

## Linguistic features of genre and method variation in translation: A computational perspective

|               |   |
|---------------|---|
| Item Type     | Chapter in book   |
| Authors       | Lapshinova-Koltunski, Ekaterina;Zampieri, Marcos  |
| Citation      | Ekaterina Lapshinova-Koltunski, Marcos Zampieri (2018). Linguistic features of genre and method variation in translation: a computational perspective. In Dominique Legallois, Thierry Charnois, Meri Larjavaara (Eds.), <i>The Grammar of Genres and Styles: From Discrete to Non-Discrete Units</i> (pp. 92–117). Berlin, Boston: De Gruyter. <a href="https://doi.org/10.1515/9783110595864-005">https://doi.org/10.1515/9783110595864-005</a> |
| DOI           | <a href="https://doi.org/10.1515/9783110595864-005">10.1515/9783110595864-005</a>   |
| Publisher     | Mouton De Gruyter   |
| Rights        | Attribution-NonCommercial-NoDerivs 3.0 United States  |
| Download date | 2026-04-14 11:50:36   |
| License       | <a href="http://creativecommons.org/licenses/by-nc-nd/3.0/us/">http://creativecommons.org/licenses/by-nc-nd/3.0/us/</a>   |
| Link to Item  | <a href="http://hdl.handle.net/2436/622010">http://hdl.handle.net/2436/622010</a>   |

# Linguistic Features of Genre and Method Variation in Translation: A Computational Perspective

Ekaterina Lapshinova-Koltunski<sup>1</sup> and Marcos Zampieri<sup>2</sup>

<sup>1</sup>Saarland University, Germany

<sup>2</sup>University of Wolverhampton, United Kingdom

## Abstract

In this paper we describe the use of text classification methods to investigate genre and method variation in an English - German translation corpus. For this purpose we use linguistically motivated features representing texts using a combination of part-of-speech tags arranged in bigrams, trigrams, and 4-grams. The classification method used in this paper is a Bayesian classifier with Laplace smoothing. We use the output of the classifiers to carry out an extensive feature analysis on the main difference between genres and methods of translation.

## 1 Introduction

In the present study, we use text classification techniques to explore variation in translation. We analyse the interplay between two dimensions influencing this variation: translation methods (human and machine translation) and text registers (e.g. fiction, political speeches, etc.). Our starting assumption is that the interplay between these dimensions is reflected in the lexico-grammar of translated texts, i.e. in their linguistic features. Our assumption here is that registers and methods represent two dimensions influencing linguistic properties of translations, and can thus be confounding in a specific task. For instance, if we want to automatically distinguish between human and machine translation, we need to exclude features which are rather register-specific as they can compromise the results of the classification.

In our previous work, see e.g. Lapshinova-Koltunski (2017), we used a set of features derived from theoretical frameworks, such as genre / register theory, e.g. Halliday and Hasan (1989); Biber (1995); Neumann (2013), or translationese studies, e.g. Baker (1993); Baroni and Bernardini (2006); Volansky et al. (2011). In the present analysis, we use a data-driven approach, which will help us to discover new language structures reflecting variation in translation. Classification techniques will help us to identify discriminative features of the two variation dimensions under analysis. For this, we train classifiers to distinguish translated texts according to either their register or method of translation, using the VARTRA corpus (Lapshinova-Koltunski, 2013), a collection of English to German translations. Our assumption is that text classification methods can level out discriminative features of different translation varieties that intuition alone cannot grasp; thus enabling us to investigate in more detail the properties of each of them. More than the classification results *per se*, we use level out interesting linguistic features that can be further used in linguistic analysis and NLP applications.

Text classification methods have been applied in a wide range of tasks such as spam detection (Medlock, 2008), native language identification (NLI) (Gebre et al., 2013) temporal text classification (Niculae et al., 2014), and the identification of lexical complexity in text (Malmasi et al.,

2016). In the aforementioned studies, researchers are interested in how well classification methods can perform or, in other words, how reliably these methods are able to attribute correct labels to a set of texts. Therefore, most researchers in text classification are concerned in exploring features and algorithms that deliver the best performance for each task. In recent works Diwersy et al. (2014); Zampieri et al. (2013), however, text classification methods were proposed to investigate language variation across corpora (e.g. diatopic and dialectal variation) using linguistically motivated features.

In this paper, we propose an approach to automatically classify translated texts regarding register and method of translation. We are interested not only in obtaining state-of-the-art classification performance, but also in leveling out interesting linguistic features from the data. The features used here constitute combinations of part-of-speech (POS) tags. These POS combinations represent, however, language patterns, e.g. a finite auxiliary verb followed by a participle represents a verbal phrase in a passive voice. Analyzing sets of the POS combinations that result from the classification experiments, we try identify those that are specific for the classes under analysis.

This paper is structured as follows: the following section presents the theoretical background, as well as the related work. Here, we also describe our previous experiments on text classification for the analysis of translation variation. In Section 3, we introduce the dataset, as well as the methodology applied for the analysis. Section 4 presents the results of the classification. We further investigate the results in Section 5, in which we concentrate on the analysis of features specific for translation methods and genres. Section 6 summarizes the findings and presents a discussion on the related issues.

## 2 Theoretical Background and Related Work

### 2.1 Theoretical Background

Translation is influenced by several factors, including the source and the target language, registers or genres a text belongs to, as well as the translation method involved. Since the present study focuses on genre and method variation, we will also base our research on the studies related to this type of variation.

Genre-specific variation of translation is related to studies within register and genre theory, e.g. Halliday and Hasan (1989), Biber (1995), which analyse contextual variation of languages. In the present paper, we use the term **genre** and not **register**, although they represent two different points of view covering the same ground, see e.g. Lee (2001), and we use the latter in our previous studies, see e.g. Lapshinova-Koltunski (2017) and Lapshinova-Koltunski and Vela (2015). Mostly, we refer to genre when speaking about a text as a member of a cultural category, about a register when we view a text as language. However, in this study we consider both as lexico-grammatical characterisations, conventionalisation and functional configuration determined by a context use.

The differences between genres can be identified through a corpus-based analysis of phonological, lexico-grammatical and textual (cohesion) features in these genres; see the studies on linguistic variation by Biber (1995) or Biber et al. (1999), and linguistic variation among genres can be traced in the distribution of these features.

Multilingual studies concern linguistic variation across languages, comparing genre and register settings specific for the languages under analysis, e.g. Biber (1995) on English, Nukulaelae Tuvaluan, Korean and Somali, and Hansen-Schirra et al. (2012) and Neumann (2013) on English and German. The latter two also consider this type of linguistic variation in translations. Other translation scholars e.g. Steiner (2004) and House (2014), also pay attention to genre and

register variation when analysing language in a multilingual context of translation. However, they either do not account for the distributions of the corresponding features, or analyse individual texts only. In the works by De Sutter et al. (2012) and Delaere and De Sutter (2013), register-related differences are also described for translated texts. Yet, these differences are identified on the level of lexical features only.

The features that are most frequently used in studies on variation in corpus-based approaches are of shallow character and include lexical density (LD), type-token-ratio (TTR), and part-of-speech (POS) proportionality. Steiner (2012) uses these features to characterise profiles of various subcorpora distinguished by language (English and German), text production type (translation and original) and eight different registers. The author defines a number of contrast types including register controlled ones which implies (1) contrasts within one register between English and German, and (2) contrasts between registers within each of the languages, see (Steiner, 2012, p. 72). In our analysis, we consider genre variation only within translations.

Applying a quantitative approach, Neumann (2013) analyses an extensive set of linguistic patterns reflecting register variation and shows the differences between the two languages under analysis. The author also demonstrates to what degree translations are adapted to the requirements of different registers, showing how both register and language typology are at work.

Kunz et al. (2017) show that register variation is also relevant for a number of textual phenomena. They analyse structural and functional subtypes of coreference, substitution, discourse connectives and ellipsis on a dataset of several registers in English and German. They are able to identify contrasts and commonalities across the two languages and registers with respect to the subtypes of all textual phenomena under analysis. The authors show that these languages differ as to the degree of variation between individual registers in the realisation of the phenomena under analysis, i.e. there is more variation in German than English. They attest the main differences in terms of preferred meaning relations: a preference for explicitly realising logico-semantic relations by discourse markers and a tendency to realise relations of identity by coreference. Interestingly, similar meaning relations are realised by different subtypes of discourse phenomena in different languages and registers.

Whereas attention is paid to genre settings in human translation analysis, they have not yet been considered much in machine translation. There exist some studies in the area of statistical machine translation (SMT) evaluation, e.g. errors in translation of new domains Irvine et al. (2013). However, the error types concern the lexical level only, as the authors operate solely with the notion of domain and not genre. Domains represent only one of the genre parameters and reflect what a text is about, i.e. its topic, and further settings are thus ignored. Although some NLP studies, e.g. those employing web resources, do argue for the importance of genre conventions, see e.g. Santini et al. (2010), genre remains out of the focus of machine translation. In the studies on adding in-domain bilingual data to the training material of SMT systems (Wu et al., 2008) or on application of in-domain comparable corpora (Irvine and Callison-Burch, 2014), again, only the notion of domain is taken into consideration.

Variation in terms of translation method has not received much attention so far. There are numerous studies in the context of NLP that address both human and machine translations (Papineni et al., 2002; Babych et al., 2004). Yet they all serve the task of automatic MT system evaluation and focus solely on translation error analysis, using human translation as a reference in the evaluation of machine translation outputs. Evaluations serves the task to prove to what extent automatically translated texts (hypothesis translations) comply with the manually translated ones (reference translations). The ranking of machine-translated texts is based on scores produced with various metrics. The metrics applied in the state-of-the-art MT evaluation are automatic and language-independent: BLEU and NIST (Doddington, 2002). However,

since they do not incorporate any linguistic features, BLEU scores need to be treated carefully, which was demonstrated by Callison-Burch et al. (2006). This fact has been advancing the development of new automatic metrics, such as METEOR (Denkowski and Lavie, 2014), Asiya (González et al., 2014) and VERTa (Comelles and Atserias, 2014). They incorporate lexical, syntactic and semantic information into their scores. The accuracy of the evaluation methods is usually proven through human evaluation. More specifically, the automatically provided scores are correlated with the human judgements which are realised by ranking MT outputs (Bojar et al., 2014; Vela and van Genabith, 2015) and others. Some of the existing metrics incorporate linguistic knowledge.

There are even more works on MT evaluation that operate with linguistically-motivated categories, e.g. Popovic and Ney (2011) or Fishel et al. (2012). However, none of them provides a comprehensive analysis of the differences between human and machine translation in terms of specific linguistically motivated features. In fact, the knowledge on the discriminative features of human and machine translation can be derived from the studies operating with machine learning procedures for MT evaluation, such as Stanojević and Simaán (2014) or Gupta et al. (2015). Corston-Oliver et al. (2001) use classifiers that learn to distinguish human translations from machine ones. These classifiers are trained with various features including lexicalised trigram perplexity, part of speech trigram perplexity and linguistic features such as branching properties of the parse, function word density, constituent length, and others. Their best results are achieved if perplexity calculations were combined with finer-grained linguistic features. Their most discriminatory features that differentiate between human and machine translations are not just word n-grams. They include the distance between pronouns, the number of second person pronoun, the number of function words, and the distance between prepositions.

Volansky et al. (2011) operate with translationese-inspired features, and are able to distinguish between manual and automatic translations in their dataset with 100% accuracy. However, the manual and automatic translations they are using have different source texts. We believe that the distinction they are able to achieve is not the distinction between translation methods, but rather between different underlying texts, since their most discriminatory features are the ones that show good performance in any text classification task (token n-grams). El-Haj et al. (2014) make use of readability as a proxy for style and analyse consistency in translation style considering how readability varies both within and between translations. They compare Arabic and English human and machine translations of the originally French novel “The Stranger” (French: *L’Étranger*). The results show that translations by humans (both male and female) are closer to each other than to automatic translations. The authors also measure closeness of translations to the original in terms of the selected measures, which should serve as an indicator of translation quality.

To the best of our knowledge, there have not been many studies published about the interplay between the two dimensions influencing translation that are in focus of our study. Kruger and van Rooy (2012) try to answer the question on the relationship between register and the features of translated language. Their hypothesis was that the translation-related features would not be strongly linked to register variation suggesting that in translated text reveal less register variation, or sensitivity to register, which is a consequence of translation-specific effects. However, their findings provide limited support for this hypothesis. They state that the distribution and prevalence of linguistic realisations of the features of translated language may vary according to register. Therefore, the concept of translated language should be more carefully analysed and defined in terms of registers (Kruger and van Rooy, 2012, p. 61–62). Jensen and McGillivray (2012) analyse the interaction between registers, source language and translators’ background on the basis of morphological features. The interaction between the dimensions of register, author and translator was also analysed by Jensen and Hareide (2013)

who use patterns of sentence alignment as features.

Thus, there is no comprehensive description of the linguistic features that represent the dimensions of translation variation. We analyse the interplay between the dimensions of genre (register) and method trying to detect specific features that reflect this interplay.

## 2.2 Previous Experiments

In our previous analyses (Zampieri and Lapshinova-Koltunski, 2015), we applied text classification methods on a set of English-German translations. We used two different sets of features: n-grams taking 1) all word forms into account and 2) semi-delexicalized text representations – all the nouns were replaced with placeholders, which represented the novelty of that approach. Our task was two-fold: (a) to discriminate between different genres (fiction, political essays, etc., a total of seven classes); and (b) to discriminate between translation methods (human professional, human student, rule-based machine and two statistical systems – five classes).

We performed several classification tasks: (1) We use word n-grams to train five translation method and seven genre classes; (2) We use delexicalised n-grams to classify four and five translation method classes and seven register classes; (3) We use delexicalised n-grams to classify between human and machine translation.

The results of the first experiment show that the classifier performs better for the distinction of genres than of translation methods (F-measure of 57.30% and 35.30% respectively). This result is not surprising, as content words (including proper nouns) are domain specific and, therefore, the classifier can better differentiate between genres that vary in their domains. The results of this experiment shows that it is important to use (semi-)delexicalised features in a dataset that represents both dimensions of variation in translation – genre and method.

The results of the second experiment show that delexicalised features reduce the performance of the genre classifier (from 57.30% to 45.40%) but increase the performance of the translation method classifier (from 35.30% to 43.10%), especially if we reduce the dataset to four classes instead of five (we concatenate both statistical machine translation outputs). The results of this experiment confirm the importance of (semi-)delexicalised features, as we achieve similar scores for the analysis of both dimensions of variation in our data.

In the last experiment, we reduce the number of translation method classes to two – human and machine. This classification is less fine-grained and represents manual and automatic procedures of translation. As expected, this experiment delivers better classifier results: up to 60.5% F-measure in distinguishing between manually and automatically translated texts.

In the last step, we performed qualitative analysis of the output features paying attention to those which turned to be most informative for the corresponding classification task. In this way, we were able to identify a set of semi-delexicalized n-grams that are discriminative for either certain genres or translation methods in our data. This step was manual and included evaluation of trigrams only, as the performance of trigram models achieved the best results in the classification task. We generated two lists of features specific either to human or machine translation, and fourteen lists of features discriminating genre pairs. The features informative for translation method included full nominal phrases that differentiated in the type of determiners (articles in human and possessives in machine translations), personal pronouns expressing coreference that differentiated in the grammatical number (singular in human and plural in machine translations), event anaphors that differentiated in the type of pronouns (demonstrative in human and personal in machine), etc. These features differed from those specific for genre identification. For instance, for the discrimination between political essays and fictional texts, discourse markers expressing different relations turned to be informative. Moreover, the lists included also features related to verbal phrases, e.g. passive vs. active voice, infinitives and

modal verbs differing in their meaning.

In this paper we base upon the results of this analysis: we use a two-member classification of translation methods (manual and automatic) and seven-member classification for genres. Moreover, we decide to fully delexicalise the features and run our experiments on delexicalised features instead of semi-delexicalised ones.

## 3 Methods

### 3.1 Data

For the purpose of our study we looked for suitable translation corpora containing different genres and methods of translation. The only corpus known to us that possesses these characteristics is VARTRA (Lapshinova-Koltunski, 2013). VARTRA comprises multiple translations from English into German. These translations were produced with five different translation methods as follows: (1) human professionals (PT1), (2) human student translators (PT2), (3) a rule-based MT system (RBMT), (4) a statistical MT system trained with a large quantity of unknown data (SMT1) and (5) a statistical MT system trained with a small amount of data (SMT2).

VARTRA contains texts from different genres, namely: political essays (ESS), fictional texts (FIC), instruction manuals (INS), popular-scientific articles (POP), letters of share-holders (SHA), prepared political speeches (SPE), and touristic leaflets (TOU). Each sub-corpus represents a translation variety, a translation setting which differs from all others in both method and genre (e.g. PT1-ESS or PT2-FIC, etc.). The corpus is tokenised, lemmatised, tagged with part-of-speech information, segmented into syntactic chunks and sentences. The annotations were obtained with Tree Tagger (Schmid, 1994).

Before classification was carried out, we split the corpus into sentences.<sup>1</sup> The length of each sentence varies between 12 and 24 tokens. This results in a dataset containing 6,200 instances.

The features used in the experiments we report in this paper were based on the combinations of POS tags arranged in form of bag-of-words (BoW), bigrams, trigrams, and 4-grams.<sup>2</sup> In Example (1), we illustrate the representation of the sentences in the corpus. (1-a) represents a sentence from the corpus, (1-b) shows the representation, where all nouns are substituted with the placeholder *PLH* resulting in what we call a semi-delexicalized text representation.

- (1) a. *Die weltweiten Herausforderungen im Bereich der Energiesicherheit erfordern ber einen Zeitraum von vielen Jahrzehnten nachhaltige Anstrengungen auf der ganzen Welt.*
- b. *Die weltweiten PLH im PLH der PLH erfordern ber einen PLH von vielen PLH nachhaltige PLH auf der ganzen PLH.*
- c. ART ADJA NN APPRART NN ART NN VVFIN APPR ART NN APPR PIAT  
       ADJA ADJA NN APPR ART ADJA NN.

This type of representation lies between fully delexicalized representations, such as the one proposed by Diwersy et al. (2014) for the study of variation in translation and diatopic variation of French texts, and the fully lexicalized representation, common in most text classification

---

<sup>1</sup>The decision to split the corpus into sentences was motivated by the amount of texts available in the VARTRA corpus. Splitting the corpus into sentences generated enough data points for text classification and made the task more challenging.

<sup>2</sup>Note that in this paper we make a clear distinction between BoW and unigrams. The BoW models used in this paper do not comprise any smoothing method, whereas the  $n$ -gram models are calculated using Laplace smoothing.

experiments, which uses all words in text without any substitution. This representation minimizes topic variation. Previous studies have shown that named entities significantly influence the performance of text classification systems (Zampieri et al., 2013; Goutte et al., 2016). We used this representation in our previous experiments that we describe in Section 2.2 above.

In (1-c), we use a fully delexicalized representation representing texts only using the POS annotation available at the VARTRA corpus. Zampieri et al. (2013) show that classification experiments using POS and morphological information as features can not only be linguistically informative, but also achieve good performance in discriminating between texts written in different Spanish varieties. Therefore, we use this representation to test whether this is also true for translated texts.

The underlying tagset used in TreeTagger is “Stuttgart/Tbinger Tagsets” (STTS)<sup>3</sup>, one of the commonly used tagsets for German. In Table 1, we illustrate a segment from the corpus with an explanation of selected tags.

**Table 1:** Illustration of the STTS-tagged corpus segment

| <b>word</b>       | <b>POS</b> | <b>category</b>     |
|-------------------|------------|---------------------|
| Die               | ART        | article             |
| weltweiten        | ADJA       | adjective           |
| Herausforderungen | NN         | common noun         |
| im                | APPRART    | preposition+article |
| Bereich           | NN         | common noun         |
| der               | ART        | article             |
| Energiesicherheit | NN         | common noun         |
| erfordern         | VVFIN      | full finite verb    |
| über              | APPR       | preposition         |
| einen             | ART        | article             |
| Zeitraum          | NN         | common noun         |

The decision to use these features was motivated by our goal of investigating translation variation influenced by both genre and method, and our aim to obtain a classification method that could perform well on different corpora by capturing structural differences between these translation varieties.

### 3.2 Algorithm

In our experiments we use a Bayesian learning algorithm similar to Naive Bayes entitled Likelihood Estimation (LE) and previously used for language identification in Zampieri and Gebre (2012, 2014). Just like Naive Bayes classifiers, LE works based on an independence assumption that the presence of a particular feature of a class is not related to the presence of any other feature. The independence assumption makes the algorithm extremely fast and a good fit for text classification tasks. Bayesian classifiers are inspired by Bayes theorem represented by the following equation:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (1)$$

Where  $P(A|B)$  is a conditional probability of  $A$  given  $B$ . Using the notation by Kibriya et al. (2004), a Naive Bayes classifier computes class probabilities for a given document and a set of

<sup>3</sup><http://www.ims.uni-stuttgart.de/forschung/ressourcen/lexika/TagSets/stts-table.html>

classes  $C$ . It assigns each document  $t_i$  to the class with the highest probability  $P(c|t_i)$ .

$$P(c|t_i) = \frac{P(t_i|c)P(c)}{P(t_i)} \quad (2)$$

LE calculates a likelihood function on smoothed n-gram language models. Smoothing is carried out using the Laplace smoothing calculated as follows:

$$P_{lap}(w_1...w_n) = \frac{C(w_1...w_n) + 1}{N + B} \quad (3)$$

The language models can contain characters and words (e.g. bigrams and trigrams), linguistically motivated features such as parts-of-speech (POS) or morphological categories such as the one used in Zampieri et al. (2013) for the study of diatopic variation. In this paper LE is used with POS tags as features.

Models are first calculated for each particular class in the dataset. Subsequently LE calculates the probability of a document belonging to a given class. In our case classes are represented either by genres or method of translation. The function that calculates the probability of a document given a class, represented by  $L$  (language model) is the following:

$$P(L|text) = \arg \max_L \sum_{i=1}^N \log P(n_i|L) + \log P(L) \quad (4)$$

Where  $N$  is the number of  $n$ -grams in the test text. The language model  $L$  with the highest probability determines the predicted class of each document.

## 4 Classification Results

In this section, we present the results obtained in various classification experiments using a Bayesian classifier. To evaluate the performance of the classifiers we used standard metrics in text classification such as precision, recall, f-measure, and accuracy. The linguistic analysis and discussion of the most important differences between both method and genre variation will be presented later in Section 5.

### 4.1 Translation Methods: Human vs. Machine

In this first experiment we investigate differences between translation methods. The VARTRA corpus divides translation methods into five categories, three representing automatic methods and two containing translations produced by humans. In Zampieri and Lapshinova-Koltunski (2015) we trained a classifier to discriminate between these five methods and we observed that variation was more prominent when comparing human vs. machine translations. For this reason we unify PT1 and PT2 into one class and RBMT, SMT1, and SMT2 into the other.

We represent texts using the POS tags as features as presented earlier in this chapter. We use a total of 600 texts for each class split in 400 documents for training and 200 for testing. Results are presented in terms of precision, recall, and f-measure in Table 2. The baseline is 50% accuracy.

In all three settings, the model performs above the expected baseline of 50.0% f-measure. The best performance is obtained using a POS trigram model (62.5% f-measure and precision).

**Table 2:** N-grams: Human x Machine

| Features | Precision | Recall | F-Measure |
|----------|-----------|--------|-----------|
| bigrams  | 60.70%    | 60.51% | 60.61%    |
| trigrams | 62.50%    | 62.50% | 62.50%    |
| 4-grams  | 57.84%    | 57.25% | 57.54%    |

## 4.2 Genres

In this section, we train a model to automatically distinguish between seven different genres represented in our dataset. For the sake of clarity we list here the genres contained in the VARTRA corpus: political essays (ESS), fictional texts (FIC), instruction manuals (INS), popular-scientific articles (POP), letters of share-holders (SHA), prepared political speeches (SPE), and touristic leaflets (TOU). All experiments in this section are binary classification settings in which the classifier is trained to discriminate between two genres at a time. The baseline for each setting is therefore 50% accuracy.

We again use POS tags as features as described in 4.1 arranged in bigrams, trigrams, and 4-grams. We use a total of 500 texts for each class split in 300 documents for training and 200 for testing. We evaluate the performance of our method in terms of accuracy and present results in Tables 3, 4, and 5.

In Table 3 using POS bigrams we observed that the best results were obtained when discriminating instruction manuals from fictional texts, 81.25% accuracy. The worst results were obtained between speech and essays, 61.25% accuracy.

Corroborating with the findings of the previous section, we observed that for genres the overall best results are obtained when using POS trigrams. The model is able to discriminate between tourism leaflets and fictional texts with impressive results of 84% accuracy.

**Table 3:** Genres Classification in Translation in Binary Settings: POS Bigrams

| Classes    | ESS | FIC    | INS    | POP    | SHA    | SPE    | TOU    |
|------------|-----|--------|--------|--------|--------|--------|--------|
| <b>ESS</b> | -   | 78.00% | 75.75% | 65.00% | 66.25% | 61.25% | 71.75% |
| <b>FIC</b> | -   | -      | 81.25% | 79.50% | 77.75% | 74.50% | 80.50% |
| <b>INS</b> | -   | -      | -      | 74.50% | 75.50% | 79.00% | 74.25% |
| <b>POP</b> | -   | -      | -      | -      | 68.25% | 67.50% | 69.00% |
| <b>SHA</b> | -   | -      | -      | -      | -      | 66.00% | 69.25% |
| <b>SPE</b> | -   | -      | -      | -      | -      | -      | 72.75% |

We observe that in the vast majority of settings presented in Table 4, results obtained using POS trigrams were higher than those using POS bigrams.

Finally, in our last setting using POS 4-grams, we observed that this set of features do not achieve the best results in distinguishing between genres. For this reason, we preset feature analysis on POS trigrams which were the features that obtained the best results in this section.

## 5 Feature Analysis

Text classification allows us to not only measure how well certain subcorpora (e.g. human and machine translations) are distinguished from each other, but also which individual features

**Table 4:** Genres Classification in Translation in Binary Settings: POS Trigrams

| Classes    | ESS | FIC    | INS    | POP    | SHA    | SPE    | TOU    |
|------------|-----|--------|--------|--------|--------|--------|--------|
| <b>ESS</b> | -   | 76.25% | 74.00% | 66.25% | 71.00% | 61.25% | 72.50% |
| <b>FIC</b> | -   | -      | 76.50% | 76.50% | 77.75% | 76.25% | 84.00% |
| <b>INS</b> | -   | -      | -      | 74.25% | 76.75% | 76.75% | 74.75% |
| <b>POP</b> | -   | -      | -      | -      | 71.00% | 67.75% | 71.75% |
| <b>SHA</b> | -   | -      | -      | -      | -      | 65.00% | 68.00% |
| <b>SPE</b> | -   | -      | -      | -      | -      | -      | 71.25% |

**Table 5:** Genres Classification in Translation in Binary Settings: POS 4-grams

| Classes    | ESS | FIC    | INS    | POP    | SHA    | SPE    | TOU    |
|------------|-----|--------|--------|--------|--------|--------|--------|
| <b>ESS</b> | -   | 73.75% | 75.25% | 69.50% | 67.00% | 65.25% | 72.50% |
| <b>FIC</b> | -   | -      | 70.20% | 75.00% | 78.25% | 76.00% | 79.25% |
| <b>INS</b> | -   | -      | -      | 68.50% | 70.50% | 74.50% | 74.50% |
| <b>POP</b> | -   | -      | -      | -      | 71.25% | 68.75% | 70.75% |
| <b>SHA</b> | -   | -      | -      | -      | -      | 67.00% | 66.25% |
| <b>SPE</b> | -   | -      | -      | -      | -      | -      | 73.00% |

contribute to this distinction. Therefore, we analyse the output features resulting from the classification in this section. The main aim here is to identify the most informative features from the delexicalized n-grams in our experiments and to interpret them in terms of linguistic categories. This step is manual and carried out by looking through the most informative features and thus discriminative for certain genres and translation methods in our translation data.

Delexicalized trigrams consist of a sequence of words and placeholders, e.g. (1) *ART NN VMFIN* (2) *. KON PPER*, etc. Intuitively, we try to recognise more categories on a more abstract level of linguistic description, i.e. category of modality expressed through modal verbs, discourse-building devices, such as discourse markers and coreference and others for the given trigrams. Thus, example (1) represents a finite clause containing a full nominal phrase, and example (2) represents a pattern related to the level of discourse: a Connector at sentence start followed by a personal pronoun that likely refers to something previously mentioned in the text. We decide for the evaluation of trigrams, as the performance of trigram models achieved the best results in both classification tasks.

## 5.1 Translation Methods

The classification results for the distinction of translation methods outputs two lists of features: (1) the list of the features specific for human translations and (2) the list of features specific for machine translation. We analyse up to the first 20 features per translation method, summarising our observations in Table 6.

As seen from the lists, both translation methods have similar types of features that differentiate them from each other: they can be classified in terms of more abstract linguistic categories, such as discourse and modality. And some of them concern the preferred typology of phrases which can be related to the style of writing: nominal vs. verbal. The differences between the features discriminating between human and machine translations are visible on a more fine-grained level. For example, if we take into account morpho-syntactic preferences

**Table 6:** Features discriminating between human and machine translations

|           | <b>human</b>  | <b>machine</b>   |
|-----------|---|--|
| Discourse | conjunct at sent.start followed by a personal referring expression<br><br>coreference at sentence start (demonstrative)<br>coreference and negation | conjunct at sent.start followed by a full NP<br><br>adverbial at clause start<br>conjunct at clause start<br>full and pronominal referring expressions at clause start |
| Modality  | -   | +  |
| Phrases   | NP connected to other phrases<br>conj linking NP<br>NP with named entities  | V2 phrases followed by a definite NP<br>conj linking VP<br>NP describing location<br>predicative adjectives<br>locative prepositional phrases                          |
| Verbs     | verb subcategorisation patterns   | apposition<br><br>V2 structure<br>imperative constructions   |

of discourse phenomena, we observe differences in the position of cohesive triggers: in machine translations, several patterns contain a punctuation mark in the first position, which means that the observed pattern often represents sentence and clause start. Human translations rather show preferences for more sentence-starting devices. Example (2) reveals the possible reasons for this observation: human translators tend to split longer sentences into several ones (2-b), whereas a machine translation system keeps the structure (2-c) as it was in the source (2-a).

- (2) a. *He used it to modernise the castle but he must have skimped on the kitchen, since 1639 it fell into the sea and carried away the cooks and all their pots.*  
 b. *Er erbeutete das vom Schiff mitgeführte Gold und benutzte es zur Modernisierung des Schlosses. Allerdings scheint er beim Ausbau der Küche etwas geknausert zu haben, denn dieser Teil des Schlosses rutschte im Jahre 1639 ins Meer ab. Die Küche wurden mitsamt Tpfen weggespllt.*  
 c. *Er benutzte es, um das Schloss zu modernisieren, aber er muss auf dem Kchentisch gespart haben, seit 1639 ins Meer fiel und trugen die Küche und alle ihre Töpfe..*

Another difference that is clearly seen from the examples in the corpus data is the preference of human translations for conjuncts, whereas for machine translations, subjuncts and adverbials seem to be more typical. This is seen in the illustration in example (3).

- (3) a. Human: *Und er wandte sich von der goldenen Dame ab und htte gar zu gern die silberne genommen...*  
 b. Machine: *Darber hinaus muss der Bohrvorgang verfeinert, so dass die tieferen*

*Aquiferen nicht durch Arsen-Lager Wasser rann von den flachen Grundwasserleiter durch die Bohrungen selbst vergiftet werden.*

The multiword conjunction *so dass* is very frequent in machine translations in our data: the total of 108 occurrences as compared to 6 occurrences in human translations. A closer look at the data reveals that all the occurrences of *so dass* are found in machine translations produced with a statistical system trained on a large amount of data.

Human translations have also a number of features related to modality expressed via modal verbs. Our additional quantitative analysis of the modal verb distributions shows that, in general, human translations contain slightly more modal verbs than the machine ones: 16.49% vs. 15.72% out of all sentences in our data. This means that although modal verbs are more frequent in human translations, linguistic patterns with modals are more distinctive for machine-translated texts.

There is a difference in adjectival constructions. Predicative adjectives that turn to be discriminative for machine translations are also more frequent in this translation variety than in the human one (24.86% vs. 23.72%).

- (4)
- a. *The roads are excellent, with miles of motorway and dual carriageway...*
  - b. *Es gibt ausgezeichnete Straen, davon ungefahr 112 km Autobahn und noch weit mehr Kilometer mit zweispurigen Fahrbahnen.*
  - c. *Die Straen sind sehr gut, mit Meilen von Autobahn und Schnellstrae...*

As seen from example (4), the machine-translated sentence (4-c) is closer to the source one (4-a) in terms of the predicative vs. attributive usage of the adjective (*are excellent – sind sehr gut*). At the same time, human translation (4-b) is closer to the source in terms of lexical choice (*excellent – ausgezeichnet*)

Another interesting difference is the prevalence of nominal structures in human translations as opposed to machine ones (46% vs. 24%) in the analysed trigram patterns. At the same time, machine-translated texts in our corpus contain more verbal phrases under their discriminative features (50% vs. 34%). We believe that this tendency is observed due to the shining though effect (Teich, 2003): German-English contrastive analyses, e.g. the one by Steiner (2012) show that German has a preference for nominal structures, whereas English is more verbal. So, if a similar preference is observed in English-to-German translations, this could be interpreted as a phenomenon of shining though.

## 5.2 Genres

Using the same strategy, we generate a list of features discriminating genre pairs. For the sake of space, we will concentrate on the analysis of two genres only: fictional texts and political speech. For the first one, we achieve the best results in the trigram classification, whereas the classification results seem to be the worst for the second.

### 5.2.1 Fictional texts

We analyse six lists of patterns that turn out to be discriminative for fiction in the six classification tasks involving fictional texts: (1) fiction vs. political essays, (2) fiction vs. instruction manuals, (3) fiction vs. popular-scientific texts, (4) fiction vs. letters-to-shareholders, (5) fiction vs. political speeches and (6) fiction vs. tourism leaflets.

First, we sort the patterns that occur more than once in the lists. The data contains several types of patterns: (a) informative in five classification tasks, i.e. member of five lists; (b) member of four lists; (c) member of three lists; (d) member of two lists; (e) member of one list

– informative on one particular classification task, e.g. fictional texts vs. instruction manuals. The most frequent patterns are considered to be the most specific ones for fictional texts, as they were informative in several classification tasks, i.e. in discriminating fictional texts from several other genres. For instance, the pattern , *ADV KOUS*, e.g. , *so dass* / , *noch bevor* / , *auch wenn* (a discourse marker that links a subordinate clause), is informative in the first five classification tasks (fiction against essays, instructions, letters-to-shareholders, political speeches and tourism texts).

Table 7 illustrates the distribution of fiction-discriminative patters. The final list comprises 170 patterns.

**Table 7:** Feature lists discriminating between fiction and other genres

|       | list membership | types |
|-------|-----------------|-------|
| (a)   | 5 lists         | 5     |
| (b)   | 4 lists         | 5     |
| (c)   | 3 lists         | 23    |
| (d)   | 2 lists         | 49    |
| (e)   | 1 list          | 88    |
| total |                 | 170   |

In the following, we analyse those that occur in most tasks (5 lists) and a subset of those that occur in one classification task only (1 list). Table 8 illustrates the five language patterns that turned to be the most informative in most classification tasks for the discrimination of fictional texts. The last column of the table provides the information on the genre, for which the result is not valid. The first three patterns are not discriminative for fictional texts, when classified against popular-scientific articles, the fourth pattern is not informative when instructional manuals are involved. And the last one is not discriminative for fictional text, when they are classified vs. letters-to-shareholders.

**Table 8:** Features informative for fiction in most classification tasks for fiction

| pattern                | example   | excluding |
|------------------------|---|-----------|
| \$, <i>ADV KOUS</i>    | , <i>so dass</i> / , <i>noch bevor</i> / , <i>auch wenn</i>         | POP       |
| <i>ADV KOUS PPER</i>   | <i>so dass er</i> / <i>auch wenn sie</i>                            | POP       |
| <i>ADV VVFIN \$</i> .  | <i>gern wiederholen ./ nebeneinander stellen</i>                    | POP       |
| \$( <i>PPER VMFIN</i>  | ( <i>sie knnen</i> / ( <i>wir wollen</i> / ( <i>es drfte</i>        | INS       |
| <i>PPER VMFIN PPER</i> | <i>sie knnen ihn</i> / <i>wir wollen uns</i> / <i>ich mchte ihm</i> | SHA       |

Most language patterns in Table 8 are discourse-related devices, i.e. discourse markers expressing conjunctive relations (*so dass*, *auch wenn*) or pronouns triggering cohesive reference (*er*, *sie*, *es*). The last two patterns contain also modal verbs which can be interpreted in terms of sentence or text modality.

**Table 9:** Distribution of \$, *ADV KOUS* across genres

| genre | TOU  | SHA  | SPE  | POP  | ESS  | INS  | FIC  |
|-------|------|------|------|------|------|------|------|
| freq  | 0.10 | 0.11 | 0.14 | 0.23 | 0.29 | 0.35 | 0.41 |

The discriminative power of features does not necessarily imply a high frequency of a particular pattern in fictional texts. Nevertheless, the distribution of the first pattern across genres in our data shows that this trigram is more frequent in fiction than in the other genres, see Table

9 (the number are normalised per 1000 per total number of trigrams). However, the numbers in the table do not reveal the reasons for this pattern not being discriminative in the classification task for fiction vs. popular-scientific texts.

In the last step, we analyse the list of language patterns discriminating fiction in one particular task that includes 88 trigrams. In Table 10, we present a summary of these patterns describing them in terms of more general linguistic categories, e.g. specific phrases or functions.

**Table 10:** Features informative for fiction in one classification task only

| feature                  | example pattern | language example                                      |
|--------------------------|-----------------|---|
| phrases with adjectives  | ADJA KON ADJA   | <i>heller und dunkler / ernste und vielschichtige</i> |
|                          | CARD ADJA NN    | <i>zwei junge Mnner / drei wunderschne Damen</i>      |
|                          | PPOSAT NN ADJD  | <i>ihre Hrner steil / ihre Lieder schwer</i>          |
| phrases with adverbs     | ADV VAFIN PPER  | <i>so ist es / dann hast du</i>                       |
|                          | ADV VVFIN ,     | <i>anders war , / unterwegs ist ,</i>                 |
|                          | ADV VVFIN PPER  | <i>jetzt kriegt er / dann grunzt er</i>               |
| coreference via pronouns | PPER ADV VVIN   | <i>sie weiter sprechen / dir nur sagen</i>            |
|                          | PPER VAFIN PPER | <i>sie hatte sie / ich habe sie</i>                   |
|                          | \$. PPER VMFIN  | <i>. Sie macht / Sie beginnen</i>                     |
| discourse markers        | KON PPER VMFIN  | <i>Aber sie mchte / und er konnte</i>                 |
|                          | KOUS PPOSAT NN  | <i>dass meine Mutter / ob ihr Brief</i>               |
|                          | VAFIN \$. KON   | <i>wrde. Aber / hatte . Und</i>                       |

As seen from the table, the patterns specific for fiction include adjective and adverb modification and elements that contribute to structuring discourse in a text. The latter are especially specific for narrative texts, which our fictional texts belong to. These observations coincide with the results of other empirical analyses on genres, e.g. those obtained by Neumann (2013). The author also points to personal pronouns and predicative adjectives as indicators of narration and casual style which are specific for fictional texts.

### 5.2.2 Political speeches

We proceed with the analysis of political speeches, which turned to be the hardest genre to identify in the classification with trigrams. The same analysis steps as we used for fictional texts are applied here.

First, we summarise the patterns that occur more than once in the lists. The political speeches contain less types of patterns than the fictional texts. An informative pattern can be a member of maximum four classification tasks only (for fictional texts, we also had five). So, we have four lists of patterns: (a) informative in four classification tasks, i.e. member of our lists; (b) member of three lists; (c) member of two lists; (d) member of one list – informative on one particular classification task, e.g. political speeches vs. fictional texts. As in the previous case with fictional texts, we consider the most frequent patterns to be the most specific ones for political speeches, as they contribute to the distinction of speeches from several other genres.

For instance, the pattern *PTKVZ \$. ART*, e.g. *vor . Das / bei . Die / weiter . Das* (sentence end followed by a sentence starting with a nominal phrase with an article), is informative in the following four classification tasks: political speeches vs. fiction, instructions, letters-to-shareholders, and tourism texts. Table 11 illustrates the distribution of speech-discriminative patterns.

The total number of patterns is smaller than that of fictional texts (156 vs. 170 pattern types). We believe that the more distinctive features a genre has, the more distinctive it is from other genres, and thus can be easily identified with automatic classification techniques.

**Table 11:** Features discriminating between political speeches and other genres

| list membership | types   |     |
|-----------------|---------|-----|
| (a)             | 4 lists | 3   |
| (b)             | 3 lists | 18  |
| (c)             | 2 lists | 49  |
| (d)             | 1 list  | 86  |
| total           |         | 156 |

Now we will have a closer look at the patterns that are informative in most tasks and some of those that are discriminative for political speeches in one classification task only (1 list).

**Table 12:** Features informative for political speeches in most classification tasks for political speeches

| pattern        | example  | excluding |
|----------------|--|-----------|
| PTKVZ \$. ART  | <i>vor . Das / bei . Die / weiter . Das</i>      | ESS, POP  |
| VVFIN PPER ADV | <i>arbeiten wir zurzeit / auch wenn sie</i>      | POP, TOU  |
| VVPP VAINF \$. | <i>gern wiederholen ./ nebeneinander stellen</i> | ESS, POP  |

Table 12 illustrates the three language patterns that turned to be the most informative for the discrimination of political speeches. The last column of the table provides the information on the two genres, for which the result is not valid. All the three patterns are not discriminative for political speeches, when classified against popular-scientific articles. The first and the last patterns are not informative, when political speeches are distinguished from political essays. And the second pattern is not discriminative in the classification against tourism texts.

In the last step, we analyse the list of language patterns discriminating political speeches in one particular task that includes 86 trigrams. In Table 13, we present a summary of these patterns describing them in terms of more general linguistic categories, e.g. specific phrases or functions.

As seen from the table, the patterns specific for political speeches include phrases with adjective, infinitive phrases and the elements that contribute to structuring discourse in a text. From the first sight, there are types of patterns that are similar to those analysed for the fictional texts. However, our qualitative analysis reveals substantial differences. The main differences is caused by the difference in the register orientation: in a narration (most of the fictional texts in the data), there is more orientation towards the content, whereas in political speeches, we observe a clear orientation of the author towards the audience. It is especially prominent in coreference-related features. Fictional texts utilise a great number of third person pronouns, whereas political speeches have much more first and second person pronouns, see example (5).

- (5) *Was passierte mit den Kindern? Wollen Sie sagen, dass Sie eine Milliarde Dollar ausgegeben haben und nicht wissen... [What happened to the children? Do you mean that you spent a billion dollars and you don't know...]*

**Table 13:** Features informative for political speeches in one classification task only

| feature                  | example pattern  | language example  |
|--------------------------|--|---|
| phrases with adjectives  | ADJA NN KOKOM<br>ADJD APPR ART<br>ADV ADJA NN          | <i>politische Themen wie /<br/>weltweite Probleme wie<br/>wichtig fr die / mglich fr die<br/>sehr geehrte Mitglieder<br/>/ ebenfalls sprunghafte<br/>Fortschritte</i>       |
| infinitive phrases       | \$, PTKZU VVINF<br>ADJD PTKZU VVINF<br>PTKZU VVINF KON | <i>, zu sprechen / , zu bekmpfen<br/>/ , zu besprechen<br/>schwer zu entscheiden /<br/>richtig zu verteidigen<br/>zu bernehmen und / zu unter-<br/>sttzen und</i>           |
| coreference via pronouns | PPER PPOSAT NN<br>PPER VVINF \$,<br>\$, VMFIN PPER     | <i>wir unsere Ziele / wir unseren<br/>Feinden / ich Ihre Fragen<br/>wir prfen , / Sie antworten , /<br/>wir erreichen ,<br/>, mssen wir / , mchte ich / ,<br/>knnen wir</i> |
| discourse markers        | \$. ADV VAFIN<br>\$. ADV VVFIN<br>VAFIN \$. KON        | <i>. Bisher haben / . Natrlich ist<br/>. Mglicherweise brauchen / :<br/>Erstens versucht<br/>wrde. Aber / hatte . Und</i>   |

This again, coincides to what was previously observed in register/genre-related analysis (Biber et al., 1999; Neumann, 2013).

## 6 Conclusion and Discussion

This paper is, to our knowledge, the first attempt to use text classification techniques to discriminate methods and genres in translations using fully delexicalized text representations and to identify their specific features and relevant systemic differences in a single study. We report results of up to 62.50% f-measure in distinguishing between human and machine translations using POS trigrams and 81.25% accuracy in discriminating between speech and essays.

The results obtained using POS tags as features was surprisingly higher than those obtained using (semi-)delexicalized representations presented in Zampieri and Lapshinova-Koltunski (2015). This seems to indicate relevant systemic differences across genres and methods of translation that algorithms relying on (morpho)-syntactic features are able to recognize.

At the same time, the results show that it is much harder to differentiate between translation methods than between different genres, even if fully delexicalized features are used. This confirms the results by Lapshinova-Koltunski (2017) which shows that if we compare the influence of the genre dimensions in translation variation is much stronger than that of translation method.

The results of our analysis can find application in both human and machine translation. In

the first case, they deliver valuable knowledge on the translation product, which is influenced by the methods used in the process and the context of text production expressed by the genre. In case of machine translation, the results will provide a method to automatically identify genres in translation data thus helping to separate out-of-genre data from a training corpus.

The resulting lists of features can also be beneficial for automatic genre classification or human vs. machine distinction tasks. The knowledge on the differences between genres that these features deliver can also help to understand main differences between texts translated by humans and with machine translation systems. This information is especially valuable for translator training. Nowadays, translator training includes courses on post-editing technologies, since the application of such technologies has increased in translation industry recently. Translators need to know where the main problems (not necessarily errors) of machine-translated texts lie and what differs them from the texts by professional translators. This knowledge increases productivity in translation process.

In our future work, we want to perform a classification task for translation method within each genres. We assume that the differences between texts that differ in translation methods can be identified better, if classification is carried out within on genre only. Moreover, this will provide us with the information on how human and machine translations differ, if one particular genre is involved.

## References

- Babych, B., Hartley, A., and Sharoff, S. (2004). Modelling legitimate translation variation for automatic evaluation of mt quality. In *Proceedings of LREC-2004*, volume Vol. 3.
- Baker, M. (1993). Corpus linguistics and translation studies: Implications and applications. *Text and technology: In honour of John Sinclair*, 233:250.
- Baroni, M. and Bernardini, S. (2006). A new approach to the study of translationese: Machine-learning the difference between original and translated text. *Literary and Linguistic Computing*, 21(3):259–274.
- Biber, D. (1995). *Dimensions of Register Variation. A Cross Linguistic Comparison*. Cambridge University Press, Cambridge.
- Biber, D., Johansson, S., Leech, G., Conrad, S., and Finegan, E. (1999). *Longman Grammar of Spoken and Written English*. Longman, Harlow.
- Bojar, O., Buck, C., Federmann, C., Haddow, B., Koehn, P., Monz, C., Post, M., and Specia, L., editors (2014). *Proceedings of the Ninth Workshop on Statistical Machine Translation*. Association for Computational Linguistics, Baltimore, Maryland, USA.
- Callison-Burch, C., Osborne, M., and Koehn, P. (2006). Re-Evaluation the Role of Bleu in Machine Translation Research. In *Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics*, pages 249–256.
- Comelles, E. and Atserias, J. (2014). VERTa Participation in the WMT14 Metrics Task. In *Proceedings of the 9th Workshop on Statistical Machine Translation (WMT)*, pages 368–375, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Corston-Oliver, S., Gamon, M., and Brockett, C. (2001). A machine learning approach to the automatic evaluation of machine translation. In *Proceedings of the 3th Annual Meeting on Association for Computational Linguistics*, pages 148–155.

- De Sutter, G., Delaere, I., and Plevoets, K. (2012). Lexical lectometry in corpus-based translation studies: Combining profile-based correspondence analysis and logistic regression modeling. In *Quantitative Methods in Corpus-based Translation Studies: a Practical Guide to Descriptive Translation Research*, volume 51, pages 325–345. John Benjamins Publishing Company, Amsterdam, The Netherlands.
- Delaere, I. and De Sutter, G. (2013). Applying a multidimensional, register-sensitive approach to visualize normalization in translated and non-translated Dutch. *Belgian Journal of Linguistics*, 27:43–60.
- Denkowski, M. and Lavie, A. (2014). Meteor Universal: Language Specific Translation Evaluation for Any Target Language. In *Proceedings of the Workshop on Statistical Machine Translation (WMT)*.
- Diwersy, S., Evert, S., and Neumann, S. (2014). A semi-supervised multivariate approach to the study of language variation. *Linguistic Variation in Text and Speech, within and across Languages*.
- Doddington, G. (2002). Automatic Evaluation of Machine Translation Quality Using N-gram Co-Occurrence Statistics. In *Proceedings of the 2nd International Conference on Human Language Technologies (HLT)*, pages 138–145.
- El-Haj, M., Rayson, P., and Hall, D. (2014). Language independent evaluation of translation style and consistency: Comparing human and machine translations of camus novel “the stranger”. In Sojka, P., Horák, A., Kopeček, I., and Pala, K., editors, *Proceedings of the 17th International Conference TSD 2014*, volume 8655 of *Lecture Notes in Computer Science*, Brno, Czech Republic. Springer.
- Fishel, M., Sennrich, R., Popovic, M., and Bojar, O. (2012). Terrorcat: a translation error categorization-based mt quality metric. In *7th Workshop on Statistical Machine Translation*.
- Gebre, B. G., Zampieri, M., Wittenburg, P., and Heskens, T. (2013). Improving native language identification with tf-idf weighting. In *Proceedings of the 8th NAACL Workshop on Innovative Use of NLP for Building Educational Applications (BEA8)*, Atlanta, USA.
- González, M., Barrón-Cedeño, A., and Màrquez, L. (2014). IPA and STOUT: Leveraging Linguistic and Source-based Features for Machine Translation Evaluation. In *Proceedings of the 9th Workshop on Statistical Machine Translation (WMT)*, pages 394–401, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Goutte, C., Léger, S., Malmasi, S., and Zampieri, M. (2016). Discriminating Similar Languages: Evaluations and Explorations. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC)*, pages 1800–1807, Portoroz, Slovenia.
- Gupta, R., Orăsan, C., and van Genabith, J. (2015). ReVal: A Simple and Effective Machine Translation Evaluation Metric Based on Recurrent Neural Networks. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Lisbon, Portugal.
- Halliday, M. and Hasan, R. (1989). *Language, context and text: Aspects of language in a social-semiotic perspective*. Oxford University Press, Oxford.

- Hansen-Schirra, S., Neumann, S., and Steiner, E. (2012). *Cross-linguistic Corpora for the Study of Translations. Insights from the Language Pair English-German*. de Gruyter, Berlin, New York.
- House, J. (2014). *Translation Quality Assessment. Past and Present*. Routledge.
- Irvine, A. and Callison-Burch, C. (2014). Using comparable corpora to adapt mt models to new domains. In *Proceedings of the ACL Workshop on Statistical Machine Translation (WMT)*.
- Irvine, A., Morgan, J., Carpuat, M., III, H. D., and Munteanu, D. S. (2013). Measuring machine translation errors in new domains. *TACL*, 1:429–440.
- Jenset, G. B. and Hareide, L. (2013). A multidimensional approach to aligned sentences in translated text. *Bergen Language and Linguistic Studies*, 3:195–210.
- Jenset, G. B. and McGillivray, B. (2012). Multivariate analyses of affix productivity in translated english. In Oakes, M. P. and Ji, M., editors, *Quantitative Methods in Corpus-Based Translation Studies*, pages 301–324. John Benjamins.
- Kibriya, A., Frank, E., Pfahringer, B., and Holmes, G. (2004). Multinomial naive bayes for text categorization revisited. In *Proceedings of the Australian Conference on Artificial Intelligence*, pages 488–499.
- Kruger, H. and van Rooy, B. (2012). Register and the Features of Translated Language. *Across Languages and Cultures*, 13(1):33–65.
- Kunz, K., Degaetano-Ortlieb, S., Lapshinova-Koltunski, E., Menzel, K., and Steiner, E. (2017). Gecco – an empirically-based comparison of english-german cohesion. In De Sutter, G., Delaere, I., and Lefer, M.-A., editors, *New Ways of Analysing Translational Behaviour in Corpus-Based Translation Studies*. Mouton de Gruyter. TILSM series.
- Lapshinova-Koltunski, E. (2013). VARTRA: A comparable corpus for analysis of translation variation. In *Proceedings of the Sixth Workshop on Building and Using Comparable Corpora*, pages 77–86, Sofia, Bulgaria. Association for Computational Linguistics.
- Lapshinova-Koltunski, E. (2017). Exploratory analysis of dimensions influencing variation in translation: The case of text register and translation method. In De Sutter, G., Delaere, I., and Lefer, M.-A., editors, *New Ways of Analysing Translational Behaviour in Corpus-Based Translation Studies*. Mouton de Gruyter. TILSM series.
- Lapshinova-Koltunski, E. and Vela, M. (2015). Measuring registerness in human and machine translation: A text classification approach. In *Proceedings of the Second Workshop on Discourse in Machine Translation*, pages 122–131, Lisbon, Portugal. Association for Computational Linguistics.
- Lee, D. Y. (2001). Genres, registers, text types, domains and styles: clarifying the concepts and navigating a path through the bnc jungle. *Technology*, 5:37–72.
- Malmasi, S., Dras, M., and Zampieri, M. (2016). LTG at SemEval-2016 Task 11: Complex Word Identification with Classifier Ensembles. In *Proceedings of SemEval*.
- Medlock, B. (2008). Investigating classification for natural language processing tasks. Technical report, University of Cambridge - Computer Laboratory.

- Neumann, S. (2013). *Contrastive Register Variation. A Quantitative Approach to the Comparison of English and German*. De Gruyter Mouton, Berlin, Boston.
- Niculae, V., Zampieri, M., Dinu, L. P., and Ciobanu, A. M. (2014). Temporal text ranking and automatic dating of texts. In *14th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2014)*. Association for Computational Linguistics.
- Papineni, K., Roukus, S., Ward, T., and Zhu, W.-J. (2002). BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318.
- Popovic, M. and Ney, H. (2011). Towards automatic error analysis of machine translation output. *Computational Linguistics*, 37(4):657–688.
- Santini, M., Mehler, A., and Sharoff, S. (2010). Riding the rough waves of genre on the web. In Mehler, A., Sharoff, S., and Santini, M., editors, *Genres on the Web: Computational Models and Empirical Studies*, pages 3–30. Springer.
- Schmid, H. (1994). Probabilistic part-of-speech tagging using decision trees. In *Proceedings of International Conference on New Methods in Language Processing*, Manchester, UK.
- Stanojević, M. and Simaán, K. (2014). BEER: BETter Evaluation as Ranking. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 414–419, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Steiner, E. (2004). *Translated Texts. Properties, Variants, Evaluations*. Peter Lang Verlag, Frankfurt/M.
- Steiner, E. (2012). A characterization of the resource based on shallow statistics. In Hansen-Schirra, S., Neumann, S., and Steiner, E., editors, *Cross-linguistic Corpora for the Study of Translations. Insights from the Language Pair English-German*. Mouton de Gruyter, Berlin, New York.
- Teich, E. (2003). *Cross-Linguistic Variation in System und Text. A Methodology for the Investigation of Translations and Comparable Texts*. Mouton de Gruyter, Berlin.
- Vela, M. and van Genabith, J. (2015). Re-assessing the WMT2013 Human Evaluation with Professional Translators Trainees. In *Proceedings of the 18th Annual Conference of the European Association for Machine Translation (EAMT)*.
- Volansky, V., Ordan, N., and Wintner, S. (2011). More human or more translated? original texts vs. human and machine translations. In *Proceedings of the 11th Bar-Ilan Symposium on the Foundations of AI With ISCOL*.
- Wu, H., Wang, H., and Zong, C. (2008). Domain adaptation for statistical machine translation with domain dictionary and monolingual corpora. In Scott, D. and Uszkoreit, H., editors, *Proceedings of the 22nd International Conference on Computational Linguistics (COLING-2008)*, pages 993–1000, Manchester, UK.
- Zampieri, M. and Gebre, B. G. (2012). Automatic identification of language varieties: The case of Portuguese. In *Proceedings of KONVENS2012*, pages 233–237, Vienna, Austria.
- Zampieri, M. and Gebre, B. G. (2014). Varclass: An open source language identification tool for language varieties. In *Proceedings of Language Resources and Evaluation (LREC)*, Reykjavik, Iceland.

- Zampieri, M., Gebre, B. G., and Diwersy, S. (2013). N-gram language models and POS distribution for the identification of Spanish varieties. In *Proceedings of TALN2013*, pages 580–587, Sable d’Olonne, France.
- Zampieri, M. and Lapshinova-Koltunski, E. (2015). Investigating genre and method variation in translation using text classification. In Sojka, P., Horák, A., Kopecek, I., and Pala, K., editors, *Text, Speech and Dialogue - 18th International Conference, TSD 2015, Plzen, Czech Republic, Proceedings*, volume 9302 of *Lecture Notes in Computer Science*, pages 41–50. Springer.