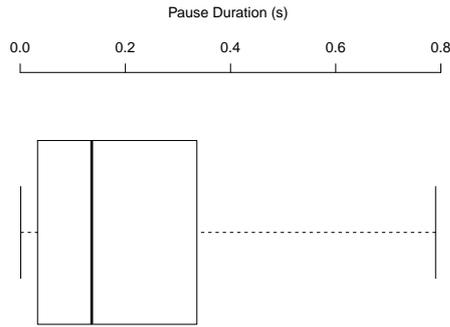


Figure 4 Boxplot for inter-keystroke pause duration in seconds (outliers not shown).

pauses between 0 and 5 seconds are used to segment the text production rhythm into “typing” and “processing” units.

In spite of the fact that the vast majority of the pauses had a duration of a few seconds, Figure 4 does not reflect their relative contribution to the total processing time. It is possible that there exist longer pauses that account for a substantial part of the segment processing time, even if they appear in a very small number (e.g. only one pause of 100 seconds accounts for the same time as one hundred pauses of 1 second). To clarify this, we generated a plot for different intervals of pause durations, summing their contributions to the total translation time. The result is represented as a weighted Pareto chart in Figure 5. Pareto charts are used to highlight the most important factor among a typically large set of them. For this purpose, bars and a line graph are used, where the frequencies of individual values are represented in descending order by bars, and the cumulative total is represented by the line. Specifically, in a Pareto chart the left vertical axis represents the frequency of occurrence, while the right vertical axis is the cumulative percentage of the total number of occurrences. In weighted Pareto charts, frequencies are multiplied by specific magnitudes such as cost or loss associated with particular events so as to better analyse their importance.⁶ The black line in the plot marks a relative frequency equal to 95%.

The plot given in Figure 5 provides valuable information about the effects of filtering pauses of a specific duration. The frequency of pauses belonging to a specific duration interval is weighted by such duration. For instance, pauses with a duration between 0 and 10 seconds (first bar in the plot) account for 58% of the processing time. Accordingly, filtering pauses of 10 seconds or more would leave 42% of the total processing time unaccounted for. As the plot shows, we need to set the threshold at 200 seconds to account for 95% of the observed total processing time.

Given these considerations, we applied two kinds of filtering to the set of inter-keystroke pauses, obtaining two new processing time measures (see also Section 3):

- **Kdur**: the total durations of *continuous typing* activity, defined as sequences of keyboard events that are at most five seconds apart.
- **Fdur**: total durations of processing, ignoring pauses of 200 seconds or more.

⁶ For a detailed introduction to this kind of analysis, see e.g. Montgomery (2004).

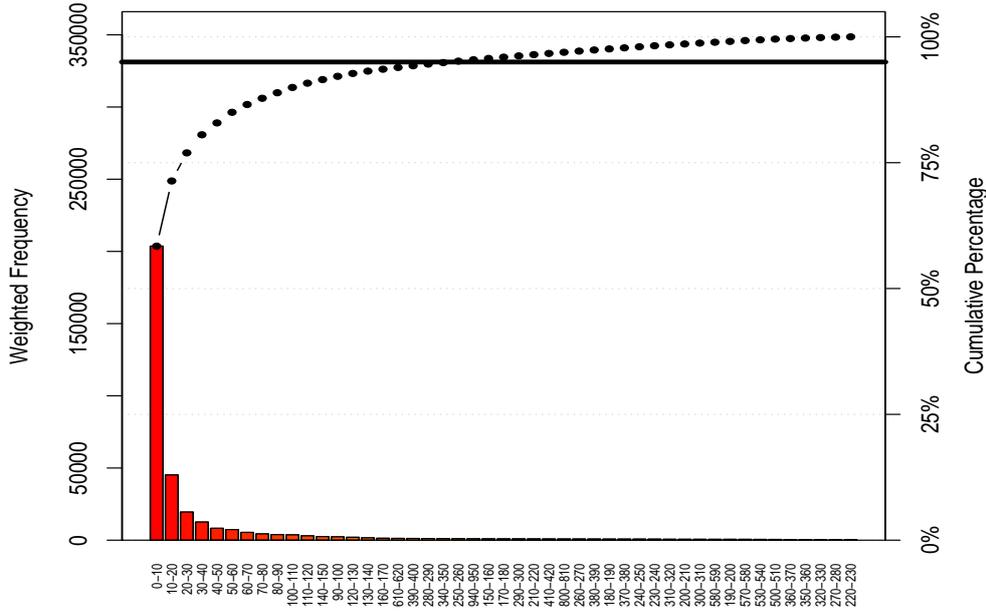


Figure 5 Weighted Pareto chart for inter-keystroke pause duration in seconds.

Table 4 shows the average segment processing times in seconds for PE, ITP, and AITP systems for three different time measurements, namely Tdur (total duration), Kdur and Fdur. PE allows for shorter processing times according to the Kdur and Fdur measures. However, for Fdur the difference between PE and ITP is small: ITP was only 5% slower. One possible explanation for the difference between PE and ITP in terms of Kdur is that ITP system users execute a higher number of short post-edit operations. Finally, Tdur values are distorted by a small number of extreme outliers in the data, as described earlier.

Table 4 Average processing time in seconds per segment in terms of Tdur, Kdur and Fdur when using PE, ITP, and AITP systems.

System	Tdur	Kdur	Fdur
PE	104.0	21.7	73.0
ITP	80.7	27.0	77.0
AITP	117.1	29.6	92.4

One important thing to take into account when analysing the average processing times is the learning curves of each system. While PE systems are typically well-known by translators, this is not the case for ITP systems. For this reason, it is also interesting to compare the processing times when working from Celer Soluciones SL (dataset 1) with those obtained when working from home (datasets 2 and 3). Since the translation editors worked from the office first, then from home, we can expect better performance with the ITP system from home after more hours of interaction (each dataset needed an average of 3.5 hours to be post-edited).

Table 5 contrasts average processing times at the office and at home in terms of Kdur and Fdur for the PE, ITP and AITP systems. As can be seen in the table, the average processing time dropped over time for all systems when the participants worked from home. The speed-up is greater for the ITP and AITP systems than for the PE system, which correlates with our assumption that editors were more familiar and thus more efficient with the conventional PE process initially.

Table 5 Average processing time in seconds per segment in terms of Kdur and Fdur when working at the office or from home using PE, ITP and AITP systems.

System	Kdur		Fdur	
	Office	Home	Office	Home
PE	27.7	19.6	88.0	67.3
ITP	35.1	24.8	94.7	71.9
AITP	37.5	27.2	111.9	87.8

5.1.2 Typing activity

Enabling interactivity has also an effect on the number of text insertions and deletions by the user. Table 6 shows the average number of manual insertions and deletions per segment using the three systems at the office, at home or for all the sessions. According to the results, the ITP system required fewer operations than the rest of the systems.

Table 6 Number of insertion and deletions operations per segment for translations generated at the office, at home or both using PE, ITP and AITP systems.

System	Office	Home	All
PE	114.9	134.6	131.3
ITP	109.6	127.2	123.6
AITP	143.2	137.0	132.6

It is important to note that these results must be interpreted in the light of the quality of the final draft produced by the translation editors, which we discuss in Section 5.2.

5.1.3 Gaze data

According to the eye-mind hypothesis of Just and Carpenter (1980), eye fixations correlate highly with what is currently processed mentally and can therefore be used as a window into instances of effortful cognitive processing. This assumption is one of the foundations of research with eye tracking: recordings of eye fixations can provide a dynamic trace of where a person’s attention is being directed.

On the basis of the eye-mind hypothesis, the average duration of gaze fixations in the source and target windows were calculated for each of the three systems in the field trial. Table 7 shows how participants exhibited a marked difference in the amount of time for which they gazed at the source and target windows. The use

Table 7 Average gaze fixations on source and target window per system.

System	PE		ITP		AITP	
	Nr.	%	Nr.	%	Nr.	%
Source window	18037	33.3	14422	26.0	16569	26.5
Target window	36193	66.7	41052	74.0	45999	73.5
Total	54230	100.0	55474	100.0	62568	100.0

of interactivity features both in ITP and AITP triggered longer gaze fixations in the target window.

Under all three system configurations the majority of gaze fixations is on the target rather than the source window. In contrast to translating from scratch, the translation editor’s task is to edit the MT output presented in the target window, and it is therefore not surprising that the primary focus is on that window. Enabling interactivity (ITP) and visualisation (AITP) shifts the focus even more into the target window. This, too, is to be expected, as the content of the target window is frequently changing in ITP mode.

5.2 Quality of edited draft translations

In this section we look at the quality of the edited draft translations for Dataset 1, which were subject to additional review after they had been submitted by the original editors. In Section 5.2.1, we use edit distance to assess the quality of the edited draft translations by comparing them against the respective reviewed, final versions. Section 5.2.2 correlates edit distance with the number of actual text modifications and revision time. (Edit distance is the minimum number of edit operations required to transform one text into the other, not the number of actual edit operations performed by the editor.)

5.2.1 Edit distance in Dataset 1

For this analysis, edit distance was computed in terms of the number of words changed. Words were chosen as the unit of editing because a word difference has typically a much closer correlation with both semantic quality and style than individual character differences. Moreover, rather than counting the absolute number of edit operations needed to transform the original text into the revised one, a relative figure (in %) is needed. This is important because the overall number of words is not the same for texts produced with the PE, ITP, and AITP systems and, without proper normalisation, differences could be due to variations in text sizes, rather than to possible quality differences. Finally, in order to ensure the estimates are true percentages, one needs to normalise by the total number of edit operations, N , including non-error matches (i.e., $N = ins + del + sub + corr$, where ins is the number of inserted words, del is the number of deleted words, sub is the number of word substitutions, and $corr$ is the number of correct words). That is, the normalised edit distance is $(ins + del + sub)/N$. Such a normalisation makes the product of the different systems fully and accurately comparable, regardless of the origin/reviewed sizes of each text.

The results of the analysis are shown in Table 8. Taking into account the 95% confidence intervals of these estimates ($\sim 1\%$), we can conclude that there is

no significant difference in the quality of the edited draft translations produced with the three assistance systems.

Table 8 Quantitative analysis of the changes introduced by the reviewers, measured as total number of words changed.

Assistance system	PE	ITP	AITP
<i>ins + del + sub</i>	286	314	307
<i>ins + del + sub + corr (=N)</i>	3082	2926	3050
Overall word changes (%)	9.3	10.7	10.1
Estimated quality (%)	90.7	89.3	89.9

When interpreting these numbers, we should also keep in mind that only Dataset 1 was analysed here. This means that the results are deduced from the translations generated while the editors were still getting used to the different systems.

5.2.2 Correlation of edit distance, revision time and text modifications

For this analysis, we counted the number of manual insertions and deletions for each of the four reviewers. Table 9 shows the average text modifications per system and reviewer R10 to R13. The table presents the average number of text modifications per segment divided by the length in characters of the segment for each of the three systems. In line with the results of Guerberof (2012), reviewers seem to follow very different reviewing styles: reviewer R10 produces the least number of text modifications, while reviewer R13 is the most eager corrector. On average reviewers insert more modifications when the post-edited text was produced with system ITP.

Table 9 Average amount of text modification (character insertions and deletions, in%) during review, per reviewer and system.

	PE	ITP	AITP	total
R10	8.9	0.8	4.8	4.8
R11	8.0	15.3	12.5	11.9
R12	9.4	9.8	8.8	9.3
R13	13.6	11.7	12.5	12.6
Total	10.0	9.4	9.7	9.7

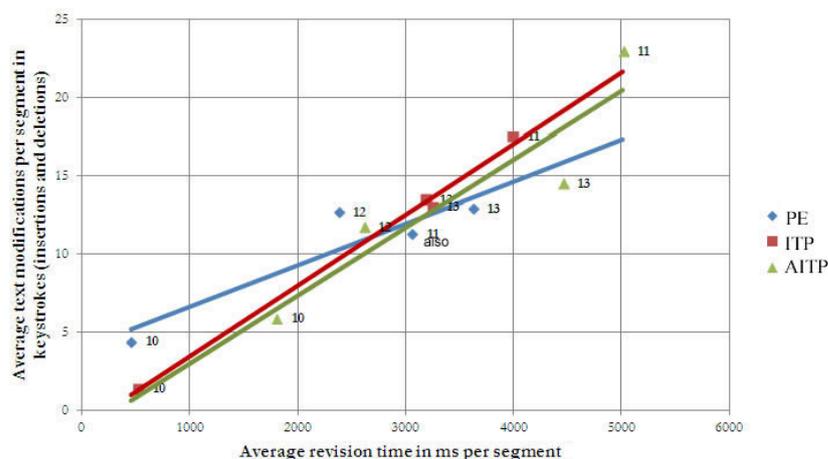
We also computed the average revision time, edit distance and number of text modifications per reviewing session, which resulted in 12 data points for each of the variables (three systems \times four reviewers). Unfortunately it was not possible to obtain reliable revision time on a segment level (which would have given many more data points) due to the fact that in the revision mode it was possible for the reviewer to read the segments, without loading them in the edit area of the workbench. As a consequence, we had to average over the entire revision session to get comparable numbers for average revision time, edit distance and number of text modifications.

Table 10 summarises correlation and significance values, and shows that there is a strong correlation between these variables, but due to the small number of data points significance is not very high.

Table 10 Correlations between keystrokes, edit distance and time in revision.

Assistance system	PE	ITP	AITP
Keystrokes vs. Time	$R^2 = .910$ $p > .081$	$R^2 = .998$ $p < .002$	$R^2 = .924$ $p > .076$
Edit dist. vs. Time	$R^2 = .740$ $p > .260$	$R^2 = .998$ $p < .002$	$R^2 = .946$ $p < .054$
Edit dist. vs. Keystrokes	$R^2 = .680$ $p > .320$	$R^2 = .999$ $p < .001$	$R^2 = .868$ $p > .132$

Figure 6 shows the correlations between text modifications and revision time. The highest correlation for all three variables can be observed in the ITP system and for the correlations between text modifications and revision time (Elming et al. 2014).

**Figure 6** Correlation between: keystrokes (insertions and deletions) vs. time in milliseconds

5.3 User feedback

User feedback was elicited from the editors in the form of questionnaires. After each session, they were asked to rate their level of overall satisfaction on a 1-5 Likert scale, where 5 corresponded the highest positive reply and 1 the lowest.

User feedback was collected regarding the following questions:

- How satisfied are you with the translations you have produced ? (Satisfaction)
- How would you rate the workbench you have just used in terms of usefulness/aids to perform a post-editing task? (Tool)
- Would you have preferred to work on your translation from scratch? (From scratch)
- Would you have preferred to work on the machine translation output without the interactivity provided by the system? (No ITP)

Table 11 summarises the feedback provided by the editors after working with each of the three systems.

Table 11 Satisfaction ratings while using PE, ITP and AITP systems

	Satisfaction			Tool			From Scratch			No ITP	
	PE	ITP	AITP	PE	ITP	AITP	PE	ITP	AITP	ITP	AITP
Editor 1	3	4	4	3	4	4	No	No	No	Yes	No
Editor 2	4	4	4	3	2	4	Yes	Yes	Yes	No	No
Editor 3	3	3	4	3	3	4	Yes	No	No	Yes	No
Editor 4	4	4	5	3	4	4	No	No	No	No	No
Editor 5	4	3	4	4	4	3	No	No	No	Yes	No
Editor 6	5	5	5	3	3	2	No	No	No	Yes	Yes
Editor 7	3	4	3	2	1	2	Yes	Yes	Yes	Yes	No
Editor 8	4	4	3	2	2	3	Yes	No	No	Yes	Yes
Editor 9	4	4	4	1	4	3	Yes	Yes	Yes	Yes	No

The results show different levels of satisfaction for the different systems. Some participants (i.e. 1, 3 and 4) seem to be more satisfied with the translations produced using interactive systems. Regarding the tool, interactive systems also are rated with a higher level of satisfaction overall, even though 7 out of 9 translators stated that they would have preferred not working with the interactivity provided by the system when using the ITP system. Their views are quite different when using AITP, since only two translators (6 and 8) continued thinking that they would have preferred to work without interactive features.

6 Related work

Improving the productivity of translators is and has been a major driver of MT research. The hope is that, in many cases, post-editing MT output will help translators to perform their work faster. Several studies were performed to evaluate the potential benefit with generally positive results. Measured reductions in translation time typically range somewhere between 18% and 34% (Flournoy and Duran 2009; Guerberof 2009; Federico et al. 2012) sometimes even reaching as high as 43% (Plitt and Masselot 2010).

The various studies differ in numerous parameters, which makes direct comparisons difficult:

- translator’s level of experience: volunteer/student (Koehn 2009a) vs. professionals (Plitt and Masselot 2010; Guerberof 2009);
- suitability of the MT system, especially when comparing older (Krings 2001) and more recent (Plitt and Masselot 2010) studies;
- the PE interface and subjects’ familiarity with it;
- language pairs and text domain;
- the way the data were collected and filtered.

Another question that has been addressed is how post-editing fares in comparison to translation from scratch in terms of final translation quality. Koehn (2009a) found that non-professionals are generally both faster and produce better translations when post-editing, a result that is consistent with later work by Plitt and Masselot (2010) as well as Green et al. (2013), showing a strong reduction in translation time and a lower error rate for professional translators as well. Interestingly, Plitt and Masselot (2010) also found that the difference between individual translators is much stronger than between language pairs and MT systems of varying

quality. Following this work, Skadiņš et al. (2011) observed (slight) negative effects of the post-editing setting for both productivity and quality for some translators but still affirm the overall helpfulness of MT suggestions.

Along with the MT systems, PE environments have developed over time, recently converging towards web-based setups (Koehn 2009b; Green et al. 2013) which integrate several aids in a single interface. Despite extensive research on Confidence Estimation for Machine Translation, such annotation has yet to be integrated. Bach et al. (2011), for example, suggested visualising word-level confidences by type size.

Besides quantity and quality, the translation process itself has been studied for many years, starting with explicit collection of translators' thoughts using Think Aloud Protocols (Krings 2001). Possible interference with the translation process quickly led to passive/indirect collection of user activity such as the logging of keystrokes and mouse movement (Langlais et al. 2000) and, more recently, gaze data (O'Brien 2009; Doherty et al. 2010; Carl 2012b). By presenting multiple languages simultaneously in an ecologically valid environment, the combination of workbench and logging functions also offers a unique opportunity to investigate broader issues of applied bilingual cognitive processing.

7 Conclusions and future work

We have presented results of a comparative evaluation of two ITP systems versus conventional PE in a professional setting. The first ("ITP") provides only text prediction; the second ("AITP") offers additional functionality designed to aid the translator/editor in the translation and post-editing process. Our results show that the ITP system accomplishes what it was designed to do: ITP minimises the number of key strokes required to generate the translations. Nevertheless, the translation time per segment was slightly higher in ITP mode than in conventional PE mode. However, depending on how processing times were measured, the difference was small: according to the Fdur measure, processing in ITP mode took only 5% longer than conventional PE. Our findings also suggest that certain types of users may benefit more from ITP as their familiarity with this mode of processing increases. In contrast, the processing time results were worse for the AITP system. This suggests that some of the advanced features that were incorporated might not be useful to increase user productivity. However, we should keep in mind that more complex systems may have a steeper learning curve. Considering that translators were already experienced post-editors, we may suspect a slight experimental bias in favour of conventional PE in the way the study was set up. To reduce the influence of familiarity vs. novelty, a respective study would have to be longitudinal, observing translator/editor performance and interaction with the tool over extended periods of time.

The analysis presented in this paper aggregates results for the different system configurations across users and text segments. A logical next step is to look in more detail at the individual editors and texts in order to determine factors that determine which mode of editing works best for whom under which conditions.

Acknowledgements This work was supported by the European Union's 7th Framework Programme (FP7/2007-2013) under grant agreement N^o 287576 (CASMACAT).

References

- Alabau V, Leiva LA, Ortiz-Martínez D, Casacuberta F (2012) User evaluation of interactive machine translation systems. In: Proceedings of the 16th Annual Conference of the European Association for Machine Translation, pp 20–23
- Alabau V, Buck C, Carl M, Casacuberta F, García-Martínez M, Germann U, González-Rubio J, Hill R, Koehn P, Leiva L, Mesa-Lao B, Ortiz-Martínez D, Saint-Amand H, Sanchis-Trilles G, Tsoukala C (2014) Casmacat: A computer-assisted translation workbench. In: Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics, pp 25–28
- Alves F, Vale D (2009) Probing the unit of translation in time: aspects of the design and development of a web application for storing, annotating, and querying translation process data. *Across Languages and Cultures* 10(2):251–273
- Bach N, Huang F, Al-Onaizan Y (2011) Goodness: A method for measuring machine translation confidence. In: Proceedings of the Annual Meeting of the Association for Computational Linguistics, pp 211–219
- Barrachina S, Bender O, Casacuberta F, Civera J, Cubel E, Khadivi S, Lagarda AL, Ney H, Tomás J, Vidal E, Vilar JM (2009) Statistical approaches to computer-assisted translation. *Computational Linguistics* 35(1):3–28
- Brown PF, Della Pietra SA, Della Pietra VJ, Mercer RL (1993) The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics* 19(2):263–311
- Callison-Burch C, Koehn P, Monz C, Post M, Soricut R, Specia L (2012) Findings of the 2012 workshop on statistical machine translation. In: Proceedings of the Seventh Workshop on Statistical Machine Translation, pp 10–51
- Carl M (2012a) The CRITT TPR-DB 1.0: A database for empirical human translation process research. In: Proceedings of the AMTA 2012 Workshop on Post-Editing Technology and Practice, pp 1–10
- Carl M (2012b) Translog-II: a program for recording user activity data for empirical reading and writing research. In: Proceedings of the Eighth International Conference on Language Resources and Evaluation, pp 4108–4112
- Carl M (2014) Produkt- und Prozesseinheiten in der CRITT Translation Process Research Database. In: Ahrens B (ed) *Translationswissenschaftliches Kolloquium III: Beiträge zur Übersetzungs- und Dolmetschwissenschaft* (Köln/Germersheim), Peter Lang, Frankfurt am Main, pp 247–266
- Carl M, Kay M (2011) Gazing and typing activities during translation : A comparative study of translation units of professional and student translators. *Meta* 56(4):952–975
- Doherty S, O’Brien S, Carl M (2010) Eye tracking as an MT evaluation technique. *Machine Translation* 24(1):1–13
- Elming J, Carl M, Balling LW (2014) Investigating user behaviour in post-editing and translation using the Casmacat workbench. In: O’Brien S, Winther Balling L, Carl M, Simard M, Specia L (eds) *Post-editing of machine translation: Processes and applications*, Cambridge Scholar Publishing, Newcastle upon Tyne, pp 147–169
- Federico M, Cattelan A, Trombetti M (2012) Measuring user productivity in machine translation enhanced computer assisted translation. In: Proceedings of the Tenth Biennial Conference of the Association for Machine Translation in the Americas

- Flournoy R, Duran C (2009) Machine translation and document localization at adobe: From pilot to production. In: Proceedings of MT Summit XII
- Green S, Heer J, Manning CD (2013) The efficacy of human post-editing for language translation. In: Proceedings of SIGCHI Conference on Human Factors in Computing Systems, pp 439–448
- Guerberof A (2009) Productivity and quality in mt post-editing. In: Proceedings of MT Summit XII-Workshop: Beyond Translation Memories: New Tools for Translators MT
- Guerberof A (2012) Productivity and quality in the post-editing of outputs from translation memories and machine translation. Ph.D. Thesis
- Just MA, Carpenter PA (1980) A theory of reading: from eye fixations to comprehension. *Psychological Review* 87(4):329
- Koehn P (2009a) A process study of computer-aided translation. *Machine Translation* 23(4):241–263
- Koehn P (2009b) A web-based interactive computer aided translation tool. In: Proceedings of ACL-IJCNLP 2009 Software Demonstrations, pp 17–20
- Krings HP (2001) Repairing texts: empirical investigations of machine translation post-editing processes, vol 5. Kent State University Press
- Lacruz I, Shreve GM, Angelone E (2012) Average pause ratio as an indicator of cognitive effort in post-editing: a case study. In: Proceedings of the AMTA 2012 Workshop on Post-Editing Technology and Practice, pp 21–30
- Langlais P, Foster G, Lapalme G (2000) Transtype: A computer-aided translation typing system. In: Proceedings of the 2000 NAACL-ANLP Workshop on Embedded Machine Translation Systems, pp 46–51
- Leiva LA, Alabau V, Vidal E (2013) Error-proof, high-performance, and context-aware gestures for interactive text edition. In: Proceedings of the 2013 annual conference extended abstracts on Human factors in computing systems, pp 1227–1232
- Montgomery D (2004) Introduction to Statistical Quality Control. Wiley
- O’Brien S (2009) Eye tracking in translation process research: methodological challenges and solutions, *Copenhagen Studies in Language*, vol 38, Samfundslitteratur, Copenhagen, pp 251–266
- Ortiz-Martínez D, Casacuberta F (2014) The new Thot toolkit for fully automatic and interactive statistical machine translation. In: Proceedings of the 14th Annual Meeting of the European Association for Computational Linguistics: System Demonstrations, pp 45–48
- Plitt M, Masselot F (2010) A productivity test of statistical machine translation post-editing in a typical localisation context. *The Prague Bulletin of Mathematical Linguistics* 93(1):7–16
- Sanchis-Trilles G, Ortiz-Martínez D, Civera J, Casacuberta F, Vidal E, Hoang H (2008) Improving interactive machine translation via mouse actions. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing, pp 485–494
- Simard M, Foster G (2013) Pepr: Post-edit propagation using phrase-based statistical machine translation. In: Proceedings of MT Summit XIV, pp 191–198
- Skadiņš R, Puriņš M, Skadiņa I, Vasiljevs A (2011) Evaluation of SMT in localization to under-resourced inflected language. In: Proceedings of the 15th International Conference of the European Association for Machine Translation, pp 35–40