
A Search Engine for Discovery of Scientific Challenges and Directions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Keeping track of scientific challenges, advances and emerging directions is a fun-
2 damental part of research. However, researchers face a flood of papers that hinders
3 discovery of important knowledge. In biomedicine, this directly impacts human
4 lives. To address this problem, we present a novel task of extraction and search of
5 scientific challenges and directions, to facilitate rapid knowledge discovery. We
6 construct and release an expert-annotated corpus of texts sampled from full-length
7 papers, labeled with novel semantic categories that generalize across many types
8 of challenges and directions. We focus on a large corpus of interdisciplinary work
9 relating to the COVID-19 pandemic, ranging from biomedicine to areas such as AI
10 and economics. We apply a model trained on our data to identify challenges and
11 directions across the corpus and build a dedicated search engine. In experiments
12 with 19 researchers and clinicians using our system, we outperform a popular
13 scientific search engine in assisting knowledge discovery. Finally, we show that
14 models trained on our resource generalize to the wider biomedical domain and
15 to AI papers, highlighting its broad utility. We make our data, model and search
16 engine publicly available.¹

17 1 Introduction

18 Success in scientific efforts hinges on identifying promising and important problems to work on,
19 developing novel and effective solutions, and formulating hypotheses and directions for further
20 exploration. Each new scientific advance helps address gaps in knowledge, including potential
21 extensions and refinements of prior results. New advances often lead to new challenges and directions.
22 With millions of scientific papers published every year, sets of challenges and potential directions
23 for addressing them grow rapidly. A striking recent example is that of literature pertaining to the
24 COVID-19 pandemic [36], which exploded in unprecedented volume with researchers from across
25 diverse fields exploring the many facets of the disease and its societal ramifications. As the pandemic
26 continues worldwide, it is especially urgent to provide scientists with tools for staying aware of
27 advances, problems, and limitations faced by fellow researchers and medical professionals, and of
28 emerging hypotheses or early indications of potential solutions.

29 Unfortunately, due to the immense scale and siloed nature of the scientific community, it can
30 be difficult for researchers to keep track of their own specialty areas, let alone discover relevant
31 knowledge in areas outside their immediate focus [12, 13, 14, 30]. This can result in poor awareness
32 of failures or limitations reported in recent studies, wasting redundant resources and leading to clinical
33 decision-making uninformed about shortcomings of interventions [5]. Disturbingly, there have been
34 many cases where problems in treatments had been reported but not picked up by sectors of the
35 clinical community [6, 31, 7] leading to higher rates of morbidity and mortality [17, 9, 33].

¹Redacted for anonymity.

36 Our goal is to bolster the ability of researchers and clinicians to **keep track of difficulties, limitations**
37 **and emerging hypotheses**. This could help clinical decision making be well-informed, accelerate
38 innovation by surfacing new opportunities to work on, inspire new research directions, and match
39 challenges with potential solutions from other communities [13]. In the face of challenging medical
40 scenarios, such as the rise of a novel virus or situations where standard treatments fail, rapidly finding
41 reports of similar challenges and directions to address them could have dramatic effect [26]. Finally,
42 at the macro level, this ability could assist policymakers and funding agencies (e.g., NIH, NSF)
43 seeking to identify important challenges and promising directions to prioritize research programs; in
44 times of crisis this process needs to be done rapidly² but demands substantial human effort.

45 To address this problem and facilitate discovery of scientific knowledge, we make the following key
46 contributions:

- 47 • **Novel Task: Extraction and Search of Scientific Challenges and Directions.** We define semantic
48 categories for ‘challenges’ and ‘directions’ that generalize across many types of difficulties,
49 limitations, flaws and hypotheses or potential indications that an issue is worthy of investigation.
50 We focus on COVID-19 literature as the main test bed for our task, as it is known to be highly
51 interdisciplinary [14] with research in many different fields (e.g., AI, climatology, engineering,
52 economics) and relates to a global emergency that urgently demands tools to help researchers and
53 clinicians keep track of challenges and new opportunities.
- 54 • **Expert-Annotated Dataset, Publicly Released.** We collect and publicly release a resource of
55 2.9K expert-annotated texts from full-length COVID-19 papers, labeled by experts for challenges
56 and directions with high inter-annotator agreement. We use the data to train multi-label sentence
57 classification models that achieve high accuracy scores. We analyze model errors, discovering that
58 contextual information can both help and harm results. Based on this finding, we explore a simple
59 technique that integrates multiple ways of encoding context.
- 60 • **Novel Scientific Search Engine For Researchers and Clinicians.** We build a novel public search
61 engine that indexes challenges and directions. We apply a model trained on our dataset and apply it
62 to the full corpus of 550K COVID-19 papers to build an index of scientific challenges and potential
63 directions. We create a search engine that allows users to search for combinations of entities (e.g.,
64 names of drugs, diseases, etc.) and retrieve challenge/direction sentences that mention them.
- 65 • **Evaluating Generality: Zero-Shot Generalization across Biomedicine and AI.** We demon-
66 strate zero-shot generalization, obtaining a high MAP of over 95% when applying the model
67 trained on COVID-19 papers to a broader corpus in the general biomedical domain, and to AI
68 papers in computer science. This indicates the potential value of our resource beyond COVID-19,
69 such as for future pandemics or crises, or for helping AI researchers handle the explosion of
70 research in this area.
- 71 • **Evaluating Utility: User Studies with Researchers.** We conduct studies measuring utility. First,
72 we evaluate the system’s ability to help researchers with diverse backgrounds discover challenges
73 and directions for a given query (e.g., directions in *drug discovery*). This could also be important
74 for researchers looking into a new area, e.g., AI researchers seeking biomedical problems (Fig.
75 1). Second, we recruit nine *medical researchers working on COVID-19* in clinical practice and
76 research. These users often require finding information on challenges and directions, during
77 research or treatment planning. In both experiments, totalling 19 researchers and over 70 distinct
78 queries, our prototype outperforms PubMed, the most widely used biomedical search tool, in both
79 quality and utility for discovery of challenges and directions.

80 2 Task Overview & Definitions

81 We present a novel task of automatically identifying sentences in papers that clearly state *scientific*
82 *challenges and directions*. We consider the multi-label classification setting, where for a given
83 sentence $X = f w_1; w_2; \dots; w_T g$ with T tokens, our goal is to output two labels $Y = f c; d g$, where
84 c and d are binary targets indicating if the sentence mentions a challenge/direction, respectively.
85 Additionally, we are also given *context* sentences surrounding X : $(X_{previous}; X_{next})$, for the previous
86 and next sentences, respectively, which could be used as further input to models. The multi-label

²[https://covid19.who.int/news-detail/covid-19-who-2020-04-08/](https://covid19.who.int/news-detail/covid-19-who-2020-04-08)

Figure 1: Overview of our system. (1) We collect expert annotations of sentences mentioning challenges and directions from across the CORD-19 corpus. (2) We train a sentence identification model on this data and apply it to the full corpus to extract high-confidence sentences. (3) We build a search engine indexing challenges and directions in COVID-19 literature, allowing users to search for entities and retrieve sentences with their contexts.

87 setting allows us to capture that in many cases, sentences refer to both challenges and directions at
88 the same time. At a high level, our labels are defined as follows.

- 89 • Challenge: A sentence mentioning a problem, difficulty, law, limitation, failure, lack of clarity, or
90 knowledge gap.
- 91 • Research direction: A sentence mentioning suggestions or needs for further research, hypotheses,
92 speculations, indications or hints that an issue is worthy of exploration.

93 Figure 1 shows examples for each category. Also, we further present the motivation for the task,
94 example annotations, and why the categories are non-trivial for both humans & machines to identify
95 in Technical Appendix §A.1 & §A.2. We note that in addition to biomedical literature discussed in
96 the Introduction, our task is related to a body of research which we cover in Technical Appendix
97 §A.3.

98 3 Data Collection and Models

99 Data Collection & Annotation. We sample 3000 sentences from the full-text papers of CORD-19
100 (180k full-text papers with 25m sentences³). Four expert annotators with biomedical and bioNLP
101 backgrounds annotated the sentences with high agreement. We create a train/dev/test stratified split of
102 40%/10%/50%, splitting by distinct papers⁴. See full details in Appendix §A.4.

103 Baseline Models. We evaluate a range of baseline models for our novel task: simple key-
104 word/sentiment based heuristics, zero-shot inference based on a language model trained for natural
105 language inference (NLI), re-tuning scientific and non-scientific language models (LMs), and re-
106 tuning the LMs where they are context-aware (including the surrounding sentences in the training).
107 We also explore two customized approaches: re-tuning using a Hierarchical Attention Network
108 (HAN) [38], and an approach where we obtain outputs using several variants (“slices”) of context-
109 aware and not-context-aware versions of a re-tuned language model and combine their logits to
110 yield a final pair of logits used for prediction. See Appendix §A.5 & §A.6 for a full description.

111 Results. As seen in Table 1, the best individual classifier by F1 is PubMedBERT with a binary-F1
112 of 0.770 and 0.766 on the challenge and direction labels, respectively. Our customized approach leads
113 to an improvement of about one F1 point for both labels over the best individual model (standard
114 error of $1:05 \cdot 10^{-4}$). See full results in Appendix §A.7 & error analysis in Appendix §A.8.

³We use a snapshot of CORD-19 from 08-02-2021.

⁴See Table §3 in Appendix §A.4.

Model	Challenge			Direction		
	P	R	F1	P	R	F1
Keyword	0.535	0.760	0.628	0.455	0.792	0.578
Sentiment	0.405	0.966	0.571	0.239	0.837	0.371
NLI	0.659	0.693	0.675	0.401	0.825	0.540
RoBERTa-lg	0.723 (0.042)	0.824 (0.046)	0.769 (0.004)	0.697 (0.065)	0.825 (0.06)	0.754 (0.004)
SciBERT	0.729 (0.023)	0.799 (0.03)	0.761 (0.007)	0.719 (0.044)	0.783 (0.043)	0.749 (0.01)
PubMedBERT	0.738 (0.018)	0.804 (0.017)	0.770(0.006)	0.755 (0.017)	0.778 (0.015)	0.766(0.006)
HAN	0.671 (0.02)	0.863 (0.03)	0.759 (0.01)	0.674 (0.04)	0.804 (0.04)	0.734 (0.001)
Ctx-slices	0.742 (0.011)	0.829 (0.012)	0.783(0.004)	0.732 (0.02)	0.82 (0.03)	0.773(0.005)

Table 1: Model Results. The PubMedBERT model ne-tuned on our multi-label classification task performs best. For the neural models we present the average over 5 training seeds where the number in parentheses is the standard deviation.

115 High precision@K We observe that for 20% recall we obtain well over 90% precision, and for 40%
 116 recall about 90% precision, demonstrating the utility of our model for a search engine application
 117 (§4) where precision at the top is important.

118 Generalization For directions, we obtain MAP and AUC of around 96% for general biomedicine
 119 papers outside CORD-19, and around 95% for AI papers we sample from the computer science
 120 domain. For challenges, we reach a MAP and AUC reach around 97-98% for biomedicine and around
 121 96% for AI. See full analysis in Appendix §A.7.

122 4 Search Engine User Studies

123 We build a search engine that indexes challenges and directions across 550k papers in CORD-
 124 19. We perform entity linking to the biomedical KB of MeSH entities [24] which allows us to
 125 partially group together all challenges or directions into “topics” referring to a specific
 126 combination of concepts. See Appendix §A.9.1 for full details.

127 Experiment I - diverse researchers We recruited ten participants with diverse research backgrounds.
 128 Each participant was given twenty queries (formulated by a domain expert). We compared the amount
 129 of challenges and directions they were able to find in a limited time frame using our system vs. the
 130 popular PubMed search system. Our system, on average, yielded 2.6 times as much challenges and
 131 3 times as much directions per query. Further details are in Appendix §A.9.2.

132 Experiment II - clinical researchers We recruited nine expert MDs at a large hospital who are
 133 involved in clinical research for COVID-19 and for their specialty areas (each have over 1000
 134 citations). Each expert completed randomly ordered search tasks (queries curated by an expert; see
 135 Appendix §A.12) using the same systems as in the previous experiment. After all search tasks were
 136 completed we use a standardized Post Study System Usability Questionnaire (PSSUQ)
 137 experts strongly preferred our search engine to PubMed (overall average of 92% versus 59%, with
 138 non-normalized scores of 6.42 vs. 4.14). Further details are in Appendix §A.9.3.

139 5 Conclusion

140 We presented methods for extracting scientific challenges and directions from scholarly papers. We
 141 collected 3K expert-labeled sentences and their contexts from COVID-19 papers, and used the dataset
 142 to fine-tune scientific language models on our multi-label sentence classification task. We found that the
 143 approach can identify challenges and directions with high precision, and that using the model trained
 144 on our dataset achieves high zero-shot generalization on general biomedical papers and AI papers in
 145 computer science. We harnessed the model to index 950K sentences and build a novel search engine
 146 that allows researchers to search for biomedical entities and retrieve sentences mentioning difficulties,
 147 limitations, hypotheses and directions. Researchers using our system found that our system provided
 148 better support than PubMed in terms of utility and relevance.

149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199

References

- [1] Michael Bada, Miriam Eckert, Donald Evans, Kristin Garcia, Krista Shipley, Dmitry Sitnikov, William A Baumgartner, K Bretonnel Cohen, Karin Verspoor, Judith A Blake, et al. Concept annotation in the craft corpus. *BMC bioinformatics* 2012.
- [2] Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* pages 3606–3611, 2019.
- [3] Olivier Bodenreider. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research* 2004.
- [4] D. Card, Peter Henderson, Urvashi Khandelwal, Robin Jia, Kyle Mahowald, and Dan Jurafsky. With little power comes great responsibility. *EMNLP*, 2020.
- [5] Iain Chalmers, Michael B Bracken, Ben Djulbegovic, Silvio Garattini, Jonathan Grant, A Metin Gülmezoglu, David W Howells, John P A Ioannidis, and Sandy Oliver. How to increase value and reduce waste when research priorities are low. *The Lancet* 383(9912):156–165, 2014. ISSN 0140-6736. doi: [https://doi.org/10.1016/S0140-6736\(13\)62229-1](https://doi.org/10.1016/S0140-6736(13)62229-1). <https://www.sciencedirect.com/science/article/pii/S0140673613622291>.
- [6] M. Clarke and S. Hopewell. Many reports of randomised trials still don't begin or end with a systematic review of the relevant evidence. *Journal of the Bahrain medical society* 24:145–148, 2013.
- [7] N. Cooper, David R. Jones, and A. Sutton. The use of systematic reviews when designing studies. *Clinical Trials*, 2:260 – 264, 2005.
- [8] Franck Dernoncourt and Ji Young Lee. Pubmed 200k rct: a dataset for sequential sentence classification in medical abstracts. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)* pages 308–313, 2017.
- [9] R. Gilbert, G. Salanti, M. Harden, and S. See. Infant sleeping position and the sudden infant death syndrome: systematic review of observational studies and historical review of recommendations from 1940 to 2002. *International journal of epidemiology* 34 4:874–87, 2005.
- [10] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *arXiv preprint arXiv:2007.15772*, 2020.
- [11] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. Don't stop pretraining: Adapt language models to domains and tasks. *ACL*, 2020.
- [12] Tom Hope, Joel Chan, Aniket Kittur, and Dafna Shahaf. Accelerating innovation through analogy mining. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* pages 235–243, 2017.
- [13] Tom Hope, Jason Portenoy, Kishore Vasan, Jonathan Borchartd, Eric Horvitz, Daniel S Weld, Marti A Hearst, and Jevin West. Scisight: Combining faceted navigation and research group detection for COVID-19 exploratory scientific search. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations* pages 135–143, 2020.
- [14] Tom Hope, Aida Amini, David Wadden, Madeleine van Zuylen, Sravanthi Parasa, Eric Horvitz, Daniel Weld, Roy Schwartz, and Hannaneh Hajishirzi. Extracting a Knowledge Base of Mechanisms from COVID-19 Papers. 2021.
- [15] Ting-Hao Huang, Chieh-Yang Huang, Chien-Kuang Cornelia Ding, Yen-Chia Hsu, and C Lee Giles. Coda-19: Using a non-expert crowd to annotate research aspects on 10,000+ abstracts in the covid-19 open research dataset. *Proceedings of the 1st Workshop on NLP for COVID-19 at ACL* 2020.
- [16] Sarthak Jain, Madeleine van Zuylen, Hannaneh Hajishirzi, and Iz Beltagy. Scirex: A challenge dataset for document-level information extraction. *arXiv preprint arXiv:2005.00512*, 2020.
- [17] Katharine Ker, Phil Edwards, Pablo Perel, Haleema Shakur, and Ian Roberts. Effect of tranexamic acid on surgical bleeding: systematic review and cumulative meta-analysis. *BMJ*, 344, 2012. doi: 10.1136/bmj.e3054. URL <https://www.bmj.com/content/344/bmj.e3054>.

- 200 [18] Jin-Dong Kim, Tomoko Ohta, Yoshimasa Tsuruoka, Yuka Tateisi, and Nigel Collier. Introduction
201 to the bio-entity recognition task at jnlpba. *Proceedings of the international joint workshop on*
202 *natural language processing in biomedicine and its applications* 2004.
- 203 [19] Jin-Dong Kim, Yue Wang, and Yamamoto Yasunori. The genia event extraction shared task, 2013
204 edition-overview. In *BioNLP Shared Task Workshop* 2013.
- 205 [20] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- 206 [21] James R Lewis. Psychometric evaluation of the pssuq using data from ve years of usability studies.
207 *International Journal of Human-Computer Interaction* 2002.
- 208 [22] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy,
209 Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural
210 language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461* 2019.
- 211 [23] Jiao Li, Yueping Sun, Robin J Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter
212 Davis, Carolyn J Mattingly, Thomas C Wiegiers, and Zhiyong Lu. Biocreative v cdr task corpus: a
213 resource for chemical disease relation extraction. *Database* 2016.
- 214 [24] Carolyn E Lipscomb. Medical subject headings (mesh). *Bulletin of the Medical Library Association*
215 88(3):265, 2000.
- 216 [25] Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel S. Weld. S2ORC: The
217 Semantic Scholar Open Research Corpus. *Proceedings of ACL* 2020. URL <https://arxiv.org/abs/1911.02782> .
- 219 [26] Christopher A Longhurst, Robert A Harrington, and Nigam H Shah. A `green button` for using
220 aggregate patient data at the point of care. *Health affairs* 33(7):1229–1235, 2014.
- 221 [27] Steven Loria. textblob documentation. Release 0.15.2, 2018.
- 222 [28] Sunil Mohan and Donghui Li. Medmentions: A large biomedical corpus annotated with umls
223 concepts. In *Automated Knowledge Base Construction (AKBC)* 18.
- 224 [29] Mark Neumann, Daniel King, Iz Beltagy, and Waleed Ammar. Scispacy: Fast and robust models for
225 biomedical natural language processing. *arXiv preprint arXiv:1902.07669* 2019.
- 226 [30] Jason Portenoy, Marissa Radensky, Jevin West, Eric Horvitz, Daniel Weld, and Tom Hope. Bridger:
227 Toward bursting scientific liter bubbles and boosting innovation via novel author discovery. *arXiv*
228 preprint arXiv:2108.05669 2021.
- 229 [31] K. Robinson and S. Goodman. A systematic examination of the citation of prior research in reports
230 of randomized, controlled trials. *Annals of Internal Medicine* 54:50 – 55, 2011.
- 231 [32] Micah Shlain, Hillel Taub-Tabib, Shoval Sadde, and Yoav Goldberg. Syntactic search by example. In
232 ACL, 2020.
- 233 [33] J. Sinclair. Meta-analysis of randomized controlled trials of antenatal corticosteroid for the prevention
234 of respiratory distress syndrome: discussion. *American journal of obstetrics and gynecology* 173 1:
235 335–44, 1995.
- 236 [34] Hillel Taub-Tabib, Micah Shlain, Shoval Sadde, Dan Lahav, Matan Eyal, Yaara Cohen, and Yoav
237 Goldberg. Interactive extractive search over biomedical corpora, 2020.
- 238 [35] Byron C. Wallace, Joël Kuiper, Aakash Sharma, Mingxi Zhu, and Iain J. Marshall. Extracting pico
239 sentences from clinical trial reports using supervised distant supervision. *Mach. Learn. Res* 17(1):
240 4572–4596, January 2016. ISSN 1532-4435.
- 241 [36] Lucy Lu Wang, Kyle Lo, Yoganand Chandrasekhar, Russell Reas, Jiangjiang Yang, Darrin Eide,
242 Kathryn Funk, Rodney Kinney, Ziyang Liu, William Merrill, et al. Cord-19: The covid-19 open
243 research dataset. *arXiv preprint arXiv:2004.10706* 2020.
- 244 [37] Thomas Wolf, Julien Chaumond, Lysandre Debut, Victor Sanh, Clement Delangue, Anthony Moi,
245 Pierric Cistac, Morgan Funtowicz, Joe Davison, Sam Shleifer, et al. Transformers: State-of-the-art
246 natural language processing. *Proceedings of the 2020 Conference on Empirical Methods in Natural*
247 *Language Processing: System Demonstrations* pages 38–45, 2020.

- 248 [38] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Proceedings of
249 the 2016 Conference of the North American Chapter of the Association for Computational Linguistics:
250 Human Language Technologies
- 251 [39] Wenpeng Yin, Jamaal Hay, and Dan Roth. Benchmarking zero-shot text classification: Datasets,
252 evaluation and entailment approach. Proceedings of the 2019 Conference on Empirical Methods
253 in Natural Language Processing and the 9th International Joint Conference on Natural Language
254 Processing (EMNLP-IJCNLP) pages 3905–3914, 2019.

255 A Technical Appendix

256 A.1 Task Overview Definitions

257 The COVID-19 corpus [6] curates literature on COVID-19 and related diseases. With many thousands
258 of papers, keeping track is generally hard, and mapping the landscape of scientific challenges and
259 directions to address them is even harder. While “grand” challenges such as designing therapies and
260 handling novel virus variants are broadly known, research focuses on more granular specific challenges,
261 e.g., difficulties in functional analysis of specific viral proteins, or shortcomings of a specific treatment
262 regime for children. Each challenge, in turn, is associated with potential directions and hypotheses.

263 As written in the main paper, we present a novel task of automatically identifying sentences in papers
264 that clearly state scientific challenges and directions. We consider the multi-label classification setting,
265 where for a given sentence $x = [w_1; w_2; \dots; w_T]$ with T tokens, our goal is to output two labels
266 $Y = [c; d]$, where c and d are binary targets indicating if the sentence mentions a challenge/direction,
267 respectively. Additionally, we are also given context sentences surrounding x : $(X_{\text{previous}}; X_{\text{next}})$,
268 for the previous and next sentences, respectively, which could be used as further input to models. The
269 multi-label setting allows us to capture that in many cases, sentences refer to both challenges and
270 directions at the same time (see Table 2). At a high level, our labels are defined as follows.

- 271 • Challenge: A sentence mentioning a problem, difficulty, law, limitation, failure, lack of clarity, or
272 knowledge gap.
- 273 • Research direction: A sentence mentioning suggestions or needs for further research, hypotheses,
274 speculations, indications or hints that an issue is worthy of exploration.

275 These categories allow us to capture important information for scientists that is not captured by exist-
276 ing resources (see §A.3). As part of data annotation we provide annotators with richer explanations
277 and examples of each label (see §A.4) to make these definitions more concrete. Figure 1 shows
278 examples for each category (also see Table 2 in Technical Appendix §A.2 for more discussion).

279 Many cases of challenges and directions are non-trivial for both humans and machines to identify.
280 We demonstrate two main types of difficulties (see more discussion in Technical Appendix §A.8) —
281 cases of potentially misleading keywords, and cases where deep domain knowledge or context may
282 be required.

- 283 • Misleading keywords. Consider the following sentence: “The 15-30 mg/L albumin concentration
284 is a critical value that could indicate kidney problems when it is repeatedly exceeded.” This text
285 mentions a diagnostic measure that is an indicator of a problem, rather than an actual problem.
286 This is one example out of many other potentially misleading cases, such as cases where a negative
287 outcome occurs to an entity we wish to harm (e.g., “the viral structural integrity is destroyed”).
- 288 • Context and domain knowledge. “BV-2 cells expressed Mac1 (CD11b) and Mac2 but were
289 negative for the oligodendrocyte marker GalC and the astrocyte marker GFAP.” Depending whether
290 this sentence contains a challenge is highly non-trivial, since it requires more context and deep
291 domain knowledge to understand whether this outcome is problematic or not.

292 A.2 More examples

293 Table 2 shows example sentences for each category. In the first row (challenge, not direction),
294 the example is purely a factual description of a certain tool. In the second row (challenge,
295 direction), the statement mentions a scientific future direction, but there is no associated challenge
296 that is explicitly mentioned. In the third row (challenge, not direction), there is a mention of a disease
297 that is difficult to diagnose, but there is no mention of a suggested hypothesis or direction. Finally, in
298 the last row (challenge, direction), a medical concern is presented alongside a scientific speculation
299 on the nature of the signaling in the immune system, therefore reflecting both a challenge and a
300 direction.

⁵While many papers discuss future directions in their concluding section, our task involves capturing all mentions of directions/hypotheses/speculations/early indications appearing throughout papers (e.g., in experimental analysis sections).

Labels	Example
Not Challenge, Not Direction	Nowadays, standard structure-based virtual screening has been routinely implemented in drug discovery to quickly prioritize potential compounds for in vitro activity tests.
Not Challenge, Direction	Future studies will focus on comparative sequence analysis between the PST isolates reported herein and global isolates of PST to determine the specific geographic origin(s) for this diverse PST population.
Challenge, Not Direction	Outbreaks attributed to acute BVDV infections in feedlot calves have been described previously, although definitive diagnosis is often difficult [18].
Challenge, Direction	Thus, both PRRs could be responsible for innate immune signaling during acute DENV infection, perhaps operating in temporally distinct fashion as in WNV infection.

Table 2: Examples of Challenges and Directions.

301 A.3 Related Work

302 In addition to biomedical literature discussed in the Introduction, our work on extracting challenges
303 and directions from scientific papers is related to a large body of research.

304 Scientific information extraction and text classification. The goal in this line of work is to extract
305 structured information from literature, such as sentence-level classification into categories including
306 objectives/methods/findings⁶ or extracting entities and relations⁷, [3, 35, 19]. Unlike previous
307 work, our labelling schema encapsulates underexplored facets such as difficulties, awards, uncertainties
308 (challenges) and suggestions, hypotheses, indications that an issue is worthy of additional exploration
309 (directions). Our coarse-grained schema covers diverse variants of challenges/directions and can help
310 generalize across the interdisciplinary COVID-19 literature [14].

311 COVID-19 IE and search tools. Recent work includes visualizing COVID-19 concepts and relations
312 [13], a syntactic search engine⁸, and a search engine for causal and functional relations⁹. Ours
313 system is focused on challenges and directions, not captured by existing tools. Recent works¹⁰
314 used crowd workers to annotate abstracts (not full-texts as in this paper) for Background, Purpose,
315 Method, Finding/Contribution. As discussed in §A.4, we found that crowd workers fail on our task,
316 even though recruited with high quality assurance standards.

317 A.4 Data Collection Procedure

318 We recruited four expert annotators with biomedical and bioNLP backgrounds to annotate sentences
319 sampled across CORD-19. Annotators were given detailed annotation guidelines⁶ and had a one-hour
320 training session for reviewing the guidelines and discussing more examples. The guidelines included
321 simple explanations of challenges and directions along with introductory examples. We sampled
322 sentences from full-text papers, aiming to capture diverse, fine-grained challenges/directions that
323 often do not appear in abstracts. The subset of full-text papers in CORD-19 numbers roughly 180K
324 papers with around 25 million sentences⁷. We also provide surrounding sentences around the target
325 sentence as context.

⁶Annotation guidelines are available in REDACTED FOR ANONYMITY.

⁷We use a snapshot of CORD-19 from 08-02-2021.

Labels	Train	Dev	Test	All
Not Challenge, Not Direction	602	146	745	1493
Not Challenge, Direction	106	25	122	253
Challenge, Not Direction	288	73	382	743
Challenge, Direction	155	40	210	405

Table 3: Distribution of labels across data splits. Splits are stratified with no overlap in papers.

326 Randomly sampling sentences for annotation is highly unlikely to lead to enough challenge/direction
327 cases. To increase this likelihood, two annotators curate 280 keywords or phrases with affinity to one
328 of the two categories. Sentences mentioning at least one keyword (lemmatized) are upsampled. For
329 example, words such as unknown, limit, however provide weak signal indicating a potential mention
330 of a challenge; words like suggest, future work, explore are weak indicators of a direction. To expand
331 the list further, annotators made use of SPIKE [34] which also has a vocabulary explorer that allows
332 browsing keywords similar to an input term. Overall, the 280 keywords covered around a third of
333 sentences in CORD-19, demonstrating their breadth. We note that for most keywords context can
334 completely change their meaning; for instance, “limit” can appear in the context of “we limit the
335 discussion” which has no relation to challenges. Our set of terms with weak correlation to the label
336 (e.g., the word may that very weakly relates to directions) favors high recall rather than precision.

337 Finally, to further increase coverage, we sampled at random roughly a quarter of sentences from
338 the remaining sentences that did not contain any of the keywords, obtaining in total 3000 sentences.
339 We filter sentences that are not in English, mostly numeric/mathematical, or that are very short/long
340 (often due to PDF parsing issues), resulting in 2894 sentences and their surrounding contexts, from
341 1786 papers.

342 Annotator agreement: 60% of the sentences were labeled by all annotators with high average
343 pairwise agreement. Following common practice we measure micro-F1 and macro-F1, treating labels
344 from one annotator as ground-truth and the other as predicted, obtaining 85% for challenges and
345 88% for directions for micro-F1, and 84% and 82% for macro-F1. Positive label proportions are
346 39.66% and 22.74% for challenges/directions, respectively. We create a train/dev/test stratified split
347 of 40%/10%/50% (Table §3), splitting by distinct papers. We opt for a large, diverse test set for model
348 evaluation [4]. The sampled sentences originate from papers published in 1108 journals.

349 A.4.1 A note on crowdsourcing.

350 We also attempted crowdsourcing to scale the collection process. However, despite multiple trials
351 and strict quality assurance, the nuanced nature of the task was found to be difficult for crowd workers,
352 especially due to false negatives.

353 A.5 Baseline models

354 The classification task at hand is a multi-label sentence classification problem, with the goal of
355 predicting whether a sentence mentions a challenge, a research direction, both, or neither. The
356 definitions of the challenge and direction categories are as described in §2. We evaluate a range of
357 baseline models we examine for our novel task.

- 358 • **Keyword-based** A simple heuristic based on the lexicon we curated for data collection (§A.4) —
359 sentences with a challenge keyword are labeled as challenge, and similarly for direction.
- 360 • **Sentiment** Challenge statements potentially have a negative tone, and directions are potentially
361 more positive. We score the sentiment of each sentence using an existing tool and classify
362 negative sentiment sentences as challenges and positive ones as directions.
- 363 • **Zero-shot inference** In zero-shot classification, models predict labels they were not trained on
364 [39]. This could be particularly relevant in emerging domains such as COVID-19, where collecting

⁸Our list of keywords is available in REDACTED FOR ANONYMITY.

⁹Final labels selected by majority vote, with ties (fewer than 100 cases) adjudicated by a member of the research team.

¹⁰Using the Appen platform <https://appen.com/>.

365 large amounts of labeled data could be prohibitive. We use a language model trained for natural
366 language inference (NLI), letting the model infer whether the input ~~is~~ ^{contains} the label name. See
367 Appendix §A.5.1 for full details.

368 • Scientific language models We also experiment with fine-tuning language models that were
369 pre-trained on scientific papers. We report results for PubMedBERT-abstract-full-text¹¹ which
370 was pre-trained on PubMed paper abstracts & full texts, and for SciBERT¹¹ trained on a corpus
371 of biomedical and computer science papers. In addition, we also experiment with a non-scientific
372 language model, RoBERTa-large, which has been shown to obtain excellent results when fine-
373 tuned on scientific texts [1]. We also experimented with other language models, with very similar
374 results.

375 For all language models we fine-tune we use the Hugging Face library¹¹. We use hyperparam-
376 eter tuning with the objective of maximizing the F1-score on the development set using grid search
377 over batch size ([8,16,32]), learning rate ([1e-5, 2e-5, 3e-5, 5e-5]) and epochs (maximal value
378 of 25 epochs). We use the Adam optimizer [20] with a dropout rate of 0.3 for all neural models,
379 using a binary cross-entropy (BCE) loss over our two labels. For the sentiment analysis model, we
380 tune its threshold on the development set.

381 See customized models at Technical Appendix §A.6.

382 A.5.1 Zero-shot baseline and variations specification

383 We use BART-MNLI-large [2], a pre-trained NLI model. We find that simply feeding in “challenge”
384 and “direction” as label names, or similar variants, performs poorly, likely due to the nuanced
385 complexity of these labels. Instead of using one name, we find that enumerating multiple variants of
386 challenges (e.g., difficulty, limitation, failure) provides better results.

387 We define the following sub-labels enumerating challenges and directions. We take the different
388 variants of challenges that we use in our definition of this label — [challenge, problem, difficulty,
389 law, limitation, failure, lack of clarity, gap of knowledge] — and similarly for directions ([direction,
390 suggestion, hypothesis, need for further research, open question, future work]). Denote the former
391 list by L_c , and the latter by L_d . For each category, we compute the probability of each L_c (L_d ,
392 respectively) and take the maximal value for each set of sub-labels, denoted by m_c and m_d . If
393 $m_c \geq 0.9$ we label the sentence as a challenge, and similarly for directions, using the same
394 threshold. Otherwise, the input is classified as negative.

395 We briefly examine a few variations on the zero-shot classification baseline, in terms of the class/label
396 names given as input, to study their effect. We use the same binary threshold of 0.9 for all variants.

- 397 • Class-name: Using only the class names, i.e., “challenge” and “direction”, rather than more
398 fine-grained label names.
- 399 • Template: Using “challenge” and “future direction” as part of a template sentence following the
400 approach in Yin et al. [39]. Specifically, “This sentence is about a challenge”, “This sentence is
401 about a future direction”.
- 402 • Concatenated: Instead of [challenge, problem, difficulty, law, limitation, failure, lack of clarity, gap
403 of knowledge] as standalone inputs, we concatenate them into one string — “challenge, problem,
404 difficulty, law, limitation, failure, lack of clarity, gap of knowledge”; the same was done for
405 directions.

406 Table 4 in Technical Appendix §A.7.2 shows the results of the zero-shot variant models.

407 A.6 Context Modelling Variants

408 We also experiment with models motivated by examination of baseline errors (see Technical Appendix
409 §A.8). Specifically, we find that adding context helps in certain cases: For example, in the sentence
410 “... the patient had an extreme elevation of procalcitonin without signs of bacterial infection which
411 was misclassified as a non-challenge, adding context helped identify the unexplained elevation as
412 problematic. However, context can also introduce noise (see Table 1). We explore different ways
413 in which the context can affect predictions — during training, and during inference. In addition to
414 simply fine-tuning PubMedBERT with full context, we explore two main customized approaches.

¹¹See our code attached in the REDACTED FOR ANONYMITY.

415 Hierarchical Attention Network (HAN) [38] Recall Section §2, where candidate sentences are
416 denoted by X and their surrounding context by $X_{\text{previous}}; X_{\text{next}}$. Denote by X_{context} the concate-
417 nation: [CLS] X_{previous} [SEP] X [SEP] X_{next} [SEP]. We compute a weighted average of [CLS]
418 and the first two [SEP] tokens using attention weights, and use this average embedding for final
419 classification. The weights are learned as part of end-to-end training.¹² While this model can
420 potentially learn to re-weight the context, it encodes the full X_{context} jointly before this weighting
421 takes place, which can lead to noise propagating early on. We thus test a different approach to help
422 mitigate this issue.

423 Context Slice + Combine Let $f_X(x)$ denote the label logits emitted from the final layer of the
424 PubMedBERT model which was fine-tuned on only, for some input text. Likewise, denote
425 by $f_{X_{\text{context}}}(x)$ the logits from PubMedBERT fine-tuned using the full context. At inference
426 time, we obtain outputs using the following variants (“slices”) for each x : (1) $l_1 = f_X(x)$, (2)
427 $l_2 = f_{X_{\text{context}}}(X_{\text{context}})$, (3) $l_3 = f_X(X_{\text{context}})$, and (4) $l_4 = f_{X_{\text{context}}}(x)$. We then average
428 (“combine”) all four, yielding a final pair of logits used for prediction. (1) and (2) are just the models
429 reported in Table 1 – feeding x as input to PubMedBERT fine-tuned on, and similarly for X_{context} .
430 (3) and (4) switch between training and inference inputs: in (3) X_{context} takes X_{context} as input during
431 inference, and in (4) x is fed as input into $f_{X_{\text{context}}}$. The reason we include these is to tease apart
432 different ways in which the context may introduce noise or signal, during training using context (3)
433 and during inference (4). We empirically find all four are in agreement in roughly 70%–83% of the
434 cases for challenges/directions; 3 out of 4 agree in 20%–11%, and the rest are tied. This suggests
435 each variant can potentially capture complementary information.

436 A.7 Results

437 Classification Results. As seen in Table 1, fine-tuned scientific language models outperform the
438 Zero-Shot model, which still does well considering it had no supervision and was pre-trained on
439 non-scientific texts. The sentiment analysis and keyword-based classifiers, both based on large lists
440 of “positive/negative” keywords, have good recall but poor precision. The best individual classifier
441 by F1 is PubMedBERT with a binary-F1 of 0.770 and 0.766 on the challenge and direction labels,
442 respectively. The HAN approach was able to increase recall substantially for problems, but at the
443 cost of reduction in precision, leading to overall inferior F1 on par with PubMedBERT+context.

444 The Slice-Combine approach leads to an improvement of about one F1 point for both labels over the
445 best individual model (standard error $\approx 10^{-4}$). In an ablation experiment we compute the
446 averaged logits of $l_1; l_2$ and $l_3; l_4$ separately and also simply ensemble four model runs of fine-tuned
447 PubMedBERT, both leading to inferior results (see these additional results in Technical Appendix
448 §A.7.1). Finally, an oracle that selects the best logit l_i for each input based on ground truth
449 labels has F1 of 0.907 and 0.896 for challenges/directions, suggesting much room for future work
450 on adaptive use of context during training and inference. See in-depth analysis of additional model
451 errors in Technical Appendix §A.8.

452 Precision@Recall Our primary focus is a novel search engine application (§4). For such applica-
453 tions, it is often more important to have high precision for top retrieved results. We examine precision
454 for a range of values of recall, shown in Figure 2. We observe that for 20% recall we obtain well over
455 90% precision, and for 40% recall about 90% precision.

456 Evaluating predictions across CORD-19 To further ensure quality, we run the PubMedBERT
457 model across all sentences in CORD-19. Out of all sentences indexed in our search engine as either
458 a challenge or a direction, we sample roughly 350 sentences, with higher sampling weight given to
459 high-confidence predictions. About 190 sentences have confidence greater than 0.5, and 130 have
460 confidence lower than 0.5, with 90 sentences with confidence within the range (0.25; 0.75). These
461 sentences are labeled by an expert annotator following the same criteria used to annotate our dataset
462 (§A.4). As shown in Figure 3, we obtain very high mean average precision (MAP) of 98% and area
463 under the precision-recall curve (AUC) of over 97% for directions, and 97% / 96% for challenges.
464 We conclude that for high-confidence challenge and direction sentences indexed in our search engine,
465 accuracy is expected to be overall considerably high. Our test set consists of considerably harder
466 examples, explaining the gap in performance (see discussion in §A.8).

¹²See Yang et al. [38] for details about the general framework.

Figure 2: Precision/Recall results for the PubMedBERT model, and the zero-shot model. Precision for PubMedBERT is high for reasonably large values of recall.

467 Zero-shot generalization to biomedicine and AI domains We perform a preliminary experiment
468 examining whether a model trained on our dataset can, with no additional training, generalize to
469 identify challenges and directions in general biomedical papers, which we sample from S2ORC, a
470 larger corpus with millions of papers [25], and also AI papers sampled from a corpus of full-text
471 computer science papers [6]. In total, we sample about 1000 sentences across the two datasets,
472 following the same procedure as described above for CORD-19 sentences: we randomly sample 100
473 papers that did not appear in the CORD-19 corpus to ensure no leakage of information from our
474 training set (we filter with a paper identifier shared by both resources). From these papers, we sample
475 sentences in the same way as above for annotation. The annotator labels 630 sentences: 430 sentences
476 with confidence scores greater than 0.5 and 200 sentences with scores lower than 0.5. 150 of those
477 with scores lower than 0.5. For AI papers, we follow the same procedure, with 300 sampled sentences.
478 CORD-19 papers are highly interdisciplinary in terms of the areas they cover [3], raising the
479 possibility of using our dataset to train models that can be applied to other domains without additional
480 costly data collection. As seen in Figure 3, identification accuracy is high in this sample too. For
481 directions, we obtain MAP and AUC of around 96% for biomedicine, and around 95% for AI. For
482 challenges, MAP and AUC reach around 97-98% for biomedicine and around 96% for AI. These
483 preliminary results focus mostly on high-confidence predictions relevant for search applications, and
484 generalization could be explored more extensively in future work.

485 A.7.1 Additional Slice-Combine Context Models Results

486 In addition to the experiments reported, we tested multiple ways to combine information from four
487 variants: (i) apply average or a median on the logits, (ii) majority voting, (iii) log-odds extremization,
488 (iv) training a router model based on the logit differences, (v) running logistic regression with the
489 embedding (final layer) of each of the four input encoders and their logits as features. Aside from
490 the simple averaging, logistic regression was a close runner-up. We explore the average weights the

Figure 3: Evaluating predictions beyond our test set. We use a model trained on our data to identify challenges and directions across COVID-19 (denoted COVID), S2ORC (general biomedical papers, denoted general biomed) and SciRex (full-text AI papers, denoted AI). Accuracy is considerably high. Zero-shot generalization over non-COVID papers, even non-biomedical papers, is encouragingly high, indicating the utility of our resource beyond COVID-19.

Variation	Challenge			Direction		
	P	R	F1	P	R	F1
NLI-Zeroshot	0.659	0.693	0.675	0.401	0.825	0.54
Class-name	0.789	0.440	0.565	0.618	0.065	0.119
Template	0.439	0.941	0.599	0.589	0.401	0.478
Concatenated	0.849	0.107	0.190	0.491	0.725	0.585

Table 4: Zero-Shot Baseline Ablation Results. We provide the baseline different variants of class descriptors for challenges and directions, respectively.

491 logistic regression assigned to the four context variants. For challenges 0.14, 0.21, 0.21, and 0.35
 492 for (1)-(4) respectively; and for directions 0.32, 0.2, 0.07 and 0.45. Interestingly this suggests that
 493 training the model end-to-end with context could be useful even when the context is not available at
 494 inference.

495 Sanity testing the Slice-Combine Context Models As a sanity check we simply ensemble 4 runs
 496 of PubMedBERT, resulting in inferior F1 of 0.772 and 0.764 for challenges/directions, which further
 497 indicates the complimentary value of our four context variants beyond simpler ensembling.

498 A.7.2 Zero-Shot Baseline Ablation Results

499 Results for the ablation experiments are in Table 4. As can be observed, for challenges, the variant
 500 reported in Table 1 achieves the best results by a large margin. For directions, the variant that
 501 concatenates class descriptors into one string does marginally better.

502 A.8 Error analysis

503 We study the cases where the best fine tuned model failed to classify sentences correctly. In order to
504 do so we randomly sampled and analyzed roughly 20% of the false positive and false negative errors
505 across both labels.

506 **Challenges** The most common error that accounts for a third of wrong predictions (both false
507 positive and negative) is that in some cases deciding whether an outcome is positive or negative
508 requires a more profound understanding of the biomedical entities involved and of the context. For
509 example, the sentence consider the sentence “*The surprising conclusion of the study was that relative*
510 *to primary rat Schwann cells undergoing myelination, only 2 cell lines expressed high levels of*
511 *mRNA coding for myelin proteins and none of the cell lines expressed all of the myelin proteins*
512 *typically expressed in myelinating Schwann cells.*” The model classifies the sentence as a challenge.
513 However, an expert who read the text concluded that the outcome is non-problematic since with
514 further downstream analysis the mentioned conclusion may represent a more accurate model for
515 future analysis of myelin gene promoters. Conversely, the sentence “*It is remarkable that the patient*
516 *had an extreme elevation of procalcitonin without signs of bacterial infection.*”, was not classified
517 as a challenge, but an expert annotator did identify the issue of having a strong biological indicator
518 for a serious condition without a clear explanation for its elevation as problematic. We note that in
519 multiple cases presenting the model with the context aids with these issues. For instance, in the above
520 example, when providing the context which includes a reference to the the risk factor, the prediction
521 flips to the right call. However, as discussed in §A.5, context can also add noise in some cases. This
522 observation led to our Context Slice+Combine approach described in §A.6.

523 The second biggest cause for false positives is sentences that provide a general description of a
524 condition rather than a challenge. An example is “*The location of the headache might vary depending*
525 *on which sinuses are affected[...]*”. Such texts are tricky since they are essentially facts about a
526 condition rather than a description of an explicit challenge (e.g., the headache may be trivial to treat).
527 To make the distinction clearer, consider the text “*Colitis is a chronic digestive disease*”. It presents
528 a definition of a disease, and not an instance of an explicit problem one needs to address.

529 The second biggest cause we observe for false negatives is sentences that mention *partial* solutions
530 that can mitigate a problem. For instance, “[...] *the cellular apoptotic process is immediately*
531 *triggered as an innate defense mechanism in response to infection, but is abruptly suppressed during*
532 *the middle stage of infection.*”. In this example a defense mechanism is mentioned that can mitigate
533 the problem, but the challenge nonetheless remains. Combined, the above error causes account for
534 roughly 2/3 of the false positives and negatives. Labelers were provided input on how to deal with
535 these issues in annotation, but since they are nuanced, models may require more data or sophistication
536 to classify them correctly.

537 **Directions** The most common cause for error in Direction is the identification of future action items
538 which are general vague suggestions as directions. For example, “*In agreement with the authors*
539 *of that study, we believe that communication between laboratory specialists and clinicians should*
540 *be intensified and improved.*”. The example suggests a policy or action that does not constitute a
541 research direction or hypothesis. This error accounts for 60% of the false positives.

542 In terms of recall errors we find that sentences that suggest a hypothesis or other more implicit
543 research directions account for roughly 50% of the errors. For example “*Persistent viral shedding*
544 *may indicate different levels of virulence, host immune response and infectiousness*”. The guidelines
545 stipulate that these should be positive since researchers need to verify these directions, and indeed in
546 most cases they are correctly identified, but some instances cause false negatives.

547 The rest of the errors in challenges or directions were anecdotal rather than systemic.

548 **A note on expected errors in the downstream tasks.** Our data contains an over-sampled propor-
549 tion of tricky keywords (e.g., the word “hard” appears in both “hard material” and “hard task”), and
550 thus we expect fewer errors in the search application task (see also Figure 3). In addition, in our
551 downstream task of search we rank results by prediction confidence, and precision for high values of
552 confidence is high even on our harder test set (Figure 2). Indeed, Figure 3 shows that predictions
553 appear to have high overall accuracy in a sampled set.

554 A.9 User Studies

555 We now explore user studies designed to evaluate our framework’s utility. First, we explore whether
556 our system can be helpful for quick discovery of challenges/directions. Second, we conduct a study
557 with nine medical researchers working on COVID-19 treatment and research. In total our studies
558 include 19 researchers and over 70 distinct search queries.

559 A.9.1 Search Engine

560 **Challenge and direction indexing** We build a search engine that indexes challenges and directions
561 across the entire COVID-19 corpus up to and including August 2021. To build the search engine
562 (see Figure 5 in the Technical Appendix for a screen capture), we first apply PubMedBERT to 550K
563 papers, totalling 29M sentences. 180K of the papers are with full text, the rest are abstracts. We
564 then clean poorly tokenized sentences, non-English sentences, very short sentences or texts with
565 latex code. We classify the remaining sentences leaving 2.2M sentences — about 950K sentences
566 with high-confidence predictions for at least one of challenge/direction and their surrounding context
567 sentences. We select high-confidence sentences by using a threshold of 0.99 for both challenges and
568 directions, using a thresholds leading to well over 90% precision at top-10% on our test set.

569 **Entity-based indexing** For each sentence in our set of 2.2M, we add another layer of indexing,
570 by extracting entities and linking them to knowledge base entries. This allows us to partially group
571 together all challenges or directions into “topics” referring to a specific fine-grained combination
572 of concepts (e.g., *AI + diagnosis + pneumonia*), and facilitate entity-centric faceted search which
573 is known to be useful in scientific exploratory search [13, 14].¹³ We extract a range of biomedical
574 entities and link them to a biomedical KB of MeSH (Medical Subject Headings) entities [24]. See
575 Technical Appendix §A.10 for full technical details.

576 In the experiments that follow, we compare our system with a strong real-world system — PubMed
577 biomedical search engine¹⁴, a leading search site that clinicians and researchers regularly peruse as
578 their go-to tool. While PubMed was not designed to find challenges and directions, no existing tool is;
579 PubMed allows users to search for entities such as MeSH terms, is supported by a KB of biomedical
580 entities used for automatic query expansion, and has many other functions — and as such is a strong
581 real-world baseline.

582 A.9.2 Challenge/Direction Exploration

583 We recruited ten participants with education and experience in medicine, microbiology, public health,
584 molecular, cellular, and developmental biology, biochemistry, chemical & biological engineering,
585 environmental science, and mathematics. Participants are paid \$50 per hour of work, comparing query
586 results from our system and PubMed. Participants were given guidelines, which include definitions
587 for research challenges and directions with simple examples.¹⁵

588 Each participant was given twenty queries, split into two sections for challenges and directions,
589 respectively. For each query, participants were asked to find as many research challenges as possible
590 in no more than 3 minutes. The total number of unique queries among the participants is 65. Some
591 examples of queries used for the challenges section include “antibodies” and “inflammation, lung”,
592 with the paired entities being searched jointly; example queries for the directions section include
593 “telemedicine” and “vaccines, technology”. All queries were curated by a domain expert.

594 As seen in Figure 4, our system yielded a greater number of challenges and directions, on average,
595 than the PubMed tool. Users found roughly 4.46 challenges and 6.43 directions per query using our
596 system compared to the 2.24 challenges and 2.03 directions per query found using PubMed (p-value
597 of .00192 for challenges and .000529 for directions using a paired t-test). For each participant we
598 included 5 challenges and 5 directions that were overlapping across all participants, in order to control
599 and compare between results for the same queries. We find that on average across users, 70.0%
600 of the query results using our system led to a strictly larger number of challenges discovered than

¹³Other forms of challenge grouping, such as with embedding-based sentence clustering, are interesting to explore in future work.

¹⁴<https://pubmed.ncbi.nlm.nih.gov/>

¹⁵Full annotation guidelines are included in REDACTED FOR ANONYMITY.

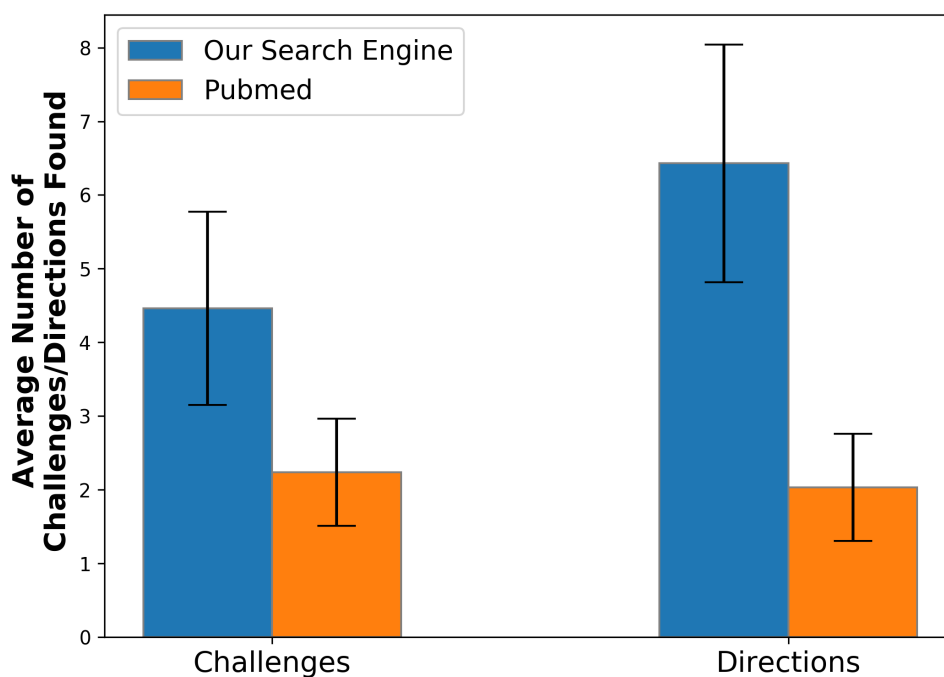


Figure 4: Study with researchers with diverse backgrounds. Participants using our search engine were able to find substantially more cases of challenges and research directions they considered useful than with PubMed. Error bars represent 90% confidence intervals.

Metric	Chal./Dir. Search	PubMed
Search	90%	48%
Utility	94%	57%
Interface	91%	68%
Overall	92%	59%

Table 5: Nine medical researchers expressed much higher satisfaction with our system (Chal./Dir.) than PubMed.

601 the respective query results using PubMed, and 22% were ties. For directions, we find a larger gap
 602 between the two systems, with 96.0% of the query results using our system yielding strictly more
 603 directions than PubMed, and 2% yielding ties.

604 A.9.3 Evaluation with Medical Researchers

605 We now report on an evaluation of our search engine performed with nine medical researchers at a
 606 large hospital.¹⁶

607 **Study.** We recruited nine expert MDs with a wide range of specialization including cardiology,
 608 pulmonary and critical care medicine, gastroenterology and general medicine who are actively
 609 involved in clinical research both for COVID-19 and specialty areas, and each have over 1000
 610 citations. Each expert completed randomly ordered search tasks (challenge/direction queries curated
 611 by an expert medical researcher; see Appendix §A.12) using both PubMed and our system. Experts
 612 using our UI viewed sentences and their contexts (previous/next sentences). In addition we also
 613 displayed metadata such as paper title, date, url. After all search tasks were completed for both
 614 systems, experts were given seven-point Likert-scale questions to judge system utility, interface, and

¹⁶In addition to the motivation discussed in the Introduction, see §Appendix A.13 for a more detailed example scenario where medical researchers need to search for challenges and directions.

615 search quality. Following [14], we use a standardized Post Study System Usability Questionnaire
616 (PSSUQ) [21], widely used in system quality research, and added questions designed to evaluate
617 search and exploration utility: *overall search accuracy, results that are not only relevant but interesting*
618 *or new, finding papers interesting to read, and ability to understand and judge each individual result*
619 *quickly without additional context.* Each question is asked twice, once for PubMed and once for our
620 system, leading to $15 \times 2 \times 6 = 180$ responses.

621 **Results.** Table 5 shows the average Likert scores (normalized to [0%,100%]) across all questions and
622 users for our system and PubMed. We group questions by three types for brevity. The results show
623 that the medical experts strongly prefer our search engine to PubMed (overall average of 92% vs.
624 59%, with non-normalized scores of 6.42 vs. 4.14). On average across all questions, the majority of
625 the nine MDs assigned our system a higher score than PubMed, at an average rate of 85% per question.
626 When considering ties, the average rate is 92%. We found that our system significantly outperformed
627 PubMed across all questions (Wilcoxon signed rank test p-value is significant at 5.409×10^{-6}). These
628 results further resonate in light of the experts' strong familiarity with PubMed and the bare-bones
629 nature of our UI.

630 **A.10 Entity-based Indexing**

631 We employ the SciSpacy library [29] to extract entities using five different NER models: one trained
632 on MedMentions [28] (a dataset with general mentions of UMLS [3] entities covering a wide range
633 of concepts), and four trained on more specialized sources (CRAFT [1], JNLPBA [18], BC5CDR
634 [23], BIONLP13CG [19]). Each entity is then automatically linked to a biomedical KB of MeSH
635 (Medical Subject Headings) entities [24] using SciSpacy's entity linking functionality that performs
636 character-trigram matching on MeSH entity names and aliases. We filter for high-confidence linked
637 entities,¹⁷ and for entities that appeared in at least 10 sentences, then selecting the top 30K unique
638 entities to be indexed by our search engine. At search time, we match user queries to MeSH aliases
639 with an autocomplete dropdown for users to select from as they type. After one entity is selected, the
640 user can search for more from a narrower list of entities that co-occur with it.

641 **A.11 Search UI**

642 Figure 5 shows a screen capture of our search user interface.

643 **A.12 Expert queries**

644 Below are examples for our expert queries:

- 645 • **Find problems/limitations/flaws** related to COVID-19 and each of (1) *hospital infections*, (2)
646 *diagnosis*, (3) *vaccines for children*, (4) *probiotics and the gastrointestinal tract*.
- 647 • **Find directions/hypotheses/potential indications** related to COVID-19 and each of (1) *mechanical*
648 *ventilators*, (2) *liver*, (3) *artificial intelligence*, (4) *drug repositioning*.

649 **A.13 Medical Search Scenario**

650 In addition to the motivation discussed in the Introduction, we briefly discuss in more detail how a
651 search engine for challenges and directions could help medical researchers when conducting literature
652 reviews. Many research ideas come to MDs with a challenge they perceive during clinical care. If
653 they are unable to find a solution to the problem based on prior experience, they then search for the
654 available scientific literature for possible guidance. When no sufficient guidance is found, research
655 projects are often commenced, starting with literature search which often involves understanding
656 and mapping out associated challenges and directions to help with formulating research questions.
657 Physicians and trainees still spend a significant amount of time doing this form of literature search
658 for fine-tuning their research question. If such a process can be simplified with automation, it could
659 potentially cut down the time and effort needed to formulate and narrow down research questions.

¹⁷Using a threshold of 0.9.

Figure 5: Screen capture of our search user interface.