

# Translated Texts Under the Lens: From Machine Translation Detection to Source Language Identification

Anonymous ACL submission

## Abstract

In this work, we tackle the problem of the detection of translated texts from different angles. On top of addressing the classic task of machine translation detection, we investigate and find the presence of common patterns across different machine translation systems as well as different source languages. Then, we show that it is possible to identify the translation systems used to produce a translated text (F1-score 88.5%) as well as the source language of the original text (F1-score 79%). We assess our tasks using Books, a new dataset we built from scratch based on excerpts of novels and the well-known Europarl dataset.

## 1 Introduction

Today, commercial machine translation systems (MTS) are used worldwide by hundreds of thousands of people for personal or working purposes. They help bridge the gap in language barriers, especially on the Web, by facilitating communication between people. However, bad actors also use these systems to massively target potential victims of email-phishing (Parmar and Jahankhani, 2021) or fake reviews of products to trick recommendation systems (Juuti et al., 2018). Thus, machine translation detectors are actively used to infer spam emails or to detect poor quality web pages (Google, Dec. 2021).

In this work we put automatic translated texts under the lens. We study the impact of the MTSs and the source language on the Machine Translation Detection (MTD) task leveraging Books, a novel dataset built from excerpt novels that we plan to release publicly. (Anonymized). We find that MTSs have common patterns that can be learned using a single MTS. Similarly we discovered that the source language of the text does not significantly impact MTD performances.

We then investigate the possibility of identifying the MTS used to produce the translation and its

source language. To explore these points, we introduce, to the best of our knowledge, two new tasks: Machine Translation Identification (MTI) and Source Language Identification (SLI). For the first task, MTI, we built a classifier that shows promising results, with an average F1-score of 88.5%. In the second task, SLI, we propose a stacked classifier able to identify the source language with an F1-score of over 79%. We also believe that this task could be helpful in forensic analysis, where malicious actors attempt to obfuscate their writing style using MTSs (Kacmarcik and Gamon, 2006; Mahmood et al., 2019).

## 2 Dataset

To assess our experiments under different settings and topic domains, we perform our study using two datasets: one extracted from novels and the other based on speech transcriptions.

The first dataset we use is *Books*, a novel dataset we introduce. To build Books, we collect 100 books originally written in 4 different languages by 100 different established writers of the XX century. In particular, we select 25 books for each of the following source languages: Italian, French, Spanish, and German. The selected books belong to several different domains and authors. Thus they have very different writing styles. From each book, we select an excerpt of approximately 10,000 characters (on average 1642.67 words per novel) and their corresponding translation from the English edition. Finally, we produce 3 more English translations for each original excerpt using the APIs of 3 state-of-the-art commercial Machine Translation Systems: Google Translate(*GT*), Microsoft Translation(*MT*), and DeepL(*DL*). At the end of the process the Books dataset is made of 400 different samples.

The second dataset we use for our experiments is Europarl (Koehn, 2005). It is a parallel corpus

081 extracted from the proceedings of the European  
 082 Parliament containing *speech transcripts* of  
 083 European parliamentarians and the corresponding  
 084 professional translations into each of the 20  
 085 European languages. The texts on this  
 086 dataset include many speech-distinctive elements  
 087 such as hesitations, broken sentences, and  
 088 repetition (Bizzoni et al., 2020). Consistently with  
 089 Books, we obtain 100 seed samples by extracting  
 090 from Europarl 25 samples for each of the 4  
 091 languages we consider. Every sample is made using  
 092 transcripts of speakers of the same source language  
 093 and contains about 10,000 characters (on average  
 094 1512.81 words per sample). We pre-process the  
 095 dataset using the tools provided by Moses (Koehn  
 096 et al., 2007). Then, we collect the parallel English  
 097 translation of each seed sample. Finally, we  
 098 translate each seed sample using the selected MTSs.  
 099 Final datasets contain 400 (100 human-translated  
 100 and 300 machine-translated) English samples.

### 101 3 Experimental Settings

102 For all the experiments, we use 60% of the dataset  
 103 as train and 40% as test. We use Python Scikit-  
 104 learn (Pedregosa et al., 2011) to implement all the  
 105 models and the feature selection techniques. We  
 106 use default parameters unless specified.

107 **Feature Description.** Tab. 3 shows the features  
 108 we use for our tasks. *Words avg* is the average  
 109 number of words for each sentence of the text,  
 110 while *Adjectives avg* is the average number of  
 111 adjectives for each text. We use the notation *Char-*  
 112 *gram (i-k)* (resp. *Word-gram (i-k)*) to indicate  
 113 all the char n-gram (resp. word n-gram) with  
 114  $n \in \{i, \dots, k\}$ . *Dist Char-gram (i-k)* are char-  
 115 grams computed over the distortion text —text  
 116 where ascii characters are replaced with a special  
 117 character (Stamatatos, 2017). *POS Word-gram(i-*  
 118 *k)* are word-grams computed over Part of Speech  
 119 (POS) tagged text. Finally, the *Type Token Ratio*  
 120 (*TTR*) is the ratio between the number of unique  
 121 words and the total number of words for a given  
 122 text. We use the Bag of Words to weight the char-  
 123 gram and word-gram, while we use Tf-Idf to weight  
 124 distortion text.

### 125 4 Machine Translation Detection

126 The *Machine Translation Detection (MTD)* task  
 127 aims to automatically detect whether a text has  
 128 been translated by a machine translation system or  
 129 is human-generated. This task was broadly studied

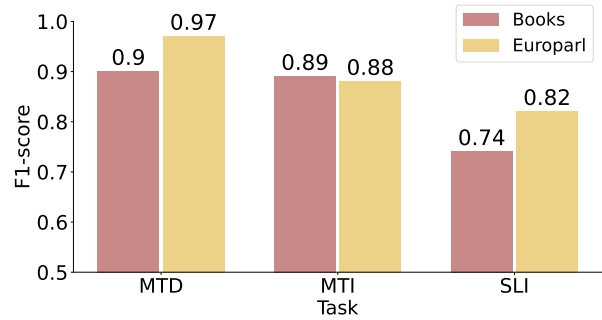


Figure 1: F1-score for the Machine Translation Detection (MTD), Machine Translation Identification (MTI) and Source Language Identification (SLI) tasks on the Books and Europarl datasets.

130 in the literature with different approaches such as  
 131 using fixed features (Aharoni et al., 2014; Li et al.,  
 132 2015), n-gram (Arase and Zhou, 2013; Popescu,  
 133 2011), coherence score (Nguyen-Son et al., 2018)  
 134 and similarity with round-trip translation (Nguyen-  
 135 Son et al., 2021). In this Section, we first want to  
 136 replicate results similar to the state-of-the-art on  
 137 our datasets. Then, we perform two experiments to  
 138 further explore the underlying patterns of machine-  
 139 translated texts.

140 For all the experiments in this section, we  
 141 use the following model. We train a Multilayer  
 142 Perceptron (Hornik et al., 1989) with a single  
 143 hidden layer made of 10 neurons and a BFGS  
 144 optimizer (Battiti and Masulli, 1990) for weights  
 145 optimization. Regarding the features, we compute  
 146 all the char n-gram with  $n \in \{1, \dots, 6\}$  and then  
 147 select the 2,500 more relevant n-gram according  
 148 to the chi-square metric (Forman et al., 2003) and  
 149 normalized with the SkLearn StandardScaler.

150 Figure 1 shows the results on Books and  
 151 Europarl datasets. We obtain a high F1-score on  
 152 both corpora (0.9 on Books and 0.97 on Europarl),  
 153 showing that our model can achieve state-of-the-  
 154 art comparable results in distinguishing machine-  
 155 translated and human-translated texts.

156 **Learning from a single MTS.** The next interesting  
 157 point we want to explore is if exists some common  
 158 pattern among the different MTSs that allow us  
 159 to detect machine-translated texts. We use only  
 160 samples translated by a single MTS, and the human  
 161 translated samples as the training set. Thus, we  
 162 repeat the experiment 3 times, once for each MTS.  
 163 Tab. 1 shows the results of this experiment for  
 164 the different datasets. As we can see, the models  
 165 achieve good results when tested on samples  
 166 produced by machine translators not present on

Train	Books	Europarl
GT	0.85	0.82
MT	0.89	0.95
DL	0.84	0.94

Table 1: F1-score for Task 1 training on a single MTS and testing on the others.

Train	Books	Europarl
IT	0.91	0.93
FR	0.85	0.74
ES	0.88	0.78
DE	0.73	0.81

Table 2: F1-score for Task 1 training on single language and testing on the others.

Feature Type	MTI		SLI	
	Books	Euro	Books	Euro
Char-gram (1-6)	318, 250	220, 593	261, 895	175, 247
Words avg	1	1	1	1
Sentence Length	1	1	1	1
Adjectives avg	1	1	1	1
Dist. Char-gram(5-8)	15, 134	12, 080	-	-
Dist. Char-gram(2-8)	-	-	13, 897	9, 522
POS Word-gram(1-6)	-	-	187, 481	145, 473
TTR	1	1	-	-
<b>All</b>	333, 389	232, 678	463, 277	330, 246

Table 3: Features used for the MTI and SLI tasks.

the training set. Interestingly, the model trained on MT achieves similar results to those obtained by training the model with the whole dataset. These results suggest that there is some common pattern among the MTSS, that the model can learn from a single MTS.

**Learning from a single language.** Since we have 4 different source languages in our dataset, we want to understand the impact they might have on the MTD task. In this experiment, we train our model using only the translation from one source language and test it against the sample produced by the other source languages and the human translated samples. Tab. 2 shows the F1-score using the different source languages. Results show that the model can learn machine translation patterns even when training only on one language, suggesting that these patterns are unrelated to the source language.

## 5 Machine Translator Identification

Results from the previous section show that there is a common pattern among the different MTSs that allow us to differentiate machine-translated text from humans-translated. In this section, we

investigate if MTS translations differ enough from each other to be able to identify which one has been used to translate a sample. Other works show that there could be potential differences between MTSs without trying to attempt to detect them. (Bhardwaj et al., 2020; Aharoni et al., 2014; Bizzoni et al., 2020) Thus, given a machine-translated text  $T'$ , our goal is to identify the MTS  $M$  that generated the text  $T'$ . We call this task *Machine Translator Identification (MTI)*. In particular, we focus on the identification of the 3 MTSs used to build the Books and Europarl datasets: *Google Translate*, *Microsoft Translation*, and *DeepL*. Given the goal of the task, for the following experiments, we use a sub-set of Europarl and Books datasets, removing from each dataset the 100 samples representing the class of human translations. For this task, we build an ensemble classifier. The first level comprises three different classifiers: a Support Vector Machine, a Logistic Regression, and a Random Tree. Then, the outputs of the classifiers are used as input to feed a hard voting layer (SkLearn VotingClassifier) for the final prediction. Tab. 3 shows the type and the number of features we use to train the three classifiers at the first level of our architecture. For all the n-gram type features, we select only the 85% most significant using SelectPercentile of SkLearn, and we standardize them with the SkLearn StandardScaler. Fig 1 reports the F1-score for the two datasets. As we can notice, our classifier performs similarly on both datasets, with an F1-score of 0.89 and 0.88 for Books and Europarl, respectively. To better understand the results, we analyze the confusion matrices of the two classifications. The confusion matrix of Books (Tab. 4) shows that GT is the hardest MTS to identify, and its misclassified samples are mostly assigned to the MT class. We found a possible explanation for these errors analyzing the BLUE score (Papineni et al., 2002) for each pair of MTS, obtaining a value of 69 for the pair GT-ML, 63 for GT-DL, and 62.9 for DL-ML. The high BLEU score between GT and MT shows that they have similar translations, leading to an erroneous classification of the GT samples. Conversely, the low similarity between the MT and DL classes could lead to the high accuracy we observe in our experiment. Finally, we obtain similar results analyzing the confusion matrix and the BLUE score for the Europarl dataset.

		Predicted		
		GT	MT	DL
Actual	GT	30	6	4
	MT	1	39	0
	DL	0	2	38

Table 4: Confusion Matrix on Books for the MTI task.

		Predicted			
		DE	ES	FR	IT
Actual	DE	25	0	5	0
	ES	2	23	4	1
	FR	0	2	27	1
	IT	1	5	9	15

Table 5: Confusion Matrix on Books for the SLI task.

		Predicted			
		DE	ES	FR	IT
Actual	DE	27	3	0	0
	ES	0	24	3	3
	FR	0	3	24	3
	IT	0	9	1	20

Table 6: Confusion Matrix on EuroParl for the SLI task.

## 6 Source Language Identification

As a final task, we propose the *Source Language Identification* (SLI). The goal of the task is, given a machine-translated text  $T'$  in a language  $L2$ , identify the source language  $L1$  of the text  $T$ . For our experiments we consider English as  $L2$  and the possible  $L1$  languages are: Italian, French, Spanish or German. This task could be considered a variation of other tasks already studied in the literature. Indeed, a similar task is the Native Language Identification (NLI), where the goal is to identify the native language  $L1$  of a person that writes a text in a second language  $L2$  (La Morgia et al., 2019; Tetreault et al., 2013). Another similar task is to determine the source language of a human-translated text (van Halteren, 2008; Koppel and Ordan, 2011). Unlike the previous study, our task focuses on identifying the source language of a text that is not written or translated by a human but by a MTS. For this task, we use the stacking ensemble technique. In particular, we stacked an AdaBoost (Freund and Schapire, 1997) model with 50 LinearSVC (Cortes and Vapnik, 1995) and a Logistic Regression (Wright, 1995) model as base estimators. Tab. 3 shows the type and the number of features we use to train the stacking classifier. For all the n-gram features, we select the top 70% according to their F-value, computed with the variance analysis (ANOVA) (St et al., 1989). Then we standardize them with a StandardScaler.

Fig 1 shows the F1-score of the model trained

and tested on both our datasets. The results suggest that identifying the source language detection is easier on EuroParl than in Books. As noted in (van Halteren, 2008) a possible reason could be that the EuroParl dataset may contain some distinctive patterns for the source language of the speaker. Instead, the Books dataset covers a wide area of topics and contains fewer clues about the speaker’s source language. Tab. 5 and 6 show the confusion matrices on the Books and EuroParl dataset. The most challenging source language to detect on both datasets is Italian, frequently misclassified as Spanish or French. German is generally better identified than the other languages except for French on the Books dataset, with five classification errors. Indeed, German has the highest F1-score among all the classes, with a value of 0.86 in Books and 0.94 in EuroParl. This is intuitive since German is a West Germanic language while the other 3 are Romance languages and have more features in common (Padró and Padró, 2004).

## 7 Conclusion and future work

In this work, we put translated text under the lens. We start evaluating the impact of MTSs and source languages on the Machine Translation Detection task. We find that MTSs have a common pattern that can be learned by a machine learning model trained with a single MTS. Moreover, the source language of the text does not significantly affect the performance of the task. Then, we introduce two new tasks: Machine Translator Identification and Source Language Identification. The goal of the first task is to identify the MTS that produced a translated text, while the second aims to identify the source language of a machine-translated text. The models we propose for both the tasks achieve an average F1-score of 88.5% and 78% respectively for MTI and SLI. Finally, we introduce Books, a novel dataset built for these tasks. Our results represent a first attempt to tackle the newly presented tasks. While we achieve good performance, we believe they could be further improved by using more advanced machine learning techniques. In our study, we perform all the analyses at the document level. In the future, it would be interesting to face the same problem at a more challenging limit, attempting to solve the tasks at the sentence level.

320  
321  
322  
323  
324  
325  
326  
327  
  
328  
329  
  
330  
331  
332  
333  
334  
335  
  
336  
337  
338  
339  
  
340  
341  
342  
343  
344  
345  
346  
  
347  
348  
349  
350  
351  
352  
353  
354  
  
355  
356  
  
357  
358  
359  
360  
  
361  
362  
363  
364  
  
365  
366  
  
367  
368  
369  
  
370  
371  
372  
373  
374

## References

Roei Aharoni, Moshe Koppel, and Yoav Goldberg. 2014. [Automatic detection of machine translated text and translation quality estimation](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 289–295, Baltimore, Maryland. Association for Computational Linguistics.

Anonymized. The resource will be available after acceptance.

Yuki Arase and Ming Zhou. 2013. [Machine translation detection from monolingual web-text](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1597–1607, Sofia, Bulgaria. Association for Computational Linguistics.

Roberto Battiti and Francesco Masulli. 1990. Bfgs optimization for faster and automated supervised learning. In *International neural network conference*, pages 757–760. Springer.

Shivendra Bhardwaj, David Alfonso Hermelo, Phillippe Langlais, Gabriel Bernier-Colborne, Cyril Goutte, and Michel Simard. 2020. [Human or neural translation?](#) In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6553–6564, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Yuri Bizzoni, Tom S Juzek, Cristina España-Bonet, Koel Dutta Chowdhury, Josef van Genabith, and Elke Teich. 2020. [How human is machine translation? comparing human and machine translations of text and speech](#). In *Proceedings of the 17th International Conference on Spoken Language Translation*, pages 280–290, Online. Association for Computational Linguistics.

Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning*, 20(3):273–297.

George Forman et al. 2003. An extensive empirical study of feature selection metrics for text classification. *J. Mach. Learn. Res.*, 3(Mar):1289–1305.

Yoav Freund and Robert E Schapire. 1997. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139.

Google. Dec. 2021. [Managing multi-regional and multilingual sites](#).

Kurt Hornik, Maxwell Stinchcombe, and Halbert White. 1989. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366.

Mika Juuti, Bo Sun, Tatsuya Mori, and N Asokan. 2018. Stay on-topic: Generating context-specific fake restaurant reviews. In *European Symposium on Research in Computer Security*, pages 132–151. Springer.

Gary Kacmarcik and Michael Gamon. 2006. [Obfuscating document stylometry to preserve author anonymity](#). In *Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions*, pages 444–451, Sydney, Australia. Association for Computational Linguistics. 375  
376  
377  
378  
379  
380

Philipp Koehn. 2005. [Europarl: A parallel corpus for statistical machine translation](#). In *Proceedings of Machine Translation Summit X: Papers*, pages 79–86, Phuket, Thailand. 381  
382  
383  
384

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. [Moses: Open source toolkit for statistical machine translation](#). In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic. Association for Computational Linguistics. 385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395

Moshe Koppel and Noam Ordan. 2011. [Translationese and its dialects](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 1318–1326, Portland, Oregon, USA. Association for Computational Linguistics. 396  
397  
398  
399  
400  
401

Massimo La Morgia, Alessandro Mei, Eugenio Nemmi, Simone Raponi, and Julinda Stefa. 2019. Nationality and geolocation-based profiling in the dark (web). *IEEE Transactions on Services Computing*. 402  
403  
404  
405

Yitong Li, Rui Wang, and Hai Zhao. 2015. [A machine learning method to distinguish machine translation from human translation](#). In *Proceedings of the 29th Pacific Asia Conference on Language, Information and Computation: Posters*, pages 354–360, Shanghai, China. 406  
407  
408  
409  
410  
411

Asad Mahmood, Faizan Ahmad, Zubair Shafiq, Padmini Srinivasan, and Fareed Zaffar. 2019. A girl has no name: Automated authorship obfuscation using mutant-x. *Proc. Priv. Enhancing Technol.*, 2019(4):54–71. 412  
413  
414  
415  
416

Hoang-Quoc Nguyen-Son, Huy H. Nguyen, Ngoc-Dung T. Tieu, Junichi Yamagishi, and Isao Echizen. 2018. [Identifying computer-translated paragraphs using coherence features](#). In *Proceedings of the 32nd Pacific Asia Conference on Language, Information and Computation*, Hong Kong. Association for Computational Linguistics. 417  
418  
419  
420  
421  
422  
423

Hoang-Quoc Nguyen-Son, Tran Thao, Seira Hidano, Ishita Gupta, and Shinsaku Kiyomoto. 2021. [Machine translated text detection through text similarity with round-trip translation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5792–5797, Online. Association for Computational Linguistics. 424  
425  
426  
427  
428  
429  
430  
431  
432

- 433 Muntsa Padró and Lluís Padró. 2004. Comparing  
434 methods for language identification. *Procesamiento*  
435 *del lenguaje natural*, 33.
- 436 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-  
437 Jing Zhu. 2002. [Bleu: a method for automatic](#)  
438 [evaluation of machine translation](#). In *Proceedings*  
439 *of the 40th Annual Meeting of the Association*  
440 *for Computational Linguistics*, pages 311–318,  
441 Philadelphia, Pennsylvania, USA. Association for  
442 Computational Linguistics.
- 443 Yogeshvar Singh Parmar and Hamid Jahankhani. 2021.  
444 Utilising machine learning against email phishing  
445 to detect malicious emails. In *Artificial Intelligence*  
446 *in Cyber Security: Impact and Implications*, pages  
447 73–102. Springer.
- 448 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel,  
449 B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,  
450 R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,  
451 D. Cournapeau, M. Brucher, M. Perrot, and  
452 E. Duchesnay. 2011. Scikit-learn: Machine learning  
453 in Python. *Journal of Machine Learning Research*,  
454 12:2825–2830.
- 455 Marius Popescu. 2011. [Studying translationese at the](#)  
456 [character level](#). In *Proceedings of the International*  
457 *Conference Recent Advances in Natural Language*  
458 *Processing 2011*, pages 634–639, Hissar, Bulgaria.  
459 Association for Computational Linguistics.
- 460 Lars St, Svante Wold, et al. 1989. Analysis of variance  
461 (anova). *Chemometrics and intelligent laboratory*  
462 *systems*, 6(4):259–272.
- 463 Efstathios Stamatatos. 2017. Authorship attribution  
464 using text distortion. In *Proceedings of the*  
465 *15th Conference of the European Chapter of the*  
466 *Association for Computational Linguistics: Volume*  
467 *1, Long Papers*, pages 1138–1149.
- 468 Joel Tetreault, Daniel Blanchard, and Aoife Cahill. 2013.  
469 A report on the first native language identification  
470 shared task. In *Proceedings of the eighth workshop*  
471 *on innovative use of NLP for building educational*  
472 *applications*, pages 48–57.
- 473 Hans van Halteren. 2008. [Source language markers](#)  
474 [in EUROPARL translations](#). In *Proceedings of the*  
475 *22nd International Conference on Computational*  
476 *Linguistics (Coling 2008)*, pages 937–944,  
477 Manchester, UK. Coling 2008 Organizing  
478 Committee.
- 479 Raymond E Wright. 1995. Logistic regression.