# Rapidly Deploying a Neural Search Engine for the COVID-19 Open Research Dataset

**Edwin Zhang,**[1] **Nikhil Gupta,**[1] **Rodrigo Nogueira,**[1] **Kyunghyun Cho,**[2,3,4] and **Jimmy Lin**[1]

[1] David R. Cheriton School of Computer Science, University of Waterloo
[2] Courant Institute of Mathematical Sciences, New York University
[3] Center for Data Science, New York University   [4] CIFAR Associate Fellow

This extended abstract represents an abridged version of Zhang et al. (2020a), posted on arXiv April 10, 2020 and concurrently submitted to this workshop. We have intentionally decided for this short piece to reflect the state of our work at that time. The latest updates on our project can be found in Zhang et al. (2020b).

The Neural Covidex is a search engine that exploits the latest neural ranking architectures to provide information access to the COVID-19 Open Research Dataset (CORD-19) curated by the Allen Institute for AI (Wang et al., 2020). It exists as part of a suite of tools we have developed to help domain experts tackle the ongoing global pandemic. We hope that improved information access capabilities to the scientific literature can inform evidence-based decision making and insight generation.

The first version of CORD-19 was released on March 13, 2020. Within a couple of weeks, our team was able to build, deploy, and share with the research community a number of open-source components that support information access to this corpus. These include: Extensions to our Anserini IR toolkit (Yang et al., 2018) and its Pyserini Python interface (Akkalyoncu Yilmaz et al., 2020) to support basic keyword search capabilities on the corpus; PyGaggle, a new library for neural text ranking that includes supervised ranking models based on T5 as well as unsupervised sentence highlighting models with BioBERT (Lee et al., 2020).

We have assembled these components into the Neural Covidex, available online at `covidex.ai`; see screenshot in Figure 1. This user interface was developed from scratch and is itself open source. Zhang et al. (2020a) described our initial efforts and shared lessons we learned along the way.

Although the application of BERT to text ranking is well known (Nogueira and Cho, 2019), we decided to deploy our latest research based on sequence-to-sequence models (Nogueira et al.,
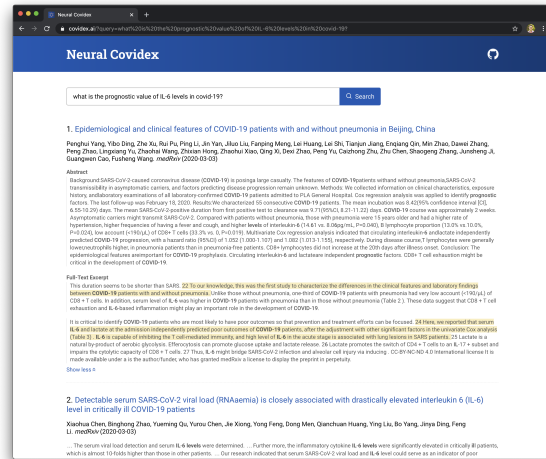


Figure 1: Screenshot of the Neural Covidex.

2020), specifically T5 (Raffel et al., 2019). This relevance classifier, which reranks BM25 results from Anserini, is fed a query $q$ and each candidate document $d$ in turn. The model is fine-tuned to produce either "true" or "false" depending on whether the document is relevant or not to the query. At inference time, we softmax the logits of the "true" and "false" tokens, and the resulting probability of the "true" token is used as the relevance score of $d$. Candidate documents are then reranked using their relevance scores. As there is no COVID-19 training data, we fine-tuned our model on the MS MARCO passage dataset (Nguyen et al., 2016), and thus our reranker operates in a zero-shot setting.

As our work pre-dated any systematic evaluation efforts by the community, at the time of submission we were unable to provide any experimental results. Since then, however, we have participated in the TREC-COVID challenge (Voorhees et al., 2020); partial results to date are reported in Zhang et al. (2020b). Nevertheless, the speed at which we were able to build and deploy the Neural Covidex is a testament to the power of open-source software and modern open-science norms.

## Acknowledgments

## References

Zeynep Akkalyoncu Yilmaz, Charles L. A. Clarke, and Jimmy Lin. 2020. A lightweight environment for learning experimental IR research practices. In *Proceedings of the 43rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2020)*, pages 2113–2116.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: a pretrained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A human generated machine reading comprehension dataset. *arXiv:1611.09268*.

Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage re-ranking with BERT. *arXiv:1901.04085*.

Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. 2020. Document ranking with a pretrained sequence-to-sequence model. *arXiv:2003.06713*.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv:1910.10683*.

Ellen Voorhees, Tasmeer Alam, Steven Bedrick, Dina Demner-Fushman, William R. Hersh, Kyle Lo, Kirk Roberts, Ian Soboroff, and Lucy Lu Wang. 2020. TREC-COVID: Constructing a pandemic information retrieval test collection. *SIGIR Forum*, 54(1).

Lucy Lu Wang, Kyle Lo, Yoganand Chandrasekhar, Russell Reas, Jiangjiang Yang, Doug Burdick, Darrin Eide, Kathryn Funk, Yannis Katsis, Rodney Kinney, Yunyao Li, Ziyang Liu, William Merrill, Paul Mooney, Dewey Murdick, Devvret Rishi, Jerry Sheehan, Zhihong Shen, Brandon Stilson, Alex Wade, Kuansan Wang, Nancy Xin Ru Wang, Chris Wilhelm, Boya Xie, Douglas Raymond, Daniel S. Weld, Oren Etzioni, and Sebastian Kohlmeier. 2020. CORD-19: The COVID-19 Open Research Dataset. *arXiv:2004.10706*.

Peilin Yang, Hui Fang, and Jimmy Lin. 2018. Anserini: reproducible ranking baselines using Lucene. *Journal of Data and Information Quality*, 10(4):Article 16.

Edwin Zhang, Nikhil Gupta, Rodrigo Nogueira, Kyunghyun Cho, and Jimmy Lin. 2020a. Rapidly deploying a neural search engine for the COVID-19 Open Research Dataset: Preliminary thoughts and lessons learned. *arXiv:2004.05125*.

Edwin Zhang, Nikhil Gupta, Raphael Tang, Xiao Han, Ronak Pradeep, Kuang Lu, Yue Zhang, Rodrigo Nogueira, Kyunghyun Cho, Hui Fang, and Jimmy Lin. 2020b. Covidex: Neural ranking models and keyword search infrastructure for the COVID-19 Open Research Dataset. *arXiv:2007.07846*.