

Pretrained Transformers for Text Ranking: BERT and Beyond

Andrew Yates,¹ Rodrigo Nogueira,² and Jimmy Lin²

¹ Max Planck Institute for Informatics, Germany

² David R. Cheriton School of Computer Science, University of Waterloo, Canada

ABSTRACT

The goal of text ranking is to generate an ordered list of texts retrieved from a corpus in response to a query. Although the most common formulation of text ranking is search, instances of the task can also be found in many natural language processing applications. This tutorial, based on a forthcoming book, provides an overview of text ranking with neural network architectures known as transformers, of which BERT is the best-known example. The combination of transformers and self-supervised pretraining has, without exaggeration, revolutionized the fields of natural language processing (NLP), information retrieval (IR), and beyond. We provide a synthesis of existing work as a single point of entry for both researchers and practitioners. Our coverage is grouped into two categories: transformer models that perform reranking in multi-stage ranking architectures and learned dense representations that perform ranking directly. Two themes pervade our treatment: techniques for handling long documents and techniques for addressing the tradeoff between effectiveness (result quality) and efficiency (query latency). Although transformer architectures and pretraining techniques are recent innovations, many aspects of their application are well understood. Nevertheless, there remain many open research questions, and thus in addition to laying out the foundations of pretrained transformers for text ranking, we also attempt to prognosticate the future.

CCS CONCEPTS

• **Information systems** → **Retrieval models and ranking.**

KEYWORDS

Multi-Stage Ranking; Learned Dense Representations

ACM Reference Format:

Andrew Yates, Rodrigo Nogueira, and Jimmy Lin. 2021. Pretrained Transformers for Text Ranking: BERT and Beyond. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*, July 11–15, 2021, Virtual Event, Canada. ACM, New York, NY, USA, 3 pages. <https://doi.org/10.1145/3404835.3462812>

1 OVERVIEW

The goal of text ranking is to generate an ordered list of texts retrieved from a corpus in response to a query for a particular task. Although the most common formulation of text ranking is

search, instances of the task can also be found in many text processing applications. This tutorial provides an overview of text ranking with neural network architectures known as transformers, of which BERT (Bidirectional Encoder Representations from Transformers) [7] is the best-known example. These models produce high quality results across many domains, tasks, and settings.

This tutorial, which is based on the preprint [21] of a forthcoming book to be published by Morgan and Claypool under the Synthesis Lectures on Human Language Technologies series, provides an overview of existing work as a single point of entry for practitioners who wish to deploy transformers for text ranking in real-world applications and researchers who wish to pursue work in this area. We cover a wide range of techniques, grouped into two categories: transformer models that perform reranking in multi-stage ranking architectures and learned dense representations that perform ranking directly. In a hands-on session we demonstrate how open-source toolkits can be used to rank documents with a variety of these approaches.

Multi-Stage Ranking Architectures. The most straightforward application of transformers to text ranking is to convert the task into a text classification problem, and then sort the texts to be ranked based on the probability that each item belongs to the relevant class. The first application of BERT to text ranking, by Nogueira and Cho [30], used BERT in exactly this manner. This *relevance classification* approach is usually deployed in a module that reranks candidate texts from an initial keyword search engine.

One key limitation of BERT is its inability to handle long input sequences and hence difficulty in ranking texts beyond a certain length (e.g., “full-length” documents such as news articles). This limitation is addressed by a number of models [1, 4, 20, 25, 30, 39], and a simple retrieve-then-rerank approach can be elaborated into a multi-stage architecture with reranker pipelines [26, 33, 37] that balance effectiveness and efficiency. On top of multi-stage ranking architectures, researchers have proposed additional innovations, including document expansion [32, 34] and term importance prediction [3, 5].

A natural question that arises is, “What’s beyond BERT?” We describe efforts to build ranking models that are faster (i.e., lower inference latency), that are better (i.e., higher ranking effectiveness), or that manifest interesting tradeoffs between effectiveness and efficiency. These include ranking models that leverage BERT variants [20], exploit knowledge distillation to train more compact student models [9], and other transformer architectures, including ground-up redesign efforts [14, 28] and adapting pretrained sequence-to-sequence models [8, 31]. These discussions set up a natural transition to ranking based on dense learned representations, the other main category of approaches we cover.

Learned Dense Representations. Arguably, the single biggest benefit brought about by modern deep learning techniques to text



This work is licensed under a Creative Commons Attribution-ShareAlike International 4.0 License.

SIGIR '21, July 11–15, 2021, Virtual Event, Canada.

© 2021 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-8037-9/21/07.

<https://doi.org/10.1145/3404835.3462812>

ranking is the move away from sparse signals, mostly limited to exact matches, to dense representations that are able to capture semantic matches to better model relevance. The potential of continuous dense representations for natural language analysis was first demonstrated nearly a decade ago with word embeddings on word analogy tasks [27]. As soon as researchers tried to build representations for any larger spans of text: phrases, sentences, paragraphs, and documents, the same issues that arise in text ranking come into focus. In fact, ranking with dense representations predates BERT by many years [6, 11, 15, 29, 38, 41].

In the context of transformers, the general setup of ranking with dense representations involves learning transformer-based encoders that convert queries and texts into dense, fixed-size vectors. In the simplest approach, ranking becomes the problem of approximate nearest neighbor (ANN) search based on some simple metric such as cosine similarity [10, 12, 13, 17, 19, 22–24, 35, 36, 40, 42]. However, recognizing that accurate ranking cannot be captured via simple metrics, researchers have explored using more complex machinery to compare dense representations [16, 18]. Here, as with multi-stage ranking architectures, limitations on text length and effectiveness–efficiency tradeoffs are important considerations. It becomes increasingly difficult to accurately capture the semantics of longer texts with fixed-sized representations, and increasingly complex comparison architectures increase latency and may necessitate reranking designs.

2 LOOKING AHEAD

Learned dense representations complement sparse (bag-of-words) term-based representations central to keyword search techniques that have dominated the landscape for more than half a century. Together, hybrid multi-stage approaches (e.g., combining both ranking and reranking) present a promising future direction.

Despite the excitement in directly ranking with dense learned representations, we anticipate that reranking transformers will remain important in the future. At a high level, there are three current approaches: *apply* existing transformer models with minimal modifications, *adapt* existing transformer models, perhaps adding additional architectural elements, and *redesign* transformer-based architectures from scratch. Which approach will prove to be most effective? The jury’s still out.

Related, in NLP we see that the GPT family [2] continues to push the frontier of larger models, more compute, and more data. For text ranking, is the simple answer to build bigger models? Probably not, since ranking has important differences with many traditional NLP tasks. But if not, what are the evolving roles of zero-shot learning, transfer learning, domain adaptation, and task-specific fine-tuning? This remains an interesting open research question.

While there are aspects of text ranking with pretrained transformers that are well understood, many promising directions await further exploration. Looking ahead, we anticipate many more exciting developments!

ACKNOWLEDGMENTS

This work was supported in part by the Canada First Research Excellence Fund and the Natural Sciences and Engineering Research Council (NSERC) of Canada.

REFERENCES

- [1] Zeynep Akkalyoncu Yilmaz, Wei Yang, Haotian Zhang, and Jimmy Lin. 2019. Cross-Domain Modeling of Sentence-Level Evidence for Document Retrieval. 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)*. Hong Kong, China, 3490–3496.
- [2] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. *arXiv:2005.14165* (2020).
- [3] Zhuyun Dai and Jamie Callan. 2019. Context-Aware Sentence/Passage Term Importance Estimation For First Stage Retrieval. *arXiv:1910.10687* (2019).
- [4] Zhuyun Dai and Jamie Callan. 2019. Deeper Text Understanding for IR with Contextual Neural Language Modeling. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2019)*. Paris, France, 985–988.
- [5] Zhuyun Dai and Jamie Callan. 2020. Context-Aware Document Term Weighting for Ad-Hoc Search. In *Proceedings of The Web Conference 2020 (WWW '20)*. 1897–1907.
- [6] Cedric De Boom, Steven Van Canneyt, Thomas Demeester, and Bart Dhoedt. 1999. Representation Learning for Very Short Texts Using Weighted Word Embedding Aggregation. *Pattern Recognition Letters* 80, C (1999), 150–156.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota, 4171–4186.
- [8] Cicero Nogueira dos Santos, Xiaofei Ma, Ramesh Nallapati, Zhiheng Huang, and Bing Xiang. 2020. Beyond [CLS] through Ranking by Generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 1722–1727.
- [9] Luyu Gao, Zhuyun Dai, and Jamie Callan. 2020. Understanding BERT Rankers Under Distillation. In *Proceedings of the 2020 ACM SIGIR International Conference on Theory of Information Retrieval (ICTIR 2020)*. 149–152.
- [10] Luyu Gao, Zhuyun Dai, Zhen Fan, and Jamie Callan. 2020. Complementing Lexical Retrieval with Semantic Residual Embedding. *arXiv:2004.13969* (2020).
- [11] Matthew Henderson, Rami Al-Rfou, Brian Strope, Yun hsuan Sung, Laszlo Lukacs, Ruiqi Guo, Sanjiv Kumar, Balint Miklos, and Ray Kurzweil. 2017. Efficient Natural Language Response Suggestion for Smart Reply. *arXiv:1705.00652* (2017).
- [12] Sebastian Hofstätter, Sophia Althammer, Michael Schröder, Mete Sertkan, and Allan Hanbury. 2020. Improving Efficient Neural Ranking Models with Cross-Architecture Knowledge Distillation. *arXiv:2010.02666* (2020).
- [13] Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently Teaching an Effective Dense Retriever with Balanced Topic Aware Sampling. In *Proceedings of the 44th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2021)*.
- [14] Sebastian Hofstätter, Markus Zlabinger, and Allan Hanbury. 2020. Interpretable & Time-Budget-Constrained Contextualization for Re-Ranking. In *Proceedings of the 24th European Conference on Artificial Intelligence (ECAI 2020)*. Santiago de Compostela, Spain.
- [15] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning Deep Structured Semantic Models for Web Search using Clickthrough Data. In *Proceedings of 22nd International Conference on Information and Knowledge Management (CIKM 2013)*. San Francisco, California, 2333–2338.
- [16] Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux, and Jason Weston. 2020. Poly-encoders: Architectures and Pre-training Strategies for Fast and Accurate Multi-sentence Scoring. In *Proceedings of the 8th International Conference on Learning Representations (ICLR 2020)*.
- [17] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 6769–6781.
- [18] Omar Khattab and Matei Zaharia. 2020. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2020)*. 39–48.
- [19] Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent Retrieval for Weakly Supervised Open Domain Question Answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy, 6086–6096.
- [20] Canjia Li, Andrew Yates, Sean MacAvaney, Ben He, and Yingfei Sun. 2020. PARADE: Passage Representation Aggregation for Document Reranking. *arXiv:2008.09093* (2020).

- [21] Jimmy Lin, Rodrigo Nogueira, and Andrew Yates. 2020. Pretrained Transformers for Text Ranking: BERT and Beyond. *arXiv:2010.06467* (2020).
- [22] Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. 2020. Distilling Dense Representations for Ranking using Tightly-Coupled Teachers. *arXiv:2010.11386* (2020).
- [23] Wenhao Lu, Jian Jiao, and Ruofei Zhang. 2020. TwinBERT: Distilling Knowledge to Twin-Structured BERT Models for Efficient Retrieval. *arXiv:2002.06275* (2020).
- [24] Sean MacAvaney, Franco Maria Nardini, Raffaele Perego, Nicola Tonello, Nazli Goharian, and Ophir Frieder. 2020. Expansion via Prediction of Importance with Contextualization. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2020)*. 1573–1576.
- [25] Sean MacAvaney, Andrew Yates, Arman Cohan, and Nazli Goharian. 2019. CEDR: Contextualized Embeddings for Document Ranking. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2019)*. Paris, France, 1101–1104.
- [26] Yoshitomo Matsubara, Thuy Vu, and Alessandro Moschitti. 2020. Reranking for Efficient Transformer-Based Answer Selection. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2020)*. 1577–1580.
- [27] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S. Corrado, and Jeff Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems 26 (NIPS 2013)*. Lake Tahoe, California, 3111–3119.
- [28] Bhaskar Mitra, Sebastian Hofstatter, Hamed Zamani, and Nick Craswell. 2020. Conformer-Kernel with Query Term Independence for Document Retrieval. *arXiv:2007.10434* (2020).
- [29] Bhaskar Mitra, Eric Nalisnick, Nick Craswell, and Rich Caruana. 2016. A Dual Embedding Space Model for Document Ranking. *arXiv:1602.01137v1* (2016).
- [30] Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage Re-ranking with BERT. *arXiv:1901.04085* (2019).
- [31] Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. 2020. Document Ranking with a Pretrained Sequence-to-Sequence Model. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. 708–718.
- [32] Rodrigo Nogueira and Jimmy Lin. 2019. From doc2query to docTTTTTquery. (2019).
- [33] Rodrigo Nogueira, Wei Yang, Kyunghyun Cho, and Jimmy Lin. 2019. Multi-Stage Document Ranking with BERT. In *arXiv:1910.14424*.
- [34] Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019. Document Expansion by Query Prediction. In *arXiv:1904.08375*.
- [35] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2020. RocketQA: An Optimized Training Approach to Dense Passage Retrieval for Open-Domain Question Answering. *arXiv:2010.08191* (2020).
- [36] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. 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)*. Hong Kong, China, 3982–3992.
- [37] Luca Soldaini and Alessandro Moschitti. 2020. The Cascade Transformer: an Application for Efficient Answer Sentence Selection. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 5697–5708.
- [38] Ledell Wu, Adam Fisch, Sumit Chopra, Keith Adams, Antoine Bordes, and Jason Weston. 2018. StarSpace: Embed All The Things!. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI 2018)*.
- [39] Zhijing Wu, Jiaxin Mao, Yiqun Liu, Jingtao Zhan, Yukun Zheng, Min Zhang, and Shaoping Ma. 2020. Leveraging Passage-Level Cumulative Gain for Document Ranking. In *Proceedings of The Web Conference 2020 (WWW '20)*. 2421–2431.
- [40] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval. *arXiv:2007.00808* (2020).
- [41] Hamed Zamani, Mostafa Dehghani, W. Bruce Croft, Erik Learned-Miller, and Jaap Kamps. 2018. From Neural Re-Ranking to Neural Ranking: Learning a Sparse Representation for Inverted Indexing. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management (CIKM 2018)*. Torino, Italy, 497–506.
- [42] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Min Zhang, and Shaoping Ma. 2020. RepBERT: Contextualized Text Embeddings for First-Stage Retrieval. *arXiv:2006.15498* (2020).