

Feasibility of BERT Embeddings For Domain-Specific Knowledge Mining

Anonymous ACL submission

Abstract

001 Extracting information from large corpora of
002 unstructured text using computational meth-
003 ods presents a challenge. Tshitoyan et al.
004 (2019) demonstrated that unsupervised math-
005 ematical word-embeddings produced by a
006 static language model could be utilized to
007 uncover ‘latent knowledge’ within a materi-
008 als science corpus. The rise of contextual-
009 ized and massively pre-trained language mod-
010 els like BERT have seen static models becom-
011 ing surpassed for most NLP tasks. Neverthe-
012 less, due to innate architectural and use dif-
013 ferences, BERT requires adaptation for knowl-
014 edge mining. This study tests the suitability
015 of BERT-derived word embeddings for knowl-
016 edge mining purposes. It utilizes a variation
017 of the approach described by Bommasani et al.
018 (2020) for creating static-equivalent vectors
019 from multiple contextualized word represen-
020 tations. It is conducted using a biomedical
021 corpus, a biomedical BERT variation and val-
022 idated using domain-specific intrinsic bench-
023 marking tools. Novel, layer-wise BERT per-
024 formance characteristics are demonstrated. A
025 key finding is that layer-wise intrinsic per-
026 formance differs for nouns and verbs. Per-
027 formance also varies according to whether a
028 word of interest belongs to BERT’s native vo-
029 cabulary or is built from sub-word represen-
030 tations: BERT-native representations perform
031 best when extracted from earlier layers, while
032 representations requiring multiple tokens per-
033 form best when extracted from the middle-to-
034 latter model layers.

035 1 Introduction

036 A vast amount of biomedical knowledge exists as
037 unstructured text within journals, books and ab-
038 stracts among other formats. This knowledge exists
039 as relationships and connections between described
040 concepts, objects and events. Information extrac-
041 tion from such corpora using supervised methods
042 requires large, manually-labelled datasets. Conse-
043 quently, these methods do not readily scale.

044 Recently, Tshitoyan et al. (2019) demonstrated
045 that known and novel relationships between en-
046 tities described within a materials science cor-
047 pus could be discovered using unsupervised, high-
048 dimensional word embeddings (Bengio et al., 2003;
049 Collobert and Weston, 2008; Collobert et al., 2011).
050 Here, the authors trained a skip-gram variant of
051 the Word2Vec neural language model (Mikolov
052 et al., 2013) on a corpus of 3.3 million materials
053 science abstracts to produce 200-dimensional em-
054 beddings for each word in the corpus vocabulary.
055 Remarkably, when the embeddings representing
056 material names (e.g. ‘Bi₂Te₃’) were ranked by
057 their cosine similarity to the representation of ‘ther-
058 moelectric,’ several novel thermoelectric conduc-
059 tors were identified. Despite the material name
060 never having appeared alongside, or within the
061 same document as ‘thermoelectric,’ the direct rela-
062 tionship between the novel material’s word repre-
063 sentation and ‘thermoelectric’ was permitted due to
064 indirect relationships between the material’s name
065 and other words/phrases such as ‘chalcogenide’
066 (chalcogenides are good thermoelectrics) and ‘band
067 gap’ (which determines thermoelectric properties)
068 within the vector space (Tshitoyan et al., 2019).
069 Venkatakrisnan et al. (2020) subsequently applied
070 the same Word2Vec skip-gram technique to an un-
071 structured text corpus of over 100 million biomed-
072 ical documents, discovering novel tissue-reservoirs
073 of the ACE2 receptor used by SARS-CoV-2 to in-
074 fect a host organism.

075 Both Tshitoyan et al. (2019) and Venkatakris-
076 nan et al. (2020) postulated that context-aware em-
077 beddings, such as those from the bidirectional en-
078 coder representation from transformers (BERT)
079 model (Devlin et al., 2018) could outperform
080 static models at these tasks. Aside from funda-
081 mentally different architecture, BERT produces
082 ‘just-in-time’ contextualized embeddings from pre-
083 tokenized sequences fed into the model individu-
084 ally. Moreover, unlike static models like Word2Vec

085 and GloVe (Pennington et al., 2014) which build
086 corpus-specific vocabularies, BERT possesses an
087 innate vocabulary of approximately 30,000 words
088 and handles extra-vocabulary words by decomposi-
089 tion into constituent sub-words. As such, a method
090 of leveraging BERT’s unique architecture and train-
091 ing on massive text corpora to ultimately yield
092 word representations capable of use in knowledge
093 mining is lacking. Bommasani et al. (2020) de-
094 scribed a method for reducing contextualized word
095 representations to static-equivalents by aggregating
096 contextualized word representations from BERT
097 over a number of contexts: These aggregated rep-
098 resentations outperformed static ones in general
099 domain intrinsic benchmarking tasks.

100 Much like static word representations, BERT-
101 derived equivalents can subsequently be adapted
102 for knowledge discovery by ranking geometric sim-
103 ilarity between represented concepts, objects or
104 processes. Nevertheless, as ‘latent knowledge’ re-
105 quires physical validation, the quality of novel lan-
106 guage model suggestions cannot easily be assessed.
107 Domain-specific intrinsic benchmarks which assess
108 semantic similarity and relatedness between word
109 representations using geometric measures (Chiu
110 et al., 2018) may be utilized as an appropriate sur-
111rogate: Higher-fidelity mathematical representations
112 of described reality are expected to approximate hu-
113 man user similarity ratings between concepts and
114 objects. This study subsequently tests the hypothe-
115 sis of both Tshitoyan et al. (2019) and Venkatakrish-
116 nan et al. (2020) that contextual language models
117 yield word representations for knowledge mining
118 that are superior to those produced by static model
119 in a biomedical domain and therefore suitable for
120 knowledge mining. Using a corpus of 500,000 ab-
121 stracts and full-text articles (Wang et al., 2020),
122 embeddings produced by a series of static models
123 are tested against aggregated contextual representa-
124 tions sampled from the corpus and passed through
125 a biomedically-trained BERT variant, and assessed
126 using domain-specific intrinsic benchmarks.

127 2 Methods

128 2.1 Dataset and Text Preprocessing

129 In response to the COVID-19 pandemic, the Coro-
130 navirus Open Research Dataset (CORD-19) was
131 released by governmental and academic institu-
132 tions. It consists of over 500,000 scholarly articles
133 (with over 200,000 full text articles and preprints)
134 and abstracts pertaining to COVID-19 (Wang et al.,

2020)¹. Corpus metadata was removed and articles
aggregated into a single file. All numbers were
replaced with a special token and selective lower-
casing was performed to preserve abbreviations.
For the Word2Vec and GloVe models, common
terms and punctuation were removed.

The BERT approach was informed by results
of an initial pilot study (see Appendix for prelimi-
nary data). Two approaches were adopted, in-
volving extracting n long or short contextual sen-
tence samples from the corpus: Long sequences
were created by splitting on periods into constituent
sentences. Short context sequences were created
by further splitting on commas into constituent
phrases. When selecting examples containing each
word, the maximum sequence length for both long
and short sequences was limited to 512, the maxi-
mum sequence length allowed by BERT, assuming
each word in the sequence is represented by a sin-
gle token.

155 2.2 BERT Approach

156 BioBERT is a variation of the original BERT model
157 which has been further pre-trained on PubMed ab-
158 stracts and PubMed Central full-text articles. It
159 outperforms general models at various biomedical
160 NLP tasks (Lee et al., 2020). The open source
161 HuggingFace (Wolf et al., 2020)² implementa-
162 tion of BioBERT v1.1 was utilized without any further
163 pre-training or fine-tuning based upon results of
164 the preliminary study (B). Depending on approach,
165 long and short sequences containing words of in-
166 terest (from the benchmarking vocabulary) were
167 selected. In order to ensure consistency between
168 static and contextual test vocabularies, as BERT is
169 able to use subword pooling for words outside its
170 native vocabulary, only words that were in both the
171 benchmarking vocabulary (see 2.4) and the CORD-
172 19 vocabulary were selected for sampling.

173 For both long and short sequence approaches,
174 $n = 500, 1000$ or 5000 samples were extracted
175 from the pre-processed corpus for tokenization.
176 Sequences were selected only if they contained
177 a single instance of the word of interest and were
178 discarded if their pre- or post-tokenized length ex-
179 ceeded 512. Here, for each word w in context c ,
180 BERT’s tokenizer will either yield a single token
181 or decompose w into k sub-word tokens, where

¹<https://www.kaggle.com/allen-institute-for-ai/CORD-19-research-challenge>

²<https://huggingface.co/>

182 $\{\mathbf{w}_c^1, \dots, \mathbf{w}_c^k\} \mapsto \mathbf{w}_c$. Tokenized sequences were
 183 then fed into the model and the sequence representa-
 184 tions were extracted from all 13 model layers. For
 185 words represented by a single 1x768 representation,
 186 this was extracted without further operations. For
 187 decomposed words, the arithmetic mean of all \mathbf{w}_c^k
 188 was taken to yield a single 1x768 representation
 189 from k sub-word representations, per context:

$$190 \quad \mathbf{w}_c = \text{mean}(\mathbf{w}_c^1, \dots, \mathbf{w}_c^k)$$

191 The arithmetic mean of the n contextual exam-
 192 ples of each word w , $\mathbf{w}_{c1}, \dots, \mathbf{w}_{cn}$ was then taken. If
 193 n examples meeting the inclusion criteria were not
 194 available, then the maximum number were taken:

$$195 \quad \mathbf{w} = \begin{cases} \text{mean}(\mathbf{w}_{c1}, \dots, \mathbf{w}_{cn}) & n = 500, 1000, 5000 \\ \text{mean}(\mathbf{w}_{c1}, \dots, \mathbf{w}_{c_{\max(n)}}) & n < 500, 1000, 5000 \end{cases}$$

196 Decision to take arithmetic mean of both sub-
 197 word representations and take either n was based
 198 on the results of [Bommasani et al. \(2020\)](#), where
 199 they found this operation outperformed other possi-
 200 ble operations (e.g. max., min., last) for both
 201 sub-word pooling and contextual aggregation (see
 202 also ([Ács et al., 2021](#))). The present approach also
 203 differed from [Bommasani et al. \(2020\)](#) who instead
 204 took the representation produced by feeding the
 205 word in isolation into the model³.

206 2.3 Static Models

207 The aggregated embeddings obtained from 2.2
 208 were compared against several static baseline mod-
 209 els including 200 and 300-dimensional Word2Vec
 210 skip-gram models, and a 300-dimensional GloVe
 211 model all trained from scratch on only CORD-
 212 19, using default hyperparameters. Addition-
 213 ally, pre-trained 200-dimensional embeddings from
 214 BioWordVec ([Zhang et al., 2019](#))⁴ were also
 215 obtained and used for benchmarking. Briefly,
 216 BioWordVec is an open set of static biomedical
 217 word vectors trained on a corpus of over 27 million
 218 articles, that additionally combine sub-word infor-
 219 mation from unlabelled biomedical text together
 220 with a biomedical controlled vocabulary.

³A single word (rather than a sequence) is an ‘unnatural’
 input for BERT, yielding a poorly-performing ‘decontextual-
 ized’ word representation (see ([Bommasani et al., 2020](#)) for
 more detail).

⁴[https://github.com/ncbi-nlp/
 BioWordVec](https://github.com/ncbi-nlp/BioWordVec)

221 2.4 Benchmarking

222 Bio-SimVerb and Bio-SimLex ([Chiu et al., 2018](#))
 223 are benchmarking resources for the biomedical do-
 224 main that offer 988 and 1000 test verb and noun
 225 pairs, respectively. These word-pairs have been
 226 extracted from 14 open biomedical ontologies and
 227 over 14,000 biomedical journals covering over 120
 228 areas of biomedicine. Additionally, some of the test
 229 word pairs are from the general domain. These re-
 230 sources address shortcomings of previous biomed-
 231 ical benchmarks such as MayoSRS ([Pakhomov
 232 et al., 2011](#)) and UMNSRS ([Pakhomov et al., 2010](#))
 233 which only test nouns, and fail to distinguish be-
 234 tween semantic relatedness and similarity ([Chiu
 235 et al., 2018](#)). The CORD-19 vocabulary covered
 236 97% of BioSimVerb and 94.43% of BioSimLex
 237 test pairs, respectively.

238 3 Results

239 3.1 Verb Benchmarks

240 The left sub-plot of Figure 3 and left column of
 241 Table 1 demonstrates the layer-wise performance
 242 of $n = 500, 1000$ and 5000 aggregated contextual-
 243 ized verb representations across all BERT layers.
 244 Performance is generally preserved regardless of
 245 sequence lengths/number of aggregated contexts.
 246 Embeddings extracted and distilled from the 7th
 247 and 8th layers performed best for all combinations.
 248 Short contexts marginally outperform longer con-
 249 texts at most layers. The best performing represen-
 250 tations for all combinations were extracted from
 251 layer 8 and distilled from 1000 contexts, though
 252 these representations did not substantially outper-
 253 form those distilled from other n . In general, repre-
 254 sentations extracted from the latter 6 layers (with
 255 the exception of layer 11) outperform the best-
 256 performing static embeddings at verb benchmark-
 257 ing.

258 3.2 Noun Benchmarks

259 The right sub-plot of 3 and right column of Ta-
 260 ble 1 demonstrates layer-wise performance of n
 261 = 500, 1000 and 5000 aggregated contextualized
 262 noun representations extracted from all BERT lay-
 263 ers. Unlike Bio-SimVerb, static models (with the
 264 exception of GloVe) outperformed aggregated noun
 265 representations from all layers and for all n . The
 266 plotted line demonstrates different morphology
 267 compared to verb benchmarks: Here, representa-
 268 tions distilled from the first 8 BERT layers out-
 269 performed those from the latter layers, increasing

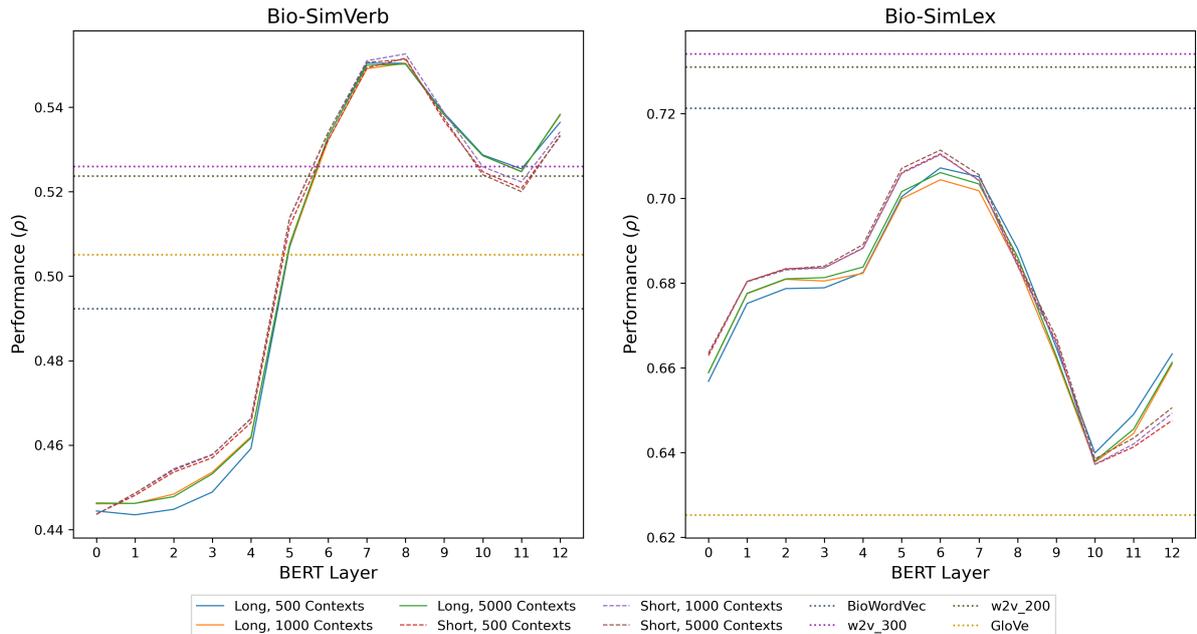


Figure 1: Layer-wise performance of BERT embeddings (0 corresponds to input layer) at both Bio-SimVerb and Bio-SimLex benchmarks. Horizontal dashed lines correspond to performance of static models.

270 until and then peaking at layer 6 before declining
 271 thereafter. Unlike verbal performance, however,
 272 embeddings extracted and distilled from shorter
 273 sequences demonstrated discernible improvement
 274 relative to those from long sequences. The best
 275 performance achieved was from $n = 1000$ embed-
 276 dings, however overall, n made little difference.

Method	Bio-SimVerb	Bio-SimLex
Long 500	0.5504 (7=8)	0.7072 (6)
Long 1000	0.5504 (8)	0.7044 (6)
Long 5000	0.5502 (8)	0.7061 (6)
Short 500	0.5516 (8)	0.7105 (6)
Short 1000	0.5526 (8)	0.7103 (6)
Short 5000	0.5513 (8)	0.7114 (6)
w2v 300	0.5260	0.7341
w2v 200	0.5237	0.7310
GloVe 300	0.5051	0.6253
BWV 200	0.4923	0.7213

Table 1: Top performing (Spearman’s ρ) distilled BERT embeddings and static embeddings. ‘Long/Short n ’ indicates embeddings distilled from n example sequences of each word. Number in brackets indicates layer. w2v200/300 = Word2Vec 200/300 dimensional embeddings. BWV = BioWordVec 200 dimensional embeddings. Bold entries indicate best overall performance.

3.3 Effect of Sub-Word Pooling

277 In an attempt to explain observed performance and
 278 explore the effects of sub-word pooling, test word-
 279 pairs from Bio-SimVerb and Bio-SimLex were sep-
 280 arated into two groups using a criteria of whether
 281 both words in a respective test pair existed in
 282 BioBERT’s native vocabulary or not. This yielded
 283 test word pairs where both were BERT-native and
 284 had a single representation, or where at least one of
 285 the words in the pair required a sub-word pooling
 286 operation before the n contextual representations
 287 could be averaged. Representations from all BERT
 288 layers using short contexts and the top perform-
 289 ing static model representations (300-dimensional
 290 Word2Vec) were then subjected to Spearman’s rank
 291 testing as per the Bio-Simverb methodology (Chiu
 292 et al., 2018), albeit using the modified, separated
 293 test pair rankings.

294 For both verbs and nouns, general and layer-wise
 295 native word performance differed to words requir-
 296 ing sub-word pooling prior to context aggregation.
 297 Moreover, it is apparent that when test-vocabulary
 298 is stratified in this regard, n has little to no bear-
 299 ing upon overall performance of the BERT embed-
 300 dings. Verbs native to BERT’s vocabulary gener-
 301 ally outperform those from the top-performing
 302 static model at all layers (see left side of Figure 2
 303 and Table 2). Performance declines steadily from
 304 layer 0-12. Performance of word representation for
 305

verbs requiring sub-word pooling is overall lower than single-token representations, demonstrating performance increase from layers 0-8. It is only at layers 7-8 where performance slightly exceeds performance of static representations.

Only BERT-native noun representations extracted from the first 4 layers demonstrated superior performance to the corresponding static embeddings. Similar to the BERT-native verbs, performance decreased from layer 0-12, reaching a trough at layer 10, before slightly increasing thereafter. Moreover, sub-word pooled noun representation performance increased substantially from layers 0-6 before declining thereafter. These representations never outperformed the static embeddings (see Figure 2 and Table 2).

Method	Bio-SimVerb	Bio-SimLex
Short 500 (S)	0.6691 (1)	0.7255 (1)
Short 500 (M)	0.4603 (8)	0.7417 (6)
Short 1000 (S)	0.6685 (1)	0.7255 (1)
Short 1000 (M)	0.4629 (8)	0.7420 (6)
Short 5000 (S)	0.6688 (1)	0.7256 (1)
Short 5000 (M)	0.4621 (8)	0.7418 (6)
Static (S)	0.5255	0.6959
Static (M)	0.4545	0.7628

Table 2: Performance of BERT embeddings aggregated from short contextual examples and with $n = 500, 1000, 5000$. **S** or **M** in brackets indicate whether representations were for words native to BERT i.e. using a single token to represent or those requiring sub-word pooling, respectively. Static representations were from a 300-dimensional Word2Vec model. Bold entries indicate best overall performance.

4 Discussion

This study demonstrates the feasibility of using BERT-derived word representations for knowledge mining purposes, however their benefit over static representations as used by Tshitoyan et al. (2019) and Venkatakrisnan et al. (2020) is less clear. Overall, this study demonstrates that relatively few, short-sequence contextual word-examples extracted from a corpus can be aggregated and utilized to yield embeddings that can outperform in the case of verbs, or approximate (in the case of nouns) those from the best performing static models trained on entire corpora. Practically this sampling-based approach may offer time and cost savings over training entire static models from scratch, when dealing with large corpora. More-

over, due to the superior performance of the BERT representations for verbs, aggregated contextualized embeddings may even be preferable when mining verb-rich text (e.g. clinical notes).

BERT’s pre-training on massive text corpora may be responsible for performance characteristics observed: Evidence is provided by the differing performance of representations innate to BERT’s vocabulary compared to representations built from multiple sub-words. It appears that n is less important than whether the word belongs to BERT’s innate vocabulary or not: For both nouns and verbs, if the word of interest is BERT-native then it is preferable to utilize representations from the earlier layers and these are superior to static ones, while if a word requires sub-word decomposition then layers 6-8 seem to be optimal. Nevertheless, more work is required to quantify the effect of multiple subwords on performance, as the split vocabulary in this study utilized a relatively imprecise criteria of $k > 1$ for test-pairs where at least one word was non-native to BERT. Moreover, though Bommasani et al. (2020) demonstrated that taking the arithmetic mean of k sub-words was the best performing method on their general-domain intrinsic benchmarking, a later study by Ács et al. (2021) showed that sub-word pooling approach mattered depending on desired downstream NLP tasks. Consequently, further exploration into both k and n parameters should be conducted.

Based on the observed performance, it might also benefit to expand BERT’s native vocabulary with domain-specific words prior to conducting pre-training. Though the pilot study showed that further-pretraining was detrimental to performance (see Appendix B), this method utilized BERT’s native vocabulary and was only tested using long sequences (which themselves underperform relative to short sequences) before being abandoned. Another consideration is that further pre-training steps are necessary to improve sub-word performance (Liu et al., 2019), which could be important for non-general domains. Moreover, as the benchmarking vocabularies incorporate both general-domain and biomedical-domain word pairs (Chiu et al., 2018), it may also be that the general domain test pairs are contributing disproportionately to performance boosts. Another area for exploration is the comparatively different layer-wise performance for BERT-native words versus extra-vocabulary words, with the former’s performance generally decreasing

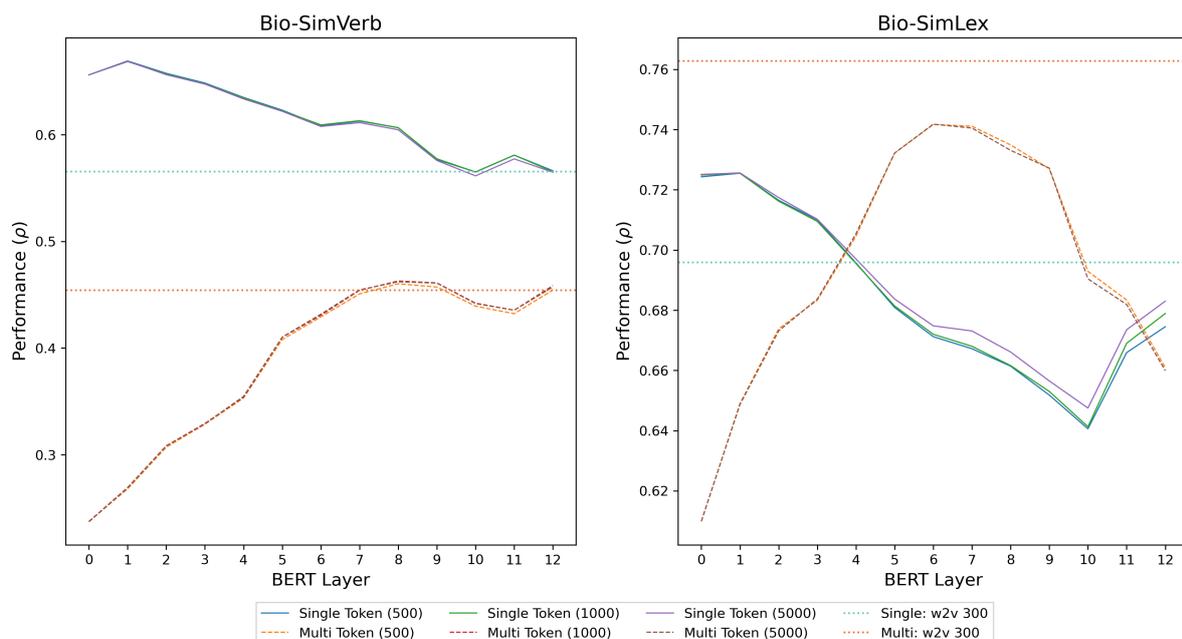


Figure 2: Layer-wise performance of BERT embeddings (0 corresponds to input layer) at both Bio-SimVerb and Bio-SimLex benchmarks. Horizontal dashed lines correspond to performance of static models.

and the latter’s increasing.

A working knowledge-mining framework utilizing BERT might consist of first extracting the vocabulary of the corpus upon which mining will be conducted and removing any irrelevant words (e.g. stop words). Then, n samples for each word in the vocabulary may be taken from the corpus and tokenized. As BERT’s attention is quadratic to the sequence (Devlin et al., 2018), and representations extracted from short sequences perform better, shorter sample sequences are desirable. Tokenized sequences can then be encoded, and representations extracted, with sub-word pooling performed if necessary. The n contextual examples of each word representation can then be averaged to yield a 1x768 dimensional representation for each word in the corpus vocabulary. It is this collection of vocabulary embeddings that can be subsequently used for mining as per Tshitoyan et al. (2019).

5 Conclusions

This study has successfully demonstrated feasibility of aggregated contextual word representations derived from BERT for biomedical knowledge mining tasks. It has also uncovered several technical and performance-related idiosyncrasies of BERT and BioBERT that require further investigation.

6 Acknowledgements

Thanks to **redacted** and **redacted** for their assistance with this study.

References

- Judit Ács, Ákos Kádár, and Andras Kornai. 2021. [Subword pooling makes a difference](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Association for Computational Linguistics, Online, pages 2284–2295. <https://doi.org/10.18653/v1/2021.eacl-main.194>.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. 2003. A neural probabilistic language model. *Journal of machine learning research* 3(Feb):1137–1155.
- Rishi Bommasani, Kelly Davis, and Claire Cardie. 2020. Interpreting pretrained contextualized representations via reductions to static embeddings. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. pages 4758–4781.
- Billy Chiu, Sampo Pyysalo, Ivan Vulić, and Anna Korhonen. 2018. Bio-simverb and bio-simlex: wide-coverage evaluation sets of word similarity in biomedicine. *BMC bioinformatics* 19(1):1–13.
- Ronan Collobert and Jason Weston. 2008. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*. pages 160–167.

445	Ronan Collobert, Jason Weston, Léon Bottou, Michael	Lucy Lu Wang, Kyle Lo, Yoganand Chandrasekhar,	501
446	Karlen, Koray Kavukcuoglu, and Pavel Kuksa.	Russell Reas, Jiangjiang Yang, Darrin Eide, Kathryn	502
447	2011. Natural language processing (almost) from	Funk, Rodney Kinney, Ziyang Liu, William Merrill,	503
448	scratch. <i>Journal of machine learning research</i>	et al. 2020. Cord-19: The covid-19 open research	504
449	12(ARTICLE):2493–2537.	dataset. <i>ArXiv</i> .	505
450	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	506
451	Kristina Toutanova. 2018. Bert: Pre-training of deep	Chaumond, Clement Delangue, Anthony Moi, Pier-	507
452	bidirectional transformers for language understand-	ric Cistac, Tim Rault, Remi Louf, Morgan Fun-	508
453	ing. <i>arXiv preprint arXiv:1810.04805</i> .	towicz, Joe Davison, Sam Shleifer, Patrick von	509
454	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim,	Platen, Clara Ma, Yacine Jernite, Julien Plu, Can-	510
455	Donghyeon Kim, Sunkyu Kim, Chan Ho So,	wen Xu, Teven Le Scao, Sylvain Gugger, Mariama	511
456	and Jaewoo Kang. 2020. Biobert: a pre-trained	Drame, Quentin Lhoest, and Alexander Rush. 2020.	512
457	biomedical language representation model for	Transformers: State-of-the-art natural language pro-	513
458	biomedical text mining. <i>Bioinformatics</i> 36(4):1234–	cessing . In <i>Proceedings of the 2020 Conference</i>	514
459	1240.	on Empirical Methods in Natural Language Pro-	515
460	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	cessing: System Demonstrations . Association for	516
461	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	Computational Linguistics, Online, pages 38–45.	517
462	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	https://doi.org/10.18653/v1/2020.emnlp-demos.6 .	518
463	Roberta: A robustly optimized bert pretraining ap-	Yijia Zhang, Qingyu Chen, Zhihao Yang, Hongfei Lin,	519
464	proach. <i>arXiv preprint arXiv:1907.11692</i> .	and Zhiyong Lu. 2019. Biowordvec, improving	520
465	Tomas Mikolov, Kai Chen, Greg Corrado, and Jef-	biomedical word embeddings with subword infor-	521
466	frey Dean. 2013. Efficient estimation of word	mation and mesh. <i>Scientific data</i> 6(1):1–9.	522
467	representations in vector space. <i>arXiv preprint</i>		
468	<i>arXiv:1301.3781</i> .		
469	Serguei Pakhomov, Bridget McInnes, Terrence Adam,		
470	Ying Liu, Ted Pedersen, and Genevieve B Melton.		
471	2010. Semantic similarity and relatedness between		
472	clinical terms: an experimental study. In <i>AMIA an-</i>		
473	<i>annual symposium proceedings</i> . American Medical In-		
474	formatics Association, volume 2010, page 572.		
475	Serguei VS Pakhomov, Ted Pedersen, Bridget McInnes,		
476	Genevieve B Melton, Alexander Ruggieri, and		
477	Christopher G Chute. 2011. Towards a framework		
478	for developing semantic relatedness reference stan-		
479	dards. <i>Journal of biomedical informatics</i> 44(2):251–		
480	265.		
481	Jeffrey Pennington, Richard Socher, and Christopher		
482	Manning. 2014. GloVe: Global vectors for word		
483	representation . In <i>Proceedings of the 2014 Confer-</i>		
484	<i>ence on Empirical Methods in Natural Language</i>		
485	<i>Processing (EMNLP)</i> . Association for Computa-		
486	tional Linguistics, Doha, Qatar, pages 1532–1543.		
487	https://doi.org/10.3115/v1/D14-1162 .		
488	Vahe Tshitoyan, John Dagdelen, Leigh Weston,		
489	Alexander Dunn, Ziqin Rong, Olga Kononova,		
490	Kristin A Persson, Gerbrand Ceder, and Anubhav		
491	Jain. 2019. Unsupervised word embeddings capture		
492	latent knowledge from materials science literature.		
493	<i>Nature</i> 571(7763):95–98.		
494	AJ Venkatakrishnan, Arjun Puranik, Akash Anand,		
495	David Zemmour, Xiang Yao, Xiaoying Wu, Ra-		
496	makrishna Chilaka, Dariusz K Murakowski, Kristo-		
497	pher Standish, Bharathwaj Raghunathan, et al. 2020.		
498	Knowledge synthesis of 100 million biomedical doc-		
499	uments augments the deep expression profiling of		
500	coronavirus receptors. <i>Elife</i> 9:e58040.		

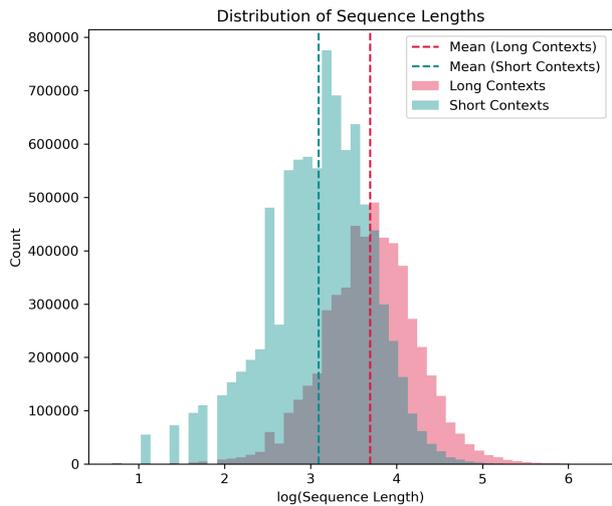


Figure 3: Distributions of log sentence lengths for long and short contextual sequences.

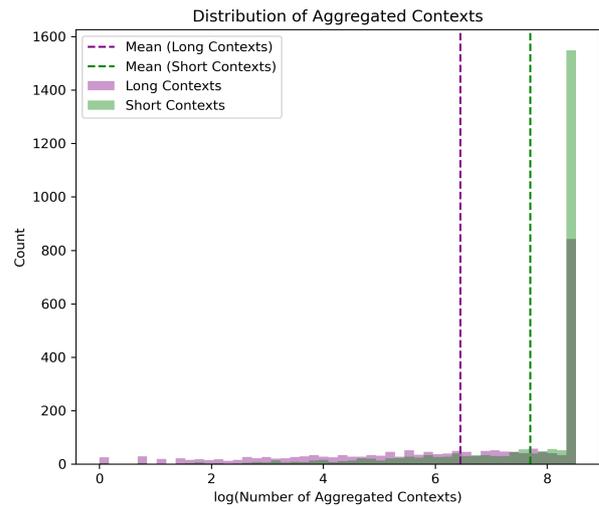


Figure 4: Distributions of log number of aggregated contexts for long and short sequence lengths. There are substantially more examples meeting $n = 5000$ for short sentences

A Corpus Sampling Characteristics

Figure 2 demonstrates the distribution of log sequence lengths for long and short sequences, respectively. Figure 3 demonstrates the distribution of log number of sequences for long and short sequences, respectively. As one of the sampling criteria was that sequences could only have a single instance of the word and had to be less than 512 words in length, together with the fact that the tokenized sequence could not exceed 512 in length, there are consequently fewer long sequence samples per word compared to short sequence examples. The mean long sequence length was 46.8 words ($\sigma = 29.5$) while the mean short sequence length was 26.3 words ($\sigma = 16.7$) (2). For long sequences, the mean number of contextual examples per word was 2330.2 ($\sigma = 2194.8$). For short sequences, the mean number of contextual examples per word was 3632.2 ($\sigma = 1948.7$) (Figure 3).

B Effect of Further Pre-Training on Word Representation Quality

An initial approach attempted was to further pre-train BioBERT using the entire CORD-19 corpus and compare performance of the contextualized word representations at Bio-SimVerb and Bio-SimLex intrinsic benchmarking tasks, with the base BioBERT model and the static baseline models. This approach used only long corpus sequences and the base BioBERT vocabulary (which itself is identical to BERT vocabulary). Pre-training was

achieved using the scripts supplied with the TensorFlow implementation of the model (<https://github.com/dmis-lab/bioBERT>) and involved first creating pre-training data using sentence examples from the corpus, before running further pre-training for 100,000 epochs. Default hyperparameters were used. For this pilot study, $n = 10, 50, 100, 500$ and 1000. The n selected examples were then all tokenized and passed through either the further-pretrained BioBERT model or the base Bio-BERT model. For either approach, representations corresponding to the word of interest were then extracted wholly (i.e. as a single 1x768 word representation, or k individual sub-word representations) and added to the list of n (explained further in 2.2). Benchmarking was performed as described in 2.4.

For the Bio-SimVerb benchmarks (Left side of Figure 4), there is a clear increase in performance by increasing n from 10 to 1000 contexts. Also apparent is that the representations extracted from the further pre-trained model underperform relative to those extracted from the base model for the same n . Biggest increases in performance are seen going from $n = 10$ to $n = 100$. Increasing n beyond this begins to demonstrate smaller performance boosts. Interestingly, best performing verb embeddings from the further pre-trained model were taken from layer 12 (see 3) while for the base model, performance peaked at embeddings extracted from

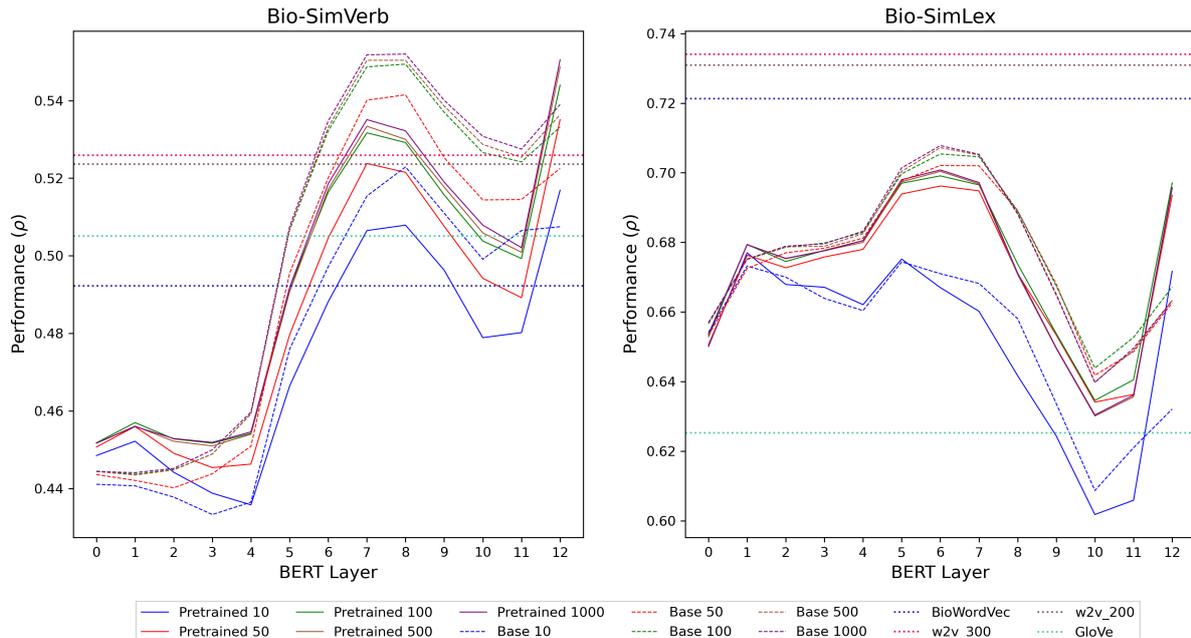


Figure 5: Layer-wise performance of BERT embeddings (0 corresponds to input layer) at both Bio-SimVerb and Bio-SimLex benchmarks. Pretrained n /Base n refer to either the further-pretrained model or the base model, respectively, followed by the n aggregated contexts. Horizontal dashed lines correspond to performance of static models.

layer 8. In some cases, embeddings taken from layer 12 of the further pre-trained model almost reached peak performance from embeddings taken from layer 8 of the base model.

For the Bio-SimLex benchmarks (Right side of Figure 4), though there was a general performance increase between representations extracted from the further-pretrained model and the base BioBERT model, it was less pronounced as it was for the verb benchmarks, with performance for the first 6 layers approximately equal before diverging thereafter. Moreover, a substantial boost is seen going from $n = 10$ to $n = 50$, becoming less pronounced as n increases. Again, performance for the representations extracted from a further pre-trained model demonstrate a trough following their maximum performance at layer 8, but increase substantially thereafter going from layer 11 to 12, though without reaching their layer 6 peak. This characteristic was not observed with the base model representations. Finally, representations from either the further-pretrained or base models did not outperform either Word2Vec 200 or 300 dimensional representations, or the BioWordVec representations.

Method	Bio-SimVerb	Bio-SimLex
Pre-Trained 10	0.5169 (12)	0.6770 (1)
Pre-Trained 50	0.5351 (12)	0.6991 (6)
Pre-Trained 500	0.5440 (12)	0.7004 (6)
Pre-Trained 1000	0.5487 (12)	0.7008 (6)
Base 10	0.5229 (8)	0.6744 (5)
Base 50	0.5415 (8)	0.7054 (6)
Base 500	0.5494 (8)	0.7072 (6)
Base 1000	0.5504 (8)	0.7078 (6)
w2v 300	0.5260	0.7341
w2v 200	0.5237	0.7310
GloVe 300	0.5051	0.6253
BWV 200	0.4923	0.7213

Table 3: Top performing (Spearman’s ρ) distilled BERT embeddings and static embeddings from pilot study. ‘Pre-Trained/Base n ’ indicates embeddings extracted from n examples taken from the distilled pre-trained or base model, respectively. Number in brackets indicates layer. w2v200/300 = Word2Vec 200/300 dimensional embeddings. BWV = BioWordVec 200 dimensional embeddings. Bold entries indicate best overall performance.