

# Towards Better Citation Intent Classification

Anonymous ACL submission

## Abstract

Accurate classification of citation intents in a scientific article provides deeper contextual understanding of and better quantifies the contributions of cited articles. This improves scientific literature platform capabilities such as search relevance, ranking and more. To our knowledge, we present the most comprehensive survey of Transformer-based language models performance on the citation intent classification task using SciCite dataset. Here, we make three recommendations. Firstly, we propose to report model performance as a distribution in contrast to a single averaged performance value. This arises from our observation that model performance is sensitive to the random seed choice resulting in wide performance variations from multiple finetuning runs. Secondly, this provides practical insights for model selection, showing the model’s best possible performance. Thus, we propose that practitioners perform multiple finetuning runs before selecting the best performing model. Thirdly, we propose a simple data augmentation to improve the distribution of model performance overall. Moving forward, we suggest exploring improvements to the finetuning and model selection process as promising future directions.

## 1 Introduction

Citations are a core part of scientific literature, providing an avenue to acknowledge the various contributions of various scientific articles. Citations are provided with specific intents, such as to provide background information, to present the use of methods or compare results from other works.

Citation intent classification is the task of identifying the intent of a specific citation. By quantifying the distribution of intents of various citations received by a particular paper, we can generate a better understanding of the nature of contributions provided by a paper beyond the citation count. For example, we can infer if a particular paper provides a useful method or result. This is useful in

many applications, such as identifying and ranking scientific articles according to the nature of their contributions.

In our paper, we will survey and evaluate the effectiveness of various state-of-the-art language models on the task of citation intent classification. In addition, we propose a data augmentation approach that improves the performance across all the surveyed models. Lastly, we discuss the implications of the results we obtained through our experiments and provide some possible directions for future work.

### 1.1 Dataset

There are various datasets of scientific literature that contain citation information. Some, such as S2ORC (Lo et al., 2020), are large but do not provide annotations for the citation intents, while smaller datasets such as ACL-ARC (Bird et al., 2008) and SciCite (Cohan et al., 2019) do provide annotations for citation intents. SciCite is an order of magnitude larger than ACL-ARC, containing 6 times the number of citations gathered from about 35 times more scientific articles from the Computer Science and Medical domains.

In this work, we focus on the citation intent classification with **SciCite** as the benchmark dataset. More specifically, the task is to correctly classify citation intents into one of three classes: Background, Method, and Result Comparison.

### 1.2 Models

In recent years, Transformer-based large language models have been the dominant approach for achieving state of the art results on NLP tasks. One reason for the effectiveness of these models is the ability to learn good language representations by pre-training on a large corpus before undergoing fine-tuning on a task-specific dataset.

In our paper, we will survey the effectiveness of a total of nine Transformer-based language models

082 on the SciCite benchmark dataset. The nine models  
083 are as follows:

- 084 1. **BERT** (Devlin et al., 2018) (Bidirectional En-  
085 coder Representations from Transformers) is  
086 a Transformer model that uses a bidirectional  
087 self-attention mechanism and pretrained on  
088 a large text corpus. It achieves high perfor-  
089 mance on many NLP tasks via transfer learn-  
090 ing.
- 091 2. **RoBERTa** (Liu et al., 2019) is based on BERT,  
092 with an improved training regime using 10  
093 times more data.
- 094 3. **DSP-RoBERTa** (Gururangan et al., 2020) is  
095 part of a family of RoBERTa models with  
096 further pretraining to adapt them for various  
097 target domains.
- 098 4. **Biomed-RoBERTa** (Lewis et al., 2020) is a  
099 RoBERTa model that is pre-trained on the  
100 full texts of about 2.7 million scientific papers  
101 from the Semantic Scholar corpus, improving  
102 the performance of the model on tasks in the  
103 biomedical domain.
- 104 5. **SciBERT** (Beltagy et al., 2019) is a BERT  
105 model pretrained on scientific text from  
106 1.14M papers in order to handle language pro-  
107 cessing tasks in scientific field.
- 108 6. **PubMedBERT** (Gu et al., 2020) is pretrained  
109 using abstracts from PubMed and full text arti-  
110 cles from PubMedCentral, achieving state-of-  
111 the-art results on various tasks in the biomed-  
112 ical domain.
- 113 7. **XLNet** (Yang et al., 2019) is a Transformer-  
114 based auto-regressive language model that re-  
115 tains the ability to learn bidirectional contexts.
- 116 8. **DeBERTa** (He et al., 2020) is a Transformer  
117 model with distangled attention mechanism  
118 and enhanced mask decoder replacing the fi-  
119 nal Softmax layer.
- 120 9. **ALBERT** (A Lite BERT) (Lan et al., 2020) is  
121 a variant of BERT that utilizes parameter shar-  
122 ing and embedding factorization to reduce the  
123 number of parameters compared to an equiva-  
124 lently sized BERT model, although at slightly  
125 lower performance.

We used the implementations and pretrained  
weights of these models through the Hugging Face  
Transformers library (HuggingFace, 2019). The  
precise implementations and checkpoints used are  
listed in the Appendix. Of these nine models, the  
some are of particular interest due to possible rele-  
vance in the corpus used for pretraining, in particu-  
lar SciBERT, DSP-RoBERTa, and PubMedBERT,  
which is pretrained on scientific literature. We also  
test some variations of a subset of the models to  
investigate the impact of model size and vocabu-  
lary casing differences between some variants of  
the models.

## 2 Related Work

Previous work on the Sci-Cite dataset was pre-  
sented in (Cohan et al., 2019), where a BiLSTM  
model with Attention and ELMo embeddings,  
along with structural scaffolds, was used to achieve  
a macro F1-score of 84.0. Following up, (Beltagy  
et al., 2019) presented fine-tuned BERT and SciB-  
ERT models that achieved with F1-score of 84.85  
and 85.49 respectively.

## 3 Methods

### 3.1 Training Settings

To be consistent with the SciBERT paper (Belt-  
agy et al., 2019) which first performed evaluation  
of BERT-base and SciBERT on SciCite dataset,  
we chose the batch size of 32 and fine-tuned all  
nine models separately with a learning rate of 1e-5.  
We run each training multiple times and report the  
mean performance.

### 3.2 Data Augmentation

We postulate that there is often a correlation be-  
tween the intent of a citation and the section in  
which the citation is located. For example, the  
section "Related Work", a citation is likely to be  
used for describing the background and positioning,  
while citation in the "Results" section is likely cited  
with the intent of result comparison. By adding  
in the parent section of the citation as additional  
context, we hypothesize that we can improve the  
context available to the model to be used for classi-  
fication. We term this data augmentation by adding  
in a **section hint**. The section hint takes the form of  
a sentence, which contains only section title, being  
added to the start of the original text. Some prepro-  
cessing is performed on the raw text of the section

173 headers available in the SciCite dataset to normal- 212  
174 ize various formatting quirks of the raw data, such 213  
175 as inconsistent capitalization schemes, inclusion of 214  
176 section numbers etc. 215

177 Example of augmentation (marked in red): 216

178 **Discussion.** More examples of contra- 217  
179 dictory results have been observed in 218  
180 bovines; some reports (Zakhartchenko 219  
181 et al., 2001; Bhuiyan et al., 2004) indi- 220  
182 cated a significant decrease in blastocyst. 221

183 We test this simple form of augmentation as an 222  
184 added experiment during our survey to see if it 223  
185 impacts the various models differently. 224

## 186 4 Results 225

187 Our experiment results are reported in **Table 1** be- 226  
188 low with macro F1-score on both original data and 227  
189 augmented data. 228

Model	RawData	Aug.Data	$\Delta$
BERT-base	84.80	84.82	+0.02
RoBERTa-base	84.01	84.59	+0.58
DSP-RoBERTa	84.90	85.61	+0.71
Biomed-RoBERTa	86.32	86.63	+0.31
SciBERT	<b>86.74</b>	<b>87.08</b>	+0.34
PubMedBERT	85.55	86.14	+0.59
XLnet-base	84.59	85.49	+0.90
DeBERTa-base	84.75	85.90	<b>+1.15</b>
ALBERT-base	84.03	84.45	+0.42

190  
191 We find that the 3 best performing models (SciB- 230  
192 ERT, Biomed-RoBERTa and PubMedBERT) are 231  
193 the ones pretrained on biomedical text, which is the 232  
194 text most relevant to the SciCite dataset. Hence, we 233  
195 have a clear indication on the benefit of pretraining 234  
196 on a relevant corpus. 235

197 DeBERTa shows the greatest improvement as a 236  
198 result of the augmented data, while also being, in 237  
199 theory, the most powerful model among those sur- 238  
200 veyed, however lacking in relevant pretraining. We 239  
201 hypothesize that DeBERTa pretrained on biomed- 240  
202 ical text would possibly be the best performing 241  
203 model. 242

204 SciBERT has the best F1-score for both the raw 243  
205 and augmented data. In fact, our result shows an 244  
206 improvement over the previous reported work, even 245  
207 without the augmented data. This is related to our 246  
208 next observation: In our experiments, we often ob- 247  
209 serve a large variance in model performance, across 248  
210 all models due to different random seeds used in 249  
211 training. We postulate that this could be due to 250

212 large variations in model performance due to dif- 213  
214 ferent random seeds, and perform further analysis 214  
in Section 5. 215

216 We can also see that the data augmentation is 215  
effective across all the models surveyed. 216

## 217 5 Analysis and Discussion 217

218 In our experiments, we observe large variations in 218  
219 model performance due to different random seeds. 219  
220 Thus, we recorded the performance of our model 220  
221 across 9 training runs per model, and plotted them 221  
222 in a box-and-whisker plot (**Figure 1**) to visualize 222  
223 the variances in performance. 223

224 We observe that the original reported F1 score 224  
225 for SciBERT (85.49) falls well within the statisti- 225  
226 cal inter-quartile range for our 9 SciBERT training 226  
227 runs, although our recorded mean and median val- 227  
228 ues are both higher. This means that the original 228  
229 reported score is not unexpected given a different 229  
230 random seed, and in our training runs with more 230  
231 random seeds, we have trained models which on 231  
232 average perform better than what was previously 232  
233 recorded. 233

234 In addition, the absolute best performing model 234  
235 (F1-score of 88.5) recorded far exceeds the best 235  
236 average score as reported in **Table 1**. In fact, the 236  
237 best performing model is Biomed-RoBERTa with a 237  
238 particular random seed, and not SciBERT, although 238  
239 it has the best score on average. 239

240 From the visualization in **Figure 1**, we can also 240  
241 see that the data augmentation is effective in at least 241  
242 one of the two following ways for every model 242  
243 we tested, resulting in an improved distribution of 243  
244 performance across multiple training runs: 244

- 245 1. To improve the absolute best performance of 245  
246 the model 246
- 247 2. To reduce the variance in performance ob- 247  
248 served during the model training 248

249 We also performed additional analysis on several 249  
250 variations of the model to determine if they have 250  
251 an impact on the task performance. 251

- 252 1. **Cased vs Uncased Models** We studied the 252  
253 difference between cased and uncased ver- 253  
254 sions of BERT-base, BERT-large and SciB- 254  
255 ERT, finding no clear correlation between the 255  
256 type of casing used and the performance of 256  
257 the model. The performance of the cased and 257  
258 uncased models are visualized in **Figure 2** in 258  
259 Appendix B. 259

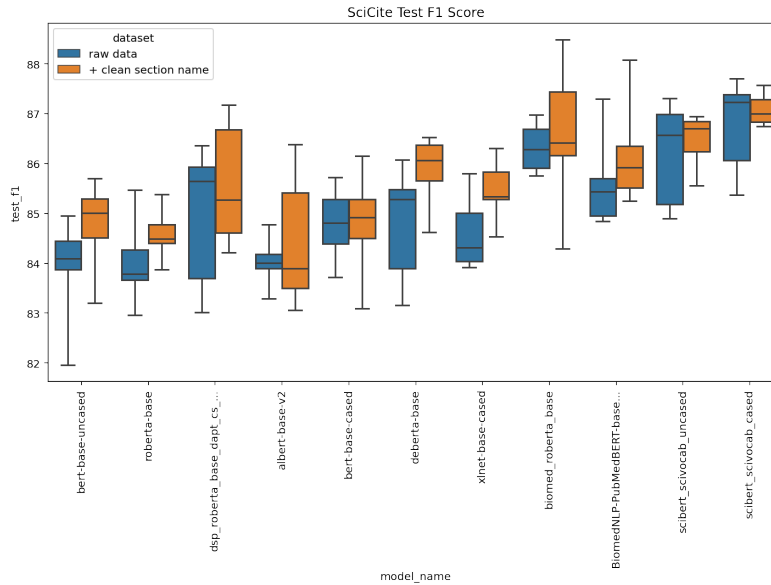


Figure 1: Spread of results.

260 **2. Base vs Large Models** We studied the differ- 289  
 261 ence between base ( 100M) and large ( 300M) 290  
 262 variants of BERT, RoBERTa, DeBERTa, XL- 291  
 263 Net and the corresponding variants of AL- 292  
 264 BERT. We observe a clear relationship be- 293  
 265 tween a larger model and improved perfor- 294  
 266 mance on the SciCite task. This also indicates 295  
 267 to us that a larger model pretrained on biomed- 296  
 268 ical tasks would likely perform better than the 297  
 269 currently available models. The performance 298  
 270 of the cased and uncased models are visual- 299  
 271 ized in **Figure 3** in Appendix B.

272 **6 Recommendations for Practitioners**

273 In this paper, we have arrived at results that al- 300  
 274 low us to propose the following recommenda- 301  
 275 tions for practitioners to achieve better practical 302  
 276 performance on the citation intent classification 303  
 task: 304

- 277 1. Train models multiple times with different ran- 305  
 278 dom seeds in order to find the best performing 306  
 279 model 307
- 280 2. Utilize data augmentation, such as the simple 308  
 281 strategy we demonstrated, as it shows meas- 309  
 282 urable improvements across all the models 310  
 283 surveyed 311
- 284 3. Report the performance score of models as a 312  
 285 distribution, rather than a singular score, in 313  
 286 order to provide a better overview of the av- 314  
 287 erage and best possible score that can result 315  
 288 from a particular model

**7 Conclusion and Future Work**

In future work, we hope to improve the practical 290  
 performance on the citation intent classification 291  
 task. We have outlined a few directions below: 292

- 293 1. Explore more methods and different permuta- 294  
 295 tions of data augmentation techniques to add 296  
 more context into the model input. 297
- 298 2. Introduce methods to perform training with 299  
 random seeds found using search algorithms 300  
 that allow us to train better performing mod- 301  
 els. 302
- 303 3. Create a better pretrained model. Based on 304  
 our observations, a DeBERTa-large model pre- 305  
 trained on biomedical texts would be a prime 306  
 candidate for the best performing model when 307  
 applied to the citation intent classification 308  
 task. 309
- 310 4. Perform our experiments on other citation in- 311  
 tent classification datasets, such as ACL-ARC, 312  
 to study the transferability of our proposed 313  
 methods. 314

315 In this short paper, we have presented a pre-  
 liminary data augmentation technique that demon-  
 strates adding more context improves the perfor-  
 mance of a model on the citation intent classifica-  
 tion task. Our method improve the performance  
 across all the models that we surveyed.

316  
317  
318  
319  
  
320  
321  
322  
323  
324  
325  
326  
327  
  
328  
329  
330  
331  
  
332  
333  
334  
335  
  
336  
337  
338  
339  
340  
  
341  
342  
343  
344  
345  
346  
347  
348  
  
349  
350  
351  
352  
  
353  
354  
  
355  
356  
357  
358  
  
359  
360  
361  
362  
363  
364  
365  
  
366  
367  
368  
369  
370

## References

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. Scibert: A pretrained language model for scientific text. *arXiv preprint arXiv:1903.10676*.

Steven Bird, Robert Dale, Bonnie J. Dorr, Bryan Gibson, Mark T. Joseph, Min-Yen Kan, Dongwon Lee, Brett Powley, Dragomir R. Radev, and Yee Fan Tan. 2008. The ACL Anthology Reference Corpus: A Reference Dataset for Bibliographic Research in Computational Linguistics. In *Proc. of the 6th International Conference on Language Resources and Evaluation Conference (LREC'08)*, pages 1755–1759.

Arman Cohan, Waleed Ammar, Madeleine Van Zuylen, and Field Cady. 2019. Structural scaffolds for citation intent classification in scientific publications. *arXiv preprint arXiv:1904.01608*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. Domain-specific language model pretraining for biomedical natural language processing.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.

HuggingFace. 2019. Pytorch transformer repository. <https://huggingface.co/>.

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations.

Patrick Lewis, Myle Ott, Jingfei Du, and Veselin Stoyanov. 2020. Pretrained language models for biomedical and clinical tasks: Understanding and extending the state-of-the-art. In *Proceedings of the 3rd Clinical Natural Language Processing Workshop*, pages 146–157, Online. Association for Computational Linguistics.

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*.

Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. S2ORC: The semantic scholar open research corpus. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4969–4983, Online. 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. *Advances in neural information processing systems*, 32.

## A Model Implementations Used

We used the PyTorch implementations of the following models from the HuggingFace Transformers library, as well as the following model checkpoints from the HuggingFace model hub:

1. **BERT** bert-base-cased and bert-base-uncased
2. **RoBERTa** roberta-base and roberta-large
3. **DSP-RoBERTa**  
dsp\_roberta\_base\_dapt\_cs\_tapt\_citation\_intent
4. **Biomed-RoBERTa**  
allenai/biomed\_roberta\_base
5. **SciBERT** allenai/scibert\_scivocab\_cased and allenai/scibert\_scivocab\_uncased
6. **PubMedBERT**  
BiomedNLP-PubMedBERT-base
7. **XLNet** xlnet-base-cased and xlnet-large-cased
8. **DeBERTa**  
microsoft/deberta-base and microsoft/deberta-large
9. **ALBERT** albert-base-v2 and albert-large-v2

The training was performed on a single V100 32GB GPU with automatic mixed precision enabled.

## B Additional diagrams

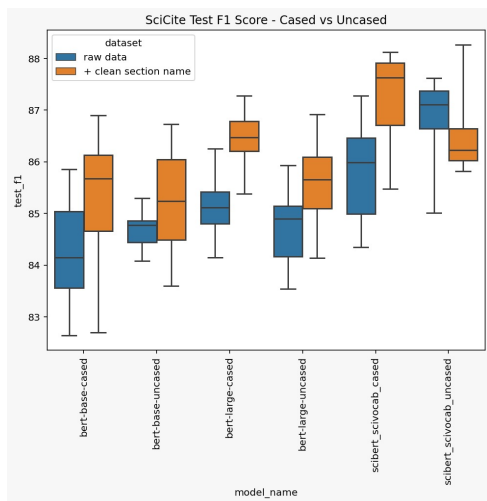


Figure 2: Spread of results.

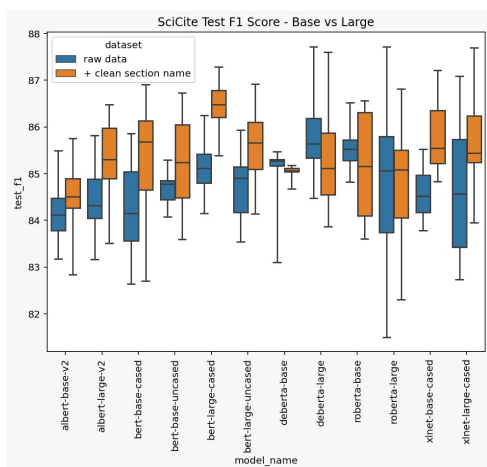


Figure 3: Spread of results.