

# Generating Bags of Words from the Sums of their Word Embeddings

No Institute Given

**Abstract.** Converting a sentence to a meaningful vector representation has uses in many NLP tasks, however very few methods allow the words to be recovered from that sentence vector representation. Being able to generate sentences from the vector representations is expected to open up many new applications. As a partway step towards this, we introduce a method for moving from sum of word embedding (SOWE) representations back to the bag of words (BOW) for the original sentences. This is done using a greedy algorithm to convert the vector to a bag of words. To our knowledge this is the first such work. It demonstrates qualitatively the ability to recreate the words from a large corpus based on its sentence embeddings. As well as practical applications for allowing classical information retrieval methods to be combined with more recent methods using the sums of word embeddings, the success of this method has theoretical implications on the degree of information maintained by the sum of embeddings representation. This lends some credence to the consideration of the SOWE as a dimensionality reduced, and meaning enhanced, data manifold for the bag of words.

## 1 Introduction

The task being tackled here is the *resynthesis* of bags of words (BOW) from sentence vector representations. In particular the generation of BOW from vectors based on the sum of the sentence’s constituent words’ embeddings (SOWE). To the knowledge of the authors, this task has not been attempted before.

The motivations for this task are the same as in the related area of sentence generation. Dinu and Baroni (2014) observe that given a sentence encodes its meaning, and the vector encodes the same meaning, then it must be possible to translate in both directions between the natural language and the vector representation. A subset of this is the unordered case (BOW), rather than true sentences, which we tackle in this paper. The success of the implementation does indicate the validity of this dual space theory, for the representations considered (where order is neglected). There are also some potential practical applications of such an implementation, often ranging around common vector space representations.

Given suitable bidirectional methods for converting between sentence vectors and bags of words, the sentence vector space can be employed as a *lingua franca* for translation between various forms of information – though with loss of word order

information. The most obvious of which is literal translation between different natural languages; however the use extends beyond this.

Several approaches have been developed for representing images and sentences in a common vector space. This is then used to select a suitable caption a list of candidates (Farhadi et al. 2010; Socher et al. 2014). Similar methods, creating a common space between images and SOWE of the keywords describing them, could be used to generate keyword descriptions using BOW resynthesis – without any need for a list. This would allow classical word-based information retrieval and indexing techniques to be applied to images.

A similar use is the replacement of vector based extractive summarisation (Kågebäck et al. 2014; Yogatama et al. 2015), with keyword based abstractive summarisation, the generation of a keyword summary from a document. The promising use of SOWE generation for all these applications is to have a separate model trained to take the source information (e.g. a picture for image description, or a cluster of sentences for abstract summarisation) as its input and train it to output a vector which is close to a target SOWE vector. This output can then be used to generate the sentence.

The method proposed in this paper has an input of a sum of word embeddings (SOWE) sentence vector, and outputs the bag of word (BOW) which it corresponds to. The input is a vector for example  $\tilde{s} = [-0.79, 1.27, 0.28, \dots, -1.29]$ , which approximates a SOWE vector, and outputs a BOW for example  $\{, : 1, \text{best}: 1, \text{it}: 2, \text{of}: 2, \text{the}: 2, \text{times}: 2, \text{was}: 2, \text{worst}: 1\}$  – the BOW for the opening line of Dickens’ *Tale of Two Cities*. Our method for BOW generation is shown in Figure 1, note that it takes as input only a word embedding vocabulary and the vector to generate the BOW from.

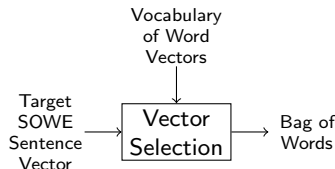


Fig. 1: The process for the regenerating sentences from SOWE-type sentence vectors.

The rest of the paper is organized into the following sections. Section 2 introduces the area, discussing in general sentence models, and prior work on generation. Section 3 explains the problem in detail and our algorithm for solving it. Section 4 described the settings used for evaluation. Section 5 discusses the results of this evaluation. The paper presents its conclusions in Section 6, including a discussion of future work.

## 2 Background

A bag of words is a classical natural language processing method for representing a text, sentence or document, commonly used in information retrieval. The text is replaced with a multiset – an unordered count of how often each word occurs

Word embeddings are vector representations of words. They have been shown to encode important syntactic and semantic properties. There are many different types of word vectors (Yin and Schtze 2015). Two of the more notable are the SkipGrams of Mikolov et al. (2013a,b) and the Global Vector word representations (GloVe) of Pennington et al. (2014). Beyond word representations are sentence vectors.

Sentence vectors represent sentences – they are often derived from word vectors. Like word vectors they can capture semantic and syntactic features. Sentence vector creation methods include the works of Le and Mikolov (2014) and Socher (2014). Far simpler than those methods, is the sum of word embeddings (SOWE). SOWE (like BOW) draws significant criticism for not only disregarding sentence structure, but disregarding word order entirely when producing the sentence vector. However this weaknesses, may be offset by the improved discrimination allowed through words directly affecting the sentence vector. It avoids the potential information loss through the indirection of more complex methods. Recent results suggest that this may allow it to be comparable overall to the more linguistically consistent embeddings when it comes to representing meaning.

White et al. (2015) found that when classifying real-world sentences into groups of semantically equivalent paraphrases, that using SOWE as the input resulted in very accurate classifications. In that work White et. al. partitioned the sentences into groups of paraphrases, then evaluated how well a linear SVM could classify unseen sentences into the class given by its meaning. They used this to evaluate a variety of different sentence embeddings techniques. They found that the classification accuracy when using SOWE as the input performed very similarly to the best performing methods – less that 0.6% worse on the harder task. From this they concluded that the mapping from the space of sentence meaning to the vector space of the SOWE, resulted in sentences with the same meaning going to distinct areas of the vector space.

Ritter et al. (2015) presented a similar task on spacial-positional meaning, which used carefully constructed artificial data, for which the meanings of the words interacted non-simply – thus theoretically favouring the more complex sentence embeddings. In their evaluation the task was classification with a Naive Bayes classifier into one of five categories of different spatial relationships. The best of the SOWE models they evaluated, outperformed the next best model by over 5%. These results suggest this simple method is still worth consideration for many sentence vector representation based tasks. SOWE is the basis of the work presented in this paper.

### 3 The Vector Selection Problem

At the core of this problem is what we will call the Vector Selection Problem, to select word vectors which sum to be closest to the target SOWE (the input). The word vectors come from a known vector vocabulary, and are selected with potential repetition. Selecting the vectors equates to selecting the words, because there is a one to one correspondence between the word embedding vectors and their words. This relies on no two words having exactly the same embeddings – which is true for all current word embedding techniques.

**Definition 1.** *The Vector Selection Problem is defined on  $(\mathcal{V}, \tilde{s}, d)$  for a finite vocabulary of vectors  $\mathcal{V}$ ,  $\mathcal{V} \subset \mathbb{R}^n$ , a target sentence vector  $\tilde{s}$ ,  $\tilde{s} \in \mathbb{R}^n$ , and any distance metric  $d$ , by:*

$$\operatorname{argmin}_{\{\forall \tilde{c} \in \mathbb{N}_0^V\}} d(\tilde{s}, \sum_{\tilde{x}_j \in \mathcal{V}} \tilde{x}_j c_j)$$

*$\tilde{x}_j$  is the vector embedding for the  $j$ th word in the vocabulary  $\tilde{x}_j \in \mathcal{V}$  and  $c_j$  is the  $j$ th element of the count vector  $\tilde{c}$  being optimised – it is the count of how many times the  $x_j$  occurs in approximation to the sum being assessed; and correspondingly it is the count of how many times the  $j$ th word from the vocabulary occurs in the bag of words. The selection problem is thus finding the right words with the right multiplicity, such that the sum of their vectors is as close to the input target vector,  $\tilde{s}$ , as possible.*

#### 3.1 NP-Hard Proof

The vector selection problem is NP-Hard. It is possible to reduce from any given instance of a *subset sum problem* to a vector selection problem. The *subset sum problem* is NP-complete (Karp 1972). It is defined: for some set of integers ( $\mathcal{S} \subset \mathbb{Z}$ ), does there exist a subset ( $\mathcal{L} \subseteq \mathcal{S}$ ) which sums to zero ( $0 = \sum_{l_i \in \mathcal{L}} l_i$ ). A suitable metric, target vector and vocabulary of vectors corresponding to the elements  $\mathcal{S}$  can be defined by a bijection; such that solving the vector selection problem will give the subset vectors corresponding a subset of  $\mathcal{S}$  with the smallest sum; which if zero indicates that the subset sum does exist, and if nonzero indicates that no such subset ( $\mathcal{L}$ ) exists. A fully detailed proof of the reduction from subset sum to the vector selection problem can be found on the first author's website.<sup>1</sup>

#### 3.2 Selection Algorithm

The algorithm proposed here to solve the selection problem is a greedy iterative processes that continues to convergence. It is a fully deterministic method, requiring no training, beyond having the word vector mapping provided. In each iteration, first a greedy search (Greedy Addition) for a path to the targeted

<sup>1</sup> [[URL removed for blinding.]]

sum point  $\tilde{s}$  is done, followed by correction with a substitution based step (n-Substitution). This process is repeated until no change is made to the path. The majority of the selection is done in the Greedy Addition step, while the n-substitution handles fine tuning.

**Greedy Addition** The greedy addition step is characterised by adding the best vector to the bag at each step (see the pseudo-code in Algorithm 1). At each step, all the vectors in the bag are summed, and then each vector in the vocabulary is added in turn to evaluate the new distance the new bag would have from the target, the bag which sums to be closest to the target becomes the current solution. This continues until there is no option to add any of the vectors without moving the sum away from the target.

Greedy Addition works surprisingly well on its own, but it is enhanced with a fine tuning step to decrease its greediness.

```

Data: the metric  $d$ 
the target sum  $\tilde{s}$ 
the vocabulary of vectors  $\mathcal{V}$ 
The current best bag of vectors  $bag_c$ : initially  $\emptyset$ 
Result: The modified  $bag_c$  which sum to be as close as greedy search can get to
the target  $\tilde{s}$ , under the metric  $d$ 

begin
   $\tilde{t} \leftarrow \sum_{x_i \in bag_c} x_i$ ;
  while true do
     $\tilde{x}^* \leftarrow \operatorname{argmin}_{x_j \in \mathcal{V}} d(\tilde{s}, \tilde{t} + \tilde{x}_j)$ ;
    if  $d(\tilde{s}, \tilde{t} + \tilde{x}^*) < d(\tilde{s}, \tilde{t})$  then
       $\tilde{t} \leftarrow \tilde{t} + \tilde{x}^*$ ;
       $bag_c \leftarrow bag_c \cup \{\tilde{x}^*\}$ ;
    else
      /* No further improving step found */
      return  $bag_c$ 
    end
  end
end

```

**Algorithm 1:** Greedy Addition. In practical implementation, the bag of vectors can be represented as list of indexes into columns of the embedding vocabulary matrix, and efficient matrix summation methods can be used.

**n-Substitution** We define a new substitution based method for fine tuning solutions called n-substitution. It can be described as considering all subbags containing up to  $n$  elements, consider replacing them with a new sub-bag of up

that size from the vocabulary, including none at all, if that would result in the overall bag getting closer to the target  $\tilde{s}$ .

The reasoning behind performing the  $n$ -substitution is to correct for greedy mistakes. Consider the 1 dimensional case where  $\mathcal{V} = 24, 25, 50$  and  $\tilde{s} = 98$ ,  $d(x, y) = |x - y|$ . Greedy addition would give  $bag_c = [50, 25, 24]$  for a norm-distance of 1, but a perfect solution is  $bag_c = [50, 24, 24]$  which is found using 1-substitution. This substitution method can be looked at as re-evaluating past decisions in light of the future decisions. In this way it lessens the greed of the addition step.

The  $n$ -substitution step has time complexity of  $O(\binom{C}{n} V^n)$  – for  $C = \sum \tilde{c}$  i.e. current cardinality of  $bag_c$ . With large vocabularies it is only practical to consider 1-substitution. With the Brown Corpus, where  $V \approx 40,000$ , it was found that 1-substitution provides a significant improvement over greedy addition alone. On a smaller trial corpora, where  $V \approx 1,000$ , 2-substitution was used and found to give further improvement. In general it is possible to initially use 1-substitution, with the overall algorithm, and if the algorithm converges to a poor solution (given the distance to the target is always known), then the selection algorithm can be retried from the converged solution, using 2-substitution and so forth. As  $n$  increases the greed decreases; at the limit the overall algorithm is not greedy at all, but is rather an exhaustive search.

## 4 Experimental Setup and Evaluations

### 4.1 Word Embeddings

GloVe representations of words are used in our evaluations (Pennington et al. 2014). There are many varieties of word embeddings which function with our algorithm. GloVe was chosen simply because of the availability of a large pre-trained vocabulary of vectors. The representations used for evaluation were pretrained on 2014 Wikipedia and Gigaword 5<sup>2</sup>. Preliminary results with SkipGrams from Mikolov et al. (2013a) suggested similar performance.

### 4.2 Corpora

The evaluation was performed on the Brown Corpus (Francis and Kucera 1979) and on a subset of the Books Corpus (Zhu et al. 2015). The Brown Corpus was sourced with samples from a 500 fictional and non-fictional works from 1961. The Books Corpus was sourced from 11,038 unpublished novels. The Books Corpus is extremely large, containing roughly 74 million sentences. After preprocessing we randomly selected 0.1% of these for evaluation.

For simplicity of evaluation, sentences containing words not found in the pretrained vector vocabulary are excluded. These were generally rare mis-spellings and unique numbers (such as serial numbers). Similarly, words which are not used in the corpus are excluded from the vector vocabulary.

<sup>2</sup> Kindly made available online at <http://nlp.stanford.edu/projects/glove/>

After the preprocessing the final corpora can be described as follows. The Brown Corpus has 42,004 sentences and a vocabulary of 40,485 words. Where-as, the Books Corpus has 66,464 sentences, and a vocabulary of 178,694 words. These corpora remain sufficiently large and complex to quantitatively evaluate the algorithm

### 4.3 Vector Selection

The Euclidean metric was used to measure how close potential solutions were to the target vector. The choice of distance metric controls how close each vector is considered to the partial sum during the greedy selection. Preliminary results on one-tenth of the Books Corpus used in the main evaluation found the city-block metric performed marginally worse than the Euclidean metric and took significantly longer to converge. The cosine similarity was not used due to issues with it not being a true distance metric.

The commonly used cosine similarity, or the linked angular distance, have an issue of zero distances between distinct points – making them not true distance metrics. For example the SOWE of “*a can can can a can*” has a zero distance under those measures to the SOWE for “*a can can*”<sup>3</sup>. That example is pathological, though a valid sentence fragment – the former referring to a canisters ability to be used to contain another can; the later to the capacity of a canister to do an undisclosed task. True metrics such as the Euclidean metric do not have this problem. Further investigation may find other better distance metrics for this step.

The Julia programming language (Bezanson et al. 2014), was used to create the implementation of the method, and the evaluation scripts for the results presented in the next section. This implementation, evaluation scripts, and the raw results are available on-line.<sup>4</sup>

## 5 Results and Discussion

Examples of the output are shown in Table 1. In that table, 8 sentences which were used for demonstration of sentence generation in Bowman et al. (2015) and Iyyer et al. (2014) have the BOW generation results shown. All examples except *a* and *f* are perfect. Example *f* is interesting as it seems that the contraction token *'re* was substituted for *are*, and *do* for *doing*. Inspections of the execution logs for running on the examples show that this was a greedy mistake that would be corrected using 2-substitution. Example *a* has for more mistakes.

The mistakes in Example *f* are seem to be related to unusual nonword tokens. Such as the three tokens with 13, 34, and 44 repetitions of the underscore character. These tokens appear in the very large Books corpus, and in the Wikipedia/Gigaword pretraining data used for word embeddings, but are generally

<sup>3</sup> The same is true for any number of repetitions of *buffalo* – each of which forms a valid sentence as noted in **tymoczko1995sweet**

<sup>4</sup> [[Authors’ URL Blinded for Review. See software/data submission]]

[illegible]



Corpus	Word Embedding Dimensions	Portion Perfect	Mean Jaccard Score	Mean Precision	Mean Recall	Mean F1 Score
Brown	50	6.3%	0.175	0.242	0.274	0.265
Brown	100	19.4%	0.374	0.440	0.530	0.477
Brown	200	44.7%	0.639	0.695	0.753	0.720
Brown	300	70.4%	0.831	0.864	0.891	0.876
Books	300	75.6%	0.891	0.912	0.937	0.923

Table 2: The performance of the BOW generation method. Note the final line is for the Books Corpus, where-as the preceding are for the Brown Corpus.

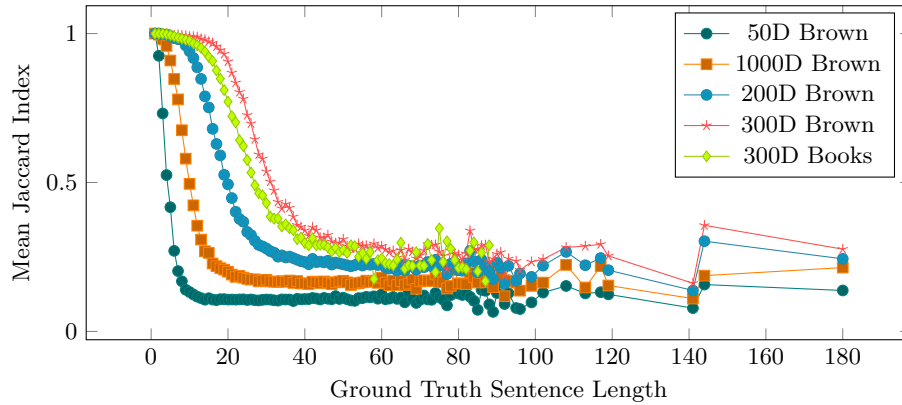


Fig. 2: The mean Jaccard index achieved during the word selection step, shown against the ground truth length of the sentence. Note that the vast majority of sentences are in the far left end of the plot. The diminishing samples is also the cause of the noise, as the ground length increases.

devoid of meaning and are used as structural elements for formatting. We theorise that because of their rarity in the pre-training data they are assigned an unusual word-embedding by GloVe. There occurrence in this example suggests that better results may be obtained by pruning the vocabulary. Either manually, or via a minimum uni-gram frequency requirement. The examples overall highlight the generally high performance the method, evaluations on the full corpora confirm this.

The Table 2 shows the quantitative performance of our method across both corpora. Five measures are reported. The most clear is the portion of exact matches – this is how often out of all the trials the method produced the exact correct bag of words. The remaining measures are all means across all the values of the measures in each trial. The Jaccard index is the portion of overlap between the reference BOW, and the output BOW – it is the cardinality of the intersection divided by that of the union. The precision is the portion of the output words

that were correct; and the recall is the portion of all correct words which were output. The  $F_1$  score is the harmonic mean of precision and recall.

The recall is higher than the precision, indicating that the method is more prone to produce additional incorrect words (lowering the precision), than to miss words out (which would lower the recall).

Initial investigation focused on the relationship between the number of dimensions in the word vector and the performance. This was carried out on the smaller Brown corpus. Results confirmed the expectation that higher dimensional embeddings allow for better generation of words. The best performing embedding size (i.e. the largest) was then used to evaluate success on the Books Corpus. The increased accuracy when using higher dimensionality embeddings remains true at all sentence lengths.

As can be seen in Figure 2 sentence length is a very significant factor in the performance of our method. As the sentences increase in length, the number of mistakes increases. However, at higher embedding dimensionality the accuracy for most sentences is high. This is because of most sentences being reasonably short. The third quartile on sentence length is 25 words for Brown, and 17 for the Books Corpus. This distribution difference is also responsible for the apparent better results on the Books Corpus, than on the Brown corpus.

While the results shown in Table 2 suggest that the Books corpus was better performing, this is due to its much shorter average sentence length. When taken as a function of the sentence length, as shown in Figure 2, performance on the Books Corpus is worse than on the Brown Corpus. It can be concluded from this result that increasing the size of the vocabulary decreases successful BOW regeneration. Books Corpus vocabulary being over four times larger results in lower performance, and the other factors the same. An important conclusion is that increase vocabulary size has less impact than increasing the sentence length or the embedding dimensionality on the method performance.

## 6 Conclusion

A method was presented for how to regenerate a bag of words, from the sum of a sentence's word embeddings. The word selection problem, of going from the sum of embeddings to the words, is NP-Hard. A greedy algorithm was found to perform well at the task, particularly for shorter sentences when high dimensional embeddings are used. It was also demonstrated that a simple probabilistic language model can be used to order the bag of words output to regenerate the original sentences.

Resynthesis degraded as sentence length increased, but remained strong with higher dimensional models up to reasonable length. It also decreased as the vocabulary size increased, however even the smaller Brown Corpus vocabulary, containing roughly 40,000 words is beyond what is suggested as the necessary as the vocabulary size for most uses (Nation 2006).

From a theoretical basis the resolvability of the selection problem shows that adding up the word vectors does preserve the information on which words were

used; particularly for higher dimensional embeddings. This shows clearly that collisions do not occur (at least with frequency) such that two unrelated sentences do not end up with the same SOWE representation.

Future work in this area would be use a stochastic language model to suggest suitable orderings for the bags of words. While this would not guarantee correct ordering every-time, we speculate that it could be used to find reasonable approximations often. Thus allowing this bag of words generation method to be used for full sentence generation, opening up a much wider range of applications.

## References

- Bezanson, Jeff et al. (2014). “Julia: A Fresh Approach to Numerical Computing”. In: arXiv: 1411.1607 [cs.MS].
- Bowman, Samuel R et al. (2015). “Generating Sentences from a Continuous Space”. In: *arXiv preprint arXiv:1511.06349*.
- Dinu, Georgiana and Marco Baroni (2014). “How to make words with vectors: Phrase generation in distributional semantics”. In: *Proceedings of ACL*, pp. 624–633.
- Farhadi, Ali et al. (2010). “Every picture tells a story: Generating sentences from images”. In: *Computer Vision–ECCV 2010*. Springer, pp. 15–29.
- Francis, W Nelson and Henry Kucera (1979). “Brown corpus manual”. In: *Brown University*.
- Iyyer, Mohit, Jordan Boyd-Graber, and Hal Daumé III (2014). “Generating Sentences from Semantic Vector Space Representations”. In: *NIPS Workshop on Learning Semantics*.
- Kågebäck, Mikael et al. (2014). “Extractive summarization using continuous vector space models”. In: *Proceedings of the 2nd Workshop on Continuous Vector Space Models and their Compositionality (CVSC)@ EACL*, pp. 31–39.
- Karp, Richard M (1972). *Reducibility among combinatorial problems*. Springer.
- Le, Quoc and Tomas Mikolov (2014). “Distributed Representations of Sentences and Documents”. In: *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pp. 1188–1196.
- Mikolov, Tomas et al. (2013a). “Efficient estimation of word representations in vector space”. In: *arXiv preprint arXiv:1301.3781*.
- Mikolov, Tomas, Wen-tau Yih, and Geoffrey Zweig (2013b). “Linguistic Regularities in Continuous Space Word Representations.” In: *HLT-NAACL*, pp. 746–751.
- Nation, I (2006). “How large a vocabulary is needed for reading and listening?” In: *Canadian Modern Language Review* 63.1, pp. 59–82.
- Pennington, Jeffrey, Richard Socher, and Christopher D. Manning (2014). “GloVe: Global Vectors for Word Representation”. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP 2014)*, pp. 1532–1543.

- Ritter, Samuel et al. (2015). “Leveraging Preposition Ambiguity to Assess Compositional Distributional Models of Semantics”. In: *The Fourth Joint Conference on Lexical and Computational Semantics*.
- Socher, Richard (2014). “Recursive Deep Learning for Natural Language Processing and Computer Vision”. PhD thesis. Stanford University.
- Socher, Richard et al. (2014). “Grounded compositional semantics for finding and describing images with sentences”. In: *Transactions of the Association for Computational Linguistics* 2, pp. 207–218.
- White, Lyndon et al. (2015). “How Well Sentence Embeddings Capture Meaning”. In: *Proc. of the Australasian Document Computing Symposium*.
- Yin, Wenpeng and Hinrich Schtze (2015). “Learning Word Meta-Embeddings by Using Ensembles of Embedding Sets”. In: eprint: 1508.04257.
- Yogatama, Dani, Fei Liu, and Noah A Smith (2015). “Extractive Summarization by Maximizing Semantic Volume”. In: *Conference on Empirical Methods in Natural Language Processing*.
- Zhu, Yukun et al. (2015). “Aligning Books and Movies: Towards Story-like Visual Explanations by Watching Movies and Reading Books”. In: *arXiv preprint arXiv:1506.06724*.