

Target-Level Sentence Simplification as Controlled Paraphrasing

Tannon Kew Sarah Ebling

Department of Computational Linguistics,

University of Zurich

{kew, ebling}@cl.uzh.ch

Abstract

Automatic text simplification aims to reduce the linguistic complexity of a text in order to make it easier to understand and more accessible. However, simplified texts are consumed by a diverse array of target audiences and what might be appropriately simplified for one group of readers may differ considerably for another. In this work we investigate a novel formulation of sentence simplification as paraphrasing with controlled decoding. This approach aims to alleviate the major burden of relying on large amounts of in-domain parallel training data, while at the same time allowing for modular and adaptive simplification. According to automatic metrics, our approach performs competitively against baselines that prove more difficult to adapt to the needs of different target audiences or require significant amounts of complex-simple parallel aligned data.

1 Introduction

Automatic text simplification (ATS) aims to reduce the linguistic complexity of a text while preserving its meaning in order to make it easier to understand and more accessible to a wider array of potential readers (Bingel and Søgaard, 2016; Sikka and Mago, 2020). These readers might include children or adults with low literacy levels, cognitive impairments, or a lack of specialist knowledge in certain topics, as well as non-native language learners (Štajner, 2021; Saggion, 2017). However, the notion of exactly what constitutes ‘simplified’ text is highly subjective and can differ considerably between different types of readers. Thus it is important to tailor solutions appropriately in order to accommodate the needs of specific target audiences.

Research on ATS, or more specifically sentence simplification (SS), has been spurred on by performance gains in neural sequence-to-sequence (seq2seq) language generation methods (Zhang and Lapata, 2017; Scarton and Specia, 2018), which

aim to do away with complex hand-crafted rules (De Belder and Moens, 2010; Siddharthan and Mandya, 2014) and improve on earlier statistical approaches (Wubben et al., 2012; Xu et al., 2016). However, fully supervised seq2seq SS approaches require a large amount of sentence-aligned parallel training data (Koehn and Knowles, 2017), which remains relatively scarce and difficult to attain.

For this reason, much work has focused on certain aspects of ATS such as lexical (Glavaš and Štajner, 2015; Kriz et al., 2018) or structural simplification (Niklaus et al., 2019; Garain et al., 2019; Narayan et al., 2017; Gao et al., 2021). Others have aimed to make better use of limited parallel complex-simple training sentences by using more sample-efficient modelling techniques that aim to predict and execute in-place edit operations (Omelianchuk et al., 2021; Dong et al., 2019). Meanwhile, despite considerable similarities between SS and paraphrasing, the task of reformulating a sentence while maintaining an equivalent meaning (Bhagat and Hovy, 2013), relatively few works have aimed to exploit paraphrases to bootstrap seq2seq-based simplification (Martin et al., 2020; Maddela et al., 2021).

We investigate this last line of work and consider an alternative framing of SS as the task of *controlled* paraphrasing. We train a large-scale paraphrase model capable of producing high-quality and diverse paraphrases and combine it with future discriminators for generation (FUDGE) (Yang and Klein, 2021) to control decoding and steer the generated paraphrase towards a specific target level for text simplification. Our experiments show that this proves to be an effective approach for generating simplified sentences for different target audiences without requiring any parallel data in the form of complex-simple aligned sentence pairs. The code and model outputs from this work are made available at <https://github.com/ZurichNLP/SimpleFUDGE>.

2 Background & Motivation

Perhaps the largest hurdle for seq2seq-based simplification is the collection of appropriately aligned complex-simple parallel data required for training robust and reliable systems (Laban et al., 2021). Furthermore, as discussed by Štajner (2021), ATS systems should be developed to support a variety of target readers and would thus benefit from modular approaches that allow for easy customisation and adaptation. Recently, however, large pre-trained generation models have continued to demonstrate impressive performance when finetuned on conditional language generation tasks (Raffel et al., 2020; Lewis et al., 2020). Along with this, there has been considerable work done on exploring ways to better control the outputs of large generative models in order to achieve certain communicative goals (Dathathri et al., 2020; Krause et al., 2021; Liu et al., 2021; Yang and Klein, 2021; Pascual et al., 2021). We see a clear link between these recent developments and the challenges associated with SS and set out to investigate a modular approach suitable for simplifying text for different target audiences and that relaxes the need for complex-simple parallel training data.

3 Method

Given a complex source sentence, our goal is to transform it into a simplified target sequence¹ that preserves its meaning. Under the traditional seq2seq framework, a target sequence $y = \{y_1, \dots, y_T\}$ can be generated autoregressively as a series of conditional probabilities over the vocabulary, whereby each target token y_i is conditioned on the source sentence $x = \{x_1, \dots, x_T\}$ and any preceding target tokens $y_{1:i-1}$,

$$P(y) = \prod_{i=1}^n P(y_i|x, y_{1:i-1}). \quad (1)$$

To ensure that the generated target sequence is appropriately simplified, we employ FUDGE (Yang and Klein, 2021), which has been shown to be effective for various controlled generation tasks. FUDGE introduces a lightweight classifier \mathcal{B} to control for a desired target attribute a during autoregressive generation with a model \mathcal{G} . In essence, it modifies the conditional probability in Equation 1 with the following Bayesian factorisation:

¹Since an appropriate simplified formulation may consist of multiple shorter sentences we refer to it as a sequence.

$$P(y_i|x, y_{1:i-1}, a) \propto P(a|y_{1:i})P(y_i|x, y_{1:i-1}). \quad (2)$$

Here, the second term is the unmodified prediction from \mathcal{G} which is combined with the conditional probability of a given all possible continuations at the current timestep i according to \mathcal{B} . For further details on FUDGE, we refer the reader to Yang and Klein (2021).

3.1 FUDGE for Target-Level Simplification

To leverage FUDGE for target-level SS, we train a classifier for *each* target level, i.e. \mathcal{B}_{Simp-l} , and combine them with the same underlying generator model \mathcal{G} . Following Yang and Klein (2021), each classifier is trained as a binary predictor on labelled subsequences of complex (Simp-0) and simple (Simp- l) texts. Since SS often involves breaking down a long complex sentence into smaller atomic sentences (Honeyfield, 1977), we train each classifier to predict labels on subsequences pertaining to consecutive sentences within a paragraph. This ensures that the classifier’s predictions do not unduly bias the generation of the end of sentence symbol ‘</S>’ after producing sentence-final punctuation.

For the underlying generator, \mathcal{G} , we fine-tune BART-large on approximately 1.4 million paraphrase sentence pairs mined from the web.² This model has no explicit knowledge of complex vs. simple language. To ensure a fair comparison to previous work, we use the exact same training data as Martin et al. (2021) and aim to keep training hyperparameters as consistent as possible (detailed in Appendix C).

In practice, combining the predictions from \mathcal{G} and \mathcal{B} relies on a single weight parameter λ . For our experiments, we derive suitable values for each target-level by sweeping over possible whole number values in the range [0,10] and selecting the best according to SARI on the validation set (see Appendix D).

4 Experimental Setup

4.1 Data

We conduct our experiments on the Newsela corpus of simplified news articles.³ In its current form, the

²In theory, given BART’s denoising autoencoding pre-training, it could also be possible to avoid fine-tuning altogether. However, initial experiments showed that the probability distribution of the off-the-shelf BART model is far too peaked for the classifier’s predictions to have any effect.

³<https://newsela.com/data/>.

	# articles	# manually aligned sentences			
		Simp-1	Simp-2	Simp-3	Simp-4
train	1,862	-	-	-	-
train	35	1,341	1,245	1,042	841
test	10	365	353	309	256
valid	5	180	163	134	87

Table 1: Newsela English corpus articles and their *manually* aligned sentences from Jiang et al. (2020) for Simp-0 to Simp-*l*.

corpus contains 1,912 English news articles that have been professionally re-written according to readability guidelines for children at multiple grade levels (Xu et al., 2015). Article versions range from Simp-0 to Simp-4, with the former referring to the original, unsimplified article, suitable for upper secondary school grades, and the latter indicating the simplest versions, suitable for lower primary school grades.⁴

While Newsela provides complex-simple alignments at the document level, it must be emphasised that this alignment is *not* a requirement for our approach and thus training examples are randomly shuffled each epoch. Nevertheless, we reason that this type of alignment is beneficial since it ensures that attribute classifiers are trained on comparable examples covering the same domains. As a consequence, each classifier must learn to distinguish between complex and simple text based on relevant characteristics, such as the lexical choices and grammatical structures for a target level, rather than exploiting potentially misleading differences in topical content (Kumar et al., 2019).

For automatic evaluation purposes, however, sentence-level alignments are a must. To this end, we make use of the manually aligned test and validation splits provided by Jiang et al. (2020). Setting aside all sentence pairs from these splits ensures that no unwanted data leakage occurs. An overview of the corpus and manually aligned sentence pairs is provided in Table 1.

4.2 Baselines

We compare our approach to two recently proposed techniques for controlled SS and a naive baseline that only uses paraphrasing.

⁴In this work, we assume that article versions (0-4) are reliable indicators of level-appropriate simplifications and provide a detailed discussion on this assumption in Appendix B.

MUSS Martin et al. (2021) leveraged large-scale paraphrase data to fine-tune BART-large in combination with the ACCESS control method for simplification (Martin et al., 2020). This method relies on four control tokens which are prepended to the source sequence and used to indicate the desired length, edit distance, lexical complexity and syntactic complexity as a ratio between the source and the output sequence. At inference time, these special tokens act as ‘control knobs’ for the simplification. Following Martin et al. (2021), we derive optimal control values through a parameter search on the validation split (see Appendix C.4).

SUPER Following Scarton and Specia (2018), we also train a level-aware supervised baseline with a single special token indicating the target level (e.g. $\langle \text{L3} \rangle = \text{Simp-3}$) prepended to each source sentence. For a fair comparison, we initialise this model from the same BART-large checkpoint as the other two models and fine-tune on the manually aligned training sentences for all Newsela levels simultaneously. This amounts to a low resource setting with a total of 4,469 training instances.

PARA In addition, we also compare to a straightforward paraphrase generated by our underlying generation model \mathcal{G} with no controls or intervention.

4.3 Evaluation Metrics

Reliably evaluating SS is an open challenge (Alva-Manchego et al., 2021). However, a range of both reference-based and reference-less automatic metrics have been proposed (Martin et al., 2018). We make use of the open-source EASSE package (Alva-Manchego et al., 2019), which implements relevant metrics such as SARI, BERTScore, Flesch-Kincaid Grade Level (FKGL) and a host of quality estimation measures for more fine-grained analysis (these metrics are detailed in Appendix A).

5 Results & Discussion

Table 2 presents the results of our experiments on the Newsela corpus. According to SARI, our primary metric, SS with FUDGE outperforms both MUSS and the supervised baseline for all target levels except for Simp-4. Our simplifications tend to exhibit lower rates of compression and higher rates of sentence splitting and additions, particularly as the degree of simplification increases. This could be considered advantageous for certain target au-

Method	SARI	BERTScore	FKGL	Comp. ratio	Sent. splits	Lev. sim.	Copies	Add prop.	Del prop.
Target Level: Simp-1			7.97	1.01	1.19	0.90	0.44	0.10	0.10
PARA	36.61	81.68	9.15	0.97	1.02	0.89	0.18	0.08	0.11
MUSS	35.69	75.95	7.75	0.81	1.00	0.84	0.01	0.07	0.24
SUPER	32.49	88.19	9.36	0.99	1.04	0.99	0.89	0.01	0.01
\mathcal{B}_{Simp-1}	36.10	80.45	8.81	0.94	1.01	0.88	0.13	0.07	0.13
Target Level: Simp-2			6.41	0.98	1.42	0.82	0.23	0.17	0.20
PARA	35.01	73.53	9.12	0.97	1.02	0.89	0.18	0.08	0.11
MUSS	36.57	65.91	7.27	0.78	1.03	0.75	0.00	0.15	0.35
SUPER	31.12	78.22	8.88	0.99	1.10	0.98	0.80	0.02	0.03
\mathcal{B}_{Simp-2}	38.32	70.75	7.42	0.96	1.25	0.84	0.08	0.12	0.17
Target Level: Simp-3			4.91	0.92	1.55	0.73	0.13	0.24	0.31
PARA	30.87	65.06	9.09	0.98	1.01	0.89	0.18	0.08	0.11
MUSS	38.05	56.03	5.19	0.62	1.01	0.68	0.00	0.12	0.45
SUPER	37.89	66.60	6.65	0.93	1.34	0.90	0.48	0.06	0.13
\mathcal{B}_{Simp-3}	39.56	61.46	6.44	1.00	1.45	0.81	0.02	0.20	0.20
Target Level: Simp-4			3.40	0.85	1.79	0.65	0.09	0.30	0.43
PARA	25.61	56.21	9.41	0.98	1.01	0.89	0.18	0.08	0.11
MUSS	39.63	51.73	5.61	0.65	1.04	0.68	0.00	0.13	0.44
SUPER	43.22	55.00	5.09	0.78	1.45	0.74	0.24	0.12	0.32
\mathcal{B}_{Simp-4}	37.03	49.60	4.60	1.02	2.14	0.76	0.00	0.28	0.28

Table 2: Target-level results on the manually aligned Newsela test set (Jiang et al., 2020). For reference-based metrics (SARI, BERTScore), where higher values are better, we highlight systems according to their performance. For FKGL and reference-less quality estimation metrics we embolden the systems that perform closest to the level-specific references (provided in the intermediary rows).

differences and in settings where the loss of too much information could be detrimental for the reader. That said, it is also possible that not all additions and sentence splits are warranted. For example, degenerate repetitions or hallucinations could potentially skew these results (see tables in Appendix G for examples). Still, the positive influence of FUDGE’s classifier is indeed most visible when comparing against the paraphrase baseline, which fails to generate more readable texts according to FKGL.

We also note that the superior performance of the fully supervised method for level Simp-4 is consistent with the findings from Spring et al. (2021), where a similar approach proved most effective for simplifying ordinary German to A1-level German, despite it being the target level with the least amount of parallel data in both studies. For other target levels, however, where the differences between source and target are perhaps more subtle, this supervised model has a strong tendency to simply copy the input sentence.

While the MUSS baseline produces decent simplifications according to SARI and FKGL, the

lower compression ratio and a higher proportion of deleted N-grams indicate that this model also tends to severely summarise the input. This information loss causes model outputs to diverge from the ground truth references and it is thus penalised heavily by BERTScore.

FUDGE’s simplification operations are performed actively during decoding by modifying the generator’s prediction logits at each timestep, with each decision being informed by the currently generated prefix and its potential continuations $y_{1:i}$. Therefore, this approach does not enforce a transformation of the input text. This is an important and desirable feature for SS as oftentimes not all parts of a sentence need to be simplified (Garbacea et al., 2021). Thus, assuming a well-trained classifier \mathcal{B} , simplification operations should occur only when appropriate given the source sentence and the generation context.⁵

Finally, using FUDGE for SS also makes use of a single hyperparameter λ which controls the

⁵In an attempt to define ‘well-trained’, we demonstrate the effect of the amount of classifier training data on simplification performance in Appendix F.

contribution from the classifier. In contrast, MUSS requires setting an appropriate value for each of the four control tokens to attain a suitable simplification. These are not only difficult to determine for each target level (see Appendix E), but the way in which these tokens interact with each other is still unclear (Martin et al., 2020), making it difficult to set these values according to any underlying intuition.

6 Conclusion & Future Work

We have explored a modular and adaptable approach to SS by reframing it as controlled paraphrasing. We used FUDGE (Yang and Klein, 2021) to steer the generation of paraphrastic sequences toward different target levels. According to automatic metrics, this approach performs competitively when compared to state-of-the-art methods. In future work we aim to conduct a more detailed analysis of the model outputs in order to better understand the qualitative differences and potential shortcomings, as well as applying this method to larger textual units beyond sentences.

Limitations

In this work we aimed to generate level-appropriate sentence simplifications without the need for parallel aligned training data which is expensive to produce and difficult to come by. Due to a lack of existing work addressing simplification for specific target levels, the number of comparable approaches and relevant evaluation data sets are inherent limitations of this work. While the approach proposed by Martin et al. (2021) is not explicitly designed to generate level-specific simplifications, we searched for the most appropriate control tokens for each target level and selected the best values according to validation performance. To this end, we used the code provided by Martin et al. (2021) but note that we did not investigate potential improvements to this parameter search. Secondly, to the best of our knowledge, the Newsela English corpus constitutes the only available resource containing simplifications for readers at specific levels. Thus, this current work does not seek to assess how well this method generalises to other domains and languages.

Finally, our system evaluation relies on automatic metrics. While we have strived to report metrics that cover the degree of simplification (SARI, FKGL) and meaning preservation (BERTScore),

we acknowledge that extensive human evaluation with the target audience is crucial in assessing the viability and validity of any system intended to adapt and rewrite text. On the one hand, degenerate outputs may lead to even more confusion and less understandability, while, on the other hand, erroneous additions and deletions could be dangerous in certain contexts (e.g. when applied to medical or legal domain texts). Therefore, a thorough qualitative analysis of the proposed approach is still required and planned for future work.

Acknowledgements

In addition to the anonymous reviewers, we would like to thank Martin Volk and Rico Sennrich for helpful comments on an earlier version of this paper.

References

- Fernando Alva-Manchego, Louis Martin, Carolina Scarton, and Lucia Specia. 2019. *EASSE: Easier automatic sentence simplification evaluation*. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations*, pages 49–54, Hong Kong, China. Association for Computational Linguistics.
- Fernando Alva-Manchego, Carolina Scarton, and Lucia Specia. 2021. *The (un)suitability of automatic evaluation metrics for text simplification*. *Computational Linguistics*, 47(4):861–889.
- Rahul Bhagat and Eduard Hovy. 2013. *Squibs: What is a paraphrase?* *Computational Linguistics*, 39(3):463–472.
- Joachim Bingel and Anders Søgaard. 2016. *Text simplification as tree labeling*. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 337–343, Berlin, Germany. Association for Computational Linguistics.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. *Plug and Play Language Models: A Simple Approach to Controlled Text Generation*. *arXiv:1912.02164 [cs]*.
- Jan De Belder and Marie-Francine Moens. 2010. *Text simplification for children*. Proceedings of the SIGIR Workshop on Accessible Search Systems, pages 19–26. ACM; New York.
- Yue Dong, Zichao Li, Mehdi Rezagholizadeh, and Jackie Chi Kit Cheung. 2019. *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*, pages 3393–3402, Florence, Italy. Association for Computational Linguistics.
- Yanjun Gao, Ting-Hao Huang, and Rebecca J. Passonneau. 2021. **ABCD: A graph framework to convert complex sentences to a covering set of simple sentences.** In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3919–3931, Online. Association for Computational Linguistics.
- Avishek Garain, Arpan Basu, Rudrajit Dawn, and Sudip Kumar Naskar. 2019. **Sentence Simplification using Syntactic Parse trees.** In *2019 4th International Conference on Information Systems and Computer Networks (ISCON)*, pages 672–676, Mathura, India. IEEE.
- Cristina Garbacea, Mengtian Guo, Samuel Carton, and Qiaozhu Mei. 2021. **Explainable Prediction of Text Complexity: The Missing Preliminaries for Text Simplification.** *arXiv:2007.15823 [cs]*.
- Goran Glavaš and Sanja Štajner. 2015. **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)*, pages 63–68, Beijing, China. Association for Computational Linguistics.
- John Honeyfield. 1977. **Simplification.** *TESOL Quarterly*, 11(4):431–440.
- Chao Jiang, Mounica Maddela, Wuwei Lan, Yang Zhong, and Wei Xu. 2020. **Neural CRF model for sentence alignment in text simplification.** In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7943–7960, Online. Association for Computational Linguistics.
- Philipp Koehn and Rebecca Knowles. 2017. **Six challenges for neural machine translation.** In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.
- Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. 2021. **GeDi: Generative discriminator guided sequence generation.** In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4929–4952, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Reno Kriz, Eleni Miltsakaki, Marianna Apidianaki, and Chris Callison-Burch. 2018. **Simplification using paraphrases and context-based lexical substitution.** 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 207–217, New Orleans, Louisiana. Association for Computational Linguistics.
- Sachin Kumar, Shuly Wintner, Noah A. Smith, and Yulia Tsvetkov. 2019. **Topics to avoid: Demoting latent confounds in text classification.** In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4153–4163, Hong Kong, China. Association for Computational Linguistics.
- Philippe Laban, Tobias Schnabel, Paul Bennett, and Marti A. Hearst. 2021. **Keep it simple: Unsupervised simplification of multi-paragraph text.** In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6365–6378, Online. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. **BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension.** In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. **DExperts: Decoding-time controlled text generation with experts and anti-experts.** In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6691–6706, Online. Association for Computational Linguistics.
- Mounica Maddela, Fernando Alva-Manchego, and Wei Xu. 2021. **Controllable text simplification with explicit paraphrasing.** In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3536–3553, Online. Association for Computational Linguistics.
- Louis Martin, Éric de la Clergerie, Benoît Sagot, and Antoine Bordes. 2020. **Controllable sentence simplification.** In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4689–4698, Marseille, France. European Language Resources Association.
- Louis Martin, Angela Fan, Éric de la Clergerie, Antoine Bordes, and Benoît Sagot. 2021. **MUSS: Multilingual Unsupervised Sentence Simplification by Mining Paraphrases.** *arXiv:2005.00352 [cs]*.

- Louis Martin, Samuel Humeau, Pierre-Emmanuel Mazaré, Éric de La Clergerie, Antoine Bordes, and Benoît Sagot. 2018. [Reference-less quality estimation of text simplification systems](#). In *Proceedings of the 1st Workshop on Automatic Text Adaptation (ATA)*, pages 29–38, Tilburg, the Netherlands. Association for Computational Linguistics.
- Shashi Narayan, Claire Gardent, Shay B. Cohen, and Anastasia Shimorina. 2017. [Split and rephrase](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 606–616, Copenhagen, Denmark. Association for Computational Linguistics.
- Christina Niklaus, Matthias Cetto, André Freitas, and Siegfried Handschuh. 2019. [DisSim: A discourse-aware syntactic text simplification framework for English and German](#). In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 504–507, Tokyo, Japan. Association for Computational Linguistics.
- Kostiantyn Omelianchuk, Vipul Raheja, and Oleksandr Skurzshanskyi. 2021. [Text Simplification by Tagging](#). In *Proceedings of the 16th Workshop on Innovative Use of NLP for Building Educational Applications*, pages 11–25, Online. Association for Computational Linguistics.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. [fairseq: A fast, extensible toolkit for sequence modeling](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.
- Damian Pascual, Beni Egressy, Clara Meister, Ryan Cotterell, and Roger Wattenhofer. 2021. [A plug-and-play method for controlled text generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3973–3997, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [Glove: Global Vectors for Word Representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer](#). *arXiv:1910.10683 [cs, stat]*.
- Horacio Saggion. 2017. *Automatic Text Simplification*. Synthesis Lectures on Human Language Technologies. Morgan & Claypool.
- Carolina Scarton 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.
- Thomas Scialom, Louis Martin, Jacopo Staiano, Éric Villemonte de la Clergerie, and Benoît Sagot. 2021. [Rethinking Automatic Evaluation in Sentence Simplification](#). *arXiv:2104.07560 [cs]*.
- Advait Siddharthan and Angrosh Mandya. 2014. [Hybrid text simplification using synchronous dependency grammars with hand-written and automatically harvested rules](#). In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 722–731, Gothenburg, Sweden. Association for Computational Linguistics.
- Punardeep Sikka and Vijay Mago. 2020. [A Survey on Text Simplification](#). *arXiv:2008.08612 [cs]*.
- Nicolas Spring, Annette Rios, and Sarah Ebling. 2021. [Exploring German multi-level text simplification](#). In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 1339–1349, Held Online. IN-COMA Ltd.
- Sanja Štajner. 2021. [Automatic Text Simplification for Social Good: Progress and Challenges](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2637–2652, Online. Association for Computational Linguistics.
- Teerapaun Tanprasert and David Kauchak. 2021. [Flesch-kincaid is not a text simplification evaluation metric](#). In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 1–14, Online. Association for Computational Linguistics.
- Sowmya Vajjala and Ivana Lučić. 2018. [OneStopEnglish corpus: A new corpus for automatic readability assessment and text simplification](#). In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 297–304, New Orleans, Louisiana. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

- Sander Wubben, Antal van den Bosch, and Emiel Krahmer. 2012. [Sentence simplification by monolingual machine translation](#). In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1015–1024, Jeju Island, Korea. Association for Computational Linguistics.
- Wei Xu, Chris Callison-Burch, and Courtney Napoles. 2015. [Problems in current text simplification research: New data can help](#). *Transactions of the Association for Computational Linguistics*, 3:283–297.
- Wei Xu, 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.
- Kevin Yang and Dan Klein. 2021. [FUDGE: Controlled text generation with future discriminators](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3511–3535, Online. Association for Computational Linguistics.
- Xingxing Zhang and Mirella Lapata. 2017. [Sentence simplification with deep reinforcement learning](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 584–594, Copenhagen, Denmark. Association for Computational Linguistics.

A Evaluation Metrics for Sentence Simplification

Simplicity SARI is intended to measure simplicity by considering N-gram overlap between the model output, source sentence and one or more reference sentences. It rewards model outputs that involve edit operations, such as deletions, additions and copies, which correspond with the provided references.

Fluency and meaning preservation BERTScore uses BERT’s contextualised representations to compute the similarity between tokens in the model output and one or more references. It has been shown to correlate better than BLEU for assessing meaning preservation and fluency in SS (Scialom et al., 2021).

Readability Flesch-Kincaid Grade Level (FKGL) is often used as a proxy for estimating text simplicity without a reference. Originally developed for grading technical materials for military personnel, it considers surface-level statistics such as word and sentence length to provide a single score. However, these scores should be interpreted carefully as it has recently been shown that this metric can be misled by degenerate and disfluent outputs (Tanprasert and Kauchak, 2021).

Quality Estimation Measures For a more fine-grained analysis of model outputs, we also report quality estimation measures which are computed between the source sentence and the model’s output. These include the compression ratio, Levenshtein similarity, average number of sentence splits performed, exact copies between source and target, and the proportion of added and deleted N-grams.

B Target Levels in Newsela

Newsela contains articles written for different graded reading levels (2-12) corresponding to primary and secondary school grades in the United States. These grades can be considered true indicators of an simplicity. However, assessing target-level simplification for all available grades is challenging due to the sheer number of them and the fact that not all of them are equally well represented in the corpus. For example, there are 1,853 articles for grade 12 but only 20 articles for grade 10 and 2 for grade 11.

To simplify our target-level analysis, we follow Xu et al. (2015) and assume that Newsela’s article version IDs (0-4) are reliable indicators of a text’s simplicity and thus adopt these as our target levels (Simp- l). Figure 1 shows how the graded reading levels are distributed over our target levels. As can be seen, the article versions do not provide a clear cut aggregation of all grades since there is a limited degree of overlap, particularly between the lower Simp- l levels. Nevertheless, a rough aggregation is discernible; Simp-4 covers grades 2 to 4, Simp-3 consists predominately of grade 5, Simp-2 contains grades 6 and 7, and Simp-1 covers grades 7 and 8. Finally, Simp-0 (the complex sources articles in our study) is mostly restricted to grade 10 and up.

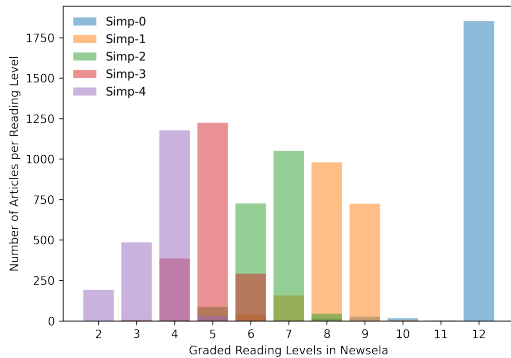


Figure 1: The distribution of graded reading levels among article versions in Newsela. For simplicity, we use as article versions our target levels for simplification.

C Details on Model Training and Inference

C.1 Resources

Model training and inference experiments were performed on NVIDIA GeForce GTX TITAN X GPUs with 12GB of memory.

C.2 Training Generation Models

For our underlying generator model \mathcal{G} and the level-aware supervised baseline, we fine-tuned BART-large using Hugging Face’s Transformers library⁶ (Wolf et al., 2020). Training parameters used for \mathcal{G} aim to replicate the settings used by Martin et al. (2021) who trained their models using Fairseq⁷

⁶<https://github.com/huggingface/transformers>.

⁷<https://github.com/facebookresearch/fairseq>.

(Ott et al., 2019). For the level-aware supervised baseline, we aimed to replicate the settings used by Spring et al. (2021) who trained their models with Sockeye⁸. Note, in contrast to the paraphrase model, the effective batch size and maximum training steps for this model are considerably smaller to account for the differences in the size of the relevant training data (1.4M paraphrase sentence pairs vs. 4k aligned simplifications).

Paraphrase Model \mathcal{G}	
hyperparameter	value
max src length	1024
max tgt length	256
eff. batch size	64
learning rate	3e-05
weight decay	0.01
optim	adamw_hf
adam betas	0.9 - 0.999
adam epsilon	1e-8
lr scheduler	polynomial
warmup steps	500
label smoothing	0.1
max steps	20000
num beams for pred	4
optim metric	loss
Level-Aware Supervised Model	
hyperparameter	value
max src length	256
max tgt length	128
eff. batch size	16
learning rate	3e-05
weight decay	0.01
optim	adamw_hf
adam betas	0.9 - 0.999
adam epsilon	1e-8
lr scheduler	polynomial
warmup steps	500
label smoothing	0.1
max steps	5000
num beams for pred	4
optim metric	rouge1

Table 3: Hyperparameters for training generation models.

C.3 Training FUDGE Classifiers

Our FUDGE classifiers \mathcal{B}_{simp-l} are unidirectional three-layer LSTM-based RNNs with hidden layer dimensionality of 512. These settings differ slightly from the original implementation by Yang and Klein (2021), who used smaller classifiers for their tasks. The embedding matrix is constructed to cover the vocabulary of the underlying

⁸<https://github.com/aws-labs/sockeye>.

ing generator model (i.e. the tokenizer is shared between \mathcal{G} and \mathcal{B}) and token embeddings are initialised using 300d pre-trained GloVe embeddings (glove-wiki-gigaword-300) (Pennington et al., 2014). For certain subwords and rare words that are OOV in GloVe, we initialise their embeddings randomly.

C.4 Inference

For all models except MUSS we run inference with beam search ($k=5$). A manual inspection of the model outputs revealed that our underlying paraphraser \mathcal{G} showed a tendency to produce repetitions in the target sequence. To counter this, we set the repetition penalty equal to 1.2 when performing inference with \mathcal{G} . All other inference hyperparameters use the default values set in Hugging Face. For each source sentence in the test set, we generate the top five model hypotheses according to the model and select the first non-empty string as the final model output.

For MUSS, we kept inference settings the same as the default set by Martin et al. (2021). The only differences are the control token values used for performing inference on each of the Newsela simplification levels, which we derive via a parameter sweep on 50 items from the respective development sets. Table 4 shows the relevant values used. Note that these values are rounded up to the nearest 0.05 inside the model.

	Comp. Ratio	Levenshtein Sim.	Word Rank Ratio	Dep. Tree Depth Ratio
Simp-1	0.30	0.99	0.54	1.45
Simp-2	0.75	0.82	0.94	0.22
Simp-3	0.52	0.85	0.45	0.62
Simp-4	0.47	0.79	0.43	0.42

Table 4: Values used for target-level inference on the Newsela English corpus with MUSS.

D Parameter Sweep for FUDGE

FUDGE has two hyperparameters which need to be set at inference time. The first is a weight λ that controls the strength of \mathcal{B} 's contribution, while the second aims to keep the cost associated with classifying all possible continuations at each decoding timestep down by only considering the best k predictions at each step (i.e. the most probable k tokens according to the model). Initial experiments

showed that λ is indeed useful for controlling the degree of simplification and finding a suitable λ is essential. Meanwhile, using different pre-selection k values (e.g. [50, 200]) had almost no effect on the resulting generation sequence when using argmax decoding techniques such as beam search. Therefore, we followed the recommendation by Yang and Klein (2021) and fixed the pre-selection $k=200$.

To get the best target-level simplifications, we searched for the optimum λ value for each combination of Newsela simplification levels and each target-level FUDGE on 50 sentences from the manually aligned validation set (Jiang et al., 2020). Figure 2 shows the resulting SARI scores. For our experiments, we selected the best scoring λ values for each simplification level and its corresponding FUDGE (i.e. plots along the diagonal). For cases where more than one possible λ delivered good results, we selected the lowest value > 0 (marked with a vertical dotted line).

It is clear from Figure 2 that cross-matching target simplification levels with FUDGES trained on a different target level would also yield good, and in some cases even better, results according to SARI (e.g. target-level Simp-2 with \mathcal{B}_{Simp-3}). We hypothesise that this is likely due to it being easier for the classifier to correctly distinguish between the positive (simple) and negative (complex) classes when the stylistic differences between simplification levels are larger. Indeed, ROC-AUC scores for each target-level classifier on the respective test sets increase from 0.67 to 0.96 going from Simp-1 to Simp-4, indicating that FUDGES trained on higher simplification levels are better at distinguishing between the classes.

E ACCESS Attributes on Newsela Corpus

While parameter sweeps are helpful, deciding on optimal attribute values for target-level simplification with ACCESS is non-trivial, especially if limited validation data is on hand. Furthermore, control values that might be suitable for one input sentence, may not be suitable for another input sentence. For example, it may not be possible to reduce an already short input sentence to half of its original length. To examine potentially suitable control values, we computed the ratio scores on source-target pairs from the manually aligned training split from Jiang et al. (2020) for all four simplification levels of the Newsela English corpus. Figure 3 shows that for most attributes the

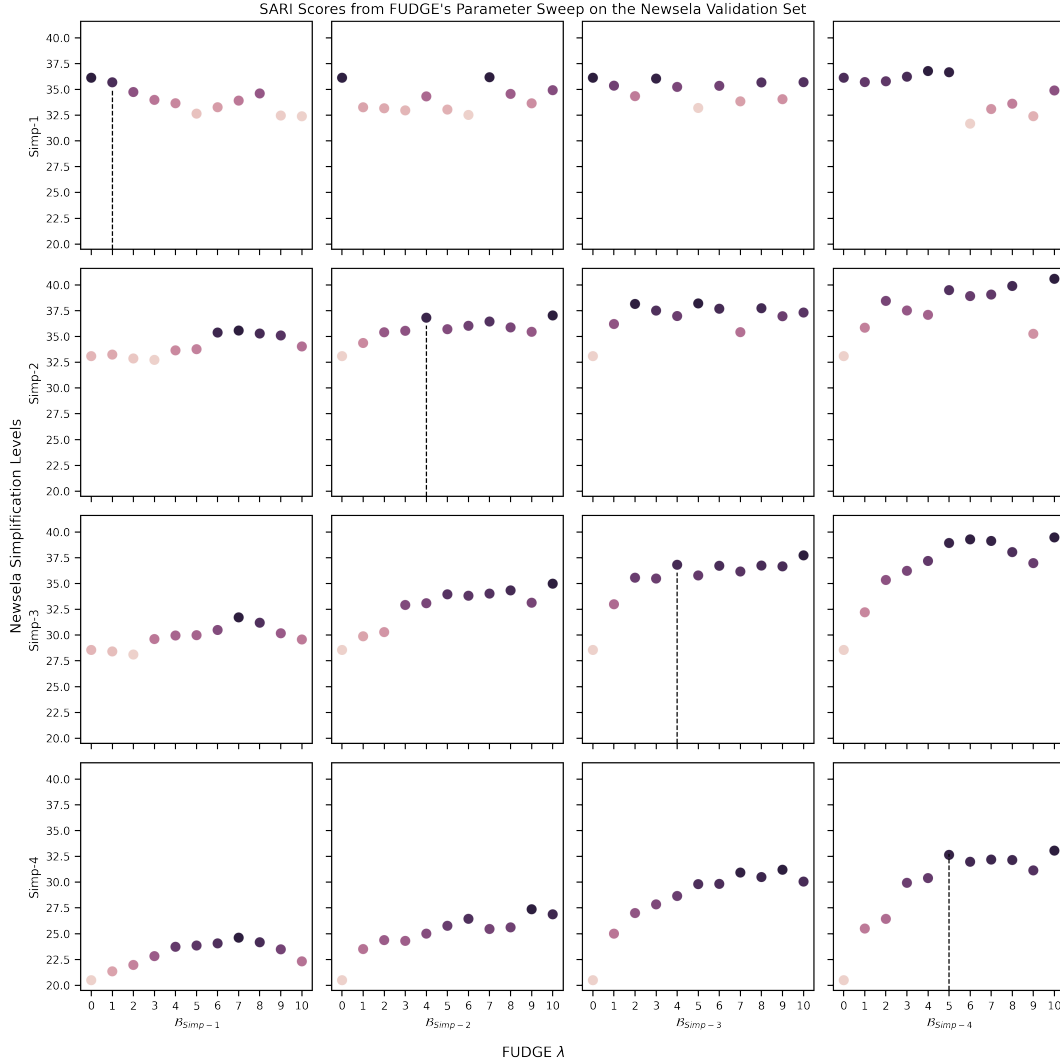


Figure 2: SARI scores from parameter sweep over different λ values for FUDGE at inference time.

largest density is on a value of 1.0, indicating no difference between the source and target. For many attribute values, the distributions are also relatively wide and flat indicating that there could be many potentially valid values, especially for the higher simplification levels (e.g. Simp-2 - Simp-4).

F Ablation Experiment

Unlike a fully-supervised seq2seq approach, FUDGE for SS does not require parallel complex-simple sentence pairs for training. Instead, it relies on contrastive instances to train its target-level classifiers. Such data is significantly easier to collect from comparable corpus resources in many languages, e.g. language learning materials (Vajjala and Lučić, 2018) or news articles produced specifi-

cally for certain target groups.⁹

However, an open question remains as to how much data is required to train a suitable classifier. While this may depend heavily on the target-level simplified text both in topical and stylistic features, we examined this question for Newsela’s Simp-4 target level. In contrast to our main experiments, here, we fix the weighted contribution from the classifier as $\lambda = 1.0$ (i.e. the minimum amount of influence). Figure 4 depicts the relationship between the amount of contrastive data used to train \mathcal{B}_{Simp-4} and the primary metrics for simplification and quality estimation.

On these plots, a strong correlation is visible between increasing the amount of contrastive data

⁹For example, Ligetil from the Danish Broadcasting Corporation (<https://www.dr.dk/ligetil/>) and Japan’s News Web Easy (<https://www3.nhk.or.jp/news/easy/>).

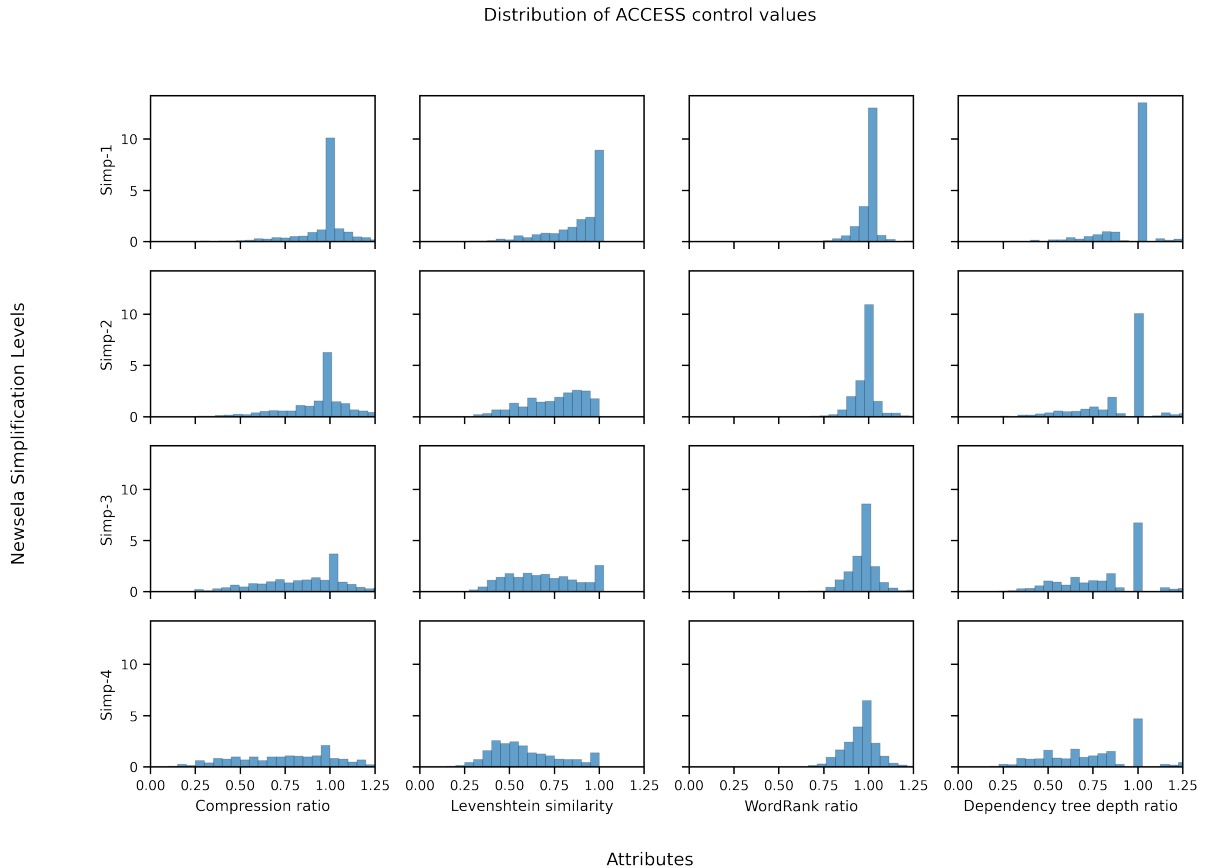


Figure 3: Density of attribute values for the four control tokens used in the ACCESS simplification method (Martin et al., 2020) and employed by MUSS (Martin et al., 2021).

and the degree to which the model simplifies the input sentences. Clearly, while more data is helpful, even small amounts of contrastive data (e.g. 500-1000 examples) can already be effective in steering the generations towards the target attribute.

G Model Output Examples

On closer inspection of the model outputs, we observed some undesirable trends among all model outputs. Firstly, the supervised approach (SUPER) tends to keep edit operations to a minimum, resulting in model outputs that are very similar to the input text. In cases where the ground truth simplification also happens to be a copy of the source, overlap metrics are unduly maximised. Secondly, MUSS tends to produce highly fluent simplifications yet these often resemble short summaries. For longer sentences, outputs can be densely packed, leading to even more complex sentences, or significant amounts of information are simply dropped from the input. It remains an open question as to whether or not this information loss is suitable for simplifying towards a target level. Finally, we also

observed that FUDGE’s model outputs (B_{simp-l}) are more susceptible to disfluencies such as redundant punctuation, subwords and phrases. This suggests that the attribute classifier can unfortunately have a negative impact on the generator in some cases.

The tables below provide randomly sampled examples of model outputs for each target-level in the Newsela English corpus. We colour parts of the simplified texts based on the edit operations applied to the source text. **Blue** indicates additions or explanations not in the source text. **Green** is used to highlight lexical and punctuation substitutions. **Yellow** shows operations on contractions (either creating or deconstructing). **Pink** indicates phrases that have been truncated or lexical deletions from the source text. **Violet** is used for larger paraphrastic segments or positionally shuffled phrases. Undesirable repetitions or hallucinations are italicised.

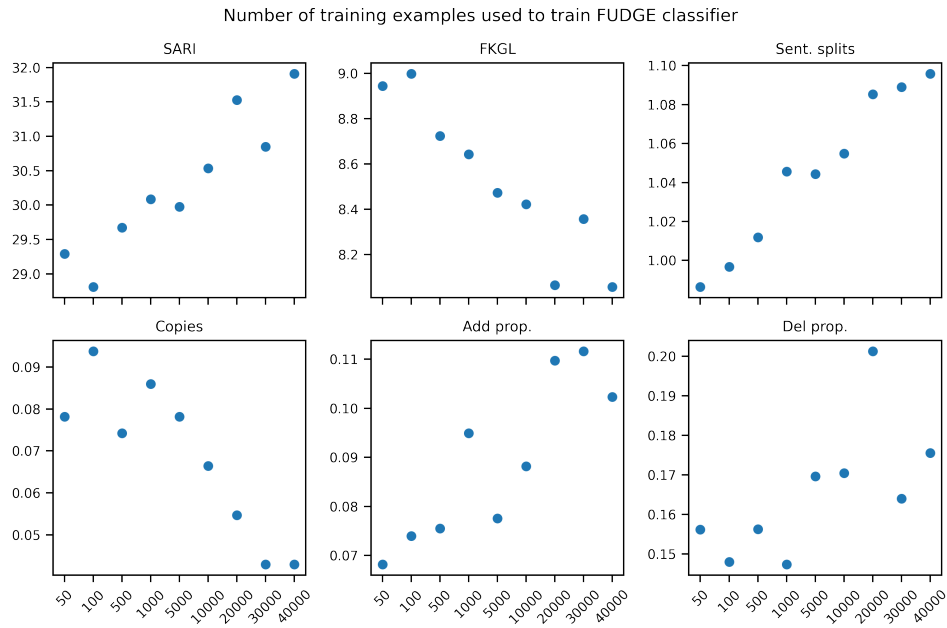


Figure 4: Automatic simplification metrics for simplification as a function of the number of contrastive complex-simple paragraphs from the Newsela English corpus used to train a FUDGE classifier \mathcal{B}_{Simp-4} .

Target Level: Simp-1

Simp-0 (source)	They found that there was indeed a pattern in their motion — though it wasn't quite what they expected.
Simp-1 (reference)	They found that there was indeed a pattern in their motion — though it was not quite what they expected.
MUSS	They found that there was a pattern in their motion - but it wasn't what they expected.
SUPER	They found that there was indeed a pattern in their motion — though it wasn't quite what they expected.
\mathcal{B}_{simp-1}	They found that there was a pattern in their movement , though it wasn't what they expected.
<hr/>	
Simp-0 (source)	The tribes helped members of the expedition through the winter, bringing them food and assisting with navigation.
Simp-1 (reference)	The tribes helped members of the expedition through the winter, bringing them food and assisting with navigation.
MUSS	The tribes helped members of the expedition by bringing them food and helping with navigation.
SUPER	The tribes helped members of the expedition through the winter, bringing them food and assisting with navigation.
\mathcal{B}_{simp-1}	The tribes helped members of the expedition through the winter, bringing food and helping with navigation.

Target Level: Simp-2

Simp-0 (source)	The feat should be possible based on Einstein's theory of general relativity, which describes gravity as the warping of spacetime by energy and matter.
Simp-2 (reference)	It should be possible based on Einstein's theory of general relativity. His famous theory describes gravity as the bending of space and time, which are one thing called "spacetime."
MUSS	Einstein's theory, called general relativity, describes gravity as the warping of spacetime by energy and matter.
SUPER	The feat should be possible based on Einstein's theory of general relativity, which describes gravity as the warping of spacetime by energy and matter.
\mathcal{B}_{simp-2}	The feat is based on Einstein's theory of general relativity. It describes gravity as the warping of spacetime by matter and energy.
Simp-0 (source)	Changes include moving the earliest date from which tribes must provide documentation of their existence into the 1930s and allowing tribes to reapply for federal recognition under certain circumstances, something previously prohibited.
Simp-2 (reference)	The revisions include establishing the 1930s as the earliest date from which tribes must provide documentation of their existence. Also, tribes may reapply for federal recognition under certain situations, something previously prohibited.
MUSS	Changes include moving the earliest date for tribes to provide documentation of existence into the 1930s, and allowing under certain circumstances tribes to reapply for federal recognition.
SUPER	Changes include moving the earliest date from which tribes must provide documentation of their existence into the 1930s and allowing tribes to reapply for federal recognition under certain circumstances, something previously prohibited.
\mathcal{B}_{simp-2}	Changes include moving the earliest date from which tribes must provide documentation of their existence. Tribes will also be allowed to reapply for federal recognition.

Target Level: Simp-3

Simp-0 (source)	Cafes, restaurants and bakeries around the nation are feeling the pain from recent outbreaks of bird flu, or avian influenza, expanding through the Midwest.
Simp-3 (reference)	The rising prices is caused by recent outbreaks of bird flu.
MUSS	Restaurants and bakeries around the nation are feeling the effects of the recent bird flu outbreak in the Midwest.
SUPER	Cafes, restaurants and bakeries around the nation are feeling the pain from recent outbreaks of bird flu. The disease is spreading through the Midwest.
\mathcal{B}_{simp-3}	Cafes, restaurants and bakeries around the nation are feeling the pain of bird flu outbreaks that have spread in the Midwest.
Simp-0 (source)	"That is, a time traveler who emerges from a Deutschian CTC enters a universe that has nothing to do with the one she exited in the future.
Simp-3 (reference)	Lloyd says a time traveler coming from Deutsch's CTC "enters a universe that has nothing to do with the one she exited."
MUSS	A Deutschian CTC traveler enters a universe that has nothing to do with the one she left.
SUPER	"That is, a time traveler who emerges from a Deutschian CTC enters a universe that has nothing to do with the one she exited in the future.
\mathcal{B}_{simp-3}	"That is, a time traveler who emerges from a Deutschian CTC enters a universe. This universe has nothing to do with the one she exited in the future. The universe she enters into has nothing to do with the one she exited in the future.

Target Level: Simp-4

Simp-0 (source)	If the particle flips the switch, the machine emits a particle—the particle—back into the CTC; if the switch isn't flipped, the machine emits nothing.
Simp-4 (reference)	If the particle flips the switch, the machine shoots a particle — the particle — back into the CTC. If the switch isn't flipped, the machine shoots out nothing.
MUSS	If the particle turns , the machine sends a particle back to the CTC. If it doesn't, it sends nothing.
SUPER	If the particle flips the switch, the machine emits a particle—the particle—back into the CTC. If the switch isn't flipped, the machine emits nothing.
\mathcal{B}_{simp-4}	If the switch is turned, the machine emits a particle-like stateBack into the CTC. If the switch is not turned , the machine emits nothing.
Simp-0 (source)	Deutsch's insight was to postulate self-consistency in the quantum realm, to insist that any particle entering one end of a CTC must emerge at the other end with identical properties.
Simp-4 (reference)	Deutsch assumes that tiny quantum particles are stable and fixed.
MUSS	In quantum theory, Deutsch insisted that any particle entering one end of a CTC must emerge at the other end with equal properties.
SUPER	Deutsch's idea was to show that any particle entering one end of a CTC must emerge at the other end <i>of a CTC must emerge at the other end</i> with identical properties.
\mathcal{B}_{simp-4}	Deutsch's idea was to postulate a very nature . He was claiming that any particle entering one end of a CTC must emerge at the other end with identical properties.