

# Simple PPDB: A Paraphrase Database for Simplification

**Ellie Pavlick**

University of Pennsylvania  
epavlick@seas.upenn.edu

**Chris Callison-Burch**

University of Pennsylvania  
ccb@cis.upenn.edu

## Abstract

We release the Simple Paraphrase Database, a subset of the Paraphrase Database (PPDB) adapted for the task of text simplification. We train a supervised model to associate simplification scores with each phrase pair, producing rankings competitive with state-of-the-art lexical simplification models. Our new simplification database contains 4.5 million paraphrase rules, making it the largest available resource for lexical simplification.

## 1 Motivation

Language is complex, and the process of reading and understanding language is difficult for many groups of people. The goal of text simplification is to rewrite text in order to make it easier to understand, for example, by children (De Belder and Moens, 2010), language learners (Petersen and Ostendorf, 2007), people with disabilities (Rello et al., 2013; Evans et al., 2014), and even by machines (Siddharthan et al., 2004). Automatic text simplification (Napoles and Dredze, 2010; Wubben et al., 2012; Xu et al., 2016) has the potential to dramatically increase access to information by making written documents available at all reading levels.

Full text simplification involves many steps, including grammatical restructuring and summarization (Feng, 2008). One of the most basic subtasks is *lexical simplification* (Specia et al., 2012)—replacing complicated words and phrases with simpler paraphrases. While there is active research in the area of lexical simplification (Coster and Kauchak, 2011a; Glavaš and Štajner, 2015; Paetzold, 2015), existing models have been by-and-large limited to single words. Often, how-

medical practitioner	→	doctor
legislative texts	→	laws
hypertension	→	high blood pressure
prevalent	→	very common
significant quantity	→	a lot
impact negatively	→	be bad

Table 1: In lexical simplification, it is often necessary to replace single words with phrases or phrases with single words. The above are examples of such lexical simplifications captured by the Simple PPDB resource.

ever, it is preferable, or even necessary to paraphrase a single complex word with multiple simpler words, or to paraphrase multiple words with a single word. For example, it is difficult to imagine a simple, single-word paraphrase of *hypertension*, but the three-word phrase *high blood pressure* is a very good simplification (Table 1). Such phrasal simplifications are overlooked by current lexical simplification models, and thus are often unavailable to the end-to-end text simplification systems that require them.

Recent research in data-driven paraphrasing has produced enormous resources containing millions of meaning-equivalent phrases (Ganitkevitch et al., 2013). Such resources capture a wide range of language variation, including the types of lexical and phrasal simplifications just described. In this work, we apply state-of-the-art machine learned models for lexical simplification in order to identify phrase pairs from the Paraphrase Database (PPDB) applicable to the task of text simplification. We introduce Simple PPDB,<sup>1</sup> a subset of the Paraphrase Database containing 4.5 million simplifying paraphrase rules. The large scale of Simple PPDB will support research into increasingly advanced methods for text simplification.

<sup>1</sup><http://www.seas.upenn.edu/~nlp/resources/simple-ppdb.tgz>

## 2 Identifying Simplification Rules

### 2.1 Paraphrase Rules

The Paraphrase Database (PPDB) is currently the largest available collection of paraphrases. Each paraphrase rule in the database has an automatically-assigned quality score between 1 and 5 (Pavlick et al., 2015). In this work, we use the PPDB-TLDR<sup>2</sup> dataset, which contains 14 million high-scoring lexical and phrasal paraphrases, and is intended to give a generally good tradeoff between precision and recall. To preprocess the data, we lemmatize all of the phrases, and remove rules which differ only by morphology, punctuation, or stop words, or which involve phrases longer than 3 words. The resulting list contains 7.5 million paraphrase rules covering 625K unique lemmatized words and phrases.

### 2.2 Lexical Simplification Model

Our goal is to build a model which can accurately identify paraphrase rules that both 1) simplify the input phrase and 2) preserve its meaning. That is, we want to avoid a model which favors “simple” words (e.g. *the*, *and*) even when they capture none of the meaning of the input phrase. We therefore train our model to make a three-way distinction between rules which simplify the input, rules which make the input less simple, and rules which generate bad paraphrases.

**Data.** We collect our training data in two phases. First, we sample 1,000 phrases from the vocabulary of the PPDB. We limit ourselves to words which also appear at least once in the Newsela corpus for text simplification (Xu et al., 2015), in order to ensure that we focus our model on the types of words for which the final resource is most likely to be applied. For each of these 1,000 words/phrases, we sample up to 10 candidate paraphrases from PPDB, stratified evenly across paraphrase quality scores. We ask workers on Amazon Mechanical Turk to rate each of the chosen paraphrase rules on a scale from 1 to 5 to indicate how well the paraphrase preserves the meaning of the original phrase. We use the same annotation design used in Pavlick et al. (2015). We have 5 workers judge each pair, omitting workers who do not provide correct answers on the embedded gold-standard pairs which we draw from WordNet. For 62% of the paraphrase rules we had

<sup>2</sup><http://paraphrase.org/#/download>

scored, the average human rating falls below 3, indicating that the meaning of the paraphrase differs substantially from that of the input. We assign all of these rules to the “bad paraphrase” class.

We take the remaining 3,758 meaning-preserving paraphrase rules (scored  $\geq 3$  in the above annotation task) and feed them into a second annotation task, in which we identify rules that simplify the input. We use the same annotation interface as in Pavlick and Nenkova (2015), which asks workers to choose which of the two phrases is simpler, or to indicate that there is no difference in complexity. We collect 7 judgements per pair and take the majority label, discarding pairs for which the majority opinion was that there was no difference. We include each rule in our training data twice, once as an instance of a “simplifying” rule, and once in the reverse direction as an instance of a “complicating” rule.

In the end, our training dataset contains 11,829 pairs, with the majority class being “bad paraphrase” (47%), and the remaining split evenly between “simplifying” and “complicating” paraphrase rules (26% each).

**Features.** We use a variety of features that have been shown in prior work to give good signal about phrases’ relative complexity. The features we include are as follows: phrase length in words and in characters, frequency according to the Google NGram corpus (Brants and Franz, 2006), number of syllables, the relative frequency of usage in Simple Wikipedia compared to normal Wikipedia (Pavlick and Nenkova, 2015), character unigrams and bigrams, POS tags, and the averaged Word2Vec word embeddings for the words in the phrase (Mikolov et al., 2013). For each phrase pair  $\langle e_1, e_2 \rangle$ , for each feature  $f$ , we include  $f(e_1)$ ,  $f(e_2)$  and  $f(e_1) - f(e_2)$ .<sup>3</sup> We also include the cosine similarity of the averaged word embeddings and the PPDB paraphrase quality score as features.

We train a multi-class logistic regression model<sup>4</sup> to predict if the application of a paraphrase rule will result in 1) simpler output, 2) more complex output, or 3) non-sense output.

**Performance.** Table 2 shows the performance of the model on cross-validation, compared to several baselines. The full model achieves 60% accuracy,

<sup>3</sup>We do not compute the difference  $f(e_1) - f(e_2)$  for sparse features, i.e. character ngrams and POS tags.

<sup>4</sup><http://scikit-learn.org/>

	Acc	Prec.
Random	47.1%	0.0%
Simple/Regular Wiki. Ratio	49.1%	47.6%
Length in Characters	51.4%	47.3%
Google Ngram Frequency	51.4%	44.2%
Number of Syllables	51.5%	45.3%
Supervised Model, W2V	54.7%	46.3%
<b>Supervised Model, Full</b>	<b>60.4%</b>	<b>52.9%</b>

Table 2: Accuracy on 10-fold cross-validation, and precision for identifying simplifying rules. Folds are constructed so that train and test vocabularies are disjoint.

5 points higher than the strongest baseline, a supervised model which uses only word embeddings as features.

### 2.3 Simple PPDB

We run the trained model described above over all 7.5 million paraphrase rules. From the predictions, we construct Simple PPDB: a list of 4.5 million simplifying paraphrase rules. A rule in Simple PPDB is represented as a triple, consisting of a syntactic category, and input phrase, and a simplified output phrase. Each rule is associated with both a paraphrase quality score from 1 to 5 (taken from PPDB 2.0), and simplification confidence score from 0 to 1.0 (our classifier’s confidence in the prediction that the rule belongs to the “simplifying” class). Note that ranking via the confidence scores of a classification model has not, to our knowledge, been explored in previous work on lexical simplification. The remainder of this paper evaluates the quality of the simplification ranking. For an evaluation of the paraphrase quality ranking, see Pavlick et al. (2015). Table 3 shows examples of some of the top ranked paraphrases according to Simple PPDB’s simplification score for several input phrases.

## 3 Evaluation

To evaluate Simple PPDB, we apply it in a setting intended to emulate the way it is likely to be used in practice. We use the Newsela Simplification Dataset (Xu et al., 2015), a corpus of manually simplified news articles. This corpus is currently the cleanest available simplification dataset and is likely to be used to train and/or evaluate the simplification systems that we envision benefitting most from Simple PPDB.

We draw a sample of 100 unique word types (“targets”) from the corpus for which Simple

PPDB has at least one candidate simplification. For each target, we take Simple PPDB’s full list of simplification rules which are of high quality according to the PPDB 2.0 paraphrase score<sup>5</sup> and which match the syntactic category of the target. On average, Simple PPDB proposes 8.8 such candidate simplifications per target.

**Comparison to existing methods.** Our baselines include three existing methods for generating lists of candidates that were proposed in prior work. The methods we test for generating lists of candidate paraphrases for a given target are: the **WordNetGenerator**, which pulls synonyms from WordNet (Devlin and Tait, 1998; Carroll et al., 1999), the **KauchakGenerator**, which generates candidates based on automatic alignments between Simple Wikipedia and normal Wikipedia (Coster and Kauchak, 2011a), and the **GlavasGenerator**, which generates candidates from nearby phrases in vector space (Glavaš and Štajner, 2015) (we use the pre-trained Word2Vec VSM (Mikolov et al., 2013)).

For each generated list, we follow Horn et al. (2014)’s supervised SVM Rank approach to rank the candidates for simplicity. We reimplement the main features of their model: namely, word frequencies according to the Google NGrams corpus (Brants and Franz, 2006) and the Simple Wikipedia corpus, and the alignment probabilities according to automatic word alignments between Wikipedia and Simple Wikipedia sentences (Coster and Kauchak, 2011b). We omit the language modeling features since our evaluation does not consider the context in which the substitution is to be applied.

All of these methods (the three generation methods and the ranker) are implemented as part of the LEXenstein toolkit (Paetzold and Specia, 2015). We use the LEXenstein implementations for the results reported here, using off-the-shelf configurations and treating each method as a black box.

**Setup.** We use each of the generate-and-rank methods to produce a ranked list of simplification candidates for each of the 100 targets drawn from the Newsela corpus. When a generation method fails to produce any candidates for a given target, we simply ignore that target for that particular method. This is to avoid giving Simple PPDB

<sup>5</sup>Heuristically, we define “high quality” as  $\geq 3.5$  for words and  $\geq 4$  for phrases.

keenly	omit	employment opportunity	remedied
most strongly	leave out	a new job	set right
deeply	delete it	opportunity	be fixed
strongly	be removed	business opportunity	be corrected
eagerly	forget about it	the job	to be resolved
very	be ignored	labour	be solved

Table 3: Examples of top-ranked simplifications proposed by Simple PPDB for several input words. Often, the best simplification for a single word is a multiword phrase, or vice-versa. These many-to-one mappings are overlooked when systems use only length or frequency as a proxy for simplicity.

an unfair advantage, since, by construction, PPDB will have full coverage of our list of 100 targets. In the end, the GlavasGenerator is evaluated over 95, the WordNetGenerator over 82, and the KauchakGenerator over 48. The results in Table 4 do not change significantly if we restrict all systems to the 48 targets which the KauchakGenerator is capable of handling. Since the GlavasGenerator is capable of producing an arbitrary number of candidates for each target, we limit the length of each of its candidate lists to be equal to the number of candidates produced by Simple PPDB for that same target.

**Human judgments.** For each of the proposed rules from all four systems, we collect human judgements on Amazon Mechanical Turk, using the same annotation interface as before. That is, we ask 7 workers to view each pair and indicate which of the two phrases is simpler, or to indicate that there is no difference. We take the majority label to be the true label for each rule. Workers show moderate agreement on the 3-way task ( $\kappa = 0.4 \pm 0.03$ ), with 14% of pairs receiving unanimous agreement and 37% receiving the same label from 6 out of 7 annotators. We note that the  $\kappa$  metric is likely a lower bound, as it punishes low agreement on pairs for which there is little difference in complexity, and thus the “correct” answer is not clear (e.g. for the pair *<matter, subject>*, 3 annotators say that *matter* is simpler, 2 say that *subject* is simpler, and 2 say there is no difference).

**Results.** Table 4 compares the different methods in terms of how well they rank simplifying rules above non-simplifying rules. Simple PPDB’s ranking of the relative simplicity achieves an averaged precision of 0.72 (0.77 P@1), compared to 0.70 (0.69 P@1) achieved by the Horn et al. (2014) system— i.e. the KauchakGenerator+SVM Ranker. We hypothesize that the performance difference between these two ranking systems is

	Avg. Prec.	P@1
Glavas+SVR	0.21	0.13
Wordnet+SVR	0.53	0.50
Kauchak+SVR	0.70	0.69
<b>Simple PPDB</b>	<b>0.72</b>	<b>0.77</b>

Table 4: Precision of relative simplification rankings of three existing lexical simplification methods compared to the Simple PPDB resource in terms of Average Precision and P@1 (both range from 0 to 1 and higher is better). All of the existing methods were evaluated using the implementations as provided in the LEXenstein toolkit.

likely due to a combination of the additional features applied in Simple PPDB’s model (e.g. word embeddings) and the difference in training data (Simple PPDB’s model was trained on 11K paraphrase pairs with trinary labels, while the Horn et al. (2014) model was trained on 500 words, each with a ranked list of paraphrases). Table 5 provides examples of the top-ranked simplification candidates proposed by each of the methods described.

<b>alarm</b>	
Glavas	<b>enrage</b> , perturb, stun
WordNet	horrify, dismay, <b>alert</b> , appal, appal
Kauchak	<b>pure</b> , <b>worry</b>
PPDB	<b>worry</b> , <b>concern</b> , <b>alert</b>
<b>genuine</b>	
Glavas	credible, <b>sort</b> , feign, <b>phoney</b> , good naturedness, <b>sincere</b> , <b>sincerely</b> , insincere, bonafide
WordNet	<b>real</b> , <b>actual</b> , unfeigned, literal, echt, <b>true</b>
Kauchak	thermal
PPDB	<b>true</b> , <b>real</b> , <b>actual</b> , <b>honest</b> , <b>sincere</b>

Table 5: Examples of candidate simplifications proposed by Simple PPDB and by three other generate-and-rank methods. Bold words were rated by humans to be simpler than the target word. Note that these candidates are judged on simplicity, not on their goodness as paraphrases.

In addition, Simple PPDB offers the largest coverage (Table 6). It has a total vocabulary of 624K unique words and phrases, and provides the largest number of potential simplifications for

	Avg. PPs per Input	Total Vocab.
Glavas+SVR	$\infty$	$\infty$
Kauchak+SVR	4.4	127K
Wordnet+SVR	6.7	155K
Simple PPDB	8.8	624K

Table 6: Overall coverage of three existing lexical simplification methods compared to the Simple PPDB resource. Glavas is marked as  $\infty$  since it generates candidates based on nearness in vector space, and in theory could generate as many words/phrases as are in the vocabulary of the vector space.

each target— for the 100 targets drawn from the Newsela corpus, PPDB provided an average of 8.8 candidates per target. The next best generator, the WordNet-based system, produces only 6.7 candidates per target on average, and has a total vocabulary of only 155K words.

## 4 Conclusion

We have described Simple PPDB, a subset of the Paraphrase Database adapted for the task of text simplification. Simple PPDB is built by applying state-of-the-art machine learned models for lexical simplification to the largest available resource of lexical and phrasal paraphrases, resulting in a web-scale resource capable of supporting research in data-driven methods for text simplification. We have shown that Simple PPDB offers substantially increased coverage of both words and multiword phrases, while maintaining high quality compared to existing methods for lexical simplification. Simple PPDB, along with the human judgements collected as part of its creation, is freely available with the publication of this paper.<sup>6</sup>

## Acknowledgments

This research was supported by a Facebook Fellowship, and by gifts from the Alfred P. Sloan Foundation, Google, and Facebook. This material is based in part on research sponsored by the NSF grant under IIS-1249516 and DARPA under number FA8750-13-2-0017 (the DEFT program). The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes. The views and conclusions contained in this publication are those of the authors and should not be interpreted as representing official policies or endorsements of DARPA and the U.S. Government.

<sup>6</sup><http://www.seas.upenn.edu/~nlp/resources/simple-ppdb.tgz>

We would especially like to thank Ani Nenkova for suggesting this line of research and for providing the initial ideas on which this work builds. We would also like to thank Courtney Napoles and Wei Xu for valuable discussions, the anonymous reviewers for thoughtful comments, and the Amazon Mechanical Turk annotators for their contributions.

## References

- Thorsten Brants and Alex Franz. 2006. Web 1T 5-gram Version 1. *Linguistic Data Consortium, Philadelphia*.
- John Carroll, Guido Minnen, Darren Pearce, Yvonne Canning, Siobhan Devlin, and John Tait. 1999. Simplifying text for language-impaired readers. In *Proceedings of EACL*, volume 99, pages 269–270.
- Will Coster and David Kauchak. 2011a. Learning to simplify sentences using Wikipedia. In *Proceedings of the Workshop on Monolingual Text-To-Text Generation*, pages 1–9, Portland, Oregon, June. Association for Computational Linguistics.
- William Coster and David Kauchak. 2011b. Simple English Wikipedia: A new text simplification task. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 665–669, Portland, Oregon, USA, June. Association for Computational Linguistics.
- Jan De Belder and Marie-Francine Moens. 2010. Text simplification for children. In *Proceedings of the SIGIR workshop on accessible search systems*, pages 19–26. ACM.
- Siobhan Devlin and John Tait. 1998. The use of a psycholinguistic database in the simplification of text for aphasic readers. *Linguistic databases*, pages 161–173.
- Richard Evans, Constantin Orasan, and Iustin Dornescu. 2014. An evaluation of syntactic simplification rules for people with autism. In *Proceedings of the 3rd Workshop on Predicting and Improving Text Readability for Target Reader Populations (PITR)*, pages 131–140.
- Lijun Feng. 2008. Text simplification: A survey. *The City University of New York, Tech. Rep.*
- Juri Ganitkevitch, Benjamin Van Durme, and Chris Callison-Burch. 2013. PPDB: The paraphrase database. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 758–764, Atlanta, Georgia, June. Association for Computational Linguistics.

- 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, July. Association for Computational Linguistics.
- Colby Horn, Cathryn Manduca, and David Kauchak. 2014. Learning a lexical simplifier using Wikipedia. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 458–463, Baltimore, Maryland, June. Association for Computational Linguistics.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Courtney Napoles and Mark Dredze. 2010. Learning Simple Wikipedia: A cogitation in ascertaining abecedarian language. In *Proceedings of the NAACL HLT 2010 Workshop on Computational Linguistics and Writing: Writing Processes and Authoring Aids*, pages 42–50, Los Angeles, CA, USA, June. Association for Computational Linguistics.
- Gustavo Paetzold and Lucia Specia. 2015. LEXenstein: A framework for lexical simplification. In *Proceedings of ACL-IJCNLP 2015 System Demonstrations*, pages 85–90, Beijing, China, July. Association for Computational Linguistics and The Asian Federation of Natural Language Processing.
- Gustavo Paetzold. 2015. Reliable lexical simplification for non-native speakers. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 9–16, Denver, Colorado, June. Association for Computational Linguistics.
- Ellie Pavlick and Ani Nenkova. 2015. Inducing lexical style properties for paraphrase and genre differentiation. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 218–224, Denver, Colorado, May–June. Association for Computational Linguistics.
- Ellie Pavlick, Pushpendre Rastogi, Juri Ganitkevitch, Benjamin Van Durme, and Chris Callison-Burch. 2015. PPDB 2.0: Better paraphrase ranking, fine-grained entailment relations, word embeddings, and style classification. 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 425–430, Beijing, China, July. Association for Computational Linguistics.
- Sarah E. Petersen and Mari Ostendorf. 2007. Text simplification for language learners: a corpus analysis. In *SLaTE*, pages 69–72. Citeseer.
- Luz Rello, Ricardo A. Baeza-Yates, and Horacio Sag-gion. 2013. The impact of lexical simplification by verbal paraphrases for people with and without dyslexia. In *CICLing (2)*, pages 501–512.
- Advaith Siddharthan, Ani Nenkova, and Kathleen McKeown. 2004. Syntactic simplification for improving content selection in multi-document summarization. In *Proceedings of the 20th international conference on Computational Linguistics*, page 896. Association for Computational Linguistics.
- Lucia Specia, Sujay Kumar Jauhar, and Rada Mihalcea. 2012. SemEval-2012 task 1: English lexical simplification. In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 347–355, Montréal, Canada, 7-8 June. 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: Long Papers-Volume 1*, pages 1015–1024. 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.