# Pretrained Transformers for Text Ranking

## BERT and Beyond

# Synthesis Lectures on Human Language Technologies

## Editor
**Graeme Hirst,** *University of Toronto*

Synthesis Lectures on Human Language Technologies is edited by Graeme Hirst of the University of Toronto. The series consists of 50- to 150-page monographs on topics relating to natural language processing, computational linguistics, information retrieval, and spoken language understanding. Emphasis is on important new techniques, on new applications, and on topics that combine two or more HLT subfields.

# Pretrained Transformers for Text Ranking

## BERT and Beyond

Jimmy Lin
University of Waterloo

Rodrigo Nogueira
University of Waterloo

Andrew Yates
University of Amsterdam and Max Planck Institute for Informatics

## 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 (NLP) applications. This book provides an overview of text ranking with neural network architectures known as transformers, of which BERT (Bidirectional Encoder Representations from Transformers) is the best-known example. The combination of transformers and self-supervised pretraining has been responsible for a paradigm shift in NLP, information retrieval (IR), and beyond.

This book provides a synthesis of existing work as a single point of entry for practitioners who wish to gain a better understanding of how to apply transformers to text ranking problems and researchers who wish to pursue work in this area. It covers a wide range of modern techniques, grouped into two high-level categories: transformer models that perform reranking in multi-stage architectures and dense retrieval techniques that perform ranking directly. Two themes pervade the book: techniques for handling long documents, beyond typical sentence-by-sentence processing in NLP, and techniques for addressing the tradeoff between effectiveness (i.e., result quality) and efficiency (e.g., query latency, model and index size). Although transformer architectures and pretraining techniques are recent innovations, many aspects of how they are applied to text ranking are relatively well understood and represent mature techniques. However, there remain many open research questions, and thus in addition to laying out the foundations of pretrained transformers for text ranking, this book also attempts to prognosticate where the field is heading.

# Contents

# Preface

BERT, unveiled by Google in October 2018, represents the culmination of a series of developments in NLP that date back many years. Together with its siblings, cousins, and intellectual descendants, these pretrained transformer models have, without exaggeration, brought about a paradigm shift in technologies that process, analyze, and otherwise manipulate human language text (and beyond). Pretrained transformers are responsible for substantial quality improvements in a range of downstream tasks, including natural language inference, syntactic and semantic analysis, text classification, question answering, summarization, and many more…and yes, text ranking.

The goal of text ranking is to generate an ordered list of texts retrieved from a corpus in response to a user's query. The most common formulation of text ranking is search, but many other NLP tasks involve aspects of text ranking as well. Despite being a relatively recent innovation, the foundations of how to apply BERT and other pretrained transformers to tackle text ranking are already quite sturdy. In many cases, improvements in effectiveness are substantial, robust, and have been widely replicated in many application scenarios. Building on these empirical foundations, this book provides a synthesis of existing work. We imagine these pages to be useful for practitioners interested in real-world solutions as well as researchers who wish to pursue work in this area.

This book covers a wide range of techniques for text ranking with pretrained transformers, grouped into two categories: transformer models that perform reranking in multi-stage architectures and dense retrieval techniques that perform ranking directly. Examples in the first category include approaches based on relevance classification, evidence aggregation from multiple segments of text, and query expansion methods. The second category involves using transformers to learn dense vector representations of texts, where ranking is formulated as comparisons between query and document vectors that take advantage of nearest neighbor search.

Two themes pervade our book: techniques for handling long documents, beyond typical sentence-by-sentence processing in NLP, and techniques for addressing the tradeoffs between effectiveness (i.e., result quality) and efficiency (e.g., query latency, model and index size). Much effort has been devoted to developing ranking models that address the mismatch between document lengths and the length limitations of input to transformers. The computational costs of inference with transformers have led to the development of alternatives and variants that aim for different tradeoffs, both within multi-stage architectures as well as with dense learned representations.

While many aspects of how to apply pretrained transformers to text ranking are well understood and can be considered mature techniques, there remain many open research questions

and unresolved issues. Thus, in addition to laying out the foundations of pretrained transformers for text ranking, this book also attempts to prognosticate where the field is heading.

It is quite remarkable that BERT is only about three years old. Taking a step back and reflecting, the field has seen an incredible amount of progress in a short amount of time. We expect even more exciting things to come!

Jimmy Lin, Rodrigo Nogueira, and Andrew Yates
August 2021

# Acknowledgments

Jimmy Lin, Rodrigo Nogueira, and Andrew Yates
August 2021