

A linguistically motivated taxonomy for Machine Translation error analysis

Ângela Costa^{1,2} · Wang Ling^{1,3,4} · Tiago Luís¹ ·
Rui Correia^{1,3,4} · Luísa Coheur^{1,3}

Received: 14 April 2014 / Accepted: 13 June 2015 / Published online: 26 June 2015
© Springer Science+Business Media Dordrecht 2015

Abstract A detailed error analysis is a fundamental step in every natural language processing task, as to be able to diagnose what went wrong will provide cues to decide which research directions are to be followed. In this paper we focus on error analysis in Machine Translation (MT). We significantly extend previous error taxonomies so that translation errors associated with Romance language specificities can be accommodated. Furthermore, based on the proposed taxonomy, we carry out an extensive analysis of the errors generated by four different systems: two mainstream online translation systems Google Translate (Statistical) and Systran (Hybrid Machine Translation), and two in-house MT systems, in three scenarios representing different challenges in the translation from English to European Portuguese. Additionally, we comment on how distinct error types differently impact translation quality.

✉ Ângela Costa
angelampcosta@gmail.com; angela@l2f.inesc-id.pt

Wang Ling
wlin@l2f.inesc-id.pt

Tiago Luís
tiago.luis@l2f.inesc-id.pt

Rui Correia
rui.correia@l2f.inesc-id.pt

Luísa Coheur
luisa.coheur@l2f.inesc-id.pt

- 1 INESC-ID, Rua Alves Redol, 9, Lisbon, Portugal
- 2 CLUNL, Avenida de Berna, 26-C, 1069-61 Lisbon, Portugal
- 3 IST, Universidade de Lisboa, Lisbon, Portugal
- 4 LTI, Carnegie Mellon University, Pittsburgh, USA

Keywords Machine translation · Error taxonomy · Error analysis · Romance languages

1 Introduction

Error analysis is the process of determining the incidence, cause and consequences of unsuccessful language (James 1998). This linguistic discipline has been applied to many research fields, such as *Foreign Language Acquisition* and *Second Language Learning and Teaching* (Corder 1967), since errors contain valuable information on the strategies that people use to acquire a language (Dulay et al. 1982) and, at the same time, allow the identification of points that need further work. In fact, according to Richards (1974), “At the level of pragmatic classroom experience, error analysis will continue to provide one means by which the teacher assesses learning and teaching and determines priorities for future effort”. More recently, error analysis has also become a focus of research in the machine translation (MT) area, where some work is dedicated to the design of taxonomies (Llitjós et al. 2005; Vilar et al. 2006; Bojar 2011) and others target error identification (Popović and Ney 2006). In this paper, we present a linguistically motivated taxonomy for translation errors that extends previous ones. Contrary to other approaches, our proposal:

- clusters different types of errors in the main areas of linguistics, allowing the precise specification of the information level needed to identify the errors and simplify possible extensions;
- allows the classification of errors that occur in Romance languages (usually ignored in previous taxonomies) and not only English;
- allows the consideration of a language’s variations.

Moreover, based on this taxonomy we perform a detailed linguistic analysis of the errors produced in the translation of English (EN) into European Portuguese (EP) texts by two mainstream online translation systems, Google Translate and Systran, and two in-house MT systems, both based on Moses technology (Koehn et al. 2007).

Google Translate¹ is the best-known translation engine, allowing automatic translation of texts in many languages; Systran² is a free online hybrid MT engine that combines rule-based and statistical MT (SMT). Moses³ is a publicly available SMT system, intensively used by MT researchers all over the world. It allows automatic training of translation models for any language pair, as long as a collection of parallel data is available, such as Europarl (Koehn 2005).

Due to the fact that SMT systems are highly dependent on the training data and thus behave differently in distinct domains, we have chosen parallel corpora with different characteristics. Therefore, we perform our experiments on a corpus that contains the one described in Costa et al. (2014), and is composed of:

¹ <http://translate.google.com>.

² <http://www.systranet.com/translate>.

³ <http://www.statmt.org/moses>.

- speech transcriptions (and respective translations of the subtitles into EP) of TED-talks;⁴
- touristic texts from the bilingual UP magazine;⁵
- TREC evaluation questions (Li and Roth 2002), translated into EP in a previous work by Costa et al. (2012).⁶

In this way, we are able to study and cover errors resulting, respectively, from speech translations, from translations within a restricted domain and also from translations over specific constructions, which is the case of questions. Moreover, the EP translations of the corpora were used to automatically evaluate the translation performed by all systems.

This paper is organised as follows: in Sect. 2 we present related work, in Sect. 3 we detail the error taxonomy, and, in Sect. 4, we describe the corpora, the tools used in our experiments and the annotation process. In Sect. 5 we analyze the errors resulting from the translations and in Sect. 6 we discuss error gravity. Finally, in Sect. 7, we highlight the main conclusions and point to future work.

2 Related work

Several studies have been developed with the goal of classifying translation errors in MT. In addition, several works focus on the identification of machine or human translation errors. Some research targets semi- or fully automatic error analysis methods, while others manually analyze these errors. In this section we survey the most significant work on these subjects.

Starting with the problem of error classification, different taxonomies have been suggested.⁷ One of the most referenced taxonomies in MT is the hierarchical classification proposed by Vilar et al. (2006). They extend the work of Llitjós et al. (2005), and split errors into five classes: “Missing Words”, “Word Order”, “Incorrect Words”, “Unknown Words” and “Punctuation Errors”. A “Missing Words” error is produced when some words in the generated sentence are missing. “Word Order” errors concern the word order of the generated sentence. This problem is solved by moving words or blocks of words within the sentence. “Incorrect Words” are found when the system is unable to find the correct translation of a given word. “Unknown Words” are “translated” simply by copying the input word to the generated sentence, without further processing. Finally, “Punctuation Errors” represent only minor disturbances, but are also considered in this taxonomy.

Inspired by the work of Vilar et al. (2006), Bojar (2011) used a similar classification that divides errors into four types: “Bad Punctuation”, “Missing Word”, “Word

⁴ <http://www.ted.com/>.

⁵ <http://upmagazine-tap.com/>.

⁶ <http://metanet4u.l2f.inesc-id.pt/>.

⁷ While not presenting a taxonomy per se, the work of Naskar et al. (2011) is interesting to note as it presents a tool for diagnostic evaluation of translation errors, DELiC4MT, available at <http://www.computing.dcu.ie/~atoral/delic4mt>, focusing mainly on linguistic checkpoints for different part-of-speech categories.

Order” and “Incorrect Words”. Basically he uses Vilar’s taxonomy, but eliminates the “Unknown Words” category.

The classification of errors by Elliott et al. (2004) was progressively developed during the analysis and manual annotation of approximately 20,000 words of MT output, translated from French into English by four systems (Systran, Reverso Promt,⁸ Compendium⁹ and SDL’s online Free Translation¹⁰). This taxonomy is slightly different, as the annotations were made according to items that a post-editor would need to amend if he/she was revising the texts to publishable quality. Error types were divided according to parts-of-speech and then sub-divided as “Inappropriate”, “Untranslated”, “Incorrect”, “Unnecessary” and “Omitted”.

At this point it is important to say that all taxonomies are influenced by the idiosyncrasies of the languages with which they are working. For instance, the work carried out by Vilar et al. (2006) concerns experiments with the language pair English–Chinese and, thus takes into consideration error types that are not relevant for European languages. For instance, working particularly with Chinese, they felt the need to introduce new types of reordering errors, as the position of the modifier changes according to the sentence construction (declarative, interrogative, subordinate/infinitival sentences). In our particular case, as we are working with translations from EN into EP, our main concern was to develop a taxonomy that captures all idiosyncrasies of Portuguese but that also works for Romance languages.

Although the purpose of this work is to classify MT errors, for the creation of our error taxonomy, we think it is also important to consider error classification studies for human errors. Regarding human translation errors, Dulay et al. (1982) suggest two major descriptive error taxonomies: the linguistic category classification (LCC) and the surface structure taxonomy (SST). LCC is based on linguistic categories (general ones, such as morphology, lexis, and grammar as well as more specific ones, such as auxiliaries, passives, and prepositions). SST focuses on the way surface structures have been altered by learners (e.g. omission, addition, misformation, and misordering). These two approaches are presented as alternative taxonomies. However, according to James (1998), there is great benefit in combining them into a single bidimensional taxonomy.

In connection with human errors, but this time errors produced by humans in a translation task, we should also mention the Multilingual eLearning in Language Engineering project¹¹ (MeLLANGE). They produced the MeLLANGE Learner Translator Corpus (Castagnoli et al. 2007) that includes work done by trainees, which was subsequently annotated for errors according to a customised error typology. While we have not done any experiments with human translation errors to date, we plan to do this in future work. However, we did keep these error types in mind when designing our taxonomy.

⁸ <http://reverso.softissimo.com/en/reverso-promt-pro>.

⁹ <http://amedida.ibit.org/compendium.php>.

¹⁰ <http://www.freetranslation.com>.

¹¹ http://corpus.leeds.ac.uk/mellange/about_mellange.html.

Considering the identification of MT errors, several automatic measures are proposed in the literature. Among these, two of the most widely used scores in SMT are BLEU (Papineni et al. 2002) and METEOR (Denkowski and Lavie 2014). BLEU scores are calculated by comparing translated segments with reference translations. Those scores are then averaged over the whole corpus to reach an estimate of the translation's overall quality. BLEU simply calculates n -gram precision without explicitly taking into account intelligibility or grammatical correctness. METEOR (Denkowski and Lavie 2014) is an automatic metric for MT evaluation that is based on a generalized concept of unigram matching between the machine-produced translation and human reference translations. It also uses other linguistic resources such as paraphrases and generally obtains better results, which is why we have chosen to use it in our automatic evaluation. However, even though automatic evaluation methods are very popular as they are quicker and less expensive than a manual evaluation, human judgements of translation performance are still more accurate. We should also add that the interpretation of these measures is not easy and they do not permit a clear identification of the engines' problems. For instance, a BLEU score of 0.20 does not allow us to precisely capture the type of error being produced by the MT system. Therefore, besides the automatic evaluation of translations, some semi-automatic error analysis has also been done. In the works described in Popović and Ney (2006) and Popović et al. (2006), errors in an English–Spanish SMT system were analysed with respect to their morphological and syntactic origin, and revealed problems in specific areas of inflectional morphology and syntactic reordering (Kirchhoff et al. 2007). A graphical user interface that automatically calculates various error measures for translation candidates and thus facilitates manual error analysis is presented in Niessen et al. (2000).

As far as manual error identification is concerned, Bojar (2011) carried out a manual evaluation of four systems: Google, PC Translator,¹² TectoMT¹³ and CU-Bojar (Bojar et al. 2009). In his work, Bojar used two techniques of manual evaluation to identify error types of these MT systems. The first technique is called “blind post-editing” and consists of an evaluation performed by two people, separately. The first annotator receives the system output and has to correct it producing an edited version; meanwhile the second annotator obtains the edited version, the source and the reference translation, and judges whether the translation is still acceptable. The second technique used was the manual annotation of the errors using a taxonomy inspired by Vilar et al. (2006).

A similar work is presented in Condon et al. (2010), but with translations to and from English to Iraqi Arabic. Errors were annotated as “Deletions”, “Insertions” or “Substitutions” for morphological classes and after they were assigned a type of error following a similar taxonomy to that proposed by Vilar et al. (2006).

Furthermore, in Fishel et al. (2012) a collection of annotated translation errors is presented, consisting of automatically produced translations and their detailed manual analysis.¹⁴ Using the collected corpora, the authors evaluated two available state-of-

¹² <http://langsoft.cz/translatorA.html>.

¹³ <http://ufal.mff.cuni.cz/tectomt>.

¹⁴ <http://terra.cl.uzh.ch/terra-corpus-collection.html>.

the-art methods of MT diagnostics and assessment: Addicter (Zeman et al. 2011)¹⁵ and Hjerson (Popović and Ney 2011).¹⁶ Addicter is an open-source tool that uses a method explicitly based on aligning the hypothesis and reference translations to devise the various error types. Hjerson decomposes the WER and PER metrics over the words in the translations with the same aim.

The Framework for Machine Translation Evaluation (FEMTI)¹⁷ is a tool created to help people that evaluate MT systems. FEMTI has two classifications incorporated: the first one consists of characteristics of the contexts where the MT systems can be applied. The second one lists the MT systems' characteristics, as well as the metrics proposed to measure them. People that use this framework have to specify the intended context for the MT system in the first classification and submit. In return, FEMTI proposes a set of characteristics that are important in that particular context, using its embedded knowledge base. All the characteristics and evaluation metrics can be changed. After this task is completed, the evaluators can print the evaluation plan and do the evaluation.

Concluding this section, we should mention the work described in Secară (2005), which presents a survey on state-of-the-art translation evaluation methods, but on a much more linguistically oriented approach. Here the focus of most of the analysed frameworks is on annotation schemes and error weighting for assessing the quality of a translated text, and on including post-editing feedback from human experts in error reduction and translation improvement.

3 Taxonomy

Error identification is not always a straightforward task. Not all errors are easy to find: some are diffused throughout the sentence or larger units of text that contains them (James 1998). Underlying the identification issue remains the problem of their classification. Our taxonomy classifies errors in terms of “the linguistic item which is affected by the error” (Dulay et al. 1982). Thus, relatively coarse categories—**Orthography**, **Lexis**, **Grammar**, **Semantic**, and **Discourse**—indicate the language level where the error is located. In the following sections we explain each one of these categories and specify the subcategories of the linguistic units where the error occurs. This description is illustrated with errors resulting from EN to EP translations. As usual, each error is identified with an asterisk, which is placed before the wrong expression.

3.1 Orthography level

Orthography level errors include all the errors concerning misuse of punctuation and misspelling of words. We divide orthography level errors into three types:

¹⁵ <https://wiki.ufal.ms.mff.cuni.cz/user:zeman:addicter>.

¹⁶ <http://www.dfki.de/~mapo02/hjerson>.

¹⁷ <http://www.issco.unige.ch:8080/cocoon/femti/st-home.html>.

punctuation, capitalization and spelling. Each incorrect use of punctuation represents a punctuation error. In the following example, a comma is erroneously inserted.

Example: Punctuation error

EN: *green tea*

EP: *chá*, verde*

Correct translation: *chá verde*

A capitalization error occurs when there is an inappropriate use of capital letters (for instance, the use of a small caption in the first letter of a proper noun). In the following example, the English sentence is correct, as the pronoun *I* is always spelt with a capital letter. Meanwhile, the Portuguese sentence does not have a subject (it is not expressed, but was previously mentioned) and probably because of this the verb was spelt with a capital letter.

Example: Capitalization error

EN: *... on time, I can console myself...*

EP: *... a tempo, *Posso consolar-me...*

Correct translation: *... a tempo, posso consolar-me...*

Finally, a spelling mistake concerns the substitution, addition or deletion of one or more letters (or an accent) to the orthography of a word.

Example: Spelling error

EN: *Basilica of the Martyrs*

EP: *Basílica dos *Mátires*

Correct translation: *Basílica dos Mártires*

Although a capitalization error could be considered a spelling mistake, we opted to provide both categories, and define them at the same abstraction level. This is due to the fact that if a capitalization error is common in natural language processing tasks, such as Automatic Speech Recognition and MT, a spelling mistake is not, as usually systems are trained with texts that do not have many spelling errors (news, parliament sessions, etc.). On the other hand, if we consider a human translation, spelling mistakes tend to be frequent, but capitalization errors are rare.¹⁸ For this reason, we have decided to keep both type of errors in the taxonomy, so that both human and MT errors can be covered.

3.2 Lexis level

Under this category we have considered all errors affecting lexical items. It should be clear that, contrary to spelling errors that respect the characters used within a word, lexis errors concern the way each word, as a whole, is translated. Thus, the following

¹⁸ Although in some languages they can be more frequent, as for instance in German, where all nouns are spelled with capital letter. This can be a problem for foreign students that do not have this particularity in their mother tongue.

types of errors at the *lexis* level are taken into account: omission, addition and untranslated. Moreover, omission and addition errors are then analysed considering the type of words they affect: (a) **content words** (or lexical words), i.e., words that carry the content or the meaning of a sentence (such as nouns (*John, room*) or adjectives (*happy, new*)); (b) **function words** (or grammatical words), that is, words that have little lexical meaning, but instead serve to express grammatical relationships with other words within a sentence (e.g. prepositions (*of, at*) and pronouns (*he, it, anybody*)).

Omission errors happen when the translation of a word present in the source text is missing in the resulting translation; an addition error represents the opposite phenomenon: the translation of a word that was not present in the source text and was added to the target text.

Example: Omission error (content word)

EN: *In his inaugural address, Barack Obama*

EP: *No seu * inaugural, Barack Obama*

Correct translation: *No seu discurso inaugural, Barack Obama*

Example: Omission error (function word)

EN: *In India*

EP: *Em Índia*

Correct translation: *Na Índia* (*Na* is the contraction of the proposition *em* and the article *a*, which was missing in the translation)

In the first example, the word *address*, a content word, was not translated and so it was missing from the sentence in Portuguese. In the second example, the missing word is a pronoun (function word). The country *India* in Portuguese is always preceded by a definite article that, in this case, is missing from the translation output.

Example: Addition error (content word)

EN: *This time I'm not going to miss*

EP: *Desta vez *correr não vou perder*

Correct translation: *Desta vez não vou perder*

Example: Addition error (function word)

EN: *highlights the work*

EP: **Já destaca-se o trabalho*

Correct translation: *destaca-se o trabalho*

These last two examples concern the addition of words to the translation output. In the first sentence, the translation engine added the verb *correr* (run) to the Portuguese sentence. Literally, the sentence was translated as *This time run I'm not going to miss*. In the second example, the added word was a function word, the adverb *já* (already). In this case, the sentence was roughly translated as *Already highlights the work*. This is not a grammatically or semantically wrong sentence; the only problem is that the word *already* was not in the source text.

Besides omitting or adding words to the translation, one other situation can occur, namely when a word is not translated (untranslated). An untranslated error

is very common in MT, because when the engine cannot find any translation candidate for a given source word, an option often chosen is to copy it to the translation output ‘as is’. This often results in a successful ‘translation’, e.g. where proper nouns are not to be translated, or between languages with many cognates. An opposite example is the following, where the MT system had no translation for the word *botany*, so this word was simply copied intact to the output sentence.

Example: Untranslated errors

EN: *in the world of botany*

EP: *no mundo da *botany*

Correct translation: *no mundo da botânica*

3.3 Grammar level

Grammar level errors are deviations in the morphological and syntactical aspects of language. On this level of analysis we identified two types of errors: misselection errors and misordering errors.

Misselection are morphological errors that the words may present. This is the case of problems at word class-level (for instance, an adjective is needed, but the translation engine returns a noun, instead), and at verbal level (tense and person). Errors of agreement (gender, number, person), and in contractions (between prepositions and articles) also fall into this category. When we have more than one of these problems in the same word we called it a ‘blend’.

Example: Misselection error (word class)

EN: *world*

EP: *mundial (worldwide)*

Correct translation: *mundo*

Example: Misselection error (verb level: tense)

EN: *Even though this is a long list*

EP: *Mesmo que esta *é uma longa lista (é (to be) should be in the subjunctive and not in the indicative mood)*

Correct translation: *Mesmo que esta seja uma longa lista*

Example: Misselection error (verb level: person)

EN: *Theater-goers can discover*

EP: *As pessoas que vão ao teatro *pode descobrir (the correct form of the verb is in the the third person plural and not in the third person singular)*

Correct translation: *As pessoas que vão ao teatro podem descobrir*

Example: Misselection error (verb level: blend)

EN: *If I go to see the Dario Fo play*

EP: *Se *vai ver a peça de Dário Fo. (ir (to go) should be in the conditional and not in the indicative mood and in the first person singular not the third person.)*

Correct translation: *Se eu for ver a peça de Dário Fo.*

Example: Misselection error (agreement: gender)

EN: *The German artist Thomas Schutte*

EP: **(A artista alemã) Thomas Schutte* (in Portuguese, like all morphologically rich languages, the pronoun and adjective have to agree in gender and number with the noun, in this case masculine and singular, not feminine)

Correct translation: *O artista alemão Thomas Schutte*

Example: Misselection error (agreement: number)

EN: *moral skills*

EP: *capacidades *moral* (both the adjective and noun have to be in the plural form)

Correct translation: *capacidades morais*

Example: Misselection error (agreement: person)

EN: *learn from our failures*

EP: *aprender com os *vossos fracassos* (the use of the possessive pronoun *vossos* is grammatically correct, but it is not in the correct person)

Correct translation: *aprender com os nossos fracassos*

Example: Misselection error (agreement: blend)

EN: *funky clothes shops*

EP: *lojas *simpático de roupa* (the adjective is not in the correct gender and number)

Correct translation: *lojas simpáticas de roupa*

Contraction problems are typical of Romance languages, such as Portuguese, as many prepositions (like *em*) have to be contracted when adjacent to an article, e.g. *na* results from the contraction of the preposition *em* (*in*) and the article *a* (*a*).¹⁹

Example: Misselection error (contraction)

EN: *in an environment*

EP: *em um ambiente*

Correct translation: *num ambiente*

Finally, Misordering errors are related with syntactic problems that the sentences may demonstrate. We should point out that a good translation does not only involve selecting the right forms to use in the right context, but also arranging them in the right order. In Portuguese, certain word classes such as adverbials and adjectives seem to be especially sensitive to misordering.

Example: Misordering error

EN: *A person is wise.*

EP: *Uma pessoa sábia *é.* (*A person wise *is.*)

Correct translation: *Uma pessoa é sábia.*

¹⁹ We should not confuse Omission errors of a function word with Misselection errors (contraction). In the first case, in the phrase *na (em + a) Índia (in India)*, the article *a* is missing, so we have an Omission error. Meanwhile, if we had a contraction problem, the sentence would be *em a Índia*, where both preposition and article were correctly selected but were not contracted as they should be (*em + a = na*).

3.4 Semantic level

Semantic errors are problems related to the meaning of the words and subsequent wrong word selection. We have individuated three different types of errors: confusion of senses, wrong choice, collocational error and idioms.

Confusion of senses is the case of a word that was translated into something representing one of its possible meanings, but, in the given context, the chosen translation is not correct.

Example: Confusion of senses errors

EN: *the authentic tea set that includes a tray, teapot and glasses* (glasses means ‘spectacles’, but it is also the plural of the noun *glass*)

EP: *um autêntico jogo de chá que inclui bandeja, bule e *óculos* (the authentic tea set that includes a tray, teapot and spectacles)

Correct translation: *um autêntico jogo de chá que inclui bandeja, bule e copos*

As far as wrong choice errors are concerned, they occur when a wrong word, without any apparent relation, is used to translate a given source word.

Example: Wrong choice errors

EN: *in the same quarter*

EP: *no mesmo *histórica* (in the same *historical)

Correct translation: *no mesmo bairro*

We should not confuse Wrong Choice with Confusion of senses. An example of the first case is the translation of *care* as *conta* (check), where there is no semantic relation between these two words. In contrast, the translation of *glasses* as *óculos* (glasses) is a predictable Confusion of senses, as the English word *glasses* can be translated into two different words in Portuguese: ‘glasses to drink’ (*copos*) and ‘glasses to see with’ (*óculos*).

Collocations are the other words any particular word normally keeps company with (James 1998). They have a compositional meaning, contrary to idioms, but the selection of their constituents is not semantically motivated. Collocational errors could be considered an instantiation of the previous error, but we have decided to take them into consideration separately. This decision was made because in the case of Confusion of senses errors, we account for single word misuse; in contrast, Collocational errors occur in blocks of words.

Example: Collocational errors

EN: *high wind*

EP: *vento alto* (literally means *tall wind*)

Correct translation: *vento forte*

Idiomatic errors concern errors in idiomatic expressions that the system does not know and translates as regular text. These expressions cannot be literally translated as their meaning it is not literal and, in many cases, the equivalent expression in the target language is very different.

Example: Idioms

EN: *kick the bucket*

EP: *dar um pontapé ao balde*

Correct translation: *esticou o pernil* (idiomatic expression that means *to die*)

3.5 Discourse level

Discourse-level errors are discursive options that are not the most expected. We consider three different situations at the Discourse level: style, variety and should not be translated.

Style errors concern a bad stylistic choice of words when translating a sentence. A typical example is the repetition of a word in a nearby context, where a synonym should have been a better option.

Example: Style errors

EN: *permission to be allowed to improvise*

EP: *autorização para ser autorizado a improvisar* (*permission to be permitted to*)

Correct translation: *permissão para ser autorizado a improvisar*

Variety errors cover the cases when the target of the translation is a certain language, but instead lexical or grammatical structures from a variety of that language are used. This is what happens, for instance, when the target of a translation is EP and Brazilian Portuguese (BP) is returned, which is very common in Google translations. With Variety errors, this taxonomy is then able to capture this phenomenon.

Example: Variety errors

EN: *in his speech*

EP: *em seu discurso*

Correct translation: *no seu discurso* (in EP, we need an article before the possessive pronoun (*seu*) that, in this case, is contracted with the preposition *em* (*em + o = no*))

Under the should not be translated category, we have considered all the word's sequences in the source language that should not be translated in the target language. In this particular case, we can find, for instance, books or film titles. In this example, we have the name of a Portuguese play originally Portuguese in the English text, but the engine tried to translate it and only succeeded in adding errors.

Example: Should not be translated errors

EN: *Havia um Menino que era Pessoa*

EP: *Havia hum Menino era Opaco Pessoa*

Correct translation: *Havia um Menino que era Pessoa*

3.6 General Overview

In Fig. 1 we resume the taxonomy previously presented.

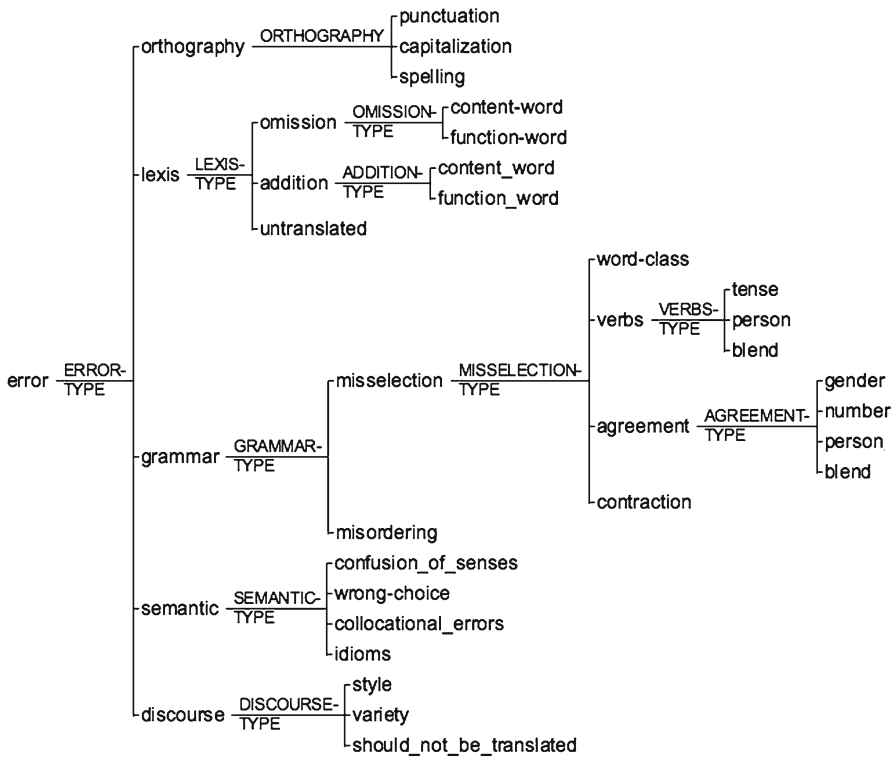


Fig. 1 Taxonomy

3.7 Comparison with other taxonomies

In this section we compare our taxonomy with the ones created for MT and described in Sect. 2, namely the ones presented in Bojar (2011) and Vilar et al. (2006), as well as with the one traditionally used for classifying students’ errors, also described in Sect. 2 and detailed in Dulay et al. (1982). From now on, we will call these taxonomies **Bojar**, **Vilar** and **Dulay**, respectively. Table 1 summarizes the following discussion.

As far as orthography level errors are concerned, **Bojar** and **Vilar** taxonomies only consider punctuation level errors. In contrast, **Dulay**’s taxonomy also considers spelling errors. Our taxonomy considers capitalization errors in addition to these types.

Concerning lexis errors, all the taxonomies agree on the proposed different types of errors, but there are significant differences in how grammar errors are captured. **Vilar** does not consider word class errors, mentioned by all the other taxonomies, and, on this level, it only classifies gender and number, while we found it important to add other types of agreement, as they are significant for Portuguese and Romance languages in general; **Bojar** does not talk about agreement problems, but **Dulay** does take this type of error into consideration. Contraction errors are not mentioned by these three taxonomies. We hypothesize that this happens because most of the taxonomies are built for English and this phenomenon is not compulsory in

Table 1 Comparison with other taxonomies

Error types	Bojar	Vilar	Dulay
Orthography			
PUNCTUATION	✓	✓	✓
CAPITALIZATION	✗	✗	✗
SPELLING	✗	✗	✓
Lexis			
OMISSION	✓	✓	✓
ADDITION	✓	✓	✓
UNTRANSLATED	✓	✓	✓
Grammar			
WORD CLASS	✓	✗	✓
VERBS	✓	✓	✓
AGREEMENT	✗	✓	✓
CONTRACTION	✗	✗	✓
MISORDERING	✓	✓	✓
Semantic			
CONFUSION OF SENSES	✓	✓	✓
WRONG CHOICE	✓	✓	✓
COLLOCATIONAL ERRORS	✗	✗	✗
IDIOMS	✗	✓	✓
Discourse			
STYLE	✗	✓	✓
VARIETY	✗	✗	✗
SHOULD NOT BE TRANSLATED	✗	✗	✗

English, as contractions are an informal option in language use. However, in several cases contractions are obligatory in Portuguese and languages like Italian, German, Spanish or French, so for this reason we consider that contraction errors must be included in an error taxonomy.

Regarding semantic errors, **Vilar** does not mention collocational errors and **Bojar** does not include idiomatic or collocational errors. **Dulay** talks about all these types of semantic errors. All of the previously mentioned authors report the existence of wrong choice of word errors. At the discourse level, style errors are mentioned by **Vilar et al. (2006)**, but not by **Bojar (2011)**. **Dulay** only takes into consideration errors of style, as **Dulay's** taxonomy was conceived to assess human errors, but not all the categories used in MT have a direct equivalent in a student's typology of mistakes. Variety errors are assessed only by our taxonomy. As previously mentioned, this type of error is very frequent in Google's translations of EP. However, the same thing happens between American English and British English, although this type of error is not considered in taxonomies built for English.

3.8 Language-dependent and -independent errors

Having discussed our proposed taxonomy, it is important to differentiate between errors that could happen in any MT task and are independent of the source and target

Table 2 Language-independent errors

Error types	Language- independent
Orthography	
PUNCTUATION	✓
CAPITALIZATION	✗
SPELLING	✓
Lexis	
OMISSION	✓
ADDITION	✓
UNTRANSLATED	✓
Grammar	
WORD CLASS	✗
VERBS	✗
AGREEMENT	✗
CONTRACTION	✗
MISORDERING	✗
Semantic	
CONFUSION OF SENSES	✓
WRONG CHOICE	✓
COLLOCATIONAL ERRORS	✓
IDIOMS	✓
Discourse	
STYLE	✓
VARIETY	✗
SHOULD NOT BE TRANSLATED	✓

language, and those errors that are language-dependent and that may not occur in every language. In Table 2 we present a resumé of the following discussion.

We cannot be sure that every existing language has punctuation, but we know that a great majority of them have. We should point out that the punctuation symbols may be different. In Greek, the question mark is written as the English semicolon and in Spanish an inverted question mark is used at the beginning of a question and the normal question mark is used at the end.

Capitalization errors should be considered language-specific, as languages like Chinese, Arabic or Korean do not have capital letters.

Considering spelling mistakes, every written language has a standard orthography, and any misuse of these rules may be marked as a spelling mistake.

As far as lexis errors are concerned, omissions, additions and untranslated words can be present in any automatic translation, independent of the language of focus. Regarding grammatical problems, this category of mistakes is language-specific and not all of the grammatical error types make sense in every language. According to Keenan and Stabler (2010), “different languages do have non-trivially different grammars: their grammatical categories are defined internal to the language and may fail to be comparable to ones used for other languages. Their rules, ways

of building complex expressions from simpler ones, may also fail to be isomorphic across languages.”

As [De Saussure \(1916\)](#) defended, language is ambiguous and polysemous by definition. By this we mean that in every language there are semantic problems that can arise in an automatic translation. For instance, any idiomatic expression has a non-literal meaning that cannot always be captured by a literal translation. Information about the context of use is necessary for it to be well interpreted and translated.

Concerning discourse errors, style errors may occur in every language, as different social contexts require an appropriate discourse. As far as variety is concerned, not all languages have a variety, like American English and British English or European Portuguese and Brazilian Portuguese, so this category should only be used for these cases.

Finally, words that should not be translated and that should be kept in the language of the source language is a language-independent problem.

4 Experimental setup

In this section we briefly describe the corpora and the tools we have used in our experiments. We also present the annotation agreement resulting from the annotation of the translation errors (according with the proposed taxonomy) in each corpus.

4.1 Corpora

As previously stated, the error analysis was carried out on a corpus of 750 sentence pairs, composed of:

- 250 pairs of sentences taken from TED talks (from now on, the TED corpus);
- 250 pairs of sentences taken from the UP magazine from TAP (henceforth, the TAP corpus);
- 250 pairs of questions taken from a corpus made available by [Li and Roth \(2002\)](#), from the TREC collection (from now on, the Questions corpus).

The TED corpus is composed of TED talk subtitles and corresponding EP translations. These were created by volunteers and are available on the TED website. As we are dealing with subtitles (and not transcriptions), content is aligned to fit the screen, and, thus some manual pre-processing was needed to connect segments in order to obtain parallel sentences.

The TAP corpus comprises 51 editions of the bilingual Portuguese national airline company magazine, divided into 2,100 files for EN and EP. It has almost 32,000 aligned sentences and a total of 724,000 Portuguese words and 730,000 English words.

The parallel corpus of Questions (EP and EN) consists of two sets of nearly 5,500 plus 500 questions each, to be used as training/test data, respectively. Details of the SMT experiments on questions can be found in [Costa et al. \(2012\)](#).

Some examples of sentences from these corpora can be found in [Table 3](#).

The Questions corpus contains short sentences and most of them start with an interrogative pronoun. The TAP corpus presents more complex grammatical structures

Table 3 Examples of sentences from the corpora

TED	The publisher bears no responsibility for return of unsolicited material and reserves the right to accept or reject any editorial and advertising material. No parts of the magazine may be reproduced without the written permission of up. The opinions expressed in this magazine are those of the authors and not necessarily those of the auditor.
TAP	They're the things you would expect: mop the floors, sweep them, empty the trash, restock the cabinets. It may be a little surprising how many things there are, but it's not surprising what they are.
Questions	Who developed the vaccination against polio? What is epilepsy? What year did the Titanic sink? Who was the first American to walk in space?

Table 4 Data used in the error analysis

Dataset	Language	Sentences	Tokens	Average sentence length
TAP	EN	250	4868	19.47
	EP	250	5521	22.08
TED	EN	250	3346	13.38
	EP	250	3894	15.58
Questions	EN	250	1856	7.42
	EP	250	2048	8.19

when compared with the TED corpus, which is influenced by its semi-spontaneous nature (some previous preparation is involved). This difference may be observed because written language tends to be more complex and intricate than speech, with longer sentences and many subordinate clauses. Spoken language tends to be full of repetitions, incomplete sentences, corrections and interruptions, which sometimes result in ungrammatical sentences.

Table 4 provides details on the number of sentences, tokens and average number of tokens per sentence that were translated. By token we understand a string of characters delimited by a white space. Therefore, not only words, but also punctuation, are tokens.

Finally, in Table 5 we can see the number of tokens per translated dataset, and the average number of tokens per sentence for each dataset.

4.2 Systems and tools

4.2.1 Machine translation systems

We tested four different systems in our evaluation: two mainstream online translation systems (Google Translate (statistical) and Systran (hybrid)), and two in-house MT systems. The online systems were run as they were²⁰ and we will denote them

²⁰ Translated on 22/10/2014.

Table 5 Number of tokens in the translated corpora on the first line, and on the second line, the average number of tokens per sentence for each dataset

Dataset	TAP	TED	Questions
Online-S	5725	3855	1955
	22.90	15.42	7.82
Online-G	5623	3956	2030
	22.49	15.82	8.12
Moses-PB	5522	3730	2068
	22.09	14.92	8.27
Moses-HPB	5507	3759	2059
	22.03	15.02	8.24

Table 6 Data used for training, tuning and testing the MT models

Dataset	Train sentences	Tuning sentences	Test sentences
TAP	4409	1000	250
TED	78,135	1000	250
Question	8,914	1000	250
Europarl	1,960,407	0	0

as Online-G and Online-S, respectively. The in-house systems were trained using Moses: the phrase-based model (Koehn et al. 2007), and the hierarchical phrase-based model (Chiang 2007), which we will henceforth denote as Moses-PB and Moses-HPB, respectively. Both systems share the same training corpora, comprised of approximately 2 million sentence pairs from Europarl (Koehn 2005). As for the in-domain corpora, we gathered the remaining sentence pairs for the TAP, Questions and TED domains after removing the held-out data, and added these into the training corpora. These contained 4409, 8904 and 78,135 sentence pairs, respectively. In total, there were 56 million tokens in English (27.32 tokens per sentence) and 58 million tokens (28.28 tokens per sentence) for Portuguese in the training set.²¹

The model was built by first running IBM model 4, with Giza++ (Och and Ney 2003), and bidirectional alignments were combined with the grow-diag-final heuristic, followed by the phrase extraction (Ling et al. 2010) for the Moses-PB model and rule extraction (Chiang 2007) for the Moses-HPB model. The parameters of the model were tuned using MERT (Och 2003) on 1,000 sentence pairs from each of the domains for this purpose. The statistics of the data used are detailed in Table 6. The splits were chosen chronologically in the order the sentences occurred in the dataset. We can see that the majority of the training data is out-of-domain (Europarl), with a relatively small in-domain parallel dataset used. Translations were also detokenized using the Moses detokenizer, and capitalized using the Portuguese capitalizer described in Batista et al. (2007).

²¹ Tokenized using the default Moses tokenizer.

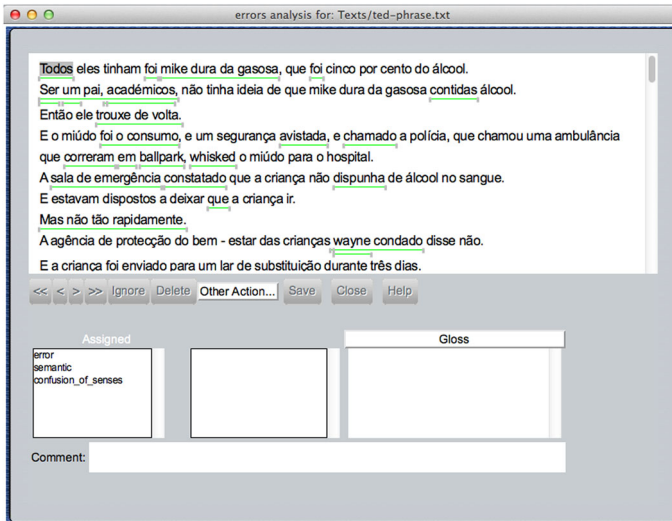


Fig. 2 UAM corpus tool

4.2.2 UAM CorpusTool

Our corpus was annotated using the UAM CorpusTool,²² a state-of-the-art environment for annotation of text corpora (see Fig. 2).

This tool allows the annotation of multiple texts using the annotation schemes previously designed. The annotation is simply done by swiping some text (clicking down and dragging to the end of the segment) and then indicating the features that are appropriate for that segment. This tool also supports a range of statistical analyses of the corpora, allowing comparisons across subsets.

4.3 Annotation agreement

For agreement purposes we compared the answers from two annotators: the linguist that developed the taxonomy and annotated the presented corpora, and another annotator with no formal linguistic instruction. The latter was given an explanation about the different types of errors in the taxonomy, and a set of annotation guidelines. Both annotators identified and classified the errors on a total of 300 sentences: 25 sentences per dataset translated by the four MT engines. Table 7 shows the agreement between the two annotators using Cohen's Kappa.

The agreement is first computed in terms of error location, i.e. whether the annotators agree where the error is placed in the sentence. Errors can be found not only in words or punctuation, but also between words (since words or punctuation marks might be missing). For each word, punctuation mark and space, we measure the agreement on a binary decision regarding the existence or non-existence of an error.

²² <http://www.wagsoft.com/CorpusTool>.

Table 7 Inter-annotator agreement

	Questions	TAP	TED
Localization	0.9717	0.9622	0.9670
L1	0.8295	0.9441	0.9140
L2	0.9216	0.9763	0.9434
L3	0.8223	0.9622	0.9662
L4	1	0.9554	0.9768

Then, we take all cases where annotators agree that there is an error, and check whether they also agree on the classification of the error, considering the first level (L1) of the taxonomy (Orthography, Lexis, Semantic, Grammar and Discourse). We repeat the same process for the second, third and fourth levels, where we gather all cases where the annotators agree on the previous level (it is possible that the annotators agree on a certain level and agree or not on the next one), and compute the agreement coefficient. As Table 7 shows, the agreement on the identification of the errors in the sentences is high for all three data sets, especially for the corpus of questions.

5 Error analysis

In the following section we analyze the errors carried out by Online-S, Online-G, Moses-PB and Moses-HPB, according to the proposed taxonomy and in the different scenarios.

5.1 Preliminary remarks

Before we begin analysing our results, there are two important issues that need to be discussed:

- We will present our results as the number of errors per dataset, but there are cases of words that contain two errors. We have calculated the number of words with two errors in the total of number of errors, and only between 3.05 and 11.18 % of the words have two errors. The only exception is Moses-HPB, which in the Question corpus contains 16.16 % of words with two errors;
- A straightforward comparison of the error types between systems is only possible at the `lexis` level. This is due to the fact that although in some situations a word may have two different error tags (for instance `misordering` and `capitalization`), words that remain untranslated or that were omitted in the target will never have errors at the `grammar` or `semantic` levels. Thus, we cannot use the number of errors to compare systems after the `lexis` level. For instance, consider that the English sentence *He was sick yesterday.* was translated as **Ele doente ontem.* (verb omitted) by system A, and as **Ele está doente ontem* by system B (verb in the wrong form). Then, system A will have a `lexis` type error and system B a `grammar`-level error. However, we cannot say that system

Table 8 Percentage of errors

System	TAP (%)	TED (%)	Questions (%)
Online-S	21	20	25
Online-G	10	14	13
Moses-PB	16	19	20
Moses-HPB	18	21	19

B has more grammatical errors than system A. That is, the number of errors can be used by each system mainly as an indicator of what the system is doing wrong.

5.2 General overview

Table 8 summarizes the percentage of errors by translation scenario relative to the number of tokens per corpus.

From these results we can observe the following:

- For the Online-S system the corpus of Questions was the most problematic document (25 % of errors). Although the syntactic form of questions is usually very simple (for instance, *What is epilepsy?*), there are problems choosing the right interrogative pronoun. For instance, Online-S translates the sentence *What is the population of Nigeria?* into **Que é a população da Nigéria?*, instead of *Qual é a população da Nigéria?*;
- Although the TAP magazine comprises long sentences,²³ this was the corpus that caused the fewest problems for the majority of the systems (10, 16 and 18 %, for Online-G, Moses-PB and Moses-HPB, respectively).

Table 9 summarises the number of errors found for each error type. In the next sections, we will discuss each specific type of error.

5.3 Lexis level errors

According to Table 10, Moses-PB and Moses-HPB performed considerably worse than Online-G and Online-S on the number of untranslated words or expressions. Although untranslated words represent the minority of lexis errors, this clearly shows a direction on Moses developers research: the translation of unknown words (that is, words never seen during training).

Regarding addition and omission errors, the system with more errors of this kind was Online-S (388 and 449, respectively) and the system with fewest addition and omission problems was Online-G with 133 and 220, respectively. We should also mention that the majority of words that were added or elided were function words and not content words.

²³ For instance, the following sentence is in the TAP corpus: *If it's Saturday, there's a play at Teatro da Trindade called Havia um Menino que era Pessoa, where theatre-goers can discover the verses the poet wrote for his nephews and nieces.*

Table 9 System error types

	TAP	TED	Questions	Errors (total)
Orthography				
Online-S	32	13	24	69
Online-G	19	12	3	34
Moses-PB	145	22	51	218
Moses-HPB	148	23	52	233
Lexis				
Online-S	441	223	219	883
Online-G	164	133	83	380
Moses-PB	247	223	136	606
Moses-HPB	315	257	128	700
Grammar				
Online-S	312	227	90	629
Online-G	157	175	72	404
Moses-PB	251	269	129	649
Moses-HPB	288	296	129	713
Semantic				
Online-S	357	288	138	783
Online-G	118	135	81	334
Moses-PB	211	182	93	486
Moses-HPB	225	184	80	489
Discourse				
Online-S	78	39	17	134
Online-G	79	84	23	186
Moses-PB	21	7	9	37
Moses-HPB	19	10	7	36

5.4 Grammar level errors

At the grammar level of errors, we have identified in Table 11 *misselection* and *misordering* errors (recall that *misselection* errors can affect verbs, agreements and contractions, while *misordering* includes word order problems).

As far as *misordering* errors are concerned, the smaller number of errors could be explained by the common subject-verb-object (SVO) structure²⁴ shared by English and Portuguese. However, although from a syntactic point of view English and Portuguese have some aspects in common, there are also differences. For instance, considering the order of the noun phrase (usually a noun, pronoun, or other noun-like word (nominal), which is optionally accompanied by a modifier such as adjectives),

²⁴ SVO is a sentence structure where the subject comes first, the verb second, and the object third, and languages may be classified according to the dominant sequence of these elements. SVO is one of the most common order in world languages.

Table 10 Lexis errors

	TAP	TED	Questions	Errors (total)
Omission				
Online-S	184	118	147	449
Online-G	104	71	45	220
Moses-PB	103	104	63	270
Moses-HPB	147	122	47	316
Addition				
Online-S	236	92	60	388
Online-G	54	53	26	133
Moses-PB	100	84	55	239
Moses-HPB	121	104	60	285
Untranslated				
Online-S	21	13	12	46
Online-G	6	9	12	27
Moses-PB	44	35	18	97
Moses-HPB	47	31	21	99

Table 11 Grammar errors

	TAP	TED	Questions	Errors (total)
Misordering				
Online-S	112	55	51	218
Online-G	36	33	31	100
Moses-PB	66	54	53	173
Moses-HPB	96	97	56	249
Misselection				
Online-S	200	172	39	411
Online-G	121	142	41	304
Moses-PB	185	215	76	476
Moses-HPB	196	199	73	464

in English the correct sequence is Adjective + Noun; meanwhile, in Portuguese, the usual order is the opposite, although in certain contexts the order Adjective + Noun is also possible. This idiosyncratic aspect of language may explain errors, such as **favorito artista* (*favourite artist*), **permanente coleção* (*permanent collection*), **alemã artista* (*german artist*), and **artística carreira* (*artistic career*). Another structure that might influence word order is the position of the direct-object and indirect-object pronouns. In Portuguese, the rule is that these pronouns should be placed after the verb, but there are many exceptions. One example of this was the sentence *para comprá-lo* (to buy it), which was wrongly translated by Online-G. In this case of infinitive construction with the preposition *para* (to), the pronoun should be placed before the verb *para o comprar* (to buy it). Meanwhile, in English, the order should be verb + pronoun (*to buy it*).

Table 12 Misselection errors

	TAP	TED	Questions	Errors (total)
Word class				
Online-S	40	20	12	72
Online-G	6	4	5	15
Moses-PB	20	31	9	52
Moses-HPB	23	29	4	56
Verbs				
Online-S	55	69	16	140
Online-G	38	79	20	137
Moses-PB	49	94	25	168
Moses-HPB	54	75	24	153
Agreement				
Online-S	96	71	7	174
Online-G	67	50	10	127
Moses-PB	107	89	42	238
Moses-HPB	106	91	43	240
Contraction				
Online-S	9	12	4	25
Online-G	10	9	6	25
Moses-PB	8	1	0	9
Moses-HPB	9	4	1	14

Considering misselection errors, Table 12 summarises the associated subtypes of errors.

The most common errors were agreement errors. We should point out that in Portuguese, according to [Cunha and Cintra \(1987\)](#), the general rule for adjectives is that they agree in gender and number with the noun they modify. In English, the agreement between the adjective and the closest noun is restricted to the words *this* and *that* (as well as *these* and *those*), as these are the only that have separate forms for singular and plural. This structural difference between languages explains many translation errors, such as *a wise man* being translated into Portuguese as **um sábio pessoa* (instead of *uma sábia pessoa*) by Moses-PB, and *the exhibition* translated as **o exposição* (instead of *a exposição*) by Online-G. In both cases the adjectives and articles have to agree in number and gender with the noun, but that does not happen, as we have a feminine noun (*pessoa*) with a masculine article (*o*), and a feminine noun (*exposição*) with a masculine article (*o*). Looking at Table 12 in more detail, we can see that there were some problems producing the correct form of the required verb, specially on the TED talks. Portuguese has a variety of tenses, aspects, and moods, as well as constructions with auxiliary verbs that makes it more grammatically complex than English. For instance, in Portuguese the verb *estar* (*to be*) is used with the present participle to indicate the present continuous aspect and the verb *ter* (*to have*) is used with the past participle for the perfect. English has a less complex tense system, and it

Table 13 Semantic errors

	TAP	TED	Questions	Errors (total)
Confusion of senses				
Online-S	293	261	124	678
Online-G	94	118	57	269
Moses-PB	156	144	73	373
Moses-HPB	179	154	67	400
Wrong choice				
Online-S	47	19	13	79
Online-G	17	7	22	46
Moses-PB	45	27	18	90
Moses-HPB	36	18	13	67
Collocational error				
Online-S	11	5	1	17
Online-G	5	5	2	12
Moses-PB	5	8	1	14
Moses-HPB	6	8	0	14
Idioms				
Online-S	5	3	0	8
Online-G	2	5	0	7
Moses-PB	5	3	0	8
Moses-HPB	4	4	0	8

is not a simple task—even for a human translator—to find the correct correspondence between the tenses of both languages.

Finally, we should take into consideration the type of *misselection* that had the fewest occurrences: *contractions*. This result was quite unexpected as this is another aspect of language where there is no congruence between English and Portuguese. In Portuguese, in several cases contractions are compulsory. For instance, the preposition *de* (*of*) can be contracted with an article and become *do* (*of* + masculine singular article), *da* (*of* + feminine singular article), *duns* (*of* + masculine plural article), or *dumas* (*of* + feminine plural article). This language rule explains errors such as *em um* (*in* + masculine singular article) by Online-G and *por a* (*by* + feminine singular article) by Moses-PB (the correct forms are *uns* and *pela*, respectively).

5.5 Semantic and discourse errors

Now taking a closer look at Table 13, we see that the *confusion of senses* error represents the majority of the semantic errors made by all engines.

Many of the *confusion of sense* errors are due to prepositions like *to*, which in Portuguese can have several translations (*para*, *a* or *de*, just to mention a few). The same happens with copular verbs like *be*, which in Portuguese can be translated into

Table 14 Discourse errors

	TAP	TED	Questions	Errors (total)
Style				
Online-S	9	8	0	17
Online-G	2	7	0	9
Moses-PB	4	2	2	8
Moses-HPB	4	6	0	10
Variety				
Online-S	31	23	6	60
Online-G	56	73	22	151
Moses-PB	0	0	0	0
Moses-HPB	1	0	0	1
Should not be translated				
Online-S	28	8	11	57
Online-G	20	3	1	24
Moses-PB	17	5	7	29
Moses-HPB	14	4	7	25

ser, *estar* or *ficar*. In the case of Moses-PB and Moses-HPB, some of these errors can be linked to the nature of Europarl. For instance, *ask* is translated as *auscultarem* (auscultate), *you* as *excelência* (excellency), *house* as *assembleia* (assembly) and *sitting* as *sessão* (session).

Considering wrong choice errors, these can produce translations with no apparent semantic explanation. An example of this is the translation of *understand* into *tradição* by Moses-HPB.

Finally, collocational errors and idioms were not a significant problem.

As far as **discourse** errors are concerned, Table 14 summarises the observed number of errors.

Clearly, translating into EP (and not BP) is where Online-G performs worse, as there are many variety errors, resulting from the fact that this engine is translating EN into BP.

5.6 BLEU and METEOR scores

We performed an automatic MT evaluation using BLEU (Papineni et al. 2002) as well as METEOR²⁵ (Denkowski and Lavie 2014), which uses other linguistic resources such as paraphrases. The results are presented in Table 15.

The MT system that had the lowest number of errors in the Questions corpus was Online-G. This is also the system with the best BLEU (58.76) and METEOR scores (72.79). Both the human and the automatic evaluation metrics agree that the system with more words tagged as errors is Online-S.

²⁵ version 1.4

Table 15 BLEU and METEOR scores achieved by the translation systems when evaluated on each test dataset

	TAP	TED	Questions
BLEU			
Online-S	17.94	19.66	32.98
Online-G	30.10	27.27	58.76
Moses-PB	45.23	26.92	42.97
Moses-HPB	44.15	27.17	43.10
METEOR			
Online-S	38.59	39.51	51.76
Online-G	49.86	48.67	72.79
Moses-PB	58.84	45.98	62.02
Moses-HPB	57.51	45.43	61.86

Table 16 BLEU and METEOR scores achieved by the translation systems trained online with the Europarl corpus

	TAP	TED	Questions
BLEU			
Moses-PB	37.50	26.59	38.60
Moses-HPB	38.27	25.61	38.60
METEOR			
Moses-PB	52.68	44.76	57.81
Moses-HPB	53.12	44.04	57.09

On the TED corpus, the best BLEU and METEOR scores are the ones from Online-G (27.27 and 48.67, respectively). According to the human evaluation this is also the system that has the lowest number of errors.

Interestingly, on the TAP corpus, we observe that our in-house SMT systems obtain a significantly better result in terms of BLEU and METEOR over Online-G, which is trained on more data, and outperforms our systems on all other datasets. However, we can observe that in terms of the total number of errors (as obtained by the human evaluation), system Online-G actually performs better than the in-house systems. In order to explain this inconsistency, we first tested whether our trained systems were over-fitting the domain using the in-domain data. This was done by training systems using only the Europarl dataset. Results are shown in Table 16, where we observe that the BLEU and METEOR scores for both Moses systems, even though having dropped drastically, are still higher than for the system Online-G (see Table 15).

This shows that the results are not caused by over-fitting the training data as the Europarl dataset is radically different from the TAP dataset. However, this could still be due to the tuning corpus, which is in-domain. We looked at the development corpora and noticed that there are many equivalent sentence pairs, such as menu items and general flight instructions, which are present in every issue of the magazine. Furthermore, many sentences are simply repetitive, such as, *Have a good flight.* and *Fancy a snack.* This happens with 53 sentence pairs (approximately 20 %) in the TAP corpus, while there are only 5 and 2 (less than 1 %) repeated sentences in the Questions and TED datasets, respectively. This allows the MT systems to tailor their output so that the translations of such content are as close as possible to the reference. This gives

a large boost in the BLEU and METEOR scores as they are biased towards finding translations that are close to the reference and not by correctness per se. Still, it is an interesting result that even in such conditions, human evaluators find more errors in Moses-PB and Moses-HPB than in Online-G. This shows the many shortcomings of these metrics, which rely on closeness to the reference or references rather than analysing their linguistic quality.

To confirm the hypothesis that the boost in scores is due to the common content in every issue of the magazine (such as flight instructions or menu items), we selected 50 sentences extracted from the magazine's main articles only. The recomputed METEOR and BLEU metrics on this subset are consistent with the previous experiments. As previously, Online-G (METEOR = 43.49; BLEU = 25.81) ranks first, while Moses-PB (METEOR = 35.31; BLEU = 22.14) and Moses-HPB (METEOR = 39.10; BLEU = 22.58) rank lower and Online-S ranks last (METEOR = 33.19; BLEU = 15.85). These scores are in line with the human annotation, which assigns better translation results to Online-G than to both Moses systems.

6 Error gravity

Having just investigated how errors occur in four different systems and three translation scenarios (Questions, TED and TAP datasets), we decided to analyze to what extent distinct error types impact translation quality. To accomplish that, we start with a subjective evaluation of the MT outputs, which consists of ranking the four translations of each sentence (Sect. 6.1). Then, by relating this rank to our taxonomy, we are able to show how the presence of each error type reflects on quality (Sect. 6.2).

6.1 Ranking translations

Using the same set of 75 sentences that were used in the experiment reported in Sect. 4.3, we carried out an evaluation similar to the one proposed in [Callison-Burch et al. \(2007\)](#). This task consists of presenting the annotators with the input sentence, the correct translation and all four MT outputs. They then decide on the order of translations based on their assessment of quality, ranking them from 1 (best) to 4 (worst). In our experiment, however, ties should only exist for translations that are exactly the same. This encouraged judges to make a decision even when facing tenuous differences such as capitalization, number or gender variations.

To report inter-annotator agreement, three annotators ranked a smaller sample of 120 translations: 10 sentences per dataset translated by the four MT engines. The agreement between the three annotators using Cohen's Kappa was 0.572, which according to [Landis and Koch \(1977\)](#) is considered a moderate agreement.

Table 17 shows the ranking of sentences per system, considering the 75 sentences (4 translations for each sentence, 300 translations in total). Online-G clearly contrasts with the other systems, having produced 50 translations that were considered the best and only 6 that were ranked in fourth place. It is important to note that the total of 90 translations that were ranked first (instead of the expected 75) results from ties that occur, for example, when multiple systems produce perfect translations.

Table 17 Number of sentences per ranking level and MT system

System	1	2	3	4
Online-S	11	31	9	24
Online-G	50	14	5	6
Moses-PB	17	16	31	11
Moses-HPB	12	14	21	28
Total	90	75	66	69

Table 18 Comparison between MT systems on 300 translations

Better than	Online-S	Online-G	Moses-PB	Moses-HPB	Total
Online-S	–	15	41	42	98
Online-G	58	–	56	60	174
Moses-PB	34	16	–	38	88
Moses-HPB	29	11	14	–	54
Total	121	42	111	140	–

With this ranking, we were able to see how many times each system was better than the other (Table 18). These results show that Online-G regularly ranks higher than the other systems. Online-S is the next best system, ranking 98 times higher than other systems. Finally, Moses-PB and Moses-HPB are the systems that ranked the lowest (88 and 54 times, respectively).

6.2 Relating error types with translation quality

Although ranking translations allows us to compare performance between systems, it is not enough to analyze error severity. If we used this ranking in a straightforward manner, we would end up grouping translations that are close to perfect with translations that contain errors that severely hinder comprehension. For instance, the translation *Quem era o Galileo?* (*Who was the Galileo?), where the only problem is the insertion of the article “the”, ranked the lowest in its group (4th place) because all the translations from the other systems were perfect. This contrasts with the case of a sentence that is placed last because it has severe comprehension problems. Consequently, if used directly, the ranking strategy described in Sect. 6.1 places the errors from these two distinct cases at the same level. In other words, it causes a discrepancy of error gravity within the same level.

To overcome this obstacle, we decided to assign a class to each set of 4 translations of the same sentence, using the following criteria:

- **All-Good:** contains sets of translations where all four of them are good, and thus the differences in their rank are caused by minor errors;
- **All-Bad:** is composed of sets of four translations where each has significant problems and errors that hinder comprehension;

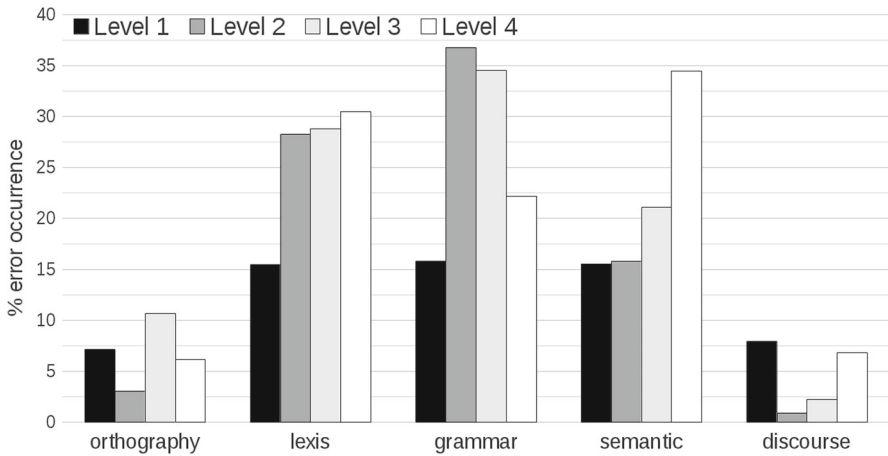


Fig. 3 Percentage of errors on mixed class

- **Mixed:** contains the remaining cases, i.e. sets of translations where some of them have only minor errors (or are perfect) while others have severe comprehension problems.

In sum, this formulation assigns a class to each set of 4 translations according to the differences in terms of quality between them. Having done that, we can now look at what happens inside each class individually, avoiding the problems that were pointed out in the aforementioned example. An error responsible for a lower rank in the All-Good class (like the article insertion in *Who was the Galileo?) does not have the same impact as an error that causes a translation to rank lower on the Mixed or All-Bad classes.

It is also important to mention that the decision to use three classes instead of a more fine-grained classification is due to the low agreement reported in the similar task of judging fluency and adequacy on a 5-point scale (Callison-Burch et al. 2007). In our task of grouping sets of translations by quality we achieved a Kappa of 0.677 (considered substantial agreement). The distribution along classes was of only 14 groups of translations assigned the class All-Good, 21 classified as All-Bad, and the majority (40) ending up classified as Mixed.

This latter class, Mixed, is the one that we are most interested to look at in detail. In this class, we can truly distinguish between ranks, since translations that ranked higher are expected to be good while translations that were placed last are likely to have severe comprehension problems. For this reason, in Fig. 3 we plotted the average percentage of each error type for the translations in the Mixed class as they occur across ranks. Note that the “Level 1” column contains all translations that ranked first in their set, independent of the system that generated them. Levels 2, 3 and 4 represent the second, third and worst translations, respectively.

We can conclude that:

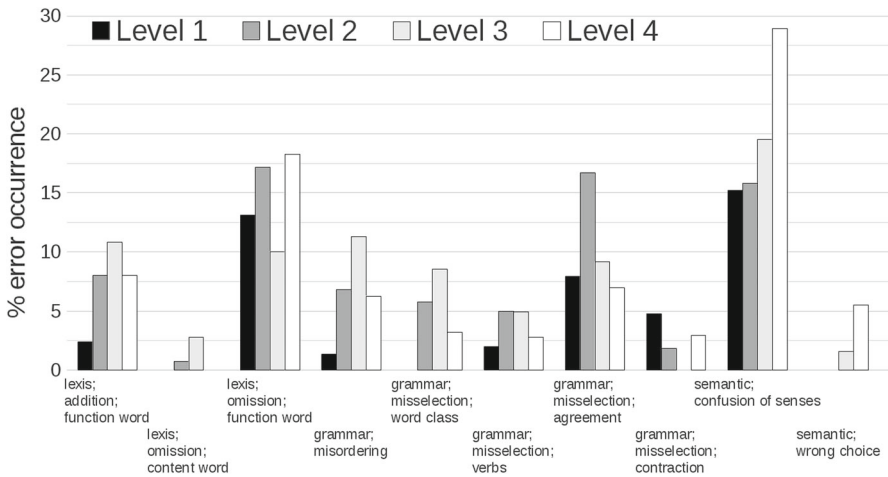


Fig. 4 Percentage of subtype of errors on mixed class

- `lexis` and `semantic` errors are clearly correlated with ranking. Results show that the more `lexis` and `semantic` errors present, the worse the translation's rank. This trend is much more accentuated for the `semantic` error type;
- lower percentages of `grammar` errors (approx. 15 %) seem to be associated with better quality translations. However, it is not clear how these errors affect translations in lower ranks;
- finally, `orthography` and `discourse` errors have the lowest percentage (< 15% in all ranks), and do not show an increasing or decreasing trend, seeming to occur at the same rate in opposing ranks.

Given some of the inconclusive results obtained in the previous figure, we decided to look further into the taxonomy and plot the subtype of errors and corresponding ranking, again in the Mixed Class (Fig. 4). For clarity, we have omitted all errors with low representation (<1 %).

When plotting error subtypes, we can see that:

- both subtypes of `semantic` errors show great correlation with ranking. The `wrong choice` subtype only occurs on translations ranked as 3 and 4. The `confusion of senses` subtype demonstrates an exponential growth as the ranking decreases, achieving almost 30 % on translations that ranked last;
- for the remaining error subtypes, both lexical and grammatical, we cannot identify a clear trend related to the ranking, with higher percentages of errors being assigned to the medial ranks.

Finally, we decided to look into what happens when all translations have severe problems. Figure 5 shows the average percentage of errors for translations in the class All-Bad. Again, we omit error types with low percentages (<1 %).

We can see that:

- contrary to what happened to the Mixed class, `confusion of senses` is not discriminative any more. Instead, it just shows a high percentage for all ranks;

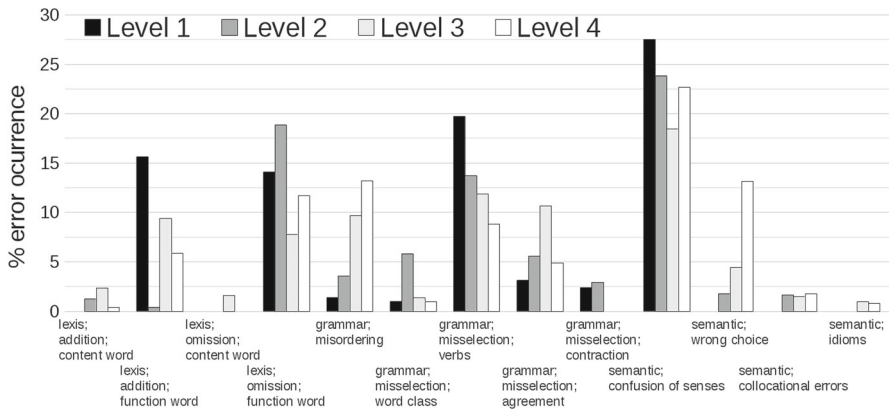


Fig. 5 Percentage of subtype of errors on all-bad class

- `wrong choice` now shows a more clear correlation;
- also opposed to what happens in the Mixed class, misordering **grammar** errors are now correlated with rank, negatively impacting comprehension;
- another difference is the occurrence of addition of content words, collocational errors, and idioms, absent from the Mixed class, and that appear always associated with lower-ranking translations.

7 Conclusions and future work

This work aims at developing a detailed taxonomy of MT errors, which extends previous taxonomies, usually focused on English errors. Therefore, the proposed taxonomy is tailored to support errors that are usually associated with morphologically richer languages, such as the Romance languages. At the basis of this taxonomy are the main areas of linguistics, we hope that it can be easily extended in order to support errors associated with specific phenomena in other Romance languages.

After establishing our taxonomy, we automatically translated sentences from three corpora, each one representing a specific translation challenge, using four different systems: two mainstream online translation systems (Google Translate (Statistical) and Systran (Rule-based)), and two in-house MT systems. Errors were manually annotated, according to the proposed taxonomy, allowing us to evaluate each system and establish some comparisons between them. For instance, we concluded that Online-G has several mistakes of *variety*. This is probably due to the fact that it translates EN into BP and not EP, as most of its training resources are BP, and it is not distinguishing between the two varieties. Regarding Moses-HPB and Moses-PB, we could find many lexis errors, specially “untranslated”, as their training corpus is limited both in size and in domain. A detailed error analysis is provided in the paper.

Regarding error gravity, we found that problems related to *confusion of senses*, *wrong choice* and *misordering* are the phenomena that most impact

translation quality, since they seem to correlate with a subjective ranking of the translations.

As far as future work is concerned, we intend to ask human translators to translate the three corpora into EP. First, we will see if we need to extend our taxonomy in order to support their mistakes. For instance, our taxonomy does not support “invented” words, which are not usual in MT, as most systems only output words that were seen during training. Nevertheless, these errors are quite commonplace in human translations. This type of error could be easily integrated into the taxonomy at the *lexis* level. Second, we will analyze the errors obtained and compare them with those committed by the MT engines.

Finally, we would like to automate some steps of our taxonomy. With some statistical learning, errors like omission, addition and words that were not translated could be automatically found and that would definitely help to make the error analysis quicker. Furthermore, having information about the most critical sentences to translate could shed some light on where the translation errors might be found.

Acknowledgments This work was partially supported by national funds through FCT - Fundação para a Ciência e a Tecnologia, under project UID/CEC/50021/2013. Ângela Costa, Wang Ling and Rui Correia are supported by PhD fellowships from FCT (SFRH/BD/85737/2012, SFRH/BD/51157/2010, SFRH/BD/51156/2010).

References

- Batista F, Mamede N, Trancoso I (2007) A lightweight on-the-fly capitalization system for automatic speech recognition. In: Proceedings of the recent advantages in natural language processing (RANLP'07), Borovets, Bulgaria
- Bojar O (2011) Analysing error types in English-Czech machine translation. *Prague Bull Math Linguist* 95:63–76
- Bojar O, Mareček D, Novák V, Popel M, Ptáček J, Rouš J, Žabokrtský Z (2009) English-Czech MT in 2008. In: Proceedings of the fourth workshop on statistical machine translation, Greece, Athens, pp 125–129
- Callison-Burch C, Forgyce C, Koehn P, Monz C, Schroeder J (2007) (meta-) evaluation of machine translation. In: Proceedings of the second workshop on statistical machine translation, Czech Republic, Prague, pp 136–158
- Castagnoli S, Ciobanu D, Kunz K, Volanschi A, Kubler, N (2007) Designing a learner translator corpus for training purposes. In: TALC7, proceedings of the 7th teaching and language corpora conference, Paris, France
- Chiang D (2007) Hierarchical phrase-based translation. *Comput Linguist* 33(2):201–228
- Condon SL, Parvaz D, Aberdeen JS, Doran C, Freeman A, Awad M (2010) Evaluation of machine translation errors in english and Iraqi Arabic. In: Proceedings of the seventh international conference on language resources and evaluation, Valletta, Malta, pp 159–168
- Corder SP (1967) The significance of learner's errors. *Int Rev Appl Linguist* 5(4):161–169
- Costa A, Luís T, Coheur L (2014) Translation errors from English to Portuguese: an annotated corpus. In: Proceedings of the ninth international conference on language resources and evaluation (LREC'14), Reykjavik, Iceland, pp 1231–1234
- Costa A, Luís T, Ribeiro J, Mendes AC, Coheur L (2012) An English-Portuguese parallel corpus of questions: translation guidelines and application in SMT. In: Proceedings of the eighth international conference on language resources and evaluation (LREC'12), Istanbul, Turkey, pp 2172–2176
- Cunha C, Cintra L (1987) *Nova Gramática do Português Contemporâneo*. Edições Sá da Costa, Lisboa
- De Saussure F (1916) *Cours de linguistique générale*. Payot, Paris

- Denkowski M, Lavie A (2014) Meteor universal: language specific translation evaluation for any target language. In: WMT 2014: proceedings of the ninth workshop on statistical machine translation, Baltimore, Maryland USA, pp 376–380
- Dulay H, Burt MK, Krashen SD (1982) *Language two*. Oxford University Press, Oxford
- Elliott D, Hartley A, Atwell E (2004) A fluency error categorization scheme to guide automated machine translation evaluation. In: Frederking RE, Taylor KB (eds) *Machine translation: from real users to research: 6th conference of the Association for Machine Translation in the Americas, AMTA 2004*, Washington, DC, vol 3265., Lecture Notes in Computer Science. Springer, Berlin, pp 64–73
- Fishel M, Bojar O, Popović M (2012) Terra: a collection of translation error-annotated corpora. In: *Proceedings of the eighth international conference on language resources and evaluation (LREC'12)*, Istanbul, Turkey, pp 7–14
- James C (1998) *Errors in language learning and use. Exploring error analysis, applied linguistics and language study*. Routledge, New York
- Keenan EL, Stabler EP (2010) Language variation and linguistic invariants. *Lingua* 120(12):2680–2685
- Kirchhoff K, Rambow O, Habash N, Diab M (2007) Semi-automatic error analysis for large-scale statistical machine translation. In: *MT Summit XI. Proceedings, Copenhagen, Denmark*, pp 289–296
- Koehn P (2005) Europarl: A Parallel Corpus for Statistical Machine Translation. In: *MT Summit X, Conference proceedings: the tenth machine translation summit, Phuket, Thailand*, pp 79–86
- Koehn P, Hoang H, Birch A, Callison-Burch C, Federico M, Bertoldi N, Cowan B, Shen W, Moran C, Zens R, Dyer C, Bojar O, Constantin A, Herbst E (2007) Moses: open source toolkit for statistical machine translation. In: *Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions, Prague, Czech Republic, Association for Computational Linguistics*, pp 177–180
- Landis RJ, Koch GG (1977) The measurement of observer agreement for categorical data. *Biometrics* 33:159–74
- Ling W, Luis T, Graca J, Coheur L, Trancoso I (2010) Towards a general and extensible phrase-extraction algorithm. In: *IWSLT '10: international workshop on spoken language translation, France, Paris*, pp 313–320
- Li X, Roth D (2002) Learning question classifiers. In: *Proceedings of the 19th international conference on computational linguistics (COLING)*, Taipei, Taiwan, pp 556–562
- Llitjós AF, Carbonell JG, Lavie A (2005) A framework for interactive and automatic refinement of transfer-based machine translation. In: *10th EAMT conference “Practical applications of machine translation”*, Budapest, Hungary, pp 87–96
- Naskar SK, Toral A, Gaspari F, Way A (2011) Framework for diagnostic evaluation of MT based on linguistic checkpoints. In: *Proceedings of machine translation summit XIII, Xiamen, China*, pp 529–536
- Niessen S, Och FJ, Leusch G, Ney H (2000) An evaluation tool for machine translation: fast evaluation for MT research. In: *LREC-2000: second international conference on language resources and evaluation. Proceedings, Athens, Greece*, pp 39–45
- Och FJ (2003) Minimum error rate training in statistical machine translation. In: *ACL-2003: 41st annual meeting of the association for computational linguistics, Sapporo, Japan*, pp 160–167
- Och FJ, Ney H (2003) A systematic comparison of various statistical alignment models. *Comput Linguist* 29:19–51
- Papinen K, Roukos S, Ward T, Zhu WJ (2002) BLEU: a method for automatic evaluation of machine translation. In: *ACL-2002: 40th annual meeting of the association for computational linguistics, Philadelphia*, pp 311–318
- Popović M, de Gispert A, Gupta D, Lambert P, Ney H, Marino JB, Federico M, Banchs R (2006) Morphosyntactic information for automatic error analysis of statistical machine translation output. In: *HLT-NAACL 2006: proceedings of the workshop on statistical machine translation, New York, NY, USA*, pp 1–6
- Popović M, Ney H (2006) Error analysis of verb inflections in Spanish translation output. In: *TC-STAR workshop on speech-to-speech translation, Barcelona, Spain*, pp 99–103
- Popović M, Ney H (2011) Towards automatic error analysis of machine translation output. *Comput Linguist* 37(4):657–688
- Richards J (1974) *Error analysis. Perspectives on second language acquisition*. Longman, London
- Secară A (2005) Translation evaluation—a state of the art survey. *eCoLoRe/MeLLANGE Workshop*. Leeds, UK, pp 39–44

- Vilar D, Xu J, D'Haro LF, Ney H (2006) Error analysis of machine translation output. In: LREC-2006: fifth international conference on language resources and evaluation. Proceedings, Genoa, Italy, pp 697–702
- Zeman D, Fishel M, Berka J, Ondřej (2011) Addicter: what is wrong with my translations? Prague Bull Math Linguist 96:79–88