

Chemical Language Understanding Benchmark

Yunsoo Kim Hyuk Ko Jane Lee Hyun Young Heo Jinyoung Yang Sungsoo Lee Kyu-hwang Lee

LG Chem

{kys.930303, kohyuk, janelee, hyheo11, yang.jy, ssoolee, amine}@lgchem.com

Abstract

In this paper, we introduce the benchmark datasets named CLUB (Chemical Language Understanding Benchmark) to facilitate NLP research in the chemical industry. We have 4 datasets consisted of text and token classification tasks. As far as we have recognized, it is one of the first examples of chemical language understanding benchmark datasets consisted of tasks for both patent and literature articles provided by industrial organization. All the datasets are internally made by chemists from scratch. Finally, we evaluate the datasets on the various language models based on BERT and RoBERTa, and demonstrate the model performs better when the domain of the pre-trained models are closer to chemistry domain. We provide baselines for our benchmark as 0.7818 in average, and we hope this benchmark is used by many researchers in both industry and academia. The CLUB can be downloaded at https://huggingface.co/datasets/bluesky333/chemical_language_understanding_benchmark.

1 Introduction

Transformer is the prevalent network architecture in natural language processing (NLP) (Vaswani *et al.*, 2017). It uses self-attention to capture each word's influence on another in a given text. Leveraging this architecture, recent advances in pre-training language models has reached state-of-the-art performances on many NLP benchmark datasets, including results that surpassed human performance (Wang *et al.*, 2019). Such advancements in language models and NLP technologies can potentially streamline and simplify the labor-intensive work for the literature and patent analysis, which are crucial in the research and development domain.

The benchmark datasets such as GLUE and SuperGLUE played a pivotal role in facilitating the advancement of NLP using language models (Wang *et al.*, 2018 and Wang *et al.*, 2019). This has inspired efforts to create benchmark datasets in the science domain as well (Yu Gu *et al.*, 2020). However, these attempts are limited within the field of biology and medicine.

In chemistry, there are few datasets available, however, as far as we know there are no benchmark datasets that include tasks for both literature articles and patents (Mysore *et al.*, 2019, Friedrich *et al.*, 2020, He *et al.*, 2021). Given the predominant reliance on patents in the chemical industry's research, especially in the early stages of product development, it is important to have datasets with patent documents to enable language models to comprehend the distinctive patent writing style, thereby performing better on tasks with patent documents.

On the other hand, academic literature often serves as the source of information that leads to new ideas for experimentation. Thus, it is critical to build a language model that understands both literature articles and patents and benchmark datasets with texts from both patents and papers for the evaluation.

In this paper, we present Chemical Language Understanding Benchmark (CLUB) to facilitate NLP research in the chemical industry, especially the language model pre-training. CLUB consists of two datasets for patents and two datasets for papers. In terms of tasks, it includes two datasets for token classification such as chemical named entity recognition, and two datasets for text classification such as patent area classification. All these datasets are internally made by chemists. We do not rely on any preexisting publicly available datasets or shared tasks. Finally, we provide the performance of various language models including the ones pre-trained with chemistry literature articles and patents as the baselines for our benchmark datasets.

Tasks	Class Group (source corpus)	Sample Type (Number)	Average token length (std)	Class name	Definition	Train	Dev	
Text CLS	PETRO- CHEMICAL (Patent)	Paragraph (2,775)	448.19 (403.81)	Household	Patents for products used in household such as PET bottles	436	120	
				Construct	Patents for products used in construction such as PVC pipes	77	25	
				Automobile	Patents for products used in automobile such as Tires	312	89	
				HouseConst	Patents for products used in household and construction	481	93	
				IndustConst	Patents for products used in industrial and construction	274	62	
				Catalyst	Patents for catalyst used for production	334	94	
				Process	Patents for production process of the products	306	72	
	RHEOLOGY (Journal)	Sentence (2,017)	55.04 (16.46)	Biodegrad_Poly	biodegradable polymer (plastic material)	553	151	
				Poly_Struc	the crystal structure of polymer which is related with mechanical properties	421	105	
				Biodgrad_Prop	biodegradable property of polymer	470	97	
				Mechanical_Prop	mechanical property of polymer	90	31	
				Rheological_Prop	rheological property of polymer which is related with polymer processability	78	19	
	Token CLS	CATALYST (Patent)	Sentence (4,663)	42.07 (14.59)	Precatalyst	Pre-catalyst form of metallocene catalyst	365	71
					Olefin	Include monomers and comonomers that participate in the synthesis of supported catalyst	947	153
Solvent					A solvent that creates a reaction environment	1,287	356	
Additive					Additives necessary for the catalyst synthesis reaction include scavengers and cocatalysts.	402	131	
Support					Support material for synthesis	417	83	
BATTERY (Journal)		Sentence (3,750)	40.73 (10.79)	Cathode_Material	Lithium compound used for cathode electrode among the components of lithium ion battery	1,411	402	
				Coating_Material	Materials coated for the purpose of improving structural stability and chemical resistance of cathode materials	1,510	359	
				Coating_Method	Method for coating the coating material on the surface of the cathode material	409	134	

Table 1: CLUB datasets for text and token classification (CLS).

85 2 Tasks

86 The CLUB Benchmark is created from scratch to
87 evaluate language models that understand the
88 fields of chemistry and materials science. The
89 benchmark dataset includes two types of tasks: text
90 classification and token classification. To evaluate
91 the representation power of the language model for
92 both patents and literature articles, each task
93 consisted of a dataset created from the patent text
94 and a dataset created from the paper text. Various
95 topics such as polymers, rheology, catalysts, and
96 batteries were selected to evaluate different fields
97 of chemistry and materials science. The detailed
98 composition of the data set is summarized in Table
99 1.

100 2.1 Text Classification

101 Text classification task is to assign a sentence or
102 document to a proper class. In this paper, we
103 present two classification datasets: RHEOLOGY
104 for sentence classification and
105 PETROCHEMICAL for document classification.
106 These datasets comprise corpora from both
107 patents and journal articles with a focus on the
108 topics of polymers, rheology, and overall
109 petrochemicals. Each dataset is available in JSON
110 format with “id”, “sentence”, and “labels” as keys.

111
112 **RHEOLOGY** sentence classification dataset
113 contains the five groups that represent the
114 polymer structures and properties, especially for
115 biodegradable polymers. It consists of 2,017
116 sentences collected from the research paper. Each
117 sentence of the RHEOLOGY classification
118 dataset is annotated by experts manually. The
119 detailed explanation of each group is presented in
120 Table 1.

121
122 **PETROCHEMICAL** dataset categorizes patents
123 into seven groups within the petrochemical
124 industry. Each group of patents accounts for
125 important parts of the industry. The petrochemical
126 industry uses **catalysts** to make the final polymer
127 products for different applications such as PET
128 bottles (**household applications**), rubber
129 (**automobile applications**), and PVC plastics
130 (**construction applications**). This production is
131 done on a factory scale, so it has its **production**
132 **process**. The seven groups consist of 5
133 applications: 1) household, 2) automobile, 3)
134 construction, 4) household & construction, and 5)
135 automobile & construction. The other two groups
136 are catalysts and processes.

137 2.2 Token Classification

138 Token classification, which includes named entity
139 recognition task, identifies tokens belonging to
140 defined classes. Considering our interests, we
141 defined the CATALYST class group and the
142 BATTERY class group as shown in Table 1. We
143 created the named entity recognition benchmark
144 dataset based on these definitions. The labeling
145 was performed by expert researchers with over
146 five years of experience in relevant fields. The
147 labeling was done in IOB format (inside, outside,
148 beginning). The labeled data was then converted
149 into JSON format with “id”, “tokens”, and “labels”
150 as keys.

151 We preprocess the token classification datasets
152 to adjust the sentence length to be less than the
153 maximum sequence length. As for named entity
154 recognition, each token has labels, and tokens that
155 come after the maximum sequence length would
156 be discarded. Thus, the model would not be able
157 to learn from those discarded tokens. We
158 minimized this issue by making the distribution of
159 the sequence length more like the gaussian
160 distribution (Appendix A).

161
162 **CATALYST** is a dataset for recognizing
163 materials involved in catalyst synthesis reactions
164 in the full text of patents. Pre-catalyst, additive,
165 olefin, solvent, and supporting material are
166 substances that participate in this reaction, and
167 these are defined as classes. “Pre-catalyst” is the
168 main substance to make the catalyst. “Additives”
169 are added to make the polymer with different
170 characteristics. “Olefin” is the monomer that
171 makes the polymer using the catalyst. “Solvent” is
172 for the polymerization of the monomer to the
173 polymer for the catalyst. “Supporting material” is
174 used to support the catalyst to do the
175 polymerization better as well as more stable.

176
177 **BATTERY** is a dataset for recognizing cathode
178 materials from literature articles related to
179 lithium-ion batteries including all-solid-state
180 batteries. There are four key components of a
181 battery: cathode material, anode material,
182 separator, and electrolyte. “Cathode material”
183 refers to the lithium compound used in the
184 positive electrode of a battery and is the most
185 important element in a battery because it has a
186 decisive effect on the energy density, power
187 output, and cycle life of the battery. This dataset
188 also has "coating material" and "coating method"
189 classes which are material and method to coat the
190 surface of the cathode material.

191 3 Dataset Statistics

192 All datasets have been divided into a training set
193 and a development set (also known as the
194 evaluation set), following an 80/20 split ratio.

195 3.1 PETROCHEMICAL dataset

196 The PETROCHEMICAL dataset is composed of
197 2,775 paragraphs. As the dataset is created with
198 titles, abstracts, and claims of patents, so it has the
199 average paragraph length of 448.19 tokens, which
200 is considerably longer than the other three datasets.
201 Also, the standard deviation for the paragraph
202 length is 403.81 tokens, which is also larger than
203 the others. For the seven classes of the dataset, the
204 respective counts of paragraphs are as follows:
205 “Household” – 556, “Construct” – 102,
206 “Automobile” – 401, “HouseConst” – 574,
207 “IndustConst” – 336, “Catalyst” – 428, and
208 “Process” – 378.

209 3.2 RHEOLOGY dataset

210 The RHEOLOGY dataset is made up of 2,017
211 sentences with an average sentence length of
212 55.03 tokens. The standard deviation of the
213 sentence length is 16.46 tokens. 704 sentences
214 were labeled as “Biodegrad_Poly” class and 526
215 sentences were labeled as “Poly_Struc”. The
216 “Biodegrad_Prop”, “Mechanical_Prop”, and
217 “Rheological_Prop” classes, which are classes
218 related to material’s properties, were labeled with
219 567, 121, and 97 sentences, respectively.

220 3.3 CATALYST dataset

221 The CATALYST dataset consists of 4,663
222 sentences. The average sentence length is 42.07
223 tokens with 14.59 tokens for standard deviation.
224 “Solvent” class was labeled the most with 1,643
225 times, followed by “Olefin” class which as labeled
226 1,100 times. “Precatalyst”, “Additive”, and
227 “Support” were labeled 436, 533, and 500 times,
228 respectively.

229 3.4 BATTERY dataset

230 The BATTERY dataset consists of 3,750
231 sentences, and the average sentence length is
232 40.73 tokens with 10.79 tokens as standard
233 deviation. The token classification breakdown
234 shows that “Cathode_Material” and
235 “Coating_Material” classes were labeled 1,813
236 times and 1,869 times, respectively. Meanwhile,
237 the “Coating_Method” class was 543 times.

238 4 Methods

239 4.1 Baseline Models

240 **BERT-Base** We use the BERT-base weights
241 released on Hugging Face model repository
242 (Devlin *et al.*, 2018). Both **cased** and **uncased**
243 versions of the model are used. We refer to each
244 version as **BERT-cased** and **BERT-uncased**
245 respectively throughout our papers. The model is
246 pre-trained with a corpus made up of BooksCorpus
247 and text parts of English Wikipedia for 1 M steps.
248 The corpus is about 16GB. The pre-training batch
249 size is 256 sequences. This model utilizes a
250 wordpiece vocabulary. The vocab size is 28,894.

251 **BioBERT** We use **BioBERT-v1.2** weights released
252 on Hugging Face model repository (Lee *et al.*,
253 2020). This is a BERT-base-cased model pre-
254 trained with PubMed abstracts from the BERT-
255 base-cased initial checkpoints. It was trained for
256 200K steps on PubMed abstracts, 270K steps on
257 PubMed Central (PMC) full texts, and another 1 M
258 steps on PubMed abstracts. The pre-training corpus
259 is about 25GB. The pre-training batch size is 192.
260 As a continued pre-trained model, it uses the same
261 vocabulary as the **BERT-base-cased** model.

262 **SciBERT** We use **sciBERT-scivocab-uncased**
263 released on Hugging Face model repository
264 (Beltagy *et al.*, 2019). This is a pre-trained BERT
265 model with 1.14 M Semantic Scholar papers,
266 which is comprised of computer science (18%) and
267 biomedical domain (82%). It differs from
268 BioBERT as it is pre-trained from scratch. The
269 papers are full texts and resulting in a corpus size
270 of 20GB. The pre-training batch size and steps are
271 unknown. It has its own wordpiece vocabulary
272 made from the pre-training corpus. The vocabulary
273 has more science terms. The vocab size is 30,990.

274 **RoBERTa** We use **RoBERTa-base** model released
275 on Hugging Face model repository (Liu *et al.*,
276 2019). It is an improvement of BERT model with a
277 larger pre-training dataset and better optimized
278 hyperparameter settings. The model is pre-trained
279 with a 160GB corpus made up of BERT pre-
280 training corpus plus News and Web contents
281 crawled. It is trained for 1 M steps. The pre-training
282 batch size is 256 sequences. The model uses byte
283 pair encoding (BPE) vocabulary, which is different
284 from BERT’s wordpiece vocabulary. The vocab
285 size is 50,000.

Task	Text classification (Accuracy)		Token classification (F1)		Average
	RHEOLOGY	PETRO-CHEMICAL	CATALYST	BATTERY	
BERT-cased	0.7970	0.8099	0.6601	0.7532	0.7550
BERT-uncased	0.7921	0.8105	0.6944	0.7571	0.7635
RoBERTa	0.7958	0.7990	0.6899	0.7658	0.7626
BioBERT	0.7978	0.8086	0.7092	0.7636	0.7698
SciBERT	0.7938	0.8045	0.7314	0.7602	0.7724
RoBERTa-PM-M3	0.7983	0.8079	0.7194	0.7815	0.7767
RoBERTa-lit	0.8017	0.8126	0.7332	0.7772	0.7811
RoBERTa-lit-pat	0.7968	0.8205	0.7323	0.7777	0.7818

Table 2: Performance of the model for the benchmark tasks. The evaluation for the text classification tasks was done using accuracy and the evaluation of the token classification tasks was done using macro-average of F1 scores. The evaluation result is the average of performances over ten runs.

RoBERTa-PM-M3 We use **RoBERTa-base-PM** weights released on Hugging Face model repository (Lewis *et al.*, 2020). It is a **RoBERTa-base** model pre-trained with a text corpus made of 27GB of PubMed abstracts, 60GB of PMC full texts, and 3.3 GB of the Medical Information Mart for Intensive Care (MIMIC-III). The model is trained for 500K steps on the corpus with a batch size of 8,192 sequences. It uses byte-pair encoding vocabulary made from the corpus, so it has a different BPE encoding vocabulary from RoBERTa-base. The vocabulary has more biomedical terms. The vocab size is 50,000.

4.2 Pre-training

For the chemistry pre-training, we gathered a large amount of chemistry patents and literature articles to train two different versions of models.

RoBERTa-lit We use **RoBERTa-PM-M3** weights as the initial checkpoint to pre-train the model with chemistry articles. We collected the abstracts of the articles using Open Academic Graphs and used the chemistry field of study to filter the ones that belong to the chemistry domain (Tang *et al.*, 2008 and Sinha *et al.*, 2015). For the filtered ones, all the abstracts were used as the training corpus. We train the model with the corpus for 1 epoch.

RoBERTa-lit-pat We use **RoBERTa-lit** weights as the initial checkpoint to pre-train the model this time with chemistry patents. We collected the

patents using USPTO BulkDownload. We filtered the chemical patents using the CPC code. For the filtered ones, abstracts, claims, and embodiment texts were used as the training corpus together with the **RoBERTa-lit**'s corpus. We train the model with the corpus for 1 epoch.

RoBERTa-lit and **RoBERTa-lit-pat** were pre-trained with NVIDIA V100 GPU and the hyperparameter setting follows the pre-training setup for **RoBERTa-PM-M3**. We also used mixed precision for training. We used the masked language model objective for the pre-training.

We expect that by pre-training the models with chemistry data, the models can learn the chemistry domain knowledge better and thus perform better on the CLUB benchmark.

4.3 Finetuning Language Models

For each dataset, we fine-tuned each models for 10 epochs with a 5e-05 learning rate on a single V100 GPU. We used 0.1 warm-up ratio, and cosine with restarts as the learning scheduler type. The training batch size was 128 and the evaluation batch size was 128. The maximum input length was 256. AdamW was used as the optimizer with a weight decay of 0.01. We used mixed precision for efficient training. We fine-tuned the model for 10 different seed initializations.

4.4 Evaluation

We evaluated all models using the accuracy for text classification tasks and the macro-average F1 score for token classification tasks. We chose the accuracy as the evaluation metric for the text classification due to its interpretability in measuring the effectiveness of the models. For token classification tasks, the use of the IOB scheme, which resulted in the "O" label being the dominant class, limited us from using the evaluation metric as text classification tasks. To provide a more balanced evaluation, we computed the F1 score of each token class excluding the "O" class, and used the macro-average of these F1 scores as the evaluation metric. For both types of tasks, the performance was averaged over ten runs with different seed initializations to reduce variance caused by randomness.

5 Results and Discussion

The performance of each model on the benchmark tasks is shown in Table 2. In general, our RoBERTa-lit-pat model outperformed the other models on average across the tasks. The result of BioBERT models pre-trained with a bio-related corpus was better than that of BERT base models, highlighting the impact of domain specific pre-training. SciBERT model pre-trained with a broad scientific literature articles performed well, especially in CATALYST task, though it still had a lower performance than RoBERTa models pre-trained with chemistry corpus. RoBERTa-PM-M3 model outperformed other models in the BATTERY task, but its overall performance was lower than that of the RoBERTa-lit-pat model.

In the text classification task, RoBERTa-lit model was the best model in the RHEOLOGY task and RoBERTa-lit-pat model score the highest in the PETROCHEMICAL task. This suggests that inclusion of patents in pre-training yields better performance in tasks with patent documents. As the PETROCHEMICAL dataset includes titles, abstracts, and representative claims of patents, the terminology used in the dataset is quite different from the terminology used in other datasets made up of literature articles. This is due to the nature of patents to protect an invention, leading them to be written in a more general manner to encompass a broader patent space.

In the CATALYST task, it was very interesting that RoBERTa-lit model, solely pre-trained on

academic papers, showed the best results in the task with patents. This task involved labeling only the embodiment section of the patent. The terminology used in the embodiment part of the patent is closer to academic language than the language used in patent claims. This could explain why a model trained only on articles could perform better in this task.

For the BATTERY task, RoBERTa-PM-M3 model had the best performance, closely followed by RoBERTa-lit-pat model. Notably RoBERTa-lit and RoBERTa-lit-pat models still showed good average performance despite only being pre-trained for one epoch. It is plausible that the performance of RoBERTa-lit-pat improves further with additional training epochs. Due to our GPU infrastructure limitations, we leave this for future work.

6 Conclusion

Chemical Language Understanding Benchmark (CLUB) is the first benchmark in the chemistry industry aimed at chemical language model evaluation with tasks for both patents and journal articles. The introduction of this benchmark is expected to catalyze research in natural language processing, particularly in information extraction, within the chemistry domain.

In the course of establishing baseline performance for the CLUB, we tested existing pre-trained models as well as our novel pre-trained models. Remarkably, the RoBERTa model pre-trained on chemical patents and literature articles, reached the highest average score, 0.7818. This performance highlights the advantage of pre-training models with a corpus closely aligned with the target domain.

Our benchmark provides a powerful tool for evaluating language models' learning capacity in the chemistry context. In addition, the tasks in our benchmark can be leveraged to accelerate the literature and patent analysis by automatically extracting information such as new chemical molecules and experiment settings.

Thus, these tasks can be the foundation of an information extraction based expert system. This system would generate structured knowledge from a large volume of papers and patents and help researchers to conduct their experiments on time without falling behind the research trends.

Our benchmark sets the foundation for future advancements in chemical language understanding.

452 It contributes to the acceleration of scientific
453 discovery in the field by integrating natural
454 language processing into chemical research and
455 development.

456 Limitations

457 Because we were doing the manual labeling with
458 experts in the field, we were only limited to two
459 types of tasks: token classification and text
460 classification. We hope to expand the benchmark to
461 include other types of tasks such as summarization,
462 question and answering, and sentence similarity in
463 the future. Sentence similarity for patents is the task
464 we aim to add for the next version because it can
465 be used to find the infringement in patents.

466 While the CLUB provides a robust benchmark for
467 evaluating language models in the context of
468 chemistry, it is not without its limitations. The
469 present version of CLUB only includes two types
470 of tasks: token classification and text classification.
471 This constraint arises primarily from the manual
472 labeling process which involved domain experts.

473 However, we aim to extend the benchmark in the
474 future to include a wider range of tasks such as
475 summarization, question answering, and sentence
476 similarity assessments. We are particularly
477 interested in the sentence similarity task for patents
478 as this could be leveraged for identifying potential
479 patent infringements.

480 Acknowledgements

481 We express our sincere gratitude to the anonymous
482 reviewers who contributed their valuable time and
483 effort to provide insightful and constructive
484 feedback on this work. Your detailed comments
485 and suggestions have greatly aided us in refining
486 our work. We would also like to extend our thanks
487 to everyone involved in the creation and
488 development of the labeled datasets. This work
489 would not have been possible without your
490 dedication and collaborative efforts. Last but not
491 the least, we are grateful for the continuous support
492 and resources provided by our institutions, which
493 have been fundamental in conducting this research.

494 References

495 Alex Wang, Amanpreet Singh, Julian Michael, Felix
496 Hill, Omer Levy, Samuel R. Bowman. 2018. [GLUE:
497 A Multi-Task Benchmark and Analysis Platform for
498 Natural Language Understanding](#). arXiv preprint
499 arXiv:1804.07461.

500 Alex Wang, Yada Pruksachatkun, Nikita Nangia,
501 Amanpreet Singh, Julian Michael, Felix Hill, Omer
502 Levy, Samuel R. Bowman. 2019. [SuperGLUE: A
503 Stickier Benchmark for General-Purpose Language
504 Understanding Systems](#). arXiv preprint
505 arXiv:1905.00537.

506 Annemarie Friedrich, Heike Adel, Federico Tomazic,
507 Johannes Hingerl, Renou Benteau, Anika
508 Marusczyk, and Lukas Lange. 2020. [The SOFC-
509 Exp Corpus and Neural Approaches to Information
510 Extraction in the Materials Science Domain](#). In
511 *Proceedings of the 58th Annual Meeting of the
512 Association for Computational Linguistics*, pages
513 1255–1268, Online. Association for Computational
514 Linguistics.

515 Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma,
516 Darrin Eide, Bo-June (Paul) Hsu, and Kuansan
517 Wang. 2015. [An Overview of Microsoft Academic
518 Service \(MAS\) and Applications](#). In *Proceedings of
519 the 24th International Conference on World Wide
520 Web (WWW '15 Companion)*. ACM, New York, NY,
521 USA, pages 243-246.

522 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
523 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
524 Kaiser, and Illia Polosukhin. 2017. [Attention is all
525 you need](#). arXiv preprint arXiv:1706.03762.

526 Iz Beltagy, Kyle Lo, and Arman Cohan. 2019.
527 [SciBERT: A Pretrained Language Model for
528 Scientific Text](#). arXiv preprint arXiv:1903.10676.

529 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
530 Kristina Toutanova. 2018. [Bert: Pre-training of deep
531 bidirectional transformers for language
532 understanding](#). arXiv preprint arXiv:1810.04805.

533 Jiayuan He, Dat Quoc Nguyen, Saber A. Akhondi,
534 Christian Druckenbrodt, Camilo Thorne, Ralph
535 Hoessel, Zubair Afzal, Zenan Zhai, Biaoyan Fang,
536 Hiyori Yoshikawa, Ameer Albahem, Lawrence
537 Cavedon, Trevor Cohn, Timothy Baldwin, and
538 Karin Verspoor. 2021. [Chemu 2020: Natural
539 language processing methods are effective for
540 information extraction from chemical patents](#).
541 *Frontiers in Research Metrics and Analytics*, 6,
542 654438.

543 Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang,
544 and Zhong Su. [ArnetMiner: Extraction and Mining
545 of Academic Social Networks](#). 2008. In *Proceedings
546 of the Fourteenth ACM SIGKDD International
547 Conference on Knowledge Discovery and Data
548 Mining (SIGKDD'2008)*. pages 990-998.

549 Jinhyuk Lee, Wonjin Yoon, Sungdong Kim,
550 Donghyeon Kim, Sunkyu Kim, Chan Ho So and
551 Jaewoo Kan. 2020. [BioBERT: a pre-trained
552 biomedical language representation model for
553 biomedical text mining](#). *Bioinformatics*, Volume 36,
554 Issue 4, February 2020, pages 1234–1240.

555 Patrick Lewis, Myle Ott, Jingfei Du, and Veselin
556 Stoyanov. 2020. [Pretrained Language Models for](#)
557 [Biomedical and Clinical Tasks: Understanding and](#)
558 [Extending the State-of-the-Art](#). In *Proceedings of the*
559 *3rd Clinical Natural Language Processing*
560 *Workshop*, pages 146–157, Online. Association for
561 Computational Linguistics.

562 Sheshera Mysore, Zach Jensen, Edward Kim, Kevin
563 Huang, Haw-Shiuan Chang, Emma Strubell, Jeffrey
564 Flanigan, Andrew McCallum, Elsa Olivetti. 2019.
565 [The Materials Science Procedural Text Corpus:](#)
566 [Annotating Materials Synthesis Procedures with](#)
567 [Shallow Semantic Structures](#). arXiv preprint
568 arXiv:1905.06939.

569 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du,
570 Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
571 Luke Zettlemoyer, and Veselin Stoyanov. 2019.
572 [RoBERTa: A Robustly Optimized BERT Pretraining](#)
573 [Approach](#). arXiv preprint arXiv:1907.11692.

574 Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto
575 Usuyama, Xiaodong Liu, Tristan Naumann,
576 Jianfeng Gao and Hoifung Poon. 2021. [Domain-](#)
577 [Specific Language Model Pretraining for](#)
578 [Biomedical Natural Language Processing](#). In *ACM*
579 *Transactions on Computing for Healthcare*, pages
580 1–23.

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602 A Adjust sentence length

603 Figure 1. shows the distribution of sentence lengths
604 in the dataset before and after the preprocessing.
605 After adjusting the sentence length, the sequence
606 length distribution follows more of a Gaussian
607 distribution than before. In the case of CATALYST
608 dataset, the number of sentences was reduced from
609 12,368 to 4,663. However, in the case of
610 BATTERY dataset, there was no change in the
611 number of the sentences. We made this
612 preprocessing to minimize the number of tokens
613 that come after the maximum sequence.

614

615

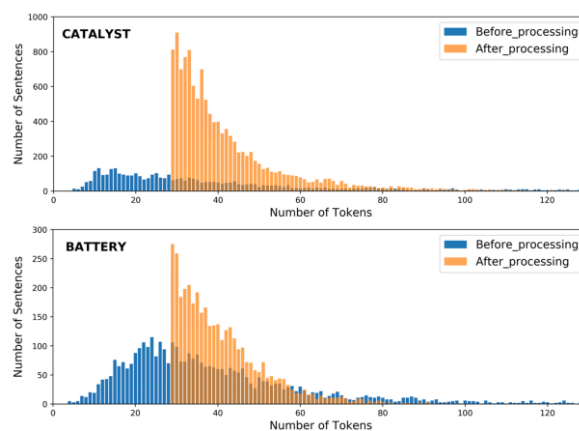


Figure 1. Distribution of sequence length before and after sentence adjustment in token classification task datasets

616

617