



This is a repository copy of *Horacio Saggion, automatic text simplification. Synthesis lectures on human language technologies, April 2017. 137 pages, ISBN:1627058680 9781627058681.*

White Rose Research Online URL for this paper:
<https://eprints.whiterose.ac.uk/155251/>

Version: Accepted Version

Article:

Scarton, C. orcid.org/0000-0002-0103-4072 (2019) Horacio Saggion, automatic text simplification. *Synthesis lectures on human language technologies*, April 2017. 137 pages, ISBN:1627058680 9781627058681. *Natural Language Engineering*, 26 (4). pp. 489-492. ISSN 1351-3249

<https://doi.org/10.1017/s1351324919000603>

This article has been published in a revised form in *Natural Language Engineering* [<https://doi.org/10.1017/S1351324919000603>]. This version is free to view and download for private research and study only. Not for re-distribution, re-sale or use in derivative works. © Cambridge University Press 2019.

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

Horacio Saggion, Automatic Text Simplification. Synthesis Lectures on Human Language Technologies, April 2017, 137 pages, ISBN:1627058680 9781627058681

Automatic text simplification (TS) is a text-to-text transformation task where the aim is to produce a simpler version of an original text. There are several important aspects to such a task:

- Audience: It is important to know the audience for which the simplified text is intended. Different audiences require different types of simplification operations. For example, the simplification requirements for a second language learner of English may be different from the requirements of a dyslexic person.
- Simplification type: Traditionally, researchers have divided TS in two classes, lexical and syntactic simplification. Lexical simplification (LS) consists of the replacement of words and phrases in the text, while SS encompasses operations at the sentence structure, such as splitting a sentence or changing it from passive to active voice. However, new state-of-the-art approaches deal with both phenomena (machine-translation (MT)-like systems, e.g.).
- Method: As with other areas in Natural Language Processing (NLP), TS can also be explored through rule-based, data-driven or hybrid approaches.

In this book, the author draws largely on his own experiences in order to address the basics of TS and to describe traditional work on this topic. In the introduction chapter, TS is defined alongside explanations of lexical and syntactic types of simplification. This chapter also includes a discussion about *how* texts are simplified by humans and *what* we can then learn in order to automatise the process. Finally, the motivation behind developing TS techniques is presented as a task with social impact that can enable different audiences to access different types of information.

Chapter 2 provides an important overview of work on Readability Assessment (RA), that is the analysis of textual complexity. Although this topic may not be always linked to TS, knowing the complexity of a text or sentence can be seen as a pre-processing step for TS or as part of its evaluation. Although previous work in TS has highlighted the problems regarding RA shallow metrics (e.g., Shardlow, 2014), further development in this area can certainly improve the results and evaluation of TS systems. Perhaps, the biggest contribution in this chapter is Section 2.3, where advanced approaches for RA are discussed.

LS is presented in Chapter 3, which highlights the work of Carroll et al. (1998) as the first approach for LS. The task of complex word identification (CWI) also appears in this chapter, where systems from the SemEval 2016 CWI shared task are detailed (Specia et al., 2012). The author also highlights that the availability of parallel data from original and simplified texts has promoted the new generation of TS systems, which could then rely on machine learning (ML) approaches. Notably, the most widely used parallel data are the Wikipedia-Simple Wikipedia (W-SW)-aligned dataset (e.g., Coster and Kauchak, 2011). LS through language modelling (also discussed in this chapter) is another approach only possible due to the acquisition of large amounts of data. Finally, an interesting topic discussed in Chapter 3 is the simplification of numerical expressions in texts. Although this may not be a trendy topic in the LS area, the author highlights the importance of this task for some audiences. In general, this chapter gives good insights about LS, including the author's own work for languages other than English (in this case, Spanish). Paetzold (2016), however, proposes a more structured way of presenting LS and a more complete survey of the topic.

Chapter 4 is dedicated to SS and is a good overview of how rule-based systems work for this task. The discussion about the work of Siddharthan (2006) and Siddharthan (2011) represents a useful survey for researchers wanting a quick reference for the topic. The detailed description of how a rule-based system for SS works, and how it can be used

through a well-known toolkit for NLP – General Architecture for Text Engineering (GATE) (Maynard et al., 2002) – makes this chapter even more useful for other researchers.

State-of-the-art approaches for TS are currently data-driven, and Chapter 5 reviews some of the work that has focused on these. TS can be viewed as an MT task, where the original document is the source and the simplified document is the target. Therefore, with the availability of parallel data for TS (such as W-SW), it is natural to explore MT-based approaches for this task (e.g. Zhu et al. (2010); Coster and Kauchak (2011)). Another approach explored by Woodsend and Lapata (2011) proposes learning TS rules from parallel data using quasi-synchronous grammars. Finally, Narayan and Gardent (2014) use a hybrid approach that relies on semantic information for split and delete operations, and an MT-based model for modelling paraphrasing.

Although TS is usually motivated using social aspects (e.g., access to information by everyone), very little work has actually resulted in tools and/or has been evaluated by its intended target audiences. Chapter 6 presents three large projects that were successful in at least one of the points above:

- PSET (e.g., Carroll et al., 1998): TS systems to adapt texts for aphasics -- English language.
- Simplext (Saggion et al., 2015): TS systems targeting people with intellectual disabilities -- Spanish language.
- PorSimples (e.g., Aluísio and Gasperin, 2010): TS systems for people with low literacy -- Portuguese language.

Chapter 7 continues to present the usefulness of TS either as a tool to help different target audiences (e.g., people with dyslexia; Rello, 2014), or as a pre-processing step to improve other NLP tasks (e.g., parsing (Jonnalagadda et al., 2009)). Work on TS for specific domains, such as the medical one (Ong et al., 2007), is also discussed in this chapter.

Finally, Chapter 8 is dedicated to presenting resources and tools available for TS, and evaluation approaches for this task. Resources for RA and parallel corpora for TS are described, including a dedicated section on the Simplext dataset and a brief discussion about the Newsela dataset. The LEXenstein toolkit for LS (Paetzold and Specia, 2015) is also presented and discussed. For evaluation, the traditional three-way human evaluation is presented, where human judges are asked to give a *Likert* score (usually from 1 to 5) in order to assess grammaticality, meaning preservation, and simplicity of automatically simplified texts. The author then correctly criticises work that uses RA metrics to evaluate simplified sentences, since such metrics are designed to work on longer texts. Automatic evaluation metrics from the MT area, such as BLEU (BiLingual Evaluation Understudy) (Papineni et al., 2002), are also discussed, as these have been largely employed by work on data-driven TS. Unfortunately, the author does not include the work of Xu et al. (2016) that presents SARI -- System output Against References and Input sentence -- as an automatic metric more adequate for TS. Similar to BLEU, SARI is a n-gram based metric that takes into account simplification references produced by humans and the original text in order to evaluate the output of a TS system. Data-driven work on TS is now mainly evaluated in terms of SARI, since BLEU is proven to be inadequate (e.g. Sulem et al., 2018a). Finally, the findings of a shared task on quality assessment for TS (QATS) (Stajner et al., 2016) are presented. The idea in QATS derives from work for quality estimation of MT (e.g. Specia et al., 2018), where ML approaches are used to build models using human assessments as labels. The ultimate goal is to create models able to generalise and automatically predict aspects such as grammaticality, meaning preservation and simplicity for unseen data points.

In summary, this book presents a useful overview of the foundation work on TS. Mainly, Chapters 2, 3 and 4 are great contributions for researchers who are either new to the topic or need to find the best references to the topics discussed. Chapters 6 and 7 also contain important information about TS, its applications and successful projects. Nevertheless, TS has been evolving at a very fast pace since this book was published. Between 2017 and

now, over 10 new approaches for TS have been proposed for the English language alone, thanks to the advance of neural deep learning techniques (e.g. Nisioi et al., 2017; Zhang and Lapata, 2017; Alva-Manchego et al., 2018; Vu et al., 2018; Guo et al., 2018; Sulem et al., 2018b; Zhao et al., 2018; Scarton and Specia, 2018; Kriz et al., 2019; Dong et al., 2019; Surya et al., 2019). Additionally, apart from SARI, SAMSA -- Simplification Automatic evaluation Measure through Semantic Annotation -- (Sulem et al., 2018c) was also proposed as a new metric for TS evaluation that uses semantic information in order to better assess sentence-level operations such as splitting. Therefore, although the reader needs to be aware of the limitations imposed on this book by the fast-growing deep learning movement in NLP and also by the growing interest in TS by the NLP community, this book nevertheless represents a useful reference of traditional work in TS.

REFERENCES

- Aluísio, Sandra Maria and Caroline Gasperin. 2010. Fostering digital inclusion and accessibility: The PorSimples project for simplification of Portuguese texts. In Proceedings of the NAACL HLT Young Investigators Workshop on Computational Approaches to Languages of the Americas, pages 46-53. Association for Computational Linguistics.
- Alva-Manchego, Fernando, Joachim Bingel, Gustavo Paetzold, Carolina Scarton, and Lucia Specia. 2017. Learning how to simplify from explicit labeling of complex-simplified text pairs. In Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 295–305, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Carroll, John, Guido Minnen, Yvonne Canning, Siobhan Devlin, and John Tait. 1998. Practical simplification of English newspaper text to assist aphasic readers. In Proceedings of AAAI-98 Workshop on Integrating Artificial Intelligence and Assistive Technology, pages 7–10, Madison, WI.
- Coster, William and David Kauchak. 2011. Learning to simplify sentences using Wikipedia. In Proceedings of the Workshop on Monolingual Text-To-Text Generation, MTTG '11, pages 1–9, ACL, Portland, Oregon. Association for Computational Linguistics.
- Dong, Yue, Zichao Li, Mehdi Rezagholizadeh, Jackie Chi Kit Cheung. 2019. EditNTS: A Neural Programmer-Interpreter Model for Sentence Simplification through Explicit Editing. Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 3393-3402, Florence, Italy. Association for Computational Linguistics.
- Guo, Han, Ramakanth Pasunuru, and Mohit Bansal. 2018. Dynamic multi-level multi-task learning for sentence simplification. In Proceedings of the 27th International Conference on Computational Linguistics, pages 462–476, Santa Fe, NM. Association for Computational Linguistics.
- Jonnalagadda, Siddhartha, Luis Tari, Jörg Hakenberg, Chitta Baral, Graciela Gonzalez. 2009. Towards Effective Sentence Simplification for Automatic Processing of Biomedical Text. Proceedings of NAACL HLT 2009: Short Papers, pages 177–180, Boulder, CO. Association for Computational Linguistics.
- Kriz, Reno João Sedoc, Marianna Apidianaki, Carolina Zheng, Gaurav Kumar, Eleni Miltsakaki, and Chris Callison-Burch. 2019. Complexity-Weighted Loss and Diverse Reranking for Sentence Simplification. Proceedings of NAACL-HLT 2019, pages 3137-3147, Florence, Italy. Association for Computational Linguistics.
- Maynard, Diana, Valentin Tablan, Hamish Cunningham, Cristian Ursu, Horacio Saggion, Kalina Bontcheva, and Yorick Wilks. Architectural elements of language engineering robustness. *Natural Language Engineering*, 8(2/3): 257-274, 2002.
- Narayan, Shashi and Claire Gardent. 2014. Hybrid simplification using deep semantics and machine translation. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 435–445, Association for Computational Linguistics, Baltimore, Maryland.
- Nisioi, Sergiu, Sanja Štajner, Simone Paolo Ponzetto, and Liviu P. Dinu. 2017. Exploring neural text simplification models. In Proceedings of the 55th Annual Meeting of the

- Association for Computational Linguistics (Volume 2: Short Papers), pages 85–91, Vancouver, Canada. Association for Computational Linguistics.
- Ong, Ethel, Jerwin Damay, Gerard Lojico, Kimberly Lu, and Dex Tarantan. 2007. Simplifying text in medical literature. *Journal of Research in Science, Computing and Engineering*, 4(1):37-47.
- Paetzold, Gustavo Henrique and Specia, Lucia. 2015. LEXenstein: A Framework for Lexical Simplification. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics, pages 85-90, Beijing, China. Association for Computational Linguistics.
- Paetzold, Gustavo Henrique. 2016. Lexical Simplification for Non-native English Speakers. PhD thesis, The University of Sheffield, Sheffield, UK.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: A method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL '02, pages 311–318, Philadelphia, Pennsylvania. Association for Computational Linguistics.
- Rello, Luz. 2014. DysWebxia. A Text Accessibility Model for People with Dyslexia. PhD thesis, Universitat Pompeu Fabra, Barcelona, Spain.
- Saggion, Horacio, Sanja Stajner, Stefan Bott, Simon Mille, Luz Rello, and Biljana Drndarevic. 2015. Making it Simplext: Implementation and evaluation of a text simplification system for Spanish. *ACM Transactions on Accessible Computing (TACCESS)*, 6(4):14.
- Scarton, Carolina and Lucia Specia. 2018. Learning Simplifications for Specific Target Audiences. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 712–718, Melbourne, Australia. Association for Computational Linguistics.
- Shardlow, Matthew. 2014. A Survey of Automated Text Simplification. *International Journal of Advanced Computer Science and Applications (IJACSA)*, Special Issue on Natural Language Processing 4(1):58–70.
- Siddharthan, Advaith. 2006. Syntactic simplification and text cohesion. *Research on Language and Computation*, 4(1):77-109.
- Siddharthan, Advaith. 2011. Text simplification using typed dependencies: A comparison of the robustness of different generation strategies. In Proceedings of the 13th European Workshop on Natural Language Generation, ENLG '11, pages 2–11, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Specia, Lucia, Carolina Scarton, Gustavo Henrique Paetzold. 2018. Quality Estimation of Machine Translation. *Synthesis Lectures on Human Language Technologies*. Morgan & Claypool Publishers, San Raphael, CA.
- Specia, Lucia, Sujay Kumar Jauhar, Rada Mihalcea. 2012. SemEval 2012 task 1: English lexical simplification. In Proceedings of the 1st Joint Conference on Lexical Computational Semantics, SemEval, pages 347-355, Montréal, Canada. Association for Computational Linguistics.
- Stajner, Sanja, Maja Popovic, and Hanna Béchara. 2016. Quality estimation for text simplification. In Proceedings of the Workshop and Shared Task on Quality Assessment for Text Simplification (QATS), Pororoz, Slovenia.
- Sulem, Elior, Omri Abend, and Ari Rappoport. 2018a. Bleu is not suitable for the evaluation of text simplification. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 738–744, Association for Computational Linguistics.
- Sulem, Elior, Omri Abend, and Ari Rappoport. 2018b. Simple and Effective Text Simplification Using Semantic and Neural Methods. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 162–173, Melbourne, Australia. Association for Computational Linguistics.
- Sulem, Elior, Omri Abend, and Ari Rappoport. 2018c. Semantic structural evaluation for text simplification. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 685-696, Association for Computational Linguistics, New Orleans, Louisiana.

- Surya, Sai, Abhijit Mishra, Anirban Laha, Parag Jain, Karthik Sankaranarayanan. 2019. Unsupervised Neural Text Simplification. Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 2058-2068, Florence, Italy. Association for Computational Linguistics.
- Vu, Tu, Baotian Hu, Tsendsuren Munkhdalai, and Hong Yu. 2018. Sentence simplification with memory-augmented neural networks. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 79-85, Association for Computational Linguistics, New Orleans, Louisiana.
- Xu, Wei, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. Optimizing statistical machine translation for text simplification. Transactions of the Association for Computational Linguistics, 4:401-415.
- Woodsend, Kristian and Mirella Lapata. 2011a. Learning to Simplify Sentences with Quasi-Synchronous Grammar and Integer Programming. In Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, pages 409-420, Edinburgh, Scotland, UK. Association for Computational Linguistics.
- Zhang, Xingxing and Mirella Lapata. 2017. Sentence simplification with deep reinforcement learning. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 595-605, Copenhagen, Denmark. Association for Computational Linguistics.
- Zhao, Sanqiang, Rui Meng, Daqing He, Saptono Andi, Parmanto Bambang. 2018. Integrating Transformer and Paraphrase Rules for Sentence Simplification. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 3164- 3173, Brussels, Belgium. Association for Computational Linguistics.
- Zhu, Zhemin, Delphine Bernhard, and Iryna Gurevych. 2010. A monolingual tree-based translation model for sentence simplification. In Proceedings of the 23rd International Conference on Computational Linguistics, COLING '10, pages 1353-1361, Stroudsburg, PA. Association for Computational Linguistics.

Dr. Carolina Scarton
Department of Computer Science, University of Sheffield, Sheffield, UK

c.scarton@sheffield.ac.uk