



**HAL**  
open science

## **Automatic Simplification of Scientific Texts: SimpleText Lab at CLEF-2022**

Liana Ermakova, Patrice Bellot, Jaap Kamps, Diana Nurbakova, Irina Ovchinnikova, Eric Sanjuan, Elise Mathurin, Sílvia Araújo, Radia Hannachi, Stéphane Huet, et al.

### ► To cite this version:

Liana Ermakova, Patrice Bellot, Jaap Kamps, Diana Nurbakova, Irina Ovchinnikova, et al.. Automatic Simplification of Scientific Texts: SimpleText Lab at CLEF-2022. *Advances in Information Retrieval*, 13186, Springer International Publishing, pp.364-373, 2022, Lecture Notes in Computer Science, 10.1007/978-3-030-99739-7\_46 . hal-03637775

**HAL Id: hal-03637775**

**<https://hal.science/hal-03637775v1>**

Submitted on 11 Apr 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Automatic Simplification of Scientific Texts: SimpleText Lab at CLEF-2022

Liana Ermakova<sup>1</sup>[0000-0002-7598-7474], Patrice Bellot<sup>2</sup>, Jaap Kamps<sup>3</sup>,  
Diana Nurbakova<sup>4</sup>, Irina Ovchinnikova<sup>5</sup>, Eric SanJuan<sup>6</sup>, Elise Mathurin<sup>1</sup>, Sílvia  
Araújo<sup>7</sup>, Radia Hannachi<sup>8</sup>, Stéphane Huet<sup>6</sup>, and Nicolas Poinso<sup>1</sup>

<sup>1</sup> Université de Bretagne Occidentale, HCTI - EA 4249, France

liana.ermakova@univ-brest.fr

<sup>2</sup> Aix Marseille Univ, Université de Toulon, CNRS, LIS, Marseille, France

<sup>3</sup> University of Amsterdam, Amsterdam, The Netherlands

<sup>4</sup> Institut National des Sciences Appliquées de Lyon, LIRIS UMR 5205 CNRS, Lyon, France

<sup>5</sup> Sechenov University, Moscow, Russia

<sup>6</sup> Avignon Université, LIA, France

<sup>7</sup> University of Minho, Portugal

<sup>8</sup> Université de Bretagne Sud, HCTI - EA 4249, France

**Abstract.** The Web and social media have become the main source of information for citizens, with the risk that users rely on shallow information in sources prioritizing commercial or political incentives rather than the correctness and informational value. Non-experts tend to avoid scientific literature due to its complex language or their lack of prior background knowledge. Text simplification promises to remove some of these barriers. The CLEF 2022 SimpleText track addresses the challenges of text simplification approaches in the context of promoting scientific information access, by providing appropriate data and benchmarks, and creating a community of NLP and IR researchers working together to resolve one of the greatest challenges of today. The track will use a corpus of scientific literature abstracts and popular science requests. It features three tasks. First, *content selection* (what is in, or out?) challenges systems to select passages to include in a simplified summary in response to a query. Second, *complexity spotting* (what is unclear?) given a passage and a query, aims to rank terms/concepts that are required to be explained for understanding this passage (definitions, context, applications). Third, *text simplification* (rewrite this!) given a query, asks to simplify passages from scientific abstracts while preserving the main content.

**Keywords:** Scientific text simplification · (Multi-document) summarization · Contextualization · Background knowledge

## 1 Introduction

Being science literate is an important ability for people. It is one of the keys for critical thinking, objective decision-making and judgment of the validity and significance of findings and arguments, which allows discerning facts from fiction. Thus, having a basic scientific knowledge may also help maintain one's health, both physiological and mental. The COVID-19 pandemic provides a good example of such a matter. Understanding

the issue itself, being aware of and applying social distancing rules and sanitary policies, choosing to use or avoid particular treatment or prevention procedures can become crucial. In the context of a pandemic, the qualified and timely information should reach everyone and be accessible. That is what motivates projects such as EasyCovid.<sup>9</sup>

However, scientific texts are often hard to understand as they require solid background knowledge and use tricky terminology. Although there were some recent efforts on text simplification (e.g. [25]), removing such understanding barriers between scientific texts and general public in an automatic manner is still an open challenge. **SimpleText Lab** brings together researchers and practitioners working on the generation of simplified summaries of scientific texts. It is a new evaluation lab that follows up the SimpleText-2021 Workshop [9]. All perspectives on automatic science popularisation are welcome, including but not limited to: Natural Language Processing (NLP), Information Retrieval (IR), Linguistics, Scientific Journalism, etc.

SimpleText provides data and benchmarks for discussion of challenges of automatic text simplification by bringing in the following tasks:

- **TASK 1: What is in (or out)?** Select passages to include in a simplified summary, given a query.
- **TASK 2: What is unclear?** Given a passage and a query, rank terms/concepts that are required to be explained for understanding this passage (definitions, context, applications,..).
- **TASK 3: Rewrite this!** Given a query, simplify passages from scientific abstracts.
- **UNSHARED TASK** We welcome any submission that uses our data!

## 2 Background

### 2.1 Content Selection

We observe an accelerating growth of scientific publication and their major impact on the society, especially in medicine (e.g. the COVID-19 pandemic) and in computer science with unprecedented use of machine learning algorithms and their societal issues (biases, explainability, etc.). Numerous initiatives try to make science understandable for everyone. Efforts have been made by scientific journalism (Nature, The Guardian, ScienceX) researchers (Papier-Maché project<sup>10</sup>), and internet forums (*Explain Like I'm 5*<sup>11</sup>). The ScienceBites<sup>12</sup> platform publishes short simple posts about individual research papers, making state-of-the-art science accessible to a wide audience. While structured abstracts are an emerging trend since they tend to be informative [12,10], non-experts are usually interested in other types of information. Popular science articles are generally much shorter than scientific publications. Thus, information selection is a crucial but understudied task in document simplification especially with regard to the target audience [39]. In many cases the information in a summary designed for an

<sup>9</sup> <https://easycovid19.org/>

<sup>10</sup> <https://papiermachesciences.org/>

<sup>11</sup> <https://www.reddit.com/r/explainlikeimfive>

<sup>12</sup> <https://sciencebites.org/>

expert in scientific domain is drastically different from that from a popularized version. Moreover, different levels of simplification, details, and explanation can be applied, e.g. for a given scientific article the Papier-Maché platform publishes two level of simplification: curiosity and advanced. Zhong et al. analyzed discourse factors related to sentence deletion on the Newsela corpus made of manually simplified sentences from news articles [39]. They found that professional editors utilize different strategies to meet the readability standards of elementary and middle schools.

The state-of-the-art in automatic summarization is achieved by deep learning models, in particular by pretrained Bidirectional Encoder Representations from Transformers (BERT) which can be used for both extractive and abstractive models [24]. It is important to study the limits of existing AI models, like GPT-2 [30] for English and CamemBERT for French [27], and how it is possible to overcome those limits. Recently, AI21 released the Jurassic-1 suite of language models, with 178B parameters for J1-Jumbo [23]. Jurassic-1 is a large AI model able to transform an existing text, e.g. in case of summarization. Multilingual T5 (mT5) is a large multilingual pretrained text-to-text transformer model developed by Google, covering 101 languages [36]. mT5 can be fine-tuned for any text-to-text generation, e.g. by using the SimpleT5 library.

## 2.2 Complexity spotting

Our analysis of the queries from different sources revealed the gap between the actual interest of the wide readership and the expectations of the journalists [28]. People are interested in biology or modern technologies as long as there is a connection with their everyday life. Thus, a simplified scientific text needs to contain references to the daily experience of people. On the one hand, the *subjective complexity of terminology* is involved when readers face concepts that go beyond their area of expertise and general knowledge, and need additional definitions or clarifications. On the other hand, the *objective complexity of terminology* is a systematic feature caused by complexity of research areas, research traditions and socio-cultural diversity. The complexity of a scientific area depends on peculiar attributes and conditions [34]. Ladyman et al. [22] suggest five such conditions: numerosity of elements, numerosity of interactions, disorder, openness, feedback. The complexity of terminology is also associated with a formal representation (*signifier*) of a term. Apart from borrowings, scientific text is rich in symbols and abbreviations (acronyms, backronyms, syllabic abbreviations, etc.) that are meant to optimize content transferring, standardize the naming of numerous elements, allow frequent interaction among them, and facilitate data processing. But readers of popularized publications expect explanations of the symbols and abbreviations.

One of the technical challenges here is thus term recognition. Robertson provided theoretical justifications of the term-weighting function IDF (inverse document frequency) in the traditional probabilistic model of information retrieval [31]. IDF shows term specificity and can be used for difficult term extraction as it is connected to the Zipf's law. WordNet [11] distance to the basic terms can be used as a measure of the term difficulty. Task-independent AI models, like GPT-2 [30], Jurassic-1 [23], Multilingual T5 [36], can be fine-tuned for the terminology extraction. It should be noted that there are tools available online such as OneClick Terms or TermoStat Web that allow us to extract intuitively mono and multiwords.

### 2.3 Text simplification

Existing works mainly focus on word/phrase-level (simplification of difficult words and constructions) [38,4,32,15,29,26] or sentence-level simplifications [40,35,7,41,5,33]. Koptient and Grabar analyzed the text transformation topology during simplification [21]. Among the most frequent transformations, they found synonymy, specification (insertion of information), generalization (deletion of information), pronominalization, substitution of adjectives by their corresponding nouns, and substitutions between singular and plural. In their further work, they proposed a rule-based system in French that combines lexical and syntactic simplification, for example, by transforming passive sentences into active sentences [19], and rating a lexicon [20]. Approaches based on rated lexicons are neither scalable nor robust to neologisms, which are frequent in scientific texts. Recent deep learning models with a large number of training parameters, like GPT-2 [30], Jurassic-1 [23], Multilingual T5 [36], can be applied for text simplification. Jurassic-1 Jumbo is the largest model publicly available with no waitlist. The AI21 studio’s playground provides ready-to-use prompts for text simplification (see De-Jargonizer). However, as Jiang et al. showed, a text simplification system depends on the quality and quantity of training data [18]. Therefore, a major step in training artificial intelligence (AI) text simplification models is the creation of high quality data.

Researchers have proposed various approaches based on expert judgment [6], readability levels [13,14], crowdsourcing [6,35,2], eye-tracking [37,16], manual annotation [17]. Traditional evaluation like comparison to the reference data by standard evaluation measures is difficult to apply as one should consider the end user (young readers, foreigners, non-experts, people with different literacy levels, people with cognitive disabilities etc.) as well as source document content.

## 3 Data set

In 2022, SimpleText’s data is two-fold: *Medicine* and *Computer Science*, as these two domains are the most popular on forums like ELI5 [28]. For both domains, we provide datasets according to our shared tasks:

- content selection relevant for non-experts;
- terminology complexity spotting in a given passage;
- simplified passages.

As in 2021, we use the Citation Network Dataset: DBLP+Citation, ACM Citation network (12th version) [1] as source of scientific documents to be simplified [8]. It contains: 4.894.083 bibliographic references published before 2020, 4.232.520 abstracts in English, 3.058.315 authors with their affiliations, and 45.565.790 ACM citations. Scientific textual content about any topic related to computer science can be extracted from this corpus together with authorship. Although we manually preselected abstracts for topics, participants also have access to use an ElasticSearch index. This Index is adequate to passage retrieval using BM25. Additional datasets have been extracted to generate Latent Dirichlet Allocation models for query expansion or train Graph Neural Networks for citation recommendation as carried out in StellarGraph<sup>13</sup> for example.

<sup>13</sup> <https://stellargraph.readthedocs.io/>

The shared datasets are: document abstract content for LDA or Word Embedding; document author relation for coauthoring analysis; document citation relation for co citation analysis; author citation relation for author impact factor analysis. These extra datasets are intended to be used to select passages by authors who are experts on the topic (highly cited by the community).

We propose 13 topics on *computer science* based on the recent  $n$  press titles from *The Guardian* enriched with keywords manually extracted from the content of the article (see Table 1). It has been checked that each keyword allows participants to extract at least 5 relevant abstracts. The use of these keywords is optional.

Query 12: Patient data from GP surgeries sold to US companies Topic 12.1: patient data
Query 13: Baffled by digital marketing? Find your way out of the maze Topic 13.1: digital marketing Topic 13.2: advertising

**Table 1.** Query examples

We selected passages that are adequate to be inserted as plain citations in the original journalistic article. The comparison of the journalistic articles with the scientific ones, as well as the analysis we carried out to choose topics, demonstrated that for non-experts the most important information is the application of an object (which problem can be solved? how to use this information/object? what are the examples?).

Text passages issued from abstracts on computer science were simplified by either a master student in Technical Writing and Translation or a pair of experts: (1) a computer scientist and (2) a professional translator, English native speaker but not specialist in computer science [8]. Each passage was discussed and rewritten multiple times until it became clear for non-computer scientists. Sentences were shortened, excluding every detail that was irrelevant or unnecessary to the comprehension of the study, and rephrased, using simpler vocabulary. If necessary, concepts were explained.

In 2022, we introduce new data based on Google Scholar and PubMed articles on muscle hypertrophy and health annotated by a master student in Technical Writing and Translation, specializing in these domains. The selected abstracts included the objectives of the study, the results and sometimes the methodology. The abstracts including only the topic of the study were excluded because of lack of information. To avoid the curse of knowledge, another master student in Technical Writing and Translation not familiar with the domain was solicited for complexity spotting.

## 4 Tasks

In 2022, SimpleText was transformed into a CLEF lab. We propose three shared tasks to help better understand the challenges as well as discuss these challenges and the way to evaluate solutions. Contributions should not exclusively rely on these shared tasks. We also welcome manual runs and runs within the unshared task.

Details on the tasks, guideline and call for contributions can be found at the SimpleText website. In this paper we just briefly introduce the planned shared tasks.

**TASK 1: What is in (or out)? Select passages to include in a simplified summary, given a query.** Based on an article from a major international newspaper general audience, this shared task aims to retrieve, from a large scientific bibliographic database with abstracts, all relevant passages to illustrate this article. Extracted passages should be adequate to be inserted as plain citations in the original paper. Sentence pooling and automatic metrics can be used to evaluate these results. The relevance of the source document can be evaluated as well as potential unresolved anaphora issues.

**TASK 2: What is unclear? Given a passage and a query, rank terms/concepts that are required to be explained for understanding this passage (definitions, context, applications,..).** The goal of this shared task is to decide which terms (up to 10) require explanation and contextualization to help a reader understand a complex scientific text — for example, with regard to a query, terms that need to be contextualized (with a definition, example and/or use-case). Terms should be ranked from 1 to 10 according to their complexity. 1 corresponds to the most difficult term, while lower ranks show that the term might be explained if there is space. Term pooling and automatic metrics (e.g. accuracy, NDCG, MSE, etc.) can be used to evaluate these results.

**TASK 3: Rewrite this! Given a query, simplify passages from scientific abstracts.** The goal of this shared task is to provide a simplified version of text passages. Participants are provided with queries and abstracts of scientific papers. The abstracts can be split into sentences as in the example. The simplified passages will be evaluated manually with use of aggregating metrics.

**UNSHARED TASK.** We welcome any submission that uses our data! This task is aimed at (but not limited to) Humanities, Social Science and Technical Communication. We encourage here manual and statistical analysis of content selection, readability and comprehensibility of simplified texts, terminology complexity analysis.

## 5 Conclusion and future work

The paper introduced the CLEF 2022 SimpleText track, containing three shared tasks and one unshared task on scientific text simplification. The created collection of simplified texts makes it possible to apply overlap metrics like ROUGE to text simplification. However, we will work on a new evaluation metric that can take into account unresolved anaphora [3] and information types. For the pilot task 2, participants will be asked to provide context for difficult terms. This context should provide a definition and take into account ordinary readers' needs to associate their particular problems with the opportunities that science provides them to solve the problems [28]. Full details about the lab can be found at the SimpleText website: <http://simpletext-project.com>. Help us to make scientific results understandable!

## 6 Acknowledgments

We thank Alain Kerhervé, University Translation Office, master students in Translation from the Université de Bretagne Occidentale, and the MaDICS research group.

## References

1. AMiner, <https://www.aminer.org/citation>
2. Alva-Manchego, F., Martin, L., Bordes, A., Scarton, C., Sagot, B., Specia, L.: Asset: A dataset for tuning and evaluation of sentence simplification models with multiple rewriting transformations. arXiv preprint arXiv:2005.00481 (2020)
3. Bellot, P., Moriceau, V., Mothe, J., SanJuan, E., Tannier, X.: INEX tweet contextualization task: Evaluation, results and lesson learned. *Inf. Process. Manage.* **52**(5), 801–819 (2016). <https://doi.org/10.1016/j.ipm.2016.03.002>, <https://doi.org/10.1016/j.ipm.2016.03.002>
4. Biran, O., Brody, S., Elhadad, N.: Putting it Simply: a Context-Aware Approach to Lexical Simplification. In: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. pp. 496–501. Association for Computational Linguistics, Portland, Oregon, USA (Jun 2011), <https://www.aclweb.org/anthology/P11-2087>
5. Chen, P., Rochford, J., Kennedy, D.N., Djamasbi, S., Fay, P., Scott, W.: Automatic Text Simplification for People with Intellectual Disabilities. In: *Artificial Intelligence Science and Technology*, pp. 725–731. WORLD SCIENTIFIC (Nov 2016). [https://doi.org/10.1142/9789813206823\\_0091](https://doi.org/10.1142/9789813206823_0091), [https://www.worldscientific.com/doi/abs/10.1142/9789813206823\\_0091](https://www.worldscientific.com/doi/abs/10.1142/9789813206823_0091)
6. De Clercq, Orphée and Hoste, Veronique and Desmet, Bart and van Oosten, Philip and De Cock, Martine and Macken, Lieve: Using the crowd for readability prediction. *NATURAL LANGUAGE ENGINEERING* **20**(3), 293–325 (2014), <http://dx.doi.org/10.1017/S1351324912000344>
7. Dong, Y., Li, Z., Rezagholizadeh, M., Cheung, J.C.K.: EditNTS: An Neural Programmer-Interpreter Model for Sentence Simplification through Explicit Editing. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. pp. 3393–3402. Association for Computational Linguistics, Florence, Italy (Jul 2019). <https://doi.org/10.18653/v1/P19-1331>, <https://www.aclweb.org/anthology/P19-1331>
8. Ermakova, L., Bellot, P., Braslavski, P., Kamps, J., Mothe, J., Nurbakova, D., Ovchinnikova, I., SanJuan, E.: Overview of SimpleText 2021 - CLEF Workshop on Text Simplification for Scientific Information Access. In: Candan, K.S., Ionescu, B., Goeuriot, L., Larsen, B., Müller, H., Joly, A., Maistro, M., Piroi, F., Faggioli, G., Ferro, N. (eds.) *Experimental IR Meets Multilinguality, Multimodality, and Interaction*. pp. 432–449. *Lecture Notes in Computer Science*, Springer International Publishing, Cham (2021). [https://doi.org/10.1007/978-3-030-85251-1\\_27](https://doi.org/10.1007/978-3-030-85251-1_27)
9. Ermakova, L., Bellot, P., Braslavski, P., Kamps, J., Mothe, J., Nurbakova, D., Ovchinnikova, I., SanJuan, E.: Text Simplification for Scientific Information Access: CLEF 2021 SimpleText Workshop. In: *Advances in Information Retrieval - 43rd European Conference on {IR} Research, {ECIR} 2021, Lucca, Italy, March 28 – April 1, 2021, Proceedings*. Lucca, Italy (2021)
10. Ermakova, L., Bordignon, F., Turenne, N., Noel, M.: Is the Abstract a Mere Teaser? Evaluating Generosity of Article Abstracts in the Environmental Sciences. *Frontiers in Research Metrics and Analytics* **3** (2018). <https://doi.org/10.3389/frma.2018.00016>, <https://www.frontiersin.org/articles/10.3389/frma.2018.00016/full>
11. Fellbaum, C. (ed.): *WordNet: An Electronic Lexical Database*. Language, Speech, and Communication, MIT Press, Cambridge, MA (1998)
12. Fontelo, P., Gavino, A., Sarmiento, R.F.: Comparing data accuracy between structured abstracts and full-text journal articles: implications in their use for informing clinical decisions. *Evidence-based medicine* **18**(6), 207–11 (2013). <https://doi.org/10.1136/eb-2013-101272>, <http://www.researchgate.net/publication/>



- [240308203\\_Comparing\\_data\\_accuracy\\_between\\_structured\\_abstracts\\_and\\_full-text\\_journal\\_articles\\_implications\\_in\\_their\\_use\\_for\\_informing\\_clinical\\_decisions](#)
13. François, T., Fairon, C.: Les apports du tal à la lisibilité du français langue étrangère. *Trait. Autom. des Langues* **54**, 171–202 (2013)
  14. Gala, N., François, T., Fairon, C.: Towards a french lexicon with difficulty measures: Nlp helping to bridge the gap between traditional dictionaries and specialized lexicons. In: *eLex-Electronic Lexicography* (2013)
  15. Glavaš, G., Štajner, S.: Simplifying Lexical Simplification: Do We Need Simplified Corpora? In: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. pp. 63–68. Association for Computational Linguistics, Beijing, China (Jul 2015). <https://doi.org/10.3115/v1/P15-2011>, <https://www.aclweb.org/anthology/P15-2011>
  16. Grabar, N., Farce, E., Sparrow, L.: Study of readability of health documents with eye-tracking approaches. In: *1st Workshop on Automatic Text Adaptation (ATA)* (2018)
  17. Grabar, N., Hamon, T.: A large rated lexicon with french medical words. In: *LREC (Language Resources and Evaluation Conference) 2016* (2016)
  18. Jiang, C., Maddela, M., Lan, W., Zhong, Y., Xu, W.: Neural CRF Model for Sentence Alignment in Text Simplification. *arXiv:2005.02324 [cs]* (Jun 2020), <http://arxiv.org/abs/2005.02324>, arXiv: 2005.02324
  19. Koptient, A., Grabar, N.: Fine-grained text simplification in French: steps towards a better grammaticality. In: *ISHIMR Proceedings of the 18th International Symposium on Health Information Management Research*. Kalmar, Sweden (Sep 2020). <https://doi.org/10.15626/ishimr.2020.xxx>, <https://hal.archives-ouvertes.fr/hal-03095247>
  20. Koptient, A., Grabar, N.: Rated Lexicon for the Simplification of Medical Texts. In: *The Fifth International Conference on Informatics and Assistive Technologies for Health-Care, Medical Support and Wellbeing HEALTHINFO 2020*. Porto, Portugal (Oct 2020), <https://hal.archives-ouvertes.fr/hal-03095275>
  21. Koptient, A., Grabar, N.: Typologie de transformations dans la simplification de textes. In: *Congrès mondial de la linguistique française*. Montpellier, France (Jul 2020), <https://hal.archives-ouvertes.fr/hal-03095235>
  22. Ladyman, J., Lambert, J., Wiesner, K.: What is a complex system? *European Journal for Philosophy of Science* **3**(1), 33–67 (Jan 2013). <https://doi.org/10.1007/s13194-012-0056-8>, <https://doi.org/10.1007/s13194-012-0056-8>
  23. Lieber, O., Sharir, O., Lentz, B., Shoham, Y.: Jurassic-1: Technical Details and Evaluation p. 9
  24. Liu, Y., Lapata, M.: Text Summarization with Pretrained Encoders. *arXiv:1908.08345 [cs]* (Sep 2019), <http://arxiv.org/abs/1908.08345>, arXiv: 1908.08345
  25. Maddela, M., Alva-Manchego, F., Xu, W.: Controllable Text Simplification with Explicit Paraphrasing. *arXiv:2010.11004 [cs]* (Apr 2021), <http://arxiv.org/abs/2010.11004>, arXiv: 2010.11004
  26. Maddela, M., Xu, W.: A Word-Complexity Lexicon and A Neural Readability Ranking Model for Lexical Simplification. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. pp. 3749–3760. Association for Computational Linguistics, Brussels, Belgium (2018). <https://doi.org/10.18653/v1/D18-1410>, <https://www.aclweb.org/anthology/D18-1410>
  27. Martin, L., Muller, B., Ortiz Suárez, P.J., Dupont, Y., Romary, L., de la Clergerie, E., Seddah, D., Sagot, B.: CamemBERT: a Tasty French Language Model. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. pp. 7203–7219. Association for Computational Linguistics, Online (2020). <https://doi.org/10.18653/v1/2020.acl-main.645>, <https://www.aclweb.org/anthology/2020.acl-main.645>

28. Ovchinnikova, I., Nurbakova, D., Ermakova, L.: What science-related topics need to be popularized? A comparative study. In: Faggioli, G., Ferro, N., Joly, A., Maistro, M., Piroi, F. (eds.) Proceedings of the Working Notes of CLEF 2021 - Conference and Labs of the Evaluation Forum, Bucharest, Romania, September 21st - to - 24th, 2021. CEUR Workshop Proceedings, vol. 2936, pp. 2242–2255. CEUR-WS.org (2021), <http://ceur-ws.org/Vol-2936/paper-203.pdf>
29. Paetzold, G., Specia, L.: Lexical Simplification with Neural Ranking. In: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers. pp. 34–40. Association for Computational Linguistics, Valencia, Spain (Apr 2017), <https://www.aclweb.org/anthology/E17-2006>
30. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I.: Language Models are Unsupervised Multitask Learners p. 24
31. Robertson, S.: Understanding inverse document frequency: on theoretical arguments for IDF. *Journal of Documentation* **60**(5), 503–520 (Jan 2004). <https://doi.org/10.1108/00220410410560582>, <https://doi.org/10.1108/00220410410560582>, publisher: Emerald Group Publishing Limited
32. Specia, L., Jauhar, S.K., Mihalcea, R.: SemEval-2012 Task 1: English Lexical Simplification. In: \*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012). pp. 347–355. Association for Computational Linguistics, Montréal, Canada (2012), <https://www.aclweb.org/anthology/S12-1046>
33. Wang, T., Chen, P., Rochford, J., Qiang, J.: Text Simplification Using Neural Machine Translation. *Proceedings of the AAAI Conference on Artificial Intelligence* **30**(1) (Mar 2016), <https://ojs.aaai.org/index.php/AAAI/article/view/9933>, number: 1
34. Wiesner, K., Ladyman, J.: Measuring complexity. arXiv:1909.13243 [nlin] (Sep 2020), <http://arxiv.org/abs/1909.13243>, arXiv: 1909.13243
35. Xu, W., Napoles, C., Pavlick, E., Chen, Q., Callison-Burch, C.: Optimizing statistical machine translation for text simplification. *Transactions of the Association for Computational Linguistics* **4**, 401–415 (2016), publisher: MIT Press
36. Xue, L., Constant, N., Roberts, A., Kale, M., Al-Rfou, R., Siddhant, A., Barua, A., Rafel, C.: mT5: A massively multilingual pre-trained text-to-text transformer. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 483–498. Association for Computational Linguistics, Online (Jun 2021). <https://doi.org/10.18653/v1/2021.naacl-main.41>, <https://aclanthology.org/2021.naacl-main.41>
37. Yaneva, V., Temnikova, I., Mitkov, R.: Accessible texts for autism: An eye-tracking study. In: Proceedings of the 17th International ACM SIGACCESS Conference on Computers & Accessibility. pp. 49–57 (2015)
38. Yatskar, M., Pang, B., Danescu-Niculescu-Mizil, C., Lee, L.: For the sake of simplicity: Unsupervised extraction of lexical simplifications from Wikipedia. In: Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics. pp. 365–368. Association for Computational Linguistics, Los Angeles, California (Jun 2010), <https://www.aclweb.org/anthology/N10-1056>
39. Zhong, Y., Jiang, C., Xu, W., Li, J.J.: Discourse Level Factors for Sentence Deletion in Text Simplification. *Proceedings of the AAAI Conference on Artificial Intelligence* **34**(05), 9709–9716 (Apr 2020). <https://doi.org/10.1609/aaai.v34i05.6520>, <https://ojs.aaai.org/index.php/AAAI/article/view/6520>, number: 05
40. Zhu, Z., Bernhard, D., Gurevych, I.: A Monolingual Tree-based Translation Model for Sentence Simplification. In: Proceedings of the 23rd International Conference on Computational

- Linguistics (Coling 2010). pp. 1353–1361. Coling 2010 Organizing Committee, Beijing, China (Aug 2010), <https://www.aclweb.org/anthology/C10-1152>
41. Štajner, S., Nisioi, S.: A Detailed Evaluation of Neural Sequence-to-Sequence Models for In-domain and Cross-domain Text Simplification. In: Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018). European Language Resources Association (ELRA), Miyazaki, Japan (May 2018), <https://www.aclweb.org/anthology/L18-1479>