

Text Simplification for Medical Wikipedia Corpora

Hoang Nguyen Hung Van
CSC580 Project Final Report

Abstract

The goal of text simplification is to transform difficult text into a version that is easier to understand and more broadly accessible. In some domains, such as healthcare and medical, fully automated approaches cannot be used since the information must be accurately preserved. In this report, I introduce the autocompletion algorithm for text simplification with medical Wikipedia corpora, which aims to assist human simplification by suggesting the next word to type when manually simplifying a text. The main part of this document is to report my experiments with four different Language Models and naive Bayes n-gram approaches. In this report, I will also show how I create the medical corpora from general Wikipedia parallel corpus for text simplification. In this report, I also show the importance of context in better assisting text simplification systems with the absolute improvement of 20.5% and 28.8%. The best model is RoBERTa, which achieves a word prediction rate of 62% on medical Wikipedia data. The project here adapted changes since my last proposal because of 3 main reasons: 1) this system is more applicable to medical and show significant results to help people with low literacy knowledge to better understand medically written documents and 2) the analysis from the project open a new research venue of using autocompletion techniques in simplifying medical documents, which are usually hard to understand mostly.

1 Introduction

Text simplification (TS) is the process of modifying the words and structure of a text while preserving the content to make the information in the text more broadly accessible (Shardlow, 2014). Most research in text simplification has focused on fully automated (Zhu et al., 2010; Coster and Kauchak, 2011; Xu et al., 2016; Zhang and Lapata, 2017;

Nishihara et al., 2019). In some domains, e.g., medicine or healthcare, using fully-automated text simplifications is not appropriate since it is critical that the information gets preserved correctly during the simplification process. Instead of fully-automated approaches, support tools such as editors are better suited to generate simplifications more efficiently and with higher quality (Kloehn et al., 2018).

Autocompletion tools suggest one or more words as the user types that could follow what has been typed so far. Autocompletion has been used in a range of applications including web queries (Cai et al., 2016), database queries (Khoussainova et al., 2010), texting (Dunlop and Crossan, 2000), and e-mail composition (Dai et al., 2019). In this document, I explore autocompletion for text simplification. In contrast to most autocomplete applications, for text simplification, in addition to the text that is being typed, I also have the additional context of the content being simplified. The work is mostly similar to interactive machine translation tools where a user translating a foreign sentence is given guidance as they type (Green et al., 2014).

In this paper, I examine the autocompletion task for sentence-level text simplification: given a difficult sentence that a user is trying to simplify and the simplification typed so far, the goal is to suggest the next word to follow what has been typed. Figure 1 shows an example difficult sentence along with the simplification that the user has typed so far. The task is to predict the next word to assist in finishing the simplification, in this case a verb like “take”, which might be continued to a partial simplification of “take place at the Chapel”.

I make three main contributions. First, I introduce the autocompletion task for sentence simplification and provide an initial analysis based on a number of recent models. Second, I show how the additional context of the difficult sentence can be

Difficult sentence	The Chapel is actively used as a place of worship and also for some concerts and college events.
Typed	Concerts and college events _____

Figure 1: An example text simplification autocompletion task. The user is simplifying the difficult sentence on top and has typed the words on the bottom so far.

integrated into these models to improve the quality of the suggestions made. Using the context of the difficult sentence significantly improves the prediction quality of the autocomplete methods. Thirdly, with the analysis of different models, I hope that it can provide precise information into how each model perform based on different categories resulting hybrid models, ensemble models, and other additional-information-utilized models that can better advance the autocompletion algorithm in medical text simplification.

2 Text Simplification Autocompletion

Given a difficult sentence that a user is trying to simplify, $d_1 d_2 \dots d_m$, and the simplification typed so far, $s_1 s_2 \dots s_i$, the autocompletion task is to suggest word s_{i+1} . I examined four recent neural models that utilize the Transformer network (Vaswani et al., 2017): RoBERTa (Liu et al., 2019), BERT (Devlin et al., 2018), XLNet (Yang et al., 2019), and GPT-2 (Radford and Wu, 2018). The four models are state-of-the-art neural language representation models that have performed well in a range of applications. I also examined the naive Bayes n-gram language model and used this as a base line for my project.

To understand the benefit of the context of the difficult sentence, I compare models that do not use context, i.e., predict only based on $s_1 s_2 \dots s_i$, and context-aware versions that incorporate the difficult sentence into the prediction.

2.1 BERT

BERT (Bidirectional Encoder Representations from Transformers) is a method for learning language representations using bidirectional training. The main advantage of BERT is that it uses a masked approach to train the model where some of the words in the training data are replaced with a [MASK] token. The model then attempts to predict the original value of the masked words based on the context provided by the other, non-masked words in the sequence. Unlike left-to-right or right-to-

left sequential models, BERT can use context both before and after the word to be predicted. BERT has been shown to produce state-of-the-art results in a wide range of generation and classification applications (Devlin et al., 2018).

I use the original BERT pre-trained model, which was trained on the BooksCorpus (Zhu et al., 2015) and English Wikipedia. To apply the model without context, I predict the masked word for the input " $s_1 s_2 \dots s_i$ [MASK]". For the context-aware version, I add the context of the difficult sentence " $d_1 d_2 \dots d_m . s_1 s_2 \dots s_i$ [MASK]". This biases the prediction to words related to those found in the encoded context from difficult sentences.

BERT is a pre-trained model designed to be fine-tuned for particular applications. For the text simplification autocompletion task, I fine-tuned BERT on a corpus of sentence-aligned difficult concatenated with the corresponding simple sentences. I used Transformer Neural Networks to fine-tune the pre-trained BERT language model¹ on this data. Since, our task is to predict the next word in the simple sentence, I mask out each word in the simple sentence portion and then predict that word.

2.2 GPT-2

Like BERT, GPT-2 is also based on the Transformer network, but uses left-to-right training and prediction. In each layer, GPT-2 has 12 independent attention mechanisms, called "heads", and the overall model contains 12 layers which can capture up to 144 different attention patterns. I use the publicly released model², which has 1.5B model parameters and is trained on web text. Since GPT-2 is a traditional left-to-right model, for the context unaware version I simply predict s_{i+1} based on $s_1 s_2 \dots s_i$. Like BERT, to incorporate the context of the difficult sentence, I prepend it to the simplified text typed so far and then predict s_{i+1} .

2.3 RoBERTa

Like BERT and GPT-2 is also based on the Transformer network. RoBERTa is RoBERTa: A Robustly Optimized BERT Pretraining Approach. In each layer, GPT-2 has 12 independent attention mechanisms, called "heads", and the overall model contains 12 layers which can capture up to 144 different attention patterns. I use the publicly released model. For the context unaware version I

¹<https://github.com/huggingface/transformers/tree/master/examples/>

²<https://github.com/openai/gpt-2>

simply predict s_{i+1} based on $s_1 s_2 \dots s_i$. Like BERT, to incorporate the context of the difficult sentence, I prepend it to the simplified text typed so far and then predict s_{i+1} .

2.4 XLNet

With the capability of modeling bidirectional contexts, denoising autoencoding based pretraining like BERT achieves better performance than pretraining approaches based on autoregressive language modeling. However, relying on corrupting the input with masks, BERT neglects dependency between the masked positions and suffers from a pretrain-finetune discrepancy. In light of these pros and cons, XLNet proposes a generalized autoregressive pretraining method that enables learning bidirectional contexts by maximizing the expected likelihood over all permutations of the factorization order and overcomes the limitations of BERT thanks to its autoregressive formulation. Furthermore, XLNet integrates ideas from Transformer-XL, the state-of-the-art autoregressive model, into pretraining. Empirically, under comparable experiment setting, XLNet outperforms BERT on 20 tasks, often by a large margin, including question answering, natural language inference, sentiment analysis, and document ranking.

3 Experiments

I introduce a new task, text simplification autocompletion, which relies on a parallel corpus of difficult and simple sentences. Given a difficult sentence and the first i words of the simple sentence, the goal is to predict the $i + 1^{th}$ simple word. I provide the first results on this task using RoBERTa, BERT, XLNet, and GPT-2, with and without context, as well as an trigram language model baseline.

3.1 Experimental setup

To evaluate the quality of the different models, I used the Simple Wikipedia parallel corpus (Kauchak, 2013), which contains 167K pairs of sentences, with one sentence from English Wikipedia and a corresponding sentence from Simple English Wikipedia. However, to emphasize my work on medical domain, I extracted around 3k pairs from this and run experiments. I used 70% of the sentence for training, 15% for development, and 15% for testing.

As an additional baseline that does not use context, I trained a trigram language model with

Difficult sentence	The Saxons built Banbury on the west bank of the River Cherwell.
Simple sentence	Banbury is part of the Cherwell district.

Figure 2: An example sentence pair from the English Wikipedia corpus.

Typed so far	Predict
Banbury	is
Banbury is	part
Banbury is part	of
Banbury is part of	the
Banbury is part of the	Cherwell
Banbury is part of the Cherwell	district

Figure 3: The resulting prediction tasks that are generated from the example in Figure 2.

Kneser-Ney smoothing using the SRILM toolkit (Stolcke, 2002). The model was trained on the simple sentences from the training portion of the dataset and predicts s_{i+1} as the word with the highest probability given the previous two words, i.e., $\argmax_{s_{i+1}} p(s_{i+1} | s_i s_{i-1})$.

The fine-tuning for four language models are done with a batch-size of 8, 8 epochs, and a learning rate of $5e^{-5}$. Early stopping was used based on the second time a decrease in the accuracy was seen.

To evaluate the models, I calculated how well the models predicted the next word in a test sentence, given the previous words. A simple test sentence of length n , $s_1 s_2 \dots s_n$, would result in $n - 1$ predictions, i.e., predict s_2 given s_1 , predict s_3 given $s_1 s_2$, etc. For example, Figure 2 shows a difficult sentence from English Wikipedia and the corresponding simplification from Simple English Wikipedia. Given this test example, I generate six prediction tasks, one for each word in the simple sentence after the first word. Figure 3 shows these six test prediction tasks. For the context-aware approaches, they also incorporated the difficult sentence. We measured the performance of a system using accuracy based on the number of predictions that exactly matched the next word in the corpus. The test corpus contained 450 sentence pairs resulting in a total of 6K individual word predictions.

3.2 Prediction performance

Table 1 shows the results for the nine different variants (trigram model, RoBERTa, BERT, XLNet, and

Model	No Context	Context-Aware
trigram	13%	–
RoBERTa	56.23%	62.4%
BERT	53.28%	50.43%
XLNet	46.2%	45.7%
GPT-2	23.2%	49%

Table 1: Accuracy for the different models on the Wikipedia test corpus of 450 sentence pairs. Context-aware approaches included the context of the difficult sentence when predicting.

GPT-2 with and without context). Both neural models significantly outperform the trigram language model; they have been trained on larger corpora and have access to more context allowing for better predictions. The RoBERTa is the best model here. One of the interesting point to point out is that RoBERTa, XLNet, BERT has a very small improvement with context while GPT-2 gains a large absolute improvement with context. In all the models, the performance is still very low and it suggests that this direction of research is open to new advancement in the future.

Table 2 shows the output of the RoBERTa model with and without context for simplifying the difficult sentence:

Each pseudostem can produce a single bunch of bananas.

The context-aware version is able to take advantage of the strong overlap between the difficult sentence and the simplified version that is being “typed”. The model without context makes reasonable predictions grammatically, but without the content priming the suggestions are poor overall.

3.3 Understanding model performance

To better understand how the models are performing and how the predictions of the models differ, I broke down the performances of the neural models by part of speech (POS), difficult sentence length, and the number of words typed so far.

POS Table 3 shows the accuracies broken down by part of speech, where the POS was automatically determined using Stanford CoreNLP (Manning et al., 2014). Due to the best performance of RoBERTa way better than the rest of the model, RoBERTa also outperforms the the other models in all of the POS-based prediction. Therefore, for the ensemble model, I need to find another way to

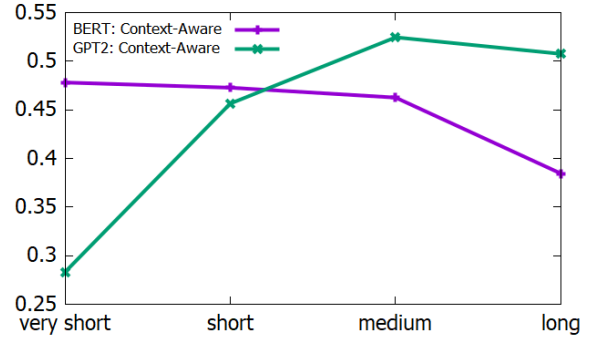


Figure 4: Accuracy for the two context-aware models based on the length of the difficult sentence: very short (≤ 5 tokens), short (6 – 15), medium (16 – 19), and long (≥ 20).

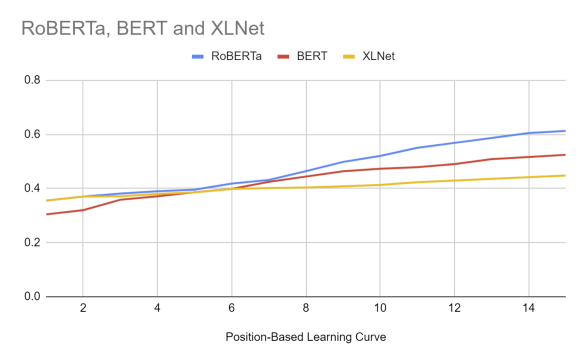


Figure 5: Accuracy for the context-aware models based on the number of words typed so far (i).

combine the models instead of hard coding based on the POS in the sentence.

Difficult sentence length Figure 4 shows the performance of the context-aware models based on the length of the difficult sentence. BERT is fairly consistent regardless of the length of the difficult sentence. Only for very long sentences does the performance drop. GPT-2 performs poorly on very short sentences, but well for other lengths. I hypothesize that the training data for GPT-2 (web text) may require more context for the more technical Wikipedia task.

Number of words typed Figure 5 shows the performance of the two context-aware models based on how many words of the simplification the model has access to, i.e., i . Early on when the sentence is first being developed, all models struggle. As more and more words are typed and more context is provided, the accuracy of both models increase. The increases starts to drop and the curves are flatten as more words given.

Accuracy@N performance Figure 4 shows the accuracy@N from RoBERTa, BERT, XLNet models on next word prediction. Accuracy@N is a

Typed so far	RoBERTa		
	No Context	Context-Aware	Actual
A	particle	pseudostem	pseudostem
A pseudostem	was	is	is
A pseudostem is	a	able	able
A pseudostem is able	to	to	to
A pseudostem is able to	create	produce	produce
A pseudostem is able to produce	a	a	a
A pseudostem is able to produce a	new	single	single
A pseudostem is able to produce a single	photon	bunch	bunch
A pseudostem is able to produce a single bunch	of	of	of
A pseudostem is able to produce a single bunch of	particles	bananas	bananas

Table 2: Sample output for simplifying the difficult sentence “Each pseudostem can produce a single bunch of bananas.” using RoBERTa with and without context. “Actual” indicates the word that should be predicted.

	RoBERTa	BERT	XLNet	GPT-2
All words	62.4	50.43	45.7	49
Nouns	60.3	48.7	45.2	51
Verbs	64	50.7	46.2	54
Adverbs	59.1	39.3	45.1	49
Adjectives	55	35.2	34.7	49
DET	76.3	68.7	51.2	51
PropNoun	25.8	21.8	17.9	34

Table 3: Accuracy of the RoBERTa, BERT, XLNet, and GPT-2 with and without context by part-of-speech on the test data.

	RoBERTa	BERT	XLNet
accuracy@2	67.2	54.5	46.9
accuracy@3	70	56.2	49.2
accuracy@4	72.1	58.0	51.3
accuracy@5	73.2	59.4	53.5
accuracy@6	73.2	59.4	53.5
accuracy@7	73.2	59.4	53.5

Table 4: Accuracy @ N of the RoBERTa, BERT, and XLNet with context on next word prediction

metric that gives a model credit as long as it can provide accurate prediction within its k suggestions. As we can see from here, this relaxing schema helps the models better assist technician because the user can pick the best word in the list of suggestion and therefore can help speed up the process and improve model performance.

Performance on predicting next K words Figure 5 show results from four neural network models in predicting the next k words. To further understand how performance affected by the more words it needs to predict, I run experiments with predicting next 1, 2, 3, and 4 words. This idea is the same with the Google email text suggestion when the model suggest multiple words at a time. As more

	RoBERTa	BERT	XLNet	GPT-2
Next 1	62.4	53.28	46.2	49
Next 2	45.1	38.7	33.5	41
Next 3	36.8	31.5	26.7	31
Next 4	31.5	24.2	21	14

Table 5: Accuracy of the RoBERTa, BERT, XLNet, and GPT-2 with context on next k prediction

words needed to predict, the performance drop significantly for the medical domain. This confirms that the task is unsolved and open for researches to further advance the autocompletion system in medical text simplification.

4 Conclusions

In this report, I have introduced a new task, text simplification autocompletion. Unlike most auto-complete tasks, for text simplification, models can be guided by the sentence that the user is simplifying. I compared a naive Bayes tri-gram language model with four recent neural models, RoBERTa, BERT, XLNet, and GPT-2, and showed how the difficult sentence could be incorporated into the prediction process. Using context resulted in significant increases in performance with the best model, RoBERTa with context, achieving a prediction accuracy of 62%, getting every other word right. I hope that this new task will allow for other interesting model adaptations to be explored. Besides, I also contribute the medical Wikipedia corpora for text simplification in medical domain,

References

- Fei Cai, Maarten De Rijke, et al. 2016. A survey of query auto completion in information retrieval. *Foundations and Trends® in Information Retrieval*.
- William Coster and David Kauchak. 2011. Learning to simplify sentences using wikipedia. In *Proceedings of the workshop on monolingual text-to-text generation*, pages 1–9. Association for Computational Linguistics.
- Andrew Dai, Benjamin Lee, Gagan Bansal, Jackie Tsay, Justin Lu, Mia Chen, Shuyuan Zhang, Tim Sohn, Yinan Wang, Yonghui Wu, et al. 2019. Gmail smart compose: Real-time assisted writing.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Mark D Dunlop and Andrew Crossan. 2000. Predictive text entry methods for mobile phones. *Personal Technologies*, 4(2-3):134–143.
- Spence Green, Sida I. Wang, Jason Chuang, Jeffrey Heer, Sebastian Schuster, and Christopher D. Manning. 2014. Human effort and machine learnability in computer aided translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1225–1236, Doha, Qatar. Association for Computational Linguistics.
- David Kauchak. 2013. Improving text simplification language modeling using unsimplified text data. In *Proceedings of the 51st annual meeting of the association for computational linguistics (volume 1: Long papers)*, pages 1537–1546.
- Nodira Khoussainova, YongChul Kwon, Magdalena Balazinska, and Dan Suciu. 2010. Snipsuggest: Context-aware autocompletion for sql. *Proceedings of the VLDB Endowment*.
- Nicholas Kloeckner, Gondy Leroy, David Kauchak, Yang Gu, Sonia Colina, Nicole P. Yuan, and Debra Revere. 2018. Improving consumer understanding of medical text: Development and validation of a new sub-simplify algorithm to automatically generate term explanations in english and spanish. *Journal of Medical Internet Research (JMIR)*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. 2014. The Stanford CoreNLP natural language processing toolkit. In *Association for Computational Linguistics (ACL) System Demonstrations*, pages 55–60.
- Daiki Nishihara, Tomoyuki Kajiwara, and Yuki Arase. 2019. Controllable text simplification with lexical constraint loss. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, pages 260–266.
- Alec Radford and Jeffrey Wu. 2018. language model and unsupervised multitask learning. *OpenAI*.
- Matthew Shardlow. 2014. A survey of automated text simplification. *International Journal of Advanced Computer Science and Applications*, 4(1):58–70.
- Andreas Stolcke. 2002. Srilm-an extensible language modeling toolkit. In *Proceedings of International Conference on Spoken Language Processing*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- 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*.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5754–5764.
- Xingxing Zhang and Mirella Lapata. 2017. Sentence simplification with deep reinforcement learning. *arXiv preprint arXiv:1703.10931*.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27.
- Zhemina Zhu, Delphine Bernhard, and Iryna Gurevych. 2010. A monolingual tree-based translation model for sentence simplification. In *Proceedings of ICCL*.