# Bias Detection and Mitigation for the Development of Fair and High-Quality Automatic Text Simplification Corpora

Victoria Muñoz-García

*University Institute for Computing Research (IUII), University of Alicante*

### Abstract

Biases embedded in text corpora pose significant challenges to the development of fair and equitable Artificial Intelligence (AI) systems. These biases, often reflective of historical inequalities and stereotypes, are inadvertently learned by Large Language Models (LLMs) during training, leading to the generation of text that can perpetuate and amplify these biases. Such biases are particularly problematic when these models are employed in real-world applications, where they can impact decision-making processes and accessibility for diverse user groups. Additionally, the complexity of generated text further exacerbates the issue for individuals with cognitive impairments, making it harder for them to understand information. In response to these challenges, ATS has emerged as a vital tool to transform complex texts into more accessible formats. Given the challenges mentioned above, the main objective of this doctoral research is to investigate and develop methodologies for the detection and mitigation of biases in training corpora used for LLMs, particularly for ATS in the tourism sector.

### Keywords

Natural language Processing, Language Models, Bias mitigation, Automatic Text Simplification,

## 1. Introduction

The continuous generation of large volumes of textual data has contributed to the proliferation of biases, which are frequently embedded in the training corpora of language models. Biases embedded in text corpora pose significant challenges to the development of fair and equitable Artificial Intelligence (AI) systems. Machine learning (ML) methods can not only reflect existing societal biases but also exacerbate them [1]. Biases and inequalities in the data may be absorbed by the algorithm and reflected in outputs when training models, which can have a significant and often harmful impact on people's lives [2]. Consequently, the development and use of high-quality corpora—characterized by fairness and explainability—are essential, as they significantly influence research outcomes. Therefore, the goal of corpus fairness is to ensure that ML models trained on these corpora do not perpetuate or amplify biases present in the data.

Moreover, the increasing complexity of information poses significant barriers to comprehension for the general public. Official communications, in particular, must be accessible to all individuals, including those with reading difficulties or cognitive impairments. However, manual text simplification is costly, as it demands considerable time and specialized expertise. Manual simplification of the existing volume of textual content is impractical [3]. This challenge underscores the necessity for automated approaches that facilitate equitable access to information. In this context, ATS seeks to transform original texts into simplified versions that are more accessible and easier to understand.

Taking these considerations into account, this PhD thesis aims to define and establish the research focus and direction of this thesis, thereby providing a structured framework and clearly defined objectives to guide the subsequent investigation. The remainder of this article is organized as follows: Section 2 presents an overview of bias detection and mitigation, as well as ATS; Section 3 outlines the research hypothesis and objectives; Section 4 details the proposed methodology; and Section 5 highlight several research questions that remain open for discussion.

## 2. Background and Related Work

Prior to outlining the research proposal, this section provides a contextual overview of the current state-of-the-art in bias detection and mitigation, as well as in Automatic Text Simplification.

### 2.1. Bias in Corpora: Implications for Fairness and Quality

One of the concerns regarding Artificial Intelligence (AI) systems lies in the fairness of their outputs, which can result in significant negative social impacts when such systems are biased or unfair. For instance, these biases may manifest as "significant prejudice towards different genders and race" [4]. To mitigate such issues at the data level, this work emphasizes corpus fairness—the principle that a corpus should be representative of the population's diversity. As defined by Mehrabi et al. [5], fairness refers to "the absence of any prejudice or favouritism toward an individual or group based on their inherent or acquired characteristics." Accordingly, this study defines a fair corpus as one that provides an accurate and balanced view of the language or phenomenon under investigation, without reinforcing existing biases.

Ensuring fairness at the corpus level directly relates to broader concerns of corpus quality, which encompasses the validity, reliability, and representativeness of the textual data. Corpus quality is determined not only by its balance and representativeness but also by the sufficiency of the data and its alignment with the target discourse domain [6] [7]. Within this framework, corpus fairness becomes a fundamental component, as it demands that the dataset avoids prejudice and adequately reflects the diversity of the target population [8]. Thus, a fair corpus is inherently a high-quality corpus—unbiased, balanced, and inclusive—ensuring that the models trained on it are less likely to perpetuate harmful stereotypes or social inequities. However, biases originate from three interrelated sources: training data bias, embedding bias, and label bias [9].

- **Training data bias** arises when the input corpus encodes historical and societal inequalities. Such biases are frequently embedded in language patterns that models internalize during pretraining, leading to the replication of stereotypes—such as associating specific professions with particular genders—or the marginalization of underrepresented groups [10, 11, 12].
- **Embedding bias** is introduced during the generation of semantic representations, where vector spaces reflect and amplify stereotypical associations. For instance, certain professions may be semantically clustered in ways that align with gender norms (e.g., aligning femininity with nursing and masculinity with medicine), thereby influencing downstream tasks in biased ways [13, 14].
- **Label bias** occurs during the fine-tuning phase, particularly in instruction tuning, where human annotators' subjective judgments and sociocultural perspectives affect the labeled data [15, 16]. Techniques such as reinforcement learning with human feedback (RLHF) further risk encoding individual annotator biases, as model alignment with human preferences may inadvertently reflect narrow value systems [17]. Mitigating these forms of bias requires comprehensive strategies, including diversifying annotator backgrounds, employing systematic fairness metrics, and establishing rigorous annotation protocols [18, 9].

Biased corpora used in training LLMs can influence not only fairness of model outputs but also the inclusiveness of simplified texts. For example, if original texts reinforce gender or cultural biases, simplification systems may perpetuate or even amplify them. Therefore, ensuring fairness in training corpora is a necessary step for developing inclusive ATS systems.

### 2.2. Automatic Text Simplification: Linguistic Levels

Promoting information accessibility is progressively more essential as the volume of textual content continues to grow, posing significant barriers to comprehension for individuals with cognitive impairments, non-native speakers, and those with limited literacy skills [19]. ATS seeks to address this issue by reducing linguistic complexity without altering the intended meaning, thereby improving the

readability and overall accessibility of written information while operating across various linguistic dimensions [3][20]: lexical, syntactic and discourse-level simplification.

**Lexical simplification (LS)** focuses on replacing complex words with simpler alternatives, typically using lexical resources or pretrained language models such as BERT [21]. The complexity of a word is often inversely related to its frequency, with general words being more polysemous and technical terms usually monosemic [22]. Therefore, the effectiveness of LS is influenced by both the topic of the text and the contextual appropriateness of the synonyms used [23]. LS typically involves five subtasks, which can be addressed either modularly or jointly in end-to-end systems [19, 24, 25, 26, 27]:

1. Complex Word or Phrase Identification (CWI/CPI), which seeks to identify terms likely to be difficult for target audiences [19, 28]; according to [25, 29], a word is considered complex if it occurs less than five times per million and meets at least one additional criterion, such as being archaic, long, borrowed, culturally uncommon, low-frequency, highly specialized, or used with an unusual meaning

2. Substitution Generation (SG), where candidate synonyms are produced using lexical databases, embeddings, or language models [25]

3. Substitution Selection (SS), which involves selecting contextually appropriate alternatives that preserve meaning [19, 25]

4. Substitution Ranking (SR), which orders synonyms by simplicity, generally informed by word frequency [26, 25]

5. Morphological generation and context adaptation to ensure grammatical integration of the substitution [26]. In cases where substitution fails, simplification may involve generalization or explanatory expansion, resulting in increased lexical and syntactic divergence from the original text [30].

**Syntactic simplification (SS)** targets structural complexity in sentences by removing or transforming difficult syntactic constructions—such as coordination, subordination, relative clauses, and passive forms—while preserving meaning [28, 20]. Techniques include sentence splitting, voice transformation, and ambiguity resolution [31, 32], as sentence length and structure have a direct impact on readability, particularly for non-native speakers and readers with cognitive impairments [33]. A set of core operations for SS has been identified, including splitting, merging, reordering, insertion, deletion, and transformations such as lexical substitution and voice alterations [32, 28, 34, 35, 36]. These operations collectively enhance sentence-level comprehensibility and often lead to discourse-level changes, thereby linking syntactic and discursive simplification [24].

At the **discourse or document level (DS)**, simplification involves modifying multi-sentence structures to improve coherence and accessibility [37]. Operations include paragraph splitting, sentence reordering, clarification of coreference chains, anaphora resolution, and title generation [20]. Since syntactic simplification can affect referential clarity and coherence, DS strategies also account for these dependencies. For example, inappropriate pronoun removal may disrupt coherence, while overuse of prepositional phrases may introduce complexity. Rule-based systems for DS have been proposed, incorporating strategies for entity replacement, substitution generation, and ranking based on referential accessibility [24]. Given that many real-world applications demand a broader contextual understanding, document-level simplification is often more applicable than sentence-level approaches [37]. Overall, the choice of linguistic level in TS—lexical, syntactic, or discursive—substantially influences the design of resources and tools employed for effective simplification.

Building on the considerations previously detailed, corpus fairness, by promoting balanced and representative data, directly contributes to the quality of linguistic resources. This is particularly relevant for Automatic Text Simplification, where the effectiveness of simplification strategies relies on the integrity and inclusiveness of the underlying corpora.

# 3. Research Description

Large Language Models (LLMs) are predominantly trained on large-scale textual corpora that often contain implicit and explicit biases. These biases are not only learned by the models but also potentially amplified, raising ethical and societal concerns.

In this research, we hypothesize that mitigating biases at the corpus level during the training phase of LLMs can significantly reduce the propagation of social and representational inequalities in downstream tasks. Specifically, we focus on the application of LLMs to the task of ATS within the touristic domain, where biased representations and inaccessible language can hinder communication. Addressing these issues jointly allows for the development of simplification systems that are not only more accessible but also fairer.

## 3.1. Objectives

This doctoral research aims to investigate and develop methods for detecting and mitigating biases in the training corpora of LLMs, with the aim of reducing biased behavior in ATS systems, particularly within the tourism sector.

### 3.1.1. Specific Objectives

Based on this objective, several sub-objectives have been defined to outline a detailed workflow:

1. Conduct a comprehensive review of current methodologies for bias detection and mitigation in corpora, with an emphasis on practices that contribute to fairness and quality in dataset creation.
2. Examining the state-of-the-art techniques in Automatic Text Simplification, including both rule-based and neural approaches.
3. Constructing high-quality and fair corpora tailored for ATS, including both general-domain and tourism-specific texts.
4. Developing a task-specific instruction corpus to enhance supervised fine-tuning for ATS.
5. Training LLMs for the ATS task in the touristic domain using the constructed corpora, with a focus on bias mitigation throughout the training pipeline.
6. Generating and publishing Spanish-language resources, including corpora, linguistic guidelines, and LLMs to support future research and development.
7. Aligning research outputs with the Sustainable Development Goals, particularly SDG 5 (Gender Equality), and SDG 10 (Reduced Inequalities).

Through the realization of these objectives, the research aims to make a significant contribution to the ethical development of LLMs and to enhance the accessibility of tourism-related content through fair and effective text simplification methodologies.

# 4. Methodology

This PhD thesis presents a methodology based on a comprehensive review of bias detection and mitigation techniques, with a focus on ATS.

The approach is divided in two main steps. The first one, consists on the data collection, which involves three different corpus:

- The refinement of the ClearSim corpus, involving data cleaning and bias mitigation, to ensure a high-quality, unbiased dataset.
- The use of a domain-specific corpus of approximately 300 million tokens, which is centered on tourism-related content.
- The creation of a simplification instruction corpus designed to train an instructed model tailored for tasks within the ATS domain.

The second step of the method consist on the training of a model that combines the bias mitigation tecnhiques, ATS task and constrained in the tourism domain. To do so, the core model for the experiments Salamandra 7b, based on a transformer architecture, will be fine-tuned with all the collected data. This methodology concludes with the development of a system that integrates bias mitigation techniques and ATS capabilities within a domain-adapted LLM framework, for the purpose of a discriminative simplified task that enhances readability without compromising accuracy.

Ethical considerations, review processes, and bias mitigation will be integrated into the development of a language model, ensuring an ethical and reliable approach to ATS in the tourism domain.

### 4.1. Related experiments and work in progress

The experiments conducted investigate gender bias in text generated by LLMs in Spanish, a language where gender is inherently embedded in its linguistic structure. The study proposes a validated methodology for quantifying this bias, which involves the creation of gendered seed-word lists, the construction of a Spanish-specific corpus with curated prompts, and the analysis of gender polarity and co-occurrences in the generated text. This methodology was tested on five state-of-the-art LLMs: GPT-3.5, GPT-4, Llama 3, Gemini 1.5, and Mixtral8x7b. The research provides a systematic framework for detecting gender bias in Spanish, revealing performance variations among models and highlighting gender disparities. The aim of this work is to enhance bias quantification methodologies and contribute to the development of more equitable AI systems.

The ongoing work focuses on the analysis and fine-tuning of the Salamandra model for the tourism domain. The Salamandra model (`Salamandra-7B-Instruct`, available at [https://huggingface.co/BSC-LT/salamandra-7b-instruct]) will initially be evaluated to establish its baseline performance. Subsequently, fine-tuning will be conducted using a corpus of simplification data in Spanish. The primary objective of the experiments is to adapt the model to the tourism domain, with particular emphasis on the language simplification classification task.

## 5. Research issues to discuss

This section addresses the challenges encountered in conducting this research that need to be taken into consideration. Future research aims to include bias detection and mitigation methods that can be applied to the Spanish language, creating high-quality corpora and enhancing the task of ATS for turism text simplification.

Our research topics cover several subjects that are open to discussion. This study focuses on three main key topics, which are outlined below, along with the questions each of these topics could encompass.

The following key questions serve as the foundation for a thorough investigation into the intersection of fairness, bias, and text simplification in this research project:

1. **Fairness**: The primary issue to be addressed concerns the considerations of fairness and quality. What criteria define fairness and quality within these contexts, and how can these criteria be effectively implemented?

2. **Bias Detection and Mitigation**: The next topic focuses on the identification and mitigation of biases. Are there comprehensive sets of features that effectively describe various types of biases in human-related data? Can attribute lists be constructed for this purpose, and are there identifiable linguistic patterns that indicate and assist in the detection of biases? Additionally, which specific patterns must be recognized and addressed? Regarding tourist-related information, how can biases be minimized, and is there a particular pattern that requires detection?

3. **Automatic Text Simplification**: From the perspective of accessibility, who is most affected by the complexity of information, and which users would benefit most from text simplification systems? Should the focus be on a specific target audience, or should a more general approach be adopted to reach a broader population?

These inquiries form the core of an in-depth exploration into fairness, bias, and text simplification within the scope of the present research.

## 6. Acknowledgments

## Declaration on Generative AI

During the preparation of this work, the authors used ChatGPT in order to: Grammar and spelling check, Paraphrase, translate and reword. After using this tool, the author reviewed and edited the content as needed and takes full responsibility for the publication's content.

## References

[1] J. P. Consuegra-Ayala, Y. Gutiérrez, Y. Almeida-Cruz, M. Palomar, Automatic annotation of protected attributes to support fairness optimization, Information Sciences (2024) 120188.

[2] T. P. Pagano, R. B. Loureiro, F. V. Lisboa, R. M. Peixoto, G. A. Guimarães, G. O. Cruz, M. M. Araujo, L. L. Santos, M. A. Cruz, E. L. Oliveira, et al., Bias and unfairness in machine learning models: a systematic review on datasets, tools, fairness metrics, and identification and mitigation methods, Big data and cognitive computing 7 (2023) 15.

[3] S. Bott, H. Saggion, Automatic simplification of spanish text for e-accessibility, in: Computers Helping People with Special Needs: 13th International Conference, ICCHP 2012, Linz, Austria, July 11-13, 2012, Proceedings, Part I 13, Springer, 2012, pp. 527–534.

[4] H. Liu, J. Dacon, W. Fan, H. Liu, Z. Liu, J. Tang, Does gender matter? towards fairness in dialogue systems, arXiv preprint arXiv:1910.10486 (2019).

[5] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, A. Galstyan, A survey on bias and fairness in machine learning, ACM computing surveys (CSUR) 54 (2021) 1–35.

[6] M. Hurtado Bodell, M. Magnusson, S. Mützel, From documents to data: A framework for total corpus quality, Socius 8 (2022) 23780231221135523.

[7] D. Miller, D. Biber, Evaluating reliability in quantitative vocabulary studies: The influence of corpus design and composition, International Journal of Corpus Linguistics 20 (2015) 30–53.

[8] H.-C. A. Intelligence, Artificial intelligence index report 2024: Public data (2024).

[9] Z. Chu, Z. Wang, W. Zhang, Fairness in large language models: A taxonomic survey, ACM SIGKDD explorations newsletter 26 (2024) 34–48.

[10] H. Kotek, R. Dockum, D. Sun, Gender bias and stereotypes in large language models, in: Proceedings of the ACM collective intelligence conference, 2023, pp. 12–24.

[11] S. Jia, T. Meng, J. Zhao, K.-W. Chang, Mitigating gender bias amplification in distribution by posterior regularization, arXiv preprint arXiv:2005.06251 (2020).

[12] T. Sun, A. Gaut, S. Tang, Y. Huang, M. ElSherief, J. Zhao, D. Mirza, E. Belding, K.-W. Chang, W. Y. Wang, Mitigating gender bias in natural language processing: Literature review, arXiv preprint arXiv:1906.08976 (2019).

[13] W.-H. Weng, A. Sellergen, A. P. Kiraly, A. D'Amour, J. Park, R. Pilgrim, S. Pfohl, C. Lau, V. Natarajan, S. Azizi, et al., An intentional approach to managing bias in general purpose embedding models, The Lancet Digital Health 6 (2024) e126–e130.

[14] O. Papakyriakopoulos, S. Hegelich, J. C. M. Serrano, F. Marco, Bias in word embeddings, in: Proceedings of the 2020 conference on fairness, accountability, and transparency, 2020, pp. 446–457.

[15] P. Haller, A. Aynetdinov, A. Akbik, Opiniongpt: Modelling explicit biases in instruction-tuned llms, in: North American Chapter of the Association for Computational Linguistics, 2023. URL: https://api.semanticscholar.org/CorpusID:261582269.

[16] M. Parmar, S. Mishra, M. Geva, C. Baral, Don't blame the annotator: Bias already starts in the annotation instructions, arXiv preprint arXiv:2205.00415 (2022).

[17] I. Itzhak, G. Stanovsky, N. Rosenfeld, Y. Belinkov, Instructed to bias: Instruction-tuned language models exhibit emergent cognitive bias, Transactions of the Association for Computational Linguistics 12 (2024) 771–785.

[18] F. Chen, L. Wang, J. Hong, J. Jiang, L. Zhou, Unmasking bias in artificial intelligence: a systematic review of bias detection and mitigation strategies in electronic health record-based models, Journal of the American Medical Informatics Association 31 (2024) 1172–1183.

[19] S. Štajner, Automatic text simplification for social good: Progress and challenges, 2021, p. 2637 – 2652. URL: https://www.scopus.com/inward/record.uri?eid=2-s2.0-85115295095&partnerID=40&md5=c9c85f46170c64633ef68e22f634c503.

[20] P. Poupet, M. Hauguel, E. Boehm, C. Roze, P. Tardy, An automated tool with human supervision to adapt difficult texts into plain language, in: S. Štajner, H. Saggio, M. Shardlow, F. Alva-Manchego (Eds.), Proceedings of the Second Workshop on Text Simplification, Accessibility and Readability, INCOMA Ltd., Shoumen, Bulgaria, Varna, Bulgaria, 2023, pp. 131–133. URL: https://aclanthology.org/2023.tsar-1.13.

[21] A. Phatak, D. W. Savage, R. Ohle, J. Smith, V. Mago, Medical text simplification using reinforcement learning (teslea): Deep learning-based text simplification approach, JMIR Medical Informatics 10 (2022). doi:10.2196/38095.

[22] E. Rolin, Q. Langlois, P. Watrin, T. François, Frenlys: A tool for the automatic simplification of french general language texts, 2021, p. 1196 – 1205. URL: https://www.scopus.com/inward/record.uri?eid=2-s2.0-85123593221&doi=10.26615%2f978-954-452-072-4_135&partnerID=40&md5=9c3e53135b71d1c106d7aba4d5391385.

[23] S. V. Pervukhina, G. V. Basenko, I. G. Ryabtseva, E. E. Sakharova, Approaches to text simplification: Can computer technologies outdo a human mind?, GEMA Online Journal of Language Studies 21 (2021).

[24] M. Romero, S. Calderon-Ramirez, M. Solis, N. Perez-Rojas, M. Chacon-Rivas, H. Saggion, Towards text simplification in spanish: A brief overview of deep learning approaches for text simplification, 2022. URL: https://www.scopus.com/inward/record.uri?eid=2-s2.0-85148476856&doi=10.1109%2fBIP56202.2022.10032482&partnerID=40&md5=8d950c2b14602f31edf62a532f29a01b.

[25] A. Todirascu, R. Wilkens, E. Rolin, T. François, D. Bernhard, N. Gala, Hector: A hybrid text simplification tool for raw texts in french, 2022, p. 4620 – 4630. URL: https://www.scopus.com/inward/record.uri?eid=2-s2.0-85144399396&partnerID=40&md5=586115bd697bdd0494280b31eca7a8fc.

[26] D. Ferres, H. Saggion, Alexsis: A dataset for lexical simplification in spanish, in: N. Calzolari, F. Bechet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, H. Mazo, H. Odijk, S. Piperidis (Eds.), LREC 2022: Thirteen International Conference on Language Resources and Evaluation, Google; S African Ctr Digital Language Resources; Vocapia Res; 3M; Emvista; Expert.ai; Grammarly; Minist Culture, Delegat Gen Langue Francaise Aux Langues France; Reg Sud Provence Alpes Cote Azur, 2022, pp. 3582–3594. 13th International Conference on Language Resources and Evaluation (LREC), Marseille, France, Jun 20-25, 2022.

[27] R. Alarcon, L. Moreno, P. Martinez, Easier corpus: A lexical simplification resource for people with cognitive impairments, PLOS ONE 18 (2023). doi:10.1371/journal.pone.0283622.

[28] S. S. Al-Thanyyan, A. M. Azmi, Automated text simplification: A survey, ACM COMPUTING SURVEYS 54 (2021). doi:10.1145/3442695.

[29] M. Shardlow, R. Evans, M. Zampieri, Predicting lexical complexity in english texts: the complex 2.0 dataset, Language Resources and Evaluation 56 (2022) 1153–1194.

[30] N. Grabar, H. Saggion, Evaluation of automatic text simplification: Where are we now, where should we go from here; [Évaluation de la simplification automatique de textes: où nous en sommes et vers où devonsnous aller], volume 1, 2022, p. 453 – 463. URL: https://www.scopus.com/inward/

record.uri?eid=2-s2.0-85149436811&partnerID=40&md5=ca0bac50f056cf185678e78a773cb723.

[31] N. Chatterjee, R. Agarwal, Depsym: A lightweight syntactic text simplification approach using dependency trees, volume 2944, 2021, p. 42 – 56. URL: https://www.scopus.com/inward/record. uri?eid=2-s2.0-85115331477&partnerID=40&md5=b3e9c726ff228f4886715035046c101a.

[32] Y. Anees, S. A. Rauf, Automatic sentence simplification in low resource setting for urdu, in: NLP4POSIMPACT 2021: the 1st Workshop on NLP for Positive Impact, Assoc Computat Linguist, 2021, pp. 60–70. 1st Workshop on Natural Language Processing for Programming (NLP4Prog), ELECTR NETWORK, AUG 05-06, 2021.

[33] L. Mucida, A. d. P. Oliveira, M. d. A. Possi, A language-independent metric for measuring text simplification that does not require a parallel corpus, volume 35, 2022. URL: https://www.scopus.com/inward/record.uri?eid=2-s2.0-85131144341&doi=10.32473%2fflairs. v35i.130608&partnerID=40&md5=279941a9151d0a2d6ba9bbd48352a325.

[34] R. Cardon, A. Bibal, R. Wilkens, D. Alfter, M. Norré, A. Müller, P. Watrin, T. François, Linguistic corpus annotation for automatic text simplification evaluation, 2022, p. 1842 – 1866. URL: https://www.scopus.com/inward/record.uri?eid=2-s2.0-85149440745&partnerID=40&md5= 4814acff05fa7912d1106695b5df26a9.

[35] A. Alshanqiti, A. Alkhodre, A. Namoun, S. Albouq, E. Nabil, A transformer seq2seq model with fast fourier transform layers for rephrasing and simplifying complex arabic text, International Journal of Advanced Computer Science and Applications 14 (2023) 888 – 898. URL: https://www.scopus.com/inward/record.uri?eid=2-s2.0-85150987857&doi=10.14569% 2fIJACSA.2023.01402101&partnerID=40&md5=866651d71d5b3653689c341102a7cba4.

[36] D. Brunato, F. Dell'Orletta, G. Venturi, Linguistically-based comparison of different approaches to building corpora for text simplification: A case study on italian, Frontiers in Psychology 13 (2022). doi:10.3389/fpsyg.2022.707630.

[37] S. Blinova, X. Zhou, M. Jaggi, C. Eickhoff, S. A. Bahrainian, SIMSUM: Document-level text simplification via simultaneous summarization, in: A. Rogers, J. Boyd-Graber, N. Okazaki (Eds.), Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, Toronto, Canada, 2023, pp. 9927–9944. URL: https://aclanthology.org/2023.acl-long.552. doi:10.18653/v1/2023.acl-long.552.