

A Survey of Model Architectures in Information Retrieval

Anonymous ACL submission

Abstract

The period from 2022 to the present has represented one of the biggest paradigm shifts in information retrieval (IR) and natural language processing (NLP). This work surveys the evolution of model architectures in IR, focusing on two key aspects: backbone models for feature extraction and end-to-end system architectures for relevance estimation. The review intentionally separates architectural considerations from training methodologies to provide a focused analysis of structural innovations in IR systems. We trace the development from traditional term-based methods to modern neural approaches, particularly highlighting the impact of transformer-based models and subsequent large language models (LLMs). We conclude with a forward-looking discussion of emerging challenges and future directions, including architectural optimizations for performance and scalability, handling of multimodal, multilingual data, and adaptation to novel application domains such as autonomous search agents that is beyond traditional search paradigms.

1 Introduction

Information Retrieval (IR) aims to retrieve relevant information sources to satisfy users' information needs. In the past decades, IR has become indispensable for efficiently and effectively accessing vast amounts of information across various applications. Beyond its traditional role, IR now also plays a critical role in assisting large language models (LLMs) to generate grounded and factual responses. Research in IR primarily centers on two key aspects: (1) extracting better query and document feature representations, and (2) developing more accurate relevance estimators. The approaches for extracting query and document features have evolved from traditional term-based methods, such as boolean logic and vector space models, to modern solutions such as dense retrieval based on pre-trained language models (Lin et al., 2022).

Relevance estimators have evolved alongside advances in feature representations. Early approaches, including probabilistic and statistical language models, computed relevance with simple similarity functions based on term-based features. Learning-to-rank (LTR) techniques later emerged, incorporating machine learning models and multi-layer neural networks for relevance estimation (Li, 2011). The success of LTR methods can be largely attributed to their extensive use of manually engineered features, derived from both statistical properties of text terms and user behavior data collected from web browsing traffic (Qin and Liu, 2013). In 2010s, a vast literature explored neural rerankers in different architectures to capture the semantic similarity between queries and documents. Then pre-trained transformers, represented by BERT (Devlin et al., 2019), quickly revolutionized the model design, leading to an era where retrieval and ranking models adopt simpler architectures for relevance estimation, such as dot product operations and MLP layer prediction heads, which operate on learned neural representations (Karpukhin et al., 2020; Nogueira et al., 2020; Lin et al., 2022).

Recent advancements of LLMs have revolutionized applied machine learning (ML) communities, including IR. One intriguing property of LLMs is that they can be used for feature extraction and relevance estimation, achieving strong performance without extensive training (Ni et al., 2022a; Nee-lakantan et al., 2022; BehnamGhader et al., 2024; Sun et al., 2023; Qin et al., 2024a, *inter alia*). The rise of LLMs in IR builds upon a rich foundation of transformer-based pre-trained language models that have evolved from earlier neural architectures. These include Transformers (Vaswani et al., 2017), Recurrent Neural Networks (RNN, Elman, 1990), Attention (Bahdanau, 2014) and pre-trained static neural representations such as Word2Vec (Mikolov, 2013) and GloVe (Pennington et al., 2014).

This work reviews the evolution of model ar-

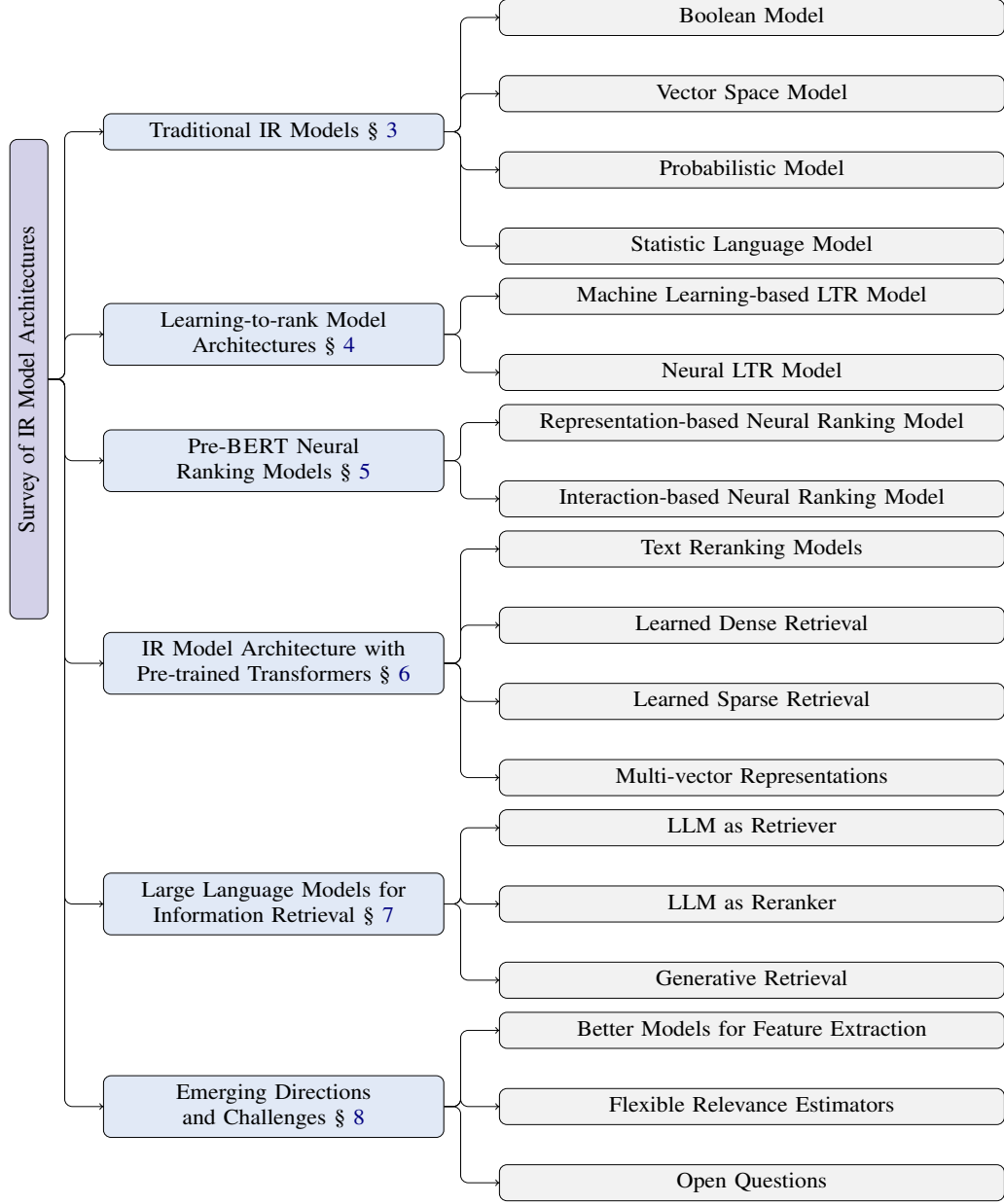


Figure 1: An overview of this survey. We focus on representative lines of works and defer details to the Appendix.

chitectures in IR (with an overview in Fig. 1). Here, the meaning of model architecture is twofold: it describes (1) backbone models for extracting query and document feature representations, and (2) end-to-end system architectures that process raw inputs, perform feature extraction, and estimate relevance. Different from prior works and surveys (Lin et al., 2022; Zhu et al., 2023), we intentionally separate our discussion of model architectures from training methodologies and deployment best practices to provide a focused architectural analysis. The shift towards neural architectures, particularly Transformer-based models, has fundamentally transformed IR by enabling rich,

contextualized representations and improved handling of complex queries. While this evolution has enhanced retrieval precision, it also presents new challenges, especially with the emergence of LLMs. These challenges include the need for architectural innovations to optimize performance and scalability, handling multimodal and multilingual data, incorporating domain-specific knowledge and understanding complex instructions. Moreover, as IR systems are increasingly integrated into diverse applications — from robotics (Xie et al., 2024), protein structure discovery (Jumper et al., 2021) to autonomous agents (Wu et al., 2023; Chen et al., 2025) that are capable of reasoning and search —

the field must evolve beyond traditional search paradigms. We conclude this survey by examining these challenges and discuss their implications for the future of IR model architectures research.

2 Background and Terminology

We focus on the classical *ad hoc* retrieval task.

Task Definition and Evaluation Given a query Q , the task is to find a ranked list of k documents, denoted as $\{D_1, D_2, \dots, D_k\}$, that exhibit the highest relevance to Q . This is achieved either by *retrieving* top- k documents from a large collection \mathcal{C} ($|\mathcal{C}| \gg |k|$), which typically comprises millions or billions of documents, or by *reranking* the top- k candidates returned by a retriever. System performance is measured using standard IR metrics such as Mean Reciprocal Rank, Recall, and normalized Discounted Cumulative Gain (nDCG).

Query and Document A *query* expresses an information need and serves as input to the *ad hoc* retrieval system. We denote *document* as the atomic unit for retrieval and ranking. Our discussions are primarily based on text-based documents, but it can also refer to a webpage or an email, depending on the actual IR application of interest.

Disentangling Model Architecture with Training Strategies Similar to other applied ML domains, the design of IR model architectures is paired with training strategies and deployment best practices. In this paper, we seek to disentangle the two and focus solely on architectures. We refer to prior surveys for a more focused review of training strategies and related topics (Schütze et al., 2008; Lin et al., 2022; Song et al., 2023).

3 Traditional IR Models

In this section, we briefly review traditional IR models prior to neural methods, with a focus on *boolean model* and *vector space model* which serve as the foundation of later development of IR models (§§ 4 to 7).¹ These models are built upon the basic unit “term” used in the representation (Nie, 2010).

Boolean Model In Boolean Models, a document D is represented by a set of terms it contains, i.e., $D = \{t_1, t_2, \dots, t_n\}$. A query Q is represented as a similar boolean expression of terms. A document is considered relevant to a query only if a logical

implication $D \rightarrow Q$ holds, i.e., the document representation logically implies the query expression.

Vector Space Model In Vector Space Models (Salton et al., 1975), the queries and documents are represented by vectors, e.g., $Q = \langle q_1, q_2, \dots, q_n \rangle$ and $D = \langle d_1, d_2, \dots, d_n \rangle$. The vector space $\mathcal{V} = \langle t_1, t_2, \dots, t_n \rangle$ is formed by all the terms the system recognizes in the documents and each element (q_i or d_i , $1 \leq i \leq n$) in the vectors represents the weight of the corresponding term in the query or the document. The weights q_i or d_i could be binary, representing presence or absence. Given the vector representations, the relevance score is estimated by a similarity function between the query Q and the document D .

4 Learning-to-Rank Model Architectures

Different from traditional IR models discussed in § 3, Learning-to-Rank (LTR) leverages the idea of supervised ML on extensively crafted numerical features (Burges et al., 2005; Burges, 2010; Qin and Liu, 2013). For each (Q_i, D_i) pair, a k -dimensional feature vector $\mathbf{x}_i \in \mathbb{R}^k$ and a relevance label y_i is provided to the ranking model f parameterized by θ . Denote the loss function as $l(\cdot)$, the ranking is trained to minimize the empirical loss on labeled training set Ψ : $\mathcal{L} = 1/|\Psi| \sum_{(\mathbf{x}_i, y_i) \in \Psi} l(f_\theta(\mathbf{x}_i), y_i)$.

Explorations in LTR models can be grouped into two directions: *ML-based models* and *neural LTR models*. Under the scope of ML models, RANKSVM (Joachims, 2006) is a pairwise LTR model based on SVM. Burges et al. (2005) studied decision trees and Wu et al. (2010) proposed LAMBDA MART based on Gradient Boosted Decision Trees (GBDT, Friedman, 2001; Ke et al., 2017). Unsurprisingly, early works also explored neural LTR models. RANKNET (Burges et al., 2005) and LAMBDA RANK (Burges et al., 2006) parameterize the LTR model with neural networks. Recent works such as GSF (Ai et al., 2019) and APPROXNDCG (Bruch et al., 2019) use multiple fully connected layers. DLCM (Ai et al., 2018a) and SETRANK (Pang et al., 2020) adopt RNN and self-attention for reranking documents. Qin et al. (2021) conducted rigorous study of benchmarking neural ranking models against GBDT-based models. See Table 1 for a list of LTR models and prior surveys on LTR techniques (Liu, 2009; Li, 2011).

LTR techniques use human-crafted numerical features and metadata like PageRank score (Brin

¹We defer the discussion of *probabilistic model* and *statistical language model* to Appx. A.

and Page, 1998) and click count as features and are still widely used in modern search systems (Google, 2019). However, it lacks the flexibility of being directly used on raw text data, and also cannot overcome the lexical mismatch problem — xthe main focus of neural ranking methods (§ 5).

5 Neural Ranking Models

Different from feature engineering of LTR (§ 4), neural ranking models utilize deep neural networks to learn feature representations directly from raw text and again use neural networks for relevance estimation.² Depending on how queries interact with documents during network processing, neural ranking models can be roughly divided into *representation-based models* and *interaction-based models* (Guo et al., 2016a).

Representation-based models can be regarded as extensions of vector space models (§ 3), which independently encode queries and documents into a latent vector space, as illustrated in Fig. 2a. The Deep Structured Semantic Model (DSSM) (Huang et al., 2013) is an early example. It utilizes word hashing and multilayer perceptrons (MLPs) to independently encode term vectors of queries and documents, enabling the computation of ranking scores based on the cosine similarity of their embeddings. Later works modify DSSM’s encoder network to better capture richer semantic and contextual information. Convolutional DSSM (Shen et al., 2014) leverages a CNN architecture to project vectors within a context window to a local contextual feature vector. A variation of DSSM replaces MLPs with a Long Short-Term Memory (LSTM) network (Hochreiter and Schmidhuber, 1997; Palangi et al., 2016; Wan et al., 2016), utilizing its memory mechanism to capture local and global context information.

Interaction-based models (Fig. 2b), on the other hand, process queries and documents jointly through neural networks. The model’s output is typically a scalar relevance score of the input query-document pair. Various network architectures have been proposed under this paradigm. MATCHPYRAMID (Pang et al., 2016) employs CNN over the interaction matrix between query and document terms. The interaction matrix is treated as an image, allowing the CNN to capture local matching

patterns through convolution and pooling operations (Hu et al., 2014). Building upon the concept of interaction-focused models, Guo et al. (2016a) highlight the importance of exact term matches in neural ranking models and proposed the Deep Relevance Matching Model (DRMM). The model constructs matching histograms for each query term to capture the distribution of matching signals across document terms. Kernel-Based Neural Ranking Model (K-NRM, Xiong et al., 2017) further advances interaction-based approaches. It employs radial basis function kernels to transform the query-document interaction matrix informative ranking features. CONV-KNRM (Dai et al., 2018) later extends it to convolutional kernels.

In addition to the development of network architecture, pre-trained embeddings (Mikolov, 2013; Pennington et al., 2014) provide semantic-based term representations to enable neural ranking models to focus on learning relevance matching patterns, improving training convergence and retrieval performance on both representation-based and interaction-based models (Levy et al., 2015).

6 IR with Pre-trained Transformers

BERT (Devlin et al., 2019) changed the research paradigm in both NLP and IR. Its success can be attributed to two factors: (1) the multi-head attention architecture (Vaswani et al., 2017) admits fine-grained, contextualized representations; (2) large-scale pre-training allows BERT to encode both semantics and world knowledge. The expressiveness of BERT has been extensively studied by prior works, e.g., Rogers et al. (2020); Tenney et al. (2019); Clark (2019). This section discusses IR architectures based on pre-trained transformers, with a focus on BERT-type encoder models.

Text Reranking Nogueira et al. (2019) first employed BERT model for reranking candidate passages from a first-stage retriever. Their model MONOBER takes as input the sequence of concatenated (Q, D) as input, and outputs a relevance score s with a linear layer on top of the BERT model (Fig. 2c). The schema has later been proved to be effective on other pre-trained encoders (Zhang et al., 2021) and encoder-decoder architectures (Nogueira et al., 2020). However, this schema faces two challenges: (1) BERT family models has a limited 512 tokens context length, making reranking long documents challenging; (2) CLS token’s single 768-dimensional representation

²We define neural information retrieval as retrieval models based on neural networks prior to pre-trained transformers. More models details are deferred to Appx. C and detailed surveys (Onal et al., 2018; Mitra et al., 2018; Xu et al., 2018)

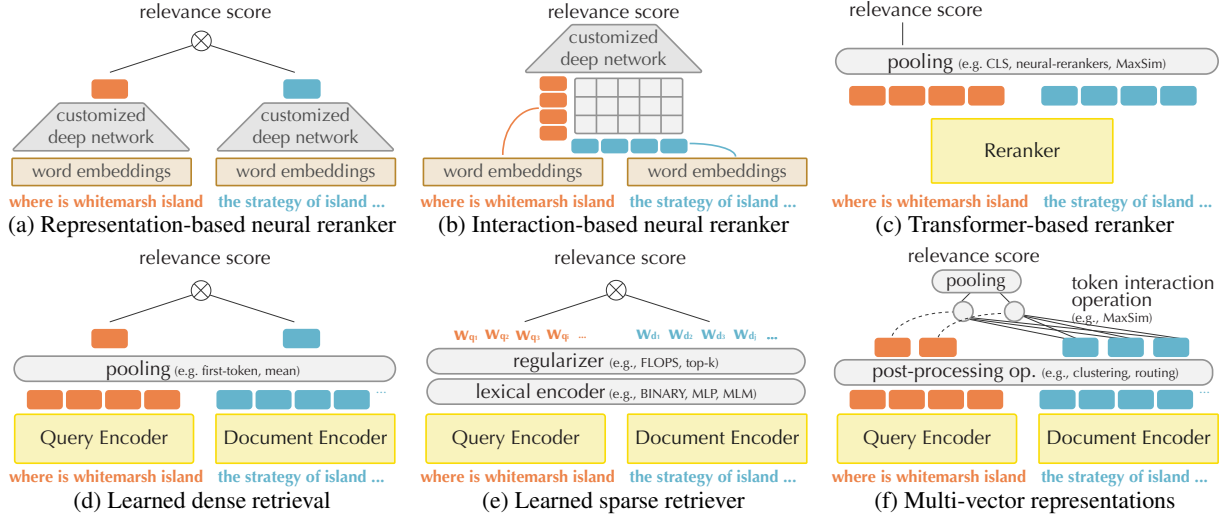


Figure 2: Illustration on neural ranking models, retriever, and reranker architectures. Brown boxes indicate uncontextualized word embeddings (e.g., Word2vec). Yellow boxes indicate pretrained Transformers (e.g., BERT).

potentially limits the expressiveness of the reranking model. Two directions have been investigated to tackle these two challenges.

In the first direction, one strategy is to segment the long document into shorter passages, score each passage individually, then aggregate the scores to get a document-level relevance score (Yilmaz et al., 2019; Dai and Callan, 2019b). In this case, the underlying model architecture remains unchanged. A different line of work seeks to perform feature-level aggregation. PARADE (Li et al., 2020) uses an additional neural network to aggregate contextualized representations from CLS tokens of passages to get the final document relevance score.

In the second direction, MacAvaney et al. (2019) discovered via CEDR that the effectiveness of reranker could be enhanced when aggregating the contextualized representations with neural ranking models such as K-NRM and DRMM. Zhang et al. (2024a) later observed that integrating token representations with late interactions could also effectively improve the reranking robustness on out-of-domain scenarios, especially for long queries.

Learned Dense Retrieval Here we use the term "bi-encoder" (Humeau et al., 2020) to refer to the model architecture commonly used for dense retrieval.³ Bi-encoder uses a backbone network to encode query Q and document D separately, then uses the encoded dense vector representations to compute the relevance score with similarity

³The term "bi-encoder" is known by many other names, such as two-tower architecture, embedding models. We refer to it as "bi-encoder" in contrast to "cross-encoder", reranker architectures that take concatenated (Q, D) as input.

functions such as dot product or cosine similarity (Xiong et al., 2020; Karpukhin et al., 2020; Gao et al., 2021c). After training, the model encodes the collection into a dense vector index, and retrieval is performed with fast nearest neighbor search techniques (Johnson et al., 2019; Malkov and Yashunin, 2016). Different from representation-based neural ranking models where distinct architectures are proposed (§ 5), existing dense retrieval methods are mostly based on pre-trained transformer language models, with variance on pooling strategy and training methodologies (detailed in Appx. D).

Learned Sparse Retrieval Similar to learned dense retrieval, learned sparse retrieval (LSR) also uses bi-encoder architecture with language models as the backbone, to transform documents into a static index for later retrieval (Zamani et al., 2018). To use a traditional inverted index for faster retrieval (Bruch et al., 2024), the query Q and document D are represented as sparse vectors whose dimensionality typically matches the vocabulary size of the backbone pre-trained transformer model (Yu et al., 2024a). Here sparsity is enforced through regularization (Formal et al., 2021b; Paria et al., 2020) and usually serves as a trade-off between effectiveness and efficiency. At a higher level, LSR can be viewed as a way to learn token importance or "impact" scores from data (Dai and Callan, 2019b; Bai et al., 2020; Mallia et al., 2021), in contrast to static formulas like BM25.

Multi-Vector Representations Learned dense retrieval’s bi-encoder architecture encodes queries and documents into single feature vectors sepa-

rately, and estimate relevance via a similarity function. This enables efficient training, indexing, and inference, but the lack of interactions between query and document terms potentially limits performance. In contrast, the "all-to-all" interaction of cross-encoder models are computationally expensive. Research has explored representing queries and documents using multiple vectors and developing corresponding relevance estimators. POLY-ENCODER (Humeau et al., 2020) computes a fixed number of vectors per query, and aggregate with softmax attention against document vectors. Due to the use of softmax operator, fast nearest neighbor search technique cannot be trivially applied. ME-BERT (Luan et al., 2021) proposes to represent documents with m vectors, and uses the maximum similarity between query vectors and document vectors to estimate relevance. COLBERT (Khattab and Zaharia, 2020) takes the multi-vector representation idea further and represents each token in query and document as a contextualized vector. Each query vector interacts with each document vector via a MaxSim operator, and the relevance score is computed by summing scalar outputs of these operators over query terms. COLBERT is trained end-to-end and achieves strong performance on public retrieval benchmarks. Numerous other studies further investigate token selection and aggregation operations, see Appx. D for details.

7 Large Language Models for IR

LLMs have exhibited proficiency in language understanding and generation, are trained to align with human preferences (OpenAI, 2023; Gemini et al., 2023; Bai et al., 2022) and able to perform complex tasks such as reasoning (Wei et al., 2022; Hurst et al., 2024; Guo et al., 2025) and planning (Song et al., 2023).⁴ In this section, we briefly discuss some works that utilize LLMs for IR tasks.

LLM as Retriever Adopting an LLM as the backbone for bi-encoder retrieval model has achieved performance improvement compared to smaller-sized models like BERT. Neelakantan et al. (2022) fine-tuned a series of off-the-shelf GPT models towards text and code representation. They empirically verified that the bi-encoder retriever’s performance can benefit from increased backbone language model capacity. Common text re-

trieval benchmarks like BEIR (Thakur et al., 2021) are currently dominated by LLM-based retrievers. A parallel line of research has explored adapting unidirectional LLM architectures into bidirectional ones to enhance representational power. LLM2VEC (BehnamGhader et al., 2024) enables bidirection and further trains LLAMA-2 (Touvron et al., 2023) with specific adaptive tasks. NV-EMBED (Lee et al., 2025) introduces a new latent attention layer and leads to improve on MTEB benchmark (Muennighoff et al., 2023).

LLM as Reranker Works discussed in § 6 have explored fine-tuning BERT-type encoder models as cross-encoder reranker. Later works further expand this paradigm to encoder-decoder models like T5 (Raffel et al., 2020) and decoder models like LLAMA (Touvron et al., 2023). Nogueira et al. (2020) fine-tuned T5 models with classification loss for passage reranking. Zhuang et al. (2023a) proposed to fine-tune T5 to produce a numerical relevance score, and optimize the model with ranking losses like RANKNET (Burgess, 2010). LISTT5 (Yoon et al., 2024) adopts the Fusion-in-Decoder architecture (Izacard and Grave, 2021) to learn a listwise reranker. RANKLLAMA (Ma et al., 2024b) fine-tunes decoder model for pointwise reranking and achieves better performance compared to T5-based rerankers. Leveraging the long-context ability of LLMs, a reranking paradigm is introduced, where LLM-based rerankers directly rerank a list of documents rather than scoring each document individually (Ma et al., 2023; Zhang et al., 2023b; Pradeep et al., 2023b,c). Instead of using raw passages, Liu et al. (2024b) used passage embeddings as input to LLMs and trained corresponding rerankers to achieve improved efficiency. Aforementioned studies still rely on labeled data and gradient updates to backbone LLMs. With the rise of instruction-following LLMs, researchers have explored using LLMs as unsupervised rerankers through prompting techniques. As this line of research does not introduce architectural changes to existing LLMs, we refer to a recent survey (Zhu et al., 2023) for further details.

Generative Retrieval Traditional IR systems follow the "index-retrieval-rank" paradigm (Schütze et al., 2008). Although effective, jointly optimizing the separate index, retrieval, and reranking modules can be challenging. A recent line of research aims to bypass the indexing step by using autoregressive language models to directly generate document

⁴In this work, we use term LLM to denote language models which are trained for text generation, including encoder-decoder models and decoder-only models.

identifiers (DocIDs). DSI (Tay et al., 2022) first constructs semantically structured DocIDs, then fine-tunes T5 models with labeled data. In the decoding phase, DSI uses all DocIDs to construct a trie and performs constrained decoding. Followup works (Wang et al., 2022b; Bevilacqua et al., 2022) further improve upon this paradigm with strategies to construct semantic DocIDs and enable robust training. A significant challenge for generative retrieval is scalability to larger corpus (Pradeep et al., 2023a). Zeng et al. (2024) utilized the residual quantization technique to construct hierarchical DocIDs and achieved comparable performance to dense retrievers on MS MARCO dataset. Generative retrieval is an active research area; see (Li et al., 2025c) for a more comprehensive review.

Remarks We note that the adoption of LLMs in IR model architectures primarily follows two main themes: (1) feature extraction, and (2) relevance estimation, as discussed in § 1. For example, LLMs’ semantic knowledge enables their strong performance in being the backbone of a retriever; and instruction-following LLMs can be directly prompted for relevance estimation. Generative retrieval and cross-encoder LLM reranking models are trained end-to-end for both feature extraction and relevance estimation. While LLMs have shown promise, several challenges and open questions remain, which leaves room for discussion (§ 8).

8 Emerging Directions and Challenges

IR systems have become crucial across diverse domains, from retrieval-augmented language modeling (Khandelwal et al., 2020a; Borgeaud et al., 2022) to applications in agents (Wu et al., 2023; Wang et al., 2024a), code generation (Wang et al., 2024c; Zhang et al., 2023a), robotics (Anwar et al., 2024), medicine (Jeong et al., 2024), and protein research (Jumper et al., 2021), *inter alia*. These developments present new challenges for IR research. Drawing from the evolution of IR architectures (§§ 3 to 7), we examine emerging trends, open problems, and potential research directions.⁵

8.1 Better Models for Feature Extraction

Scaling has been a winning recipe for modern neural networks (Kaplan et al., 2020; Hoffmann et al., 2022; Dehghani et al., 2023; Fang et al., 2024; Shao et al., 2024, *inter alia*). As IR moves toward

compute-intensive practices, we identify key areas for model improvement:

- **Data & training efficiency** Current transformer-based IR models demand extensive training data (Fang et al., 2024), making them impractical for many real-world applications. Developing architectures that can learn effectively from limited data remains crucial. Additionally, models should support parallel processing and low precision training to reduce costs and accelerate convergence (Nvidia, 2021; Fishman et al., 2024; Liu et al., 2024a).
- **Inference optimization** Real-time applications like conversational search (Mo et al., 2024b) and agent-based systems (Yao et al., 2023) require efficient handling of variable-length queries, necessitating advanced compression and optimization techniques for both retriever backbones and index structures (Dettmers and Zettlemoyer, 2023; Kumar et al., 2024; Warner et al., 2024; Bruch et al., 2024; Xu et al., 2025a, *inter alia*).
- **Multimodality & multilinguality** Future IR systems must handle diverse content types including images (Ma et al., 2024a), audio (Pusateri et al., 2024), structured data (Tan et al., 2024; Edge et al., 2024) as well as multilinguality beyond English (Zhang et al., 2023c; Enevoldsen et al., 2025). Recent advances in multimodal, multilingual retrieval (Ma et al., 2024a; Wei et al., 2025; Yu et al., 2024b; Huang et al., 2024, *inter alia*) and structured data processing (Li et al., 2023d, 2024) have demonstrated promises.
- **Transformer alternatives** While transformers have dominated recent IR research, their quadratic complexity in attention computation remains a significant bottleneck. Recent advances in linear RNNs (Peng et al., 2023, 2024; Qin et al., 2024b), state space models (Gu and Dao, 2024; Dao and Gu, 2024), and linear attention (Katharopoulos et al., 2020; Yang et al., 2024) offer alternatives with theoretical linear complexity. Although preliminary studies (Xu et al., 2025b) show limited gains compared to optimized transformers, developing efficient alternatives architectures for transformers could revolutionize large-scale information retrieval.

Strong foundation models have proven crucial for IR performance (Neelakantan et al., 2022; Ma et al., 2024b). As IR applications expand, developing foundation models that balance computational efficiency with robust performance across tasks and

⁵See Appx. F for an extended discussion *w.r.t.* deployment challenges, robustness, autonomous search agents etc.

modalities emerges as a key research priority.

8.2 Flexible Relevance Estimators

As discussed in § 6, cross-encoders provide complex non-linear relevance estimation but are computationally expensive. In contrast, bi-encoder architectures used in dense and sparse retrieval rely on linear similarity functions (e.g., inner product) to enable fast retrieval through nearest neighbor search and inverted indexing. Balancing complex relevance matching and scalable retrieval remains challenging. COLBERT (Khattab and Zaharia, 2020) addresses this by using document representation matrices with MaxSim operations, while recent work (Killingback et al., 2025) explores Hypernetworks (Ha et al., 2022) to generate query-specific neural networks for relevance estimation. The design of flexible yet scalable relevance estimators remains an active research direction.

8.3 Open Questions

The integration of IR systems into other research domains presents new challenges. We discuss key implications for future IR modeling research.

End “User” of Retrieval While traditional IR systems focus on providing search results to humans, retrieval is increasingly used to support ML models, particularly LLMs, in tasks such as generation (Gao et al., 2023), reasoning (Yao et al., 2024; Islam et al., 2024), and planning (Song et al., 2023). This shifting paradigm raises questions about task formulation, evaluation, and system optimization:

- Current IR research is grounded in human information-seeking behavior (Wilson, 2000). When the end user becomes another ML model, we must reconsider how to define and assess *relevance*. This question suggests a need for flexible, data-efficient models that are adaptable to various downstream tasks.
- Traditional IR metrics, which are designed for human-centric evaluation, may not align with downstream task performance in retrieval-augmented systems. Future IR models should support end-to-end system optimization rather than focusing solely on ranking metrics.

Autonomous Search Agent Complex tasks require retrieving long-tail knowledge using lengthy, complex queries (Soudani et al., 2024; Su et al., 2024), demanding retrieval models capable of instruction following (Weller et al., 2024a; Ravfogel et al., 2024) and reasoning (Su et al., 2024).

Existing attempts can be divided into two directions. One line of works focuses on training retrievers and rerankers that are capable of reasoning. They propose data pipelines to synthesize training data (Oh et al., 2024; Weller et al., 2024b; Shao et al., 2025, *inter alia*) and leverage strong backbone language models such as LLAMA (Dubey et al., 2024) and QWEN (Yang et al., 2025). Another line of works treats search/retrieval as an integral component of LLM reasoning chain. They consider search/retrieval system as a static tool that can be called via large reasoning model (LRM), and instead focus on improving LRM’s capability to use search tool and synthesize search results. The LRM can decide where, when and how to conduct search, and the search results subsequently influence LRM’s further reasoning and decision making (Nakano et al., 2021; Tang et al., 2025; He et al., 2025; Chen et al., 2025; Li et al., 2025a).

Despite the exciting progress, key limitations remain in building instruction following and reasoning-capable retrieval systems. For example, to enable retrievers’ reasoning capability requires strong backbone models (e.g., 7B scale), which are often infeasible for production systems. Even larger models (e.g., 32B scale) augmented with retrieval and trained via expensive reinforcement learning (Jin et al., 2025; Chen et al., 2025) still underperform simpler baselines with query decomposition and chain-of-thought prompting (Khot et al., 2023; Trivedi et al., 2023). A key open question is how to endow retrievers with strong reasoning capabilities using lightweight, scalable models. Another open challenge lies in the joint optimization of retrievers and language models within a unified, reasoning-aware framework. Lastly, the human factors of applying such autonomous search agents remains to be studied.

9 Conclusions and Closing Thoughts

Information retrieval modeling has evolved from simple term matching to complex neural networks and LLM-driven approaches, significantly improving search capabilities. Key challenges ahead include balancing computational efficiency with performance, handling diverse data types, maintaining faithfulness and trustworthiness, and integrating with emerging technologies like autonomous agents. These challenges drive opportunities for developing more adaptive, efficient, scalable and intelligent retrieval systems.

Limitations

This survey examines the evolution of IR models, with particular emphasis on challenges arising from LLMs and their implications for future architectures. Due to space constraints, we focus on representative works rather than providing an exhaustive review, with supplementary discussions of interdisciplinary research included in the Appendix.

Ethical Considerations

As this paper solely reviews existing IR developments and future research directions, we believe it presents no direct ethical concerns.

References

Qingyao Ai, Keping Bi, Jiafeng Guo, and W Bruce Croft. 2018a. Learning a deep listwise context model for ranking refinement. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 135–144.

Qingyao Ai, Keping Bi, Cheng Luo, Jiafeng Guo, and W Bruce Croft. 2018b. Unbiased learning to rank with unbiased propensity estimation. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 385–394.

Qingyao Ai, Jiaxin Mao, Yiqun Liu, and W. Bruce Croft. 2018c. [Unbiased learning to rank: Theory and practice](#). In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM '18*, page 2305–2306, New York, NY, USA. Association for Computing Machinery.

Qingyao Ai, Xuanhui Wang, Sebastian Bruch, Nadav Golbandi, Michael Bendersky, and Marc Najork. 2019. Learning groupwise multivariate scoring functions using deep neural networks. In *Proceedings of the 2019 ACM SIGIR international conference on theory of information retrieval*, pages 85–92.

Bang An, Shiyue Zhang, and Mark Dredze. 2025. [RAG LLMs are not safer: A safety analysis of retrieval-augmented generation for large language models](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5444–5474, Albuquerque, New Mexico. Association for Computational Linguistics.

Abrar Anwar, John Welsh, Joydeep Biswas, Soha Pouya, and Yan Chang. 2024. Remembr: Building and reasoning over long-horizon spatio-temporal memory for robot navigation. *arXiv preprint arXiv:2409.13682*.

Akari Asai, Timo Schick, Patrick Lewis, Xilun Chen, Gautier Izacard, Sebastian Riedel, Hannaneh Hajishirzi, and Wen-tau Yih. 2023. [Task-aware retrieval with instructions](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3650–3675, Toronto, Canada. Association for Computational Linguistics.

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024a. [Self-RAG: Learning to retrieve, generate, and critique through self-reflection](#). In *The Twelfth International Conference on Learning Representations*.

Akari Asai, Zexuan Zhong, Danqi Chen, Pang Wei Koh, Luke Zettlemoyer, Hannaneh Hajishirzi, and Wen-tau Yih. 2024b. Reliable, adaptable, and attributable language models with retrieval. *arXiv preprint arXiv:2403.03187*.

Ricardo Baeza-Yates, Berthier Ribeiro-Neto, et al. 1999. *Modern information retrieval*, volume 463. ACM press New York.

Dzmitry Bahdanau. 2014. Neural machine translation by jointly learning to align and translate. *The Twelfth International Conference on Learning Representations*.

Yang Bai, Xiaoguang Li, Gang Wang, Chaoliang Zhang, Lifeng Shang, Jun Xu, Zhaowei Wang, Fangshan Wang, and Qun Liu. 2020. Sparterm: Learning term-based sparse representation for fast text retrieval. *arXiv preprint arXiv:2010.00768*.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.

Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. 2024. [LLM2vec: Large language models are secretly powerful text encoders](#). In *First Conference on Language Modeling*.

Adam Berger and John Lafferty. 1999. Information retrieval as statistical translation. In *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 222–229.

Michele Bevilacqua, Giuseppe Ottaviano, Patrick Lewis, Scott Yih, Sebastian Riedel, and Fabio Petroni. 2022. [Autoregressive search engines: Generating substrings as document identifiers](#). In *Advances in Neural Information Processing Systems*.

Asia J Biega, Krishna P Gummadi, and Gerhard Weikum. 2018. Equity of attention: Amortizing individual fairness in rankings. In *The 41st international acm sigir conference on research & development in information retrieval*, pages 405–414.

767	Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, and	Chia-Yuan Chang, Zhimeng Jiang, Vineeth Rakesh,	821
768	Rodrigo Nogueira. 2022. Inpars: Unsupervised	Menghai Pan, Chin-Chia Michael Yeh, Guanchu	822
769	dataset generation for information retrieval. In <i>Pro-</i>	Wang, Mingzhi Hu, Zhichao Xu, Yan Zheng, Ma-	823
770	<i>ceedings of the 45th International ACM SIGIR Con-</i>	hashweta Das, et al. 2024. Main-rag: Multi-agent fil-	824
771	<i>ference on Research and Development in Information</i>	tering retrieval-augmented generation. <i>arXiv preprint</i>	825
772	<i>Retrieval</i> , pages 2387–2392.	<i>arXiv:2501.00332</i> .	826
773	Sebastian Borgeaud, Arthur Mensch, Jordan Hoff-	Yin-Wen Chang, Cho-Jui Hsieh, Kai-Wei Chang,	827
774	mann, Trevor Cai, Eliza Rutherford, Katie Milli-	Michael Ringgaard, and Chih-Jen Lin. 2010. Train-	828
775	can, George Bm Van Den Driessche, Jean-Baptiste	ing and testing low-degree polynomial data mappings	829
776	Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022.	via linear svm . <i>Journal of Machine Learning Re-</i>	830
777	Improving language models by retrieving from tril-	<i>search</i> , 11(48):1471–1490.	831
778	lions of tokens. In <i>International conference on ma-</i>	Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou,	832
779	<i>chine learning</i> , pages 2206–2240. PMLR.	Chenzheng Zhu, Haofen Wang, Jeff Z Pan, Wen	833
780	Leonid Boytsov, Preksha Patel, Vivek Sourabh, Riddhi	Zhang, Huajun Chen, Fan Yang, et al. 2025. Learn-	834
781	Nisar, Sayani Kundu, Ramya Ramanathan, and Eric	ing to reason with search for llms via reinforcement	835
782	Nyberg. 2024. Inpars-light: Cost-effective unsuper-	learning. <i>arXiv preprint arXiv:2503.19470</i> .	836
783	vised training of efficient rankers. <i>Transactions on</i>	Hyung Won Chung, Le Hou, Shayne Longpre, Barret	837
784	<i>Machine Learning Research</i> .	Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi	838
785	Sergey Brin and Lawrence Page. 1998. The anatomy of	Wang, Mostafa Dehghani, Siddhartha Brahma, et al.	839
786	a large-scale hypertextual web search engine . <i>Com-</i>	2024. Scaling instruction-finetuned language models.	840
787	<i>puter Networks</i> , 30:107–117.	<i>Journal of Machine Learning Research</i> , 25(70):1–53.	841
788	Sebastian Bruch, Franco Maria Nardini, Cosimo Rulli,	Kevin Clark. 2019. What does bert look at? an analysis	842
789	and Rossano Venturini. 2024. Efficient inverted in-	of bert’s attention. <i>arXiv preprint arXiv:1906.04341</i> .	843
790	indexes for approximate retrieval over learned sparse	Charles LA Clarke and Laura Dietz. 2024. Llm-based	844
791	representations. In <i>Proceedings of the 47th Inter-</i>	relevance assessment still can’t replace human rele-	845
792	<i>national ACM SIGIR Conference on Research and</i>	levance assessment. <i>arXiv preprint arXiv:2412.17156</i> .	846
793	<i>Development in Information Retrieval</i> , pages 152–	Gordon V. Cormack, Charles L A Clarke, and Stefan	847
794	162.	Buettcher. 2009. Reciprocal rank fusion outperforms	848
795	Sebastian Bruch, Masrour Zoghi, Michael Bendersky,	condorcet and individual rank learning methods . In	849
796	and Marc Najork. 2019. Revisiting approximate met-	<i>Proceedings of the 32nd International ACM SIGIR</i>	850
797	ric optimization in the age of deep neural networks.	<i>Conference on Research and Development in Infor-</i>	851
798	In <i>Proceedings of the 42nd international ACM SIGIR</i>	<i>mation Retrieval</i> , SIGIR ’09, page 758–759, New	852
799	<i>conference on research and development in informa-</i>	York, NY, USA. Association for Computing Machin-	853
800	<i>tion retrieval</i> , pages 1241–1244.	ery.	854
801	Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier,	Sunhao Dai, Yuqi Zhou, Liang Pang, Weihao Liu, Xi-	855
802	Matt Deeds, Nicole Hamilton, and Greg Hullender.	aolin Hu, Yong Liu, Xiao Zhang, Gang Wang, and	856
803	2005. Learning to rank using gradient descent. In	Jun Xu. 2024. Neural retrievers are biased towards	857
804	<i>Proceedings of the 22nd international conference on</i>	llm-generated content. In <i>Proceedings of the 30th</i>	858
805	<i>Machine learning</i> , pages 89–96.	<i>ACM SIGKDD Conference on Knowledge Discovery</i>	859
806	Christopher Burges, Robert Ragno, and Quoc Le. 2006.	<i>and Data Mining</i> , pages 526–537.	860
807	Learning to rank with nonsmooth cost functions. <i>Ad-</i>	Zhuyun Dai and Jamie Callan. 2019a. Context-aware	861
808	<i>vances in neural information processing systems</i> , 19.	sentence/passage term importance estimation for first	862
809	Christopher JC Burges. 2010. From ranknet to lamb-	stage retrieval. <i>arXiv preprint arXiv:1910.10687</i> .	863
810	darank to lambdamart: An overview. <i>Learning</i> ,	Zhuyun Dai and Jamie Callan. 2019b. Deeper text un-	864
811	11(23-581):81.	derstanding for ir with contextual neural language	865
812	Tianyu Cao, Neel Bhandari, Akhila Yerukola, Akari	modeling . <i>Proceedings of the 42nd International</i>	866
813	Asai, and Maarten Sap. 2025. Out of style: Rag’s	<i>ACM SIGIR Conference on Research and Develop-</i>	867
814	fragility to linguistic variation. <i>arXiv preprint</i>	<i>ment in Information Retrieval</i> .	868
815	<i>arXiv:2504.08231</i> .	Zhuyun Dai, Chenyan Xiong, Jamie Callan, and	869
816	Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and	Zhiyuan Liu. 2018. Convolutional neural networks	870
817	Hang Li. 2007. Learning to rank: from pairwise	for soft-matching n-grams in ad-hoc search . In <i>Pro-</i>	871
818	approach to listwise approach. In <i>Proceedings of the</i>	<i>ceedings of the Eleventh ACM International Confer-</i>	872
819	<i>24th international conference on Machine learning</i> ,	<i>ence on Web Search and Data Mining</i> , WSDM ’18,	873
820	pages 129–136.	page 126–134, New York, NY, USA. Association for	874
		Computing Machinery.	875

876	Zhuyun Dai, Vincent Y Zhao, Ji Ma, Yi Luan, Jianmo	Nova DasSarma, Dawn Drain, Deep Ganguli, Zac	932
877	Ni, Jing Lu, Anton Bakalov, Kelvin Guu, Keith Hall,	Hatfield-Dodds, Danny Hernandez, Andy Jones,	933
878	and Ming-Wei Chang. 2023. Promptagator: Few-	Jackson Kernion, Liane Lovitt, Kamal Ndousse,	934
879	shot dense retrieval from 8 examples. In <i>The Eleventh</i>	Dario Amodei, Tom Brown, Jack Clark, Jared Ka-	935
880	<i>International Conference on Learning Representa-</i>	plan, Sam McCandlish, and Chris Olah. 2021. A	936
881	<i>tions</i> .	mathematical framework for transformer circuits.	937
882	Tri Dao and Albert Gu. 2024. Transformers are SSMS:	<i>Transformer Circuits Thread</i> . https://transformer-	938
883	Generalized models and efficient algorithms through	circuits.pub/2021/framework/index.html .	939
884	structured state space duality. In <i>International Con-</i>	Jeffrey L. Elman. 1990. Finding structure in time . <i>Cog-</i>	940
885	<i>ference on Machine Learning (ICML)</i> .	<i>nitive Science</i> , 14(2):179–211.	941
886	Jeffrey Dean et al. 2009. Challenges in building large-	Kenneth Enevoldsen, Isaac Chung, Imene Kerboua,	942
887	scale information retrieval systems. In <i>Keynote of</i>	Márton Kardos, Ashwin Mathur, David Stap,	943
888	<i>the 2nd ACM international conference on web search</i>	Jay Gala, Wissam Siblini, Dominik Krzemiński,	944
889	<i>and data mining (WSDM)</i> , volume 10.	Genta Indra Winata, et al. 2025. Mmteb: Mas-	945
890	Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Pi-	sive multilingual text embedding benchmark. <i>arXiv</i>	946
891	otr Padlewski, Jonathan Heek, Justin Gilmer, An-	<i>preprint arXiv:2502.13595</i> .	947
892	dreas Peter Steiner, Mathilde Caron, Robert Geirhos,	Guglielmo Faggioli, Laura Dietz, Charles LA Clarke,	948
893	Ibrahim Alabdulmohsin, et al. 2023. Scaling vision	Gianluca Demartini, Matthias Hagen, Claudia Hauff,	949
894	transformers to 22 billion parameters. In <i>Internat-</i>	Noriko Kando, Evangelos Kanoulas, Martin Potthast,	950
895	<i>ional Conference on Machine Learning</i> , pages 7480–	Benno Stein, et al. 2023. Perspectives on large lan-	951
896	7512. PMLR.	guage models for relevance judgment. In <i>Proceeed-</i>	952
897	Gelei Deng, Yi Liu, Kailong Wang, Yuekang Li, Tian-	<i>ings of the 2023 ACM SIGIR International Confer-</i>	953
898	wei Zhang, and Yang Liu. 2024. Pandora: Jailbreak	<i>ence on Theory of Information Retrieval</i> , pages 39–	954
899	gpts by retrieval augmented generation poisoning.	50.	955
900	<i>arXiv preprint arXiv:2402.08416</i> .	Guglielmo Faggioli, Laura Dietz, Charles LA Clarke,	956
901	Zehang Deng, Yongjian Guo, Changzhou Han, Wan-	Gianluca Demartini, Matthias Hagen, Claudia Hauff,	957
902	lun Ma, Junwu Xiong, Sheng Wen, and Yang Xiang.	Noriko Kando, Evangelos Kanoulas, Martin Potthast,	958
903	2025. Ai agents under threat: A survey of key secu-	Benno Stein, et al. 2024. Who determines what is	959
904	rity challenges and future pathways. <i>ACM Comput-</i>	relevant? humans or ai? why not both? <i>Communica-</i>	960
905	<i>ing Surveys</i> , 57(7):1–36.	<i>tions of the ACM</i> , 67(4):31–34.	961
906	Tim Dettmers and Luke Zettlemoyer. 2023. The case for	Yan Fang, Jingtao Zhan, Qingyao Ai, Jiaxin Mao, Wei-	962
907	4-bit precision: k-bit inference scaling laws. In <i>In-</i>	hang Su, Jia Chen, and Yiqun Liu. 2024. Scaling	963
908	<i>ternational Conference on Machine Learning</i> , pages	laws for dense retrieval. In <i>Proceedings of the 47th</i>	964
909	7750–7774. PMLR.	<i>International ACM SIGIR Conference on Research</i>	965
910	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	<i>and Development in Information Retrieval</i> , pages	966
911	Kristina Toutanova. 2019. BERT: Pre-training of	1339–1349.	967
912	deep bidirectional transformers for language under-	Maxim Fishman, Brian Chmiel, Ron Banner, and Daniel	968
913	standing . In <i>Proceedings of the 2019 Conference of</i>	Soudry. 2024. Scaling fp8 training to trillion-token	969
914	<i>the North American Chapter of the Association for</i>	llms. <i>arXiv preprint arXiv:2409.12517</i> .	970
915	<i>Computational Linguistics: Human Language Tech-</i>	Thibault Formal, Stéphane Clinchant, Hervé Déjean,	971
916	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	and Carlos Lassance. 2024. Splate: Sparse late in-	972
917	4171–4186, Minneapolis, Minnesota. Association for	teraction retrieval. In <i>Proceedings of the 47th Inter-</i>	973
918	Computational Linguistics.	<i>national ACM SIGIR Conference on Research and</i>	974
919	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	<i>Development in Information Retrieval</i> , pages 2635–	975
920	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	2640.	976
921	Akhil Mathur, Alan Schelten, Amy Yang, Angela	Thibault Formal, Carlos Lassance, Benjamin Pi-	977
922	Fan, et al. 2024. The llama 3 herd of models. <i>arXiv</i>	wowarski, and Stéphane Clinchant. 2021a. Splade	978
923	<i>preprint arXiv:2407.21783</i> .	v2: Sparse lexical and expansion model for informa-	979
924	Darren Edge, Ha Trinh, Newman Cheng, Joshua	tion retrieval . <i>arXiv preprint</i> .	980
925	Bradley, Alex Chao, Apurva Mody, Steven Truitt,	Thibault Formal, Benjamin Piwowarski, and Stéphane	981
926	and Jonathan Larson. 2024. From local to global: A	Clinchant. 2021b. SPLADE: Sparse Lexical and	982
927	graph rag approach to query-focused summarization.	Expansion Model for First Stage Ranking , page	983
928	<i>arXiv preprint arXiv:2404.16130</i> .	2288–2292. Association for Computing Machinery,	984
929	Nelson Elhage, Neel Nanda, Catherine Olsson, Tom	New York, NY, USA.	985
930	Henighan, Nicholas Joseph, Ben Mann, Amanda	Yoav Freund, Raj Iyer, Robert E Schapire, and Yoram	986
931	Askell, Yuntao Bai, Anna Chen, Tom Conerly,	Singer. 2003. An efficient boosting algorithm for	987

988	combining preferences. <i>Journal of machine learning research</i> , 4(Nov):933–969.	1043
989		1044
990	Jerome H Friedman. 2001. Greedy function approx-	1045
991	imation: a gradient boosting machine. <i>Annals of</i>	1046
992	<i>statistics</i> , pages 1189–1232.	1047
993	Norbert Fuhr. 1992. Probabilistic models in information	1048
994	retrieval. <i>The computer journal</i> , 35(3):243–255.	1049
995	Luyu Gao and Jamie Callan. 2021. Condenser: a pre-	1050
996	training architecture for dense retrieval . In <i>Proceed-</i>	1051
997	<i>ings of the 2021 Conference on Empirical Methods</i>	1052
998	<i>in Natural Language Processing</i> , pages 981–993,	
999	Online and Punta Cana, Dominican Republic. Asso-	
1000	ciation for Computational Linguistics.	
1001	Luyu Gao, Zhuyun Dai, and Jamie Callan. 2020. Under-	
1002	standing bert rankers under distillation. In <i>Proceed-</i>	
1003	<i>ings of the 2020 ACM SIGIR on International Con-</i>	
1004	<i>ference on Theory of Information Retrieval</i> , pages	
1005	149–152.	
1006	Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021a.	
1007	COIL: Revisit exact lexical match in information	
1008	retrieval with contextualized inverted list . In <i>Pro-</i>	
1009	<i>ceedings of the 2021 Conference of the North Amer-</i>	
1010	<i>ican Chapter of the Association for Computational</i>	
1011	<i>Linguistics: Human Language Technologies</i> , pages	
1012	3030–3042, Online. Association for Computational	
1013	Linguistics.	
1014	Luyu Gao, Zhuyun Dai, Tongfei Chen, Zhen Fan, Ben-	
1015	jamin Van Durme, and Jamie Callan. 2021b. Comple-	
1016	ment lexical retrieval model with semantic residual	
1017	embeddings. In <i>Advances in Information Retrieval:</i>	
1018	<i>43rd European Conference on IR Research, ECIR</i>	
1019	<i>2021, Virtual Event, March 28–April 1, 2021, Pro-</i>	
1020	<i>ceedings, Part I 43</i> , pages 146–160. Springer.	
1021	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021c.	
1022	SimCSE: Simple contrastive learning of sentence em-	
1023	beddings . In <i>Proceedings of the 2021 Conference</i>	
1024	<i>on Empirical Methods in Natural Language Process-</i>	
1025	<i>ing</i> , pages 6894–6910, Online and Punta Cana, Do-	
1026	minican Republic. Association for Computational	
1027	Linguistics.	
1028	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,	
1029	Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen	
1030	Wang. 2023. Retrieval-augmented generation for	
1031	large language models: A survey. <i>arXiv preprint</i>	
1032	<i>arXiv:2312.10997</i> .	
1033	Yunfan Gao, Yun Xiong, Yijie Zhong, Yuxi Bi, Ming	
1034	Xue, and Haofen Wang. 2025. Synergizing rag and	
1035	reasoning: A systematic review. <i>arXiv preprint</i>	
1036	<i>arXiv:2504.15909</i> .	
1037	Gemini, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste	
1038	Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk,	
1039	Andrew M Dai, Anja Hauth, Katie Millican, et al.	
1040	2023. Gemini: a family of highly capable multi-	
1041	modal models. <i>arXiv preprint arXiv:2312.11805</i> .	
1042	Google. 2024. Grounding with google search .	
	AI Blog Google. 2019. Understanding searches better	1043
	than ever before .	1044
	Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen,	1045
	Yujiu Yang, Nan Duan, and Weizhu Chen. 2023.	1046
	CRITIC: Large language models can self-correct	1047
	with tool-interactive critiquing . In <i>Second Agent</i>	1048
	<i>Learning in Open-Endedness Workshop</i> .	1049
	Albert Gu and Tri Dao. 2024. Mamba: Linear-time	1050
	sequence modeling with selective state spaces . In	1051
	<i>First Conference on Language Modeling</i> .	1052
	Xinyan Guan, Jiali Zeng, Fandong Meng, Chunlei Xin,	1053
	Yaojie Lu, Hongyu Lin, Xianpei Han, Le Sun, and	1054
	Jie Zhou. 2025. Deeprag: Thinking to retrieval step	1055
	by step for large language models. <i>arXiv preprint</i>	1056
	<i>arXiv:2502.01142</i> .	1057
	Michael Günther, Jackmin Ong, Isabelle Mohr, Alaed-	1058
	dine Abdessalem, Tanguy Abel, Mohammad Kalim	1059
	Akram, Susana Guzman, Georgios Mastrapas, Saba	1060
	Sturua, Bo Wang, et al. 2023. Jina embeddings 2:	1061
	8192-token general-purpose text embeddings for long	1062
	documents. <i>arXiv preprint arXiv:2310.19923</i> .	1063
	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,	1064
	Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,	1065
	Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: In-	1066
	centivizing reasoning capability in llms via reinforce-	1067
	ment learning. <i>arXiv preprint arXiv:2501.12948</i> .	1068
	Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce	1069
	Croft. 2016a. A deep relevance matching model for	1070
	ad-hoc retrieval . In <i>Proceedings of the 25th ACM In-</i>	1071
	<i>ternational on Conference on Information and Knowl-</i>	1072
	<i>edge Management, CIKM '16</i> , page 55–64, New	1073
	York, NY, USA. Association for Computing Machin-	1074
	ery.	1075
	Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce	1076
	Croft. 2016b. Semantic matching by non-linear word	1077
	transportation for information retrieval . In <i>Proceed-</i>	1078
	<i>ings of the 25th ACM International on Conference</i>	1079
	<i>on Information and Knowledge Management, CIKM</i>	1080
	<i>'16</i> , page 701–710, New York, NY, USA. Association	1081
	for Computing Machinery.	1082
	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang,	1083
	Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xi-	1084
	angliang Zhang. 2024. Large language model based	1085
	multi-agents: A survey of progress and challenges.	1086
	<i>arXiv preprint arXiv:2402.01680</i> .	1087
	Ashim Gupta, Rishanth Rajendhran, Nathan Stringham,	1088
	Vivek Srikumar, and Ana Marasovic. 2024. Whispers	1089
	of doubt amidst echoes of triumph in NLP robustness .	1090
	In <i>Proceedings of the 2024 Conference of the North</i>	1091
	<i>American Chapter of the Association for Computa-</i>	1092
	<i>tional Linguistics: Human Language Technologies</i>	1093
	<i>(Volume 1: Long Papers)</i> , pages 5533–5590, Mexico	1094
	City, Mexico. Association for Computational Lin-	1095
	guistics.	1096
	Suchin Gururangan, Ana Marasović, Swabha	1097
	Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,	1098

1099	and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8342–8360, Online. Association for Computational Linguistics.	1153
1100		1154
1101		1155
1102		
1103		
1104		
1105	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasu-	
1106	pat, and Mingwei Chang. 2020. Retrieval augmented	
1107	language model pre-training. In <i>International confer-</i>	
1108	<i>ence on machine learning</i> , pages 3929–3938. PMLR.	
1109	David Ha, Andrew M Dai, and Quoc V Le. 2022. Hyper-	
1110	networks. In <i>International Conference on Learning</i>	
1111	<i>Representations</i> .	
1112	Yichen He, Guanhua Huang, Peiyuan Feng, Yuan Lin,	
1113	Yuchen Zhang, Hang Li, et al. 2025. Pasa: An	
1114	llm agent for comprehensive academic paper search.	
1115	<i>arXiv preprint arXiv:2501.10120</i> .	
1116	S Hochreiter and J Schmidhuber. 1997. Long short-term	
1117	memory. <i>Neural Computation MIT-Press</i> .	
1118	Jordan Hoffmann, Sebastian Borgeaud, Arthur Men-	
1119	sch, Elena Buchatskaya, Trevor Cai, Eliza Ruther-	
1120	ford, Diego de Las Casas, Lisa Anne Hendricks,	
1121	Johannes Welbl, Aidan Clark, et al. 2022. Train-	
1122	ing compute-optimal large language models. <i>arXiv</i>	
1123	<i>preprint arXiv:2203.15556</i> .	
1124	Sebastian Hofstätter, Sophia Althammer, Michael	
1125	Schröder, Mete Sertkan, and Allan Hanbury. 2020a.	
1126	Improving efficient neural ranking models with cross-	
1127	architecture knowledge distillation. <i>arXiv preprint</i>	
1128	<i>arXiv:2010.02666</i> .	
1129	Sebastian Hofstätter, Omar Khattab, Sophia Althammer,	
1130	Mete Sertkan, and Allan Hanbury. 2022. Introducing	
1131	neural bag of whole-words with colberter: Context-	
1132	tualized late interactions using enhanced reduction.	
1133	In <i>Proceedings of the 31st ACM International Confer-</i>	
1134	<i>ence on Information & Knowledge Management</i> ,	
1135	pages 737–747.	
1136	Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong	
1137	Yang, Jimmy Lin, and Allan Hanbury. 2021. Ef-	
1138	ficiently teaching an effective dense retriever with	
1139	balanced topic aware sampling. In <i>Proceedings of</i>	
1140	<i>the 44th International ACM SIGIR Conference on</i>	
1141	<i>Research and Development in Information Retrieval</i> ,	
1142	pages 113–122.	
1143	Sebastian Hofstätter, Hamed Zamani, Bhaskar Mitra,	
1144	Nick Craswell, and Allan Hanbury. 2020b. Local	
1145	self-attention over long text for efficient document	
1146	retrieval. In <i>Proceedings of the 43rd International</i>	
1147	<i>ACM SIGIR Conference on Research and Develop-</i>	
1148	<i>ment in Information Retrieval</i> , pages 2021–2024.	
1149	Sebastian Hofstätter, Markus Zlabinger, and Allan	
1150	Hanbury. 2020c. Interpretable & time-budget-	
1151	constrained contextualization for re-ranking. In	
1152	<i>ECAI 2020</i> , pages 513–520. IOS Press.	
	Jeremy Howard and Sebastian Ruder. 2018. Universal	1153
	language model fine-tuning for text classification.	1154
	<i>arXiv preprint arXiv:1801.06146</i> .	1155
	Baotian Hu, Zhengdong Lu, Hang Li, and Qingcai Chen.	1156
	2014. Convolutional neural network architectures for	1157
	matching natural language sentences. In <i>Proceedings</i>	1158
	<i>of the 28th International Conference on Neural In-</i>	1159
	<i>formation Processing Systems - Volume 2, NIPS'14</i> ,	1160
	page 2042–2050, Cambridge, MA, USA. MIT Press.	1161
	Yunhai Hu, Yilun Zhao, Chen Zhao, and Arman Cohan.	1162
	2025. Mcts-rag: Enhancing retrieval-augmented gen-	1163
	eration with monte carlo tree search. <i>arXiv preprint</i>	1164
	<i>arXiv:2503.20757</i> .	1165
	Ziniu Hu, Yang Wang, Qu Peng, and Hang Li. 2019. Un-	1166
	biased lambdamart: an unbiased pairwise learning-to-	1167
	rank algorithm. In <i>The World Wide Web Conference</i> ,	1168
	pages 2830–2836.	1169
	Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia,	1170
	David Zhang, Philip Pronin, Janani Padmanab-	1171
	han, Giuseppe Ottaviano, and Linjun Yang. 2020.	1172
	Embedding-based retrieval in facebook search. In	1173
	<i>Proceedings of the 26th ACM SIGKDD International</i>	1174
	<i>Conference on Knowledge Discovery & Data Mining</i> ,	1175
	pages 2553–2561.	1176
	Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng,	1177
	Alex Acero, and Larry Heck. 2013. Learning deep	1178
	structured semantic models for web search using	1179
	clickthrough data . In <i>Proceedings of the 22nd ACM</i>	1180
	<i>International Conference on Information & Knowl-</i>	1181
	<i>edge Management, CIKM '13</i> , page 2333–2338, New	1182
	York, NY, USA. Association for Computing Machin-	1183
	ery.	1184
	Zhiqi Huang, Puxuan Yu, Shauli Ravfogel, and James	1185
	Allan. 2024. Language concept erasure for language-	1186
	invariant dense retrieval . In <i>Proceedings of the</i>	1187
	<i>2024 Conference on Empirical Methods in Natural</i>	1188
	<i>Language Processing</i> , pages 13261–13273, Miami,	1189
	Florida, USA. Association for Computational Lin-	1190
	guistics.	1191
	Samuel Humeau, Kurt Shuster, Marie-Anne Lachaux,	1192
	and Jason Weston. 2020. Poly-encoders: Architec-	1193
	tures and pre-training strategies for fast and accurate	1194
	multi-sentence scoring . In <i>International Conference</i>	1195
	<i>on Learning Representations</i> .	1196
	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam	1197
	Perelman, Aditya Ramesh, Aidan Clark, AJ Os-	1198
	trow, Akila Welihinda, Alan Hayes, Alec Radford,	1199
	et al. 2024. Gpt-4o system card. <i>arXiv preprint</i>	1200
	<i>arXiv:2410.21276</i> .	1201
	Shayekh Bin Islam, Md Asib Rahman, K S M Tozammel	1202
	Hossain, Enamul Hoque, Shafiq Joty, and Md Rizwan	1203
	Parvez. 2024. Open-RAG: Enhanced retrieval aug-	1204
	mented reasoning with open-source large language	1205
	models . In <i>Findings of the Association for Compu-</i>	1206
	<i>tational Linguistics: EMNLP 2024</i> , pages 14231–	1207
	14244, Miami, Florida, USA. Association for Com-	1208
	putational Linguistics.	1209

- Gautier Izacard and Edouard Grave. 2021. [Leveraging passage retrieval with generative models for open domain question answering](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 874–880, Online. Association for Computational Linguistics.
- Alon Jacovi and Yoav Goldberg. 2020. [Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness?](#) In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4198–4205, Online. Association for Computational Linguistics.
- Minbyul Jeong, Jiwoong Sohn, Mujeen Sung, and Jaewoo Kang. 2024. Improving medical reasoning through retrieval and self-reflection with retrieval-augmented large language models. *Bioinformatics*, 40(Supplement_1):i119–i129.
- Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7969–7992.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*.
- Thorsten Joachims. 2006. [Training linear svms in linear time](#). In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’06, page 217–226, New York, NY, USA. Association for Computing Machinery.
- Thorsten Joachims, Laura Granka, Bing Pan, Helene Hembrooke, and Geri Gay. 2017a. Accurately interpreting clickthrough data as implicit feedback. In *Acm Sigir Forum*, volume 51, pages 4–11. Acm New York, NY, USA.
- Thorsten Joachims, Adith Swaminathan, and Tobias Schnabel. 2017b. Unbiased learning-to-rank with biased feedback. In *Proceedings of the tenth ACM international conference on web search and data mining*, pages 781–789.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547.
- Armand Joulin, Piotr Bojanowski, Tomáš Mikolov, Hervé Jégou, and Édouard Grave. 2018. Loss in translation: Learning bilingual word mapping with a retrieval criterion. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2979–2984.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. 2021. Highly accurate protein structure prediction with alphafold. *nature*, 596(7873):583–589.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Vladimir Karpukhin, Barlas Öguz, Sewon Min, Patrick Lewis, Ledell Yu Wu, Sergey Edunov, Danqi Chen, and Wen tau Yih. 2020. [Dense passage retrieval for open-domain question answering](#). *ArXiv*, abs/2004.04906.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. 2020. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pages 5156–5165. PMLR.
- Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020a. [Generalization through memorization: Nearest neighbor language models](#). In *International Conference on Learning Representations*.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020b. [Generalization through memorization: Nearest neighbor language models](#). In *International Conference on Learning Representations*.
- 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*, pages 39–48.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2023. [Decomposed prompting: A modular approach for solving complex tasks](#). In *The Eleventh International Conference on Learning Representations*.
- Julian Killingback, Hansi Zeng, and Hamed Zamani. 2025. Hypencoder: Hypernetworks for information retrieval. *arXiv preprint arXiv:2502.05364*.
- Weize Kong, Jeffrey M Dudek, Cheng Li, Mingyang Zhang, and Michael Bendersky. 2023. Sparseembed: Learning sparse lexical representations with contextual embeddings for retrieval. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2399–2403.

1321	Donald H Kraft and Duncan A Buell. 1983. Fuzzy sets	1378
1322	and generalized boolean retrieval systems. <i>International journal of man-machine studies</i> , 19(1):45–56.	1379
1323		1380
1324	Tanishq Kumar, Zachary Ankner, Benjamin F Spector,	1381
1325	Blake Bordelon, Niklas Muennighoff, Mansheej Paul,	
1326	Cengiz Pehlevan, Christopher Ré, and Aditi Raghunathan. 2024. Scaling laws for precision. <i>arXiv preprint arXiv:2411.04330</i> .	1382
1327		1383
1328		1384
1329	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	1385
1330	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient	
1331	memory management for large language model serving with pagedattention. In <i>Proceedings of the 29th Symposium on Operating Systems Principles</i> , pages	1386
1332	611–626.	1387
1333		1388
1334	Yann LeCun, Bernhard Boser, John S Denker, Donnie	1389
1335	Henderson, Richard E Howard, Wayne Hubbard, and	1390
1336	Lawrence D Jackel. 1989. Backpropagation applied	
1337	to handwritten zip code recognition. <i>Neural computation</i> , 1(4):541–551.	1391
1338		1392
1339		1393
1340		
1341	Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick	
1342	Haffner. 1998. Gradient-based learning applied to	
1343	document recognition. <i>Proceedings of the IEEE</i> ,	
1344	86(11):2278–2324.	
1345	Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan	
1346	Raiman, Mohammad Shoenybi, Bryan Catanzaro, and	
1347	Wei Ping. 2025. NV-embed: Improved techniques for training LLMs as generalist embedding models .	
1348	In <i>The Thirteenth International Conference on Learning Representations</i> .	
1349		
1350		
1351	Jinhyuk Lee, Zhuyun Dai, Xiaoqi Ren, Blair Chen,	
1352	Daniel Cer, Jeremy R Cole, Kai Hui, Michael Boratko, Rajvi Kapadia, Wen Ding, et al. 2024. Gecko:	
1353	Versatile text embeddings distilled from large language models. <i>arXiv preprint arXiv:2403.20327</i> .	
1354		
1355		
1356	Kenton Lee, Ming-Wei Chang, and Kristina Toutanova.	
1357	2019. Latent retrieval for weakly supervised open domain question answering . <i>ArXiv</i> , abs/1906.00300.	
1358		
1359	Omer Levy, Yoav Goldberg, and Ido Dagan. 2015. Improving distributional similarity with lessons learned from word embeddings . <i>Transactions of the Association for Computational Linguistics</i> , 3:211–225.	
1360		
1361		
1362		
1363	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	
1364	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	
1365	Veselin Stoyanov, and Luke Zettlemoyer. 2020a. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	
1366		
1367		
1368		
1369		
1370		
1371		
1372	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	
1373	Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020b. Retrieval-augmented generation	
1374	for knowledge-intensive nlp tasks. <i>Advances in neural information processing systems</i> , 33:9459–9474.	
1375		
1376		
1377		
	Canjia Li, Andrew Yates, Sean MacAvaney, Ben He, and	
	Yingfei Sun. 2020. Parade: Passage representation aggregation for document reranking . <i>ACM Transactions on Information Systems</i> , 42:1 – 26.	
	Chaofan Li, Zheng Liu, Shitao Xiao, and Yingxia	
	Shao. 2023a. Making large language models a better foundation for dense retrieval. <i>arXiv preprint arXiv:2312.15503</i> .	
	Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii	
	Khizbullin, and Bernard Ghanem. 2023b. Camel: Communicative agents for "mind" exploration of large language model society. <i>Advances in Neural Information Processing Systems</i> , 36:51991–52008.	
	Hang Li. 2011. Learning for ranking aggregation. In <i>Learning to Rank for Information Retrieval and Natural Language Processing</i> , pages 33–35. Springer.	
	Millicent Li, Tongfei Chen, Benjamin Van Durme, and	
	Patrick Xia. 2024. Multi-field adaptive retrieval. <i>arXiv preprint arXiv:2410.20056</i> .	
	Minghan Li, Sheng-Chieh Lin, Xueguang Ma, and	
	Jimmy Lin. 2023c. Slim: Sparsified late interaction for multi-vector retrieval with inverted indexes . In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '23</i> , page 1954–1959, New York, NY, USA. Association for Computing Machinery.	
	Sen Li, Fuyu Lv, Taiwei Jin, Guli Lin, Keping Yang,	
	Xiaoyi Zeng, Xiao-Ming Wu, and Qianli Ma. 2021. Embedding-based product retrieval in taobao search. In <i>Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining</i> , pages 3181–3189.	
	Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang,	
	Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025a. Search-o1: Agentic search-enhanced large reasoning models. <i>arXiv preprint arXiv:2501.05366</i> .	
	Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yutao Zhu, Yongkang Wu, Ji-Rong Wen, and Zhicheng Dou. 2025b. Webthinker: Empowering large reasoning models with deep research capability. <i>arXiv preprint arXiv:2504.21776</i> .	
	Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yuyao Zhang, Peitian	
	Zhang, Yutao Zhu, and Zhicheng Dou. 2025c. From matching to generation: A survey on generative information retrieval . <i>ACM Trans. Inf. Syst.</i> , 43(3).	
	Xinze Li, Zhenghao Liu, Chenyan Xiong, Shi Yu,	
	Yu Gu, Zhiyuan Liu, and Ge Yu. 2023d. Structure-aware language model pretraining improves dense retrieval on structured data. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 11560–11574.	

1431	Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long,	Sean MacAvaney, Sergey Feldman, Nazli Goharian,	1486
1432	Pengjun Xie, and Meishan Zhang. 2023e. Towards	Doug Downey, and Arman Cohan. 2022. ABNIRML:	1487
1433	general text embeddings with multi-stage contrastive	Analyzing the behavior of neural IR models . <i>Trans-</i>	1488
1434	learning. <i>arXiv preprint arXiv:2308.03281</i> .	<i>actions of the Association for Computational Linguis-</i>	1489
		<i>tics</i> , 10:224–239.	1490
1435	Jimmy Lin and Xueguang Ma. 2021. A few brief notes	Sean MacAvaney, Andrew Yates, Arman Cohan, and	1491
1436	on deepimpact, coil, and a conceptual framework	Nazli Goharian. 2019. Cedr: Contextualized em-	1492
1437	for information retrieval techniques. <i>arXiv preprint</i>	beddings for document ranking. In <i>Proceedings of</i>	1493
1438	<i>arXiv:2106.14807</i> .	<i>the 42nd international ACM SIGIR conference on</i>	1494
1439	Jimmy Lin, Rodrigo Nogueira, and Andrew Yates. 2022.	<i>research and development in information retrieval</i> ,	1495
1440	<i>Pretrained transformers for text ranking: Bert and</i>	pages 1101–1104.	1496
1441	<i>beyond</i> . Springer Nature.		
1442	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	1497
1443	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	1498
1444	Deng, Chenyu Zhang, Chong Ruan, et al. 2024a.	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	1499
1445	Deepseek-v3 technical report. <i>arXiv preprint</i>	et al. 2023. Self-refine: Iterative refinement with	1500
1446	<i>arXiv:2412.19437</i> .	self-feedback. <i>Advances in Neural Information Pro-</i>	1501
		<i>cessing Systems</i> , 36:46534–46594.	1502
1447	Qi Liu, Bo Wang, Nan Wang, and Jiaxin Mao. 2024b.	Yury Malkov and Dmitry A. Yashunin. 2016. Efficient	1503
1448	Leveraging passage embeddings for efficient listwise	and robust approximate nearest neighbor search us-	1504
1449	reranking with large language models. In <i>THE WEB</i>	ing hierarchical navigable small world graphs . <i>IEEE</i>	1505
1450	<i>CONFERENCE 2025</i> .	<i>Transactions on Pattern Analysis and Machine Intel-</i>	1506
1451	Tie-Yan Liu. 2009. Learning to rank for information	<i>ligence</i> , 42:824–836.	1507
1452	retrieval . <i>Found. Trends Inf. Retr.</i> , 3(3):225–331.		
1453	Yinhan Liu. 2019. Roberta: A robustly opti-	Antonio Mallia, Omar Khattab, Torsten Suel, and Nicola	1508
1454	mized bert pretraining approach. <i>arXiv preprint</i>	Tonellotto. 2021. Learning passage impacts for in-	1509
1455	<i>arXiv:1907.11692</i> , 364.	verted indexes. In <i>Proceedings of the 44th Inter-</i>	1510
1456	Yu-An Liu, Ruqing Zhang, Jiafeng Guo, and Maarten	<i>national ACM SIGIR Conference on Research and</i>	1511
1457	de Rijke. 2025. Robust information retrieval. In	<i>Development in Information Retrieval</i> , pages 1723–	1512
1458	<i>Proceedings of the Eighteenth ACM International</i>	1727.	1513
1459	<i>Conference on Web Search and Data Mining</i> , pages	Tomas Mikolov. 2013. Efficient estimation of word	1514
1460	1008–1011.	representations in vector space. <i>arXiv preprint</i>	1515
		<i>arXiv:1301.3781</i> , 3781.	1516
1461	Yi Luan, Jacob Eisenstein, Kristina Toutanova, and	David RH Miller, Tim Leek, and Richard M Schwartz.	1517
1462	Michael Collins. 2021. Sparse, dense, and attentional	1999. A hidden markov model information retrieval	1518
1463	representations for text retrieval. <i>Transactions of the</i>	system. In <i>Proceedings of the 22nd annual inter-</i>	1519
1464	<i>Association for Computational Linguistics</i> , 9:329–	<i>national ACM SIGIR conference on Research and</i>	1520
1465	345.	<i>development in information retrieval</i> , pages 214–221.	1521
1466	Xueguang Ma, Sheng-Chieh Lin, Minghan Li, Wenhu	Bhaskar Mitra, Nick Craswell, et al. 2018. An introduc-	1522
1467	Chen, and Jimmy Lin. 2024a. Unifying multimodal	tion to neural information retrieval. <i>Foundations and</i>	1523
1468	retrieval via document screenshot embedding . In <i>Pro-</i>	<i>Trends® in Information Retrieval</i> , 13(1):1–126.	1524
1469	<i>ceedings of the 2024 Conference on Empirical Meth-</i>	Bhaskar Mitra, Sebastian Hofstätter, Hamed Zamani,	1525
1470	<i>ods in Natural Language Processing</i> , pages 6492–	and Nick Craswell. 2021. Improving transformer-	1526
1471	6505, Miami, Florida, USA. Association for Compu-	kernel ranking model using conformer and query	1527
1472	tational Linguistics.	term independence. In <i>Proceedings of the 44th Inter-</i>	1528
1473	Xueguang Ma, Kai Sun, Ronak Pradeep, and Jimmy Lin.	<i>national ACM SIGIR Conference on Research and</i>	1529
1474	2021. A replication study of dense passage retriever.	<i>Development in Information Retrieval</i> , pages 1697–	1530
1475	<i>arXiv preprint arXiv:2104.05740</i> .	1702.	1531
1476	Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and	Fengran Mo, Abbas Ghaddar, Kelong Mao, Mehdi Reza-	1532
1477	Jimmy Lin. 2024b. Fine-tuning llama for multi-	gholizadeh, Boxing Chen, Qun Liu, and Jian-Yun Nie.	1533
1478	stage text retrieval. In <i>Proceedings of the 47th Inter-</i>	2024a. Chiq: Contextual history enhancement for	1534
1479	<i>national ACM SIGIR Conference on Research and</i>	improving query rewriting in conversational search.	1535
1480	<i>Development in Information Retrieval</i> , pages 2421–	<i>arXiv preprint arXiv:2406.05013</i> .	1536
1481	2425.		
1482	Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and	Fengran Mo, Kelong Mao, Ziliang Zhao, Hongjin	1537
1483	Jimmy Lin. 2023. Zero-shot listwise document	Qian, Haonan Chen, Yiruo Cheng, Xiaoxi Li, Yu-	1538
1484	reranking with a large language model . <i>Preprint</i> ,	tao Zhu, Zhicheng Dou, and Jian-Yun Nie. 2024b.	1539
1485	<i>arXiv:2305.02156</i> .	A survey of conversational search. <i>arXiv preprint</i>	1540
		<i>arXiv:2410.15576</i> .	1541

Fengran Mo, Kelong Mao, Yutao Zhu, Yihong Wu, Kaiyu Huang, and Jian-Yun Nie. 2023. *Convqqr: Generative query reformulation for conversational search*. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4998–5012. 1542 1543 1544 1545 1546 1547

Fengran Mo, Bole Yi, Kelong Mao, Chen Qu, Kaiyu Huang, and Jian-Yun Nie. 2024c. *Convsgd: Session data generation for conversational search*. In *Companion Proceedings of the ACM on Web Conference 2024*, pages 1634–1642. 1548 1549 1550 1551 1552

Marco Morik, Ashudeep Singh, Jessica Hong, and Thorsten Joachims. 2020. Controlling fairness and bias in dynamic learning-to-rank. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, pages 429–438. 1553 1554 1555 1556 1557 1558

Niklas Muennighoff. 2022. *Sgpt: Gpt sentence embeddings for semantic search*. *arXiv preprint arXiv:2202.08904*. 1559 1560 1561

Niklas Muennighoff, Hongjin SU, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. 2025. *Generative representational instruction tuning*. In *The Thirteenth International Conference on Learning Representations*. 1562 1563 1564 1565 1566

Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. *MTEB: Massive text embedding benchmark*. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2014–2037, Dubrovnik, Croatia. Association for Computational Linguistics. 1567 1568 1569 1570 1571 1572

Haya Nachimovsky, Moshe Tennenholtz, and Oren Kurland. 2025. A multi-agent perspective on modern information retrieval. *arXiv preprint arXiv:2502.14796*. 1573 1574 1575 1576

Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. *Webgpt: Browser-assisted question-answering with human feedback*. *arXiv preprint arXiv:2112.09332*. 1577 1578 1579 1580 1581 1582

Eric T. Nalisnick, Bhaskar Mitra, Nick Craswell, and Rich Caruana. 2016. *Improving document ranking with dual word embeddings*. *Proceedings of the 25th International Conference Companion on World Wide Web*. 1583 1584 1585 1586 1587

Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. 2022. Text and code embeddings by contrastive pre-training. *arXiv preprint arXiv:2201.10005*. 1588 1589 1590 1591 1592

Jianmo Ni, Gustavo Hernandez Abrego, Noah Constant, Ji Ma, Keith Hall, Daniel Cer, and Yinfei Yang. 2022a. *Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models*. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1864–1874, Dublin, Ireland. Association for Computational Linguistics. 1593 1594 1595 1596 1597

Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith Hall, Ming-Wei Chang, and Yinfei Yang. 2022b. *Large dual encoders are generalizable retrievers*. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9844–9855, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. 1600 1601 1602 1603 1604 1605 1606 1607

Jian-Yun Nie. 2010. *Cross-language information retrieval*. Morgan & Claypool Publishers. 1608 1609

Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. 2020. *Document ranking with a pretrained sequence-to-sequence model*. *Preprint*, arXiv:2003.06713. 1610 1611 1612

Rodrigo Nogueira, Wei Yang, Kyunghyun Cho, and Jimmy Lin. 2019. Multi-stage document ranking with bert. *arXiv preprint arXiv:1910.14424*. 1613 1614 1615

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)*, pages 1722–1727, Online. Association for Computational Linguistics. 1616 1617 1618 1619 1620 1621 1622

Zach Nussbaum, John Xavier Morris, Andriy Mulyar, and Brandon Duderstadt. 2025. *Nomic embed: Training a reproducible long context text embedder*. *Transactions on Machine Learning Research*. Reproducibility Certification. 1623 1624 1625 1626 1627

Team Nvidia. 2021. Accelerating inference with sparsity using the nvidia ampere architecture and nvidia tensorrt. In *NVIDIA Technical Blog*. 1628 1629 1630

Hanseok Oh, Hyunji Lee, Seonghyeon Ye, Haebin Shin, Hansol Jang, Changwook Jun, and Minjoon Seo. 2024. *Instructir: A benchmark for instruction following of information retrieval models*. *arXiv preprint arXiv:2402.14334*. 1631 1632 1633 1634 1635

Kezban Dilek Onal, Ye Zhang, Ismail Sengor Altinogode, Md Mustafizur Rahman, Pinar Karagoz, Alex Braylan, Brandon Dang, Heng-Lu Chang, Henna Kim, Quinten McNamara, et al. 2018. Neural information retrieval: At the end of the early years. *Information Retrieval Journal*, 21:111–182. 1636 1637 1638 1639 1640 1641

Harrie Oosterhuis. 2021. Computationally efficient optimization of plackett-luce ranking models for relevance and fairness. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1023–1032. 1642 1643 1644 1645 1646 1647

OpenAI. 2022. *Introducing chatgpt*. 1648

OpenAI. 2023. *Gpt-4 technical report*. *arXiv preprint arXiv:2303.08774*. 1649 1650

OpenAI. 2024. *Introducing chatgpt search*. 1651

1652	Hamid Palangi, Li Deng, Yelong Shen, Jianfeng Gao,	Ronak Pradeep, Kai Hui, Jai Gupta, Adam Lelkes, Hon-	1708
1653	Xiaodong He, Jianshu Chen, Xinying Song, and	glei Zhuang, Jimmy Lin, Donald Metzler, and Vinh	1709
1654	Rabab Ward. 2016. Deep sentence embedding using	Tran. 2023a. How does generative retrieval scale to	1710
1655	long short-term memory networks: analysis and ap-	millions of passages? In <i>Proceedings of the 2023</i>	1711
1656	plication to information retrieval . <i>IEEE/ACM Trans.</i>	<i>Conference on Empirical Methods in Natural Lan-</i>	1712
1657	<i>Audio, Speech and Lang. Proc.</i> , 24(4):694–707.	<i>guage Processing</i> , pages 1305–1321, Singapore. As-	1713
		sociation for Computational Linguistics.	1714
1658	Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Shengx-	Ronak Pradeep, Rodrigo Nogueira, and Jimmy Lin.	1715
1659	ian Wan, and Xueqi Cheng. 2016. Text matching as	2021. The expando-mono-duo design pattern for	1716
1660	image recognition. In <i>Proceedings of the Thirtieth</i>	text ranking with pretrained sequence-to-sequence	1717
1661	<i>AAAI Conference on Artificial Intelligence, AAAI’16</i> ,	models. <i>arXiv preprint arXiv:2101.05667</i> .	1718
1662	page 2793–2799. AAAI Press.		
1663	Liang Pang, Jun Xu, Qingyao Ai, Yanyan Lan, Xueqi	Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy	1719
1664	Cheng, and Jirong Wen. 2020. Setrank: Learning a	Lin. 2023b. Rankvicuna: Zero-shot listwise docu-	1720
1665	permutation-invariant ranking model for information	ment reranking with open-source large language	1721
1666	retrieval. In <i>Proceedings of the 43rd international</i>	models . <i>Preprint</i> , arXiv:2309.15088.	1722
1667	<i>ACM SIGIR conference on research and development</i>		
1668	<i>in information retrieval</i> , pages 499–508.	Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy	1723
		Lin. 2023c. Rankzephyr: Effective and robust	1724
1669	Biswajit Paria, Chih-Kuan Yeh, Ian E.H. Yen, Ning	zero-shot listwise reranking is a breeze! <i>Preprint</i> ,	1725
1670	Xu, Pradeep Ravikumar, and Barnabás Póczos. 2020.	arXiv:2312.02724.	1726
1671	Minimizing flops to learn efficient sparse represen-		
1672	tations . In <i>International Conference on Learning</i>	Ernest Pusateri, Anmol Walia, Anirudh Kashi, Bortik	1727
1673	<i>Representations</i> .	Bandyopadhyay, Nadia Hyder, Sayantan Mahinder,	1728
		Raviteja Anantha, Daben Liu, and Sashank Gondala.	1729
1674	Gourab K Patro, Arpita Biswas, Niloy Ganguly, Kr-	2024. Retrieval augmented correction of named	1730
1675	ishna P Gummadi, and Abhijnan Chakraborty. 2020.	entity speech recognition errors. <i>arXiv preprint</i>	1731
1676	Fairrec: Two-sided fairness for personalized recom-	<i>arXiv:2409.06062</i> .	1732
1677	mendations in two-sided platforms. In <i>Proceedings</i>		
1678	<i>of the web conference 2020</i> , pages 1194–1204.	Tao Qin and Tie-Yan Liu. 2013. Introducing LETOR	1733
		4.0 datasets . <i>CoRR</i> , abs/1306.2597.	1734
1679	Bo Peng, Eric Alcaide, Quentin Anthony, Alon Al-	Tao Qin, Tie-Yan Liu, and Hang Li. 2010. A general	1735
1680	balak, Samuel Arcadinho, Stella Biderman, Huanqi	approximation framework for direct optimization of	1736
1681	Cao, Xin Cheng, Michael Chung, Leon Derczynski,	information retrieval measures. <i>Information retrieval</i> ,	1737
1682	Xingjian Du, Matteo Grella, Kranthi Gv, Xuzheng	13:375–397.	1738
1683	He, Haowen Hou, Przemyslaw Kazienko, Jan Ko-		
1684	con, Jiaming Kong, Bartłomiej Koptyra, Hayden	Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang,	1739
1685	Lau, Jiaju Lin, Krishna Sri Ipsit Mantri, Ferdinand	Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu	1740
1686	Mom, Atsushi Saito, Guangyu Song, Xiangru Tang,	Liu, Donald Metzler, Xuanhui Wang, and Michael	1741
1687	Johan Wind, Stanisław Woźniak, Zhenyuan Zhang,	Bendersky. 2024a. Large language models are effec-	1742
1688	Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. 2023.	tive text rankers with pairwise ranking prompting . In	1743
1689	RWKV: Reinventing RNNs for the transformer era .	<i>Findings of the Association for Computational Lin-</i>	1744
1690	In <i>Findings of the Association for Computational</i>	<i>guistics: NAACL 2024</i> , pages 1504–1518, Mexico	1745
1691	<i>Linguistics: EMNLP 2023</i> , pages 14048–14077, Sin-	City, Mexico. Association for Computational Lin-	1746
1692	gapore. Association for Computational Linguistics.	guistics.	1747
1693	Bo Peng, Daniel Goldstein, Quentin Anthony, Alon	Zhen Qin, Le Yan, Honglei Zhuang, Yi Tay, Rama Ku-	1748
1694	Albalak, Eric Alcaide, Stella Biderman, Eugene	mar Pasumarthi, Xuanhui Wang, Michael Bendersky,	1749
1695	Cheah, Xingjian Du, Teddy Ferdinan, Haowen Hou,	and Marc Najork. 2021. Are neural rankers still out-	1750
1696	et al. 2024. Eagle and finch: Rwkv with matrix-	performed by gradient boosted decision trees? In	1751
1697	valued states and dynamic recurrence. <i>arXiv preprint</i>	<i>International Conference on Learning Representa-</i>	1752
1698	<i>arXiv:2404.05892</i> .	<i>tions</i> .	1753
1699	Jeffrey Pennington, Richard Socher, and Christopher	Zhen Qin, Songlin Yang, and Yiran Zhong. 2024b. Hi-	1754
1700	Manning. 2014. GloVe: Global vectors for word	erarchically gated recurrent neural network for se-	1755
1701	representation . In <i>Proceedings of the 2014 Confer-</i>	quence modeling. <i>Advances in Neural Information</i>	1756
1702	<i>ence on Empirical Methods in Natural Language Pro-</i>	<i>Processing Systems</i> , 36.	1757
1703	<i>cessing (EMNLP)</i> , pages 1532–1543, Doha, Qatar.		
1704	Association for Computational Linguistics.	Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang	1758
1705	Jay M Ponte and W Bruce Croft. 1998. A language	Ren, Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng	1759
1706	modeling approach to information retrieval. In <i>Proc.</i>	Wang. 2020. Rocketqa: An optimized training ap-	1760
1707	<i>SIGIR 1998</i> , pages 275–281.	proach to dense passage retrieval for open-domain	1761
		question answering . In <i>North American Chapter of</i>	1762
		<i>the Association for Computational Linguistics</i> .	1763

1764	Tadeusz Radecki. 1979. Fuzzy set theoretical approach to document retrieval. <i>Information Processing & Management</i> , 15(5):247–259.	
1765		
1766		
1767	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners .	
1768		
1769		
1770	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	
1771		
1772		
1773		
1774		
1775		
1776	Razieh Rahimi, Youngwoo Kim, Hamed Zamani, and James Allan. 2021. Explaining documents’ relevance to search queries. <i>arXiv preprint arXiv:2111.01314</i> .	
1777		
1778		
1779	Ori Ram, Liat Bezalel, Adi Zicher, Yonatan Belinkov, Jonathan Berant, and Amir Globerson. 2023. What are you token about? dense retrieval as distributions over the vocabulary . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2481–2498, Toronto, Canada. Association for Computational Linguistics.	
1780		
1781		
1782		
1783		
1784		
1785		
1786		
1787	Shauli Ravfogel, Valentina Pyatkin, Amir David Nissan Cohen, Avshalom Manevich, and Yoav Goldberg. 2024. Description-based text similarity . In <i>First Conference on Language Modeling</i> .	
1788		
1789		
1790		
1791	Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, QiaoQiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. RocketQAv2: A joint training method for dense passage retrieval and passage re-ranking . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 2825–2835, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
1792		
1793		
1794		
1795		
1796		
1797		
1798		
1799	Yi Ren, Hongyan Tang, and Siwen Zhu. 2022. Unbiased learning to rank with biased continuous feedback. In <i>Proceedings of the 31st ACM International Conference on Information & Knowledge Management</i> , pages 1716–1725.	
1800		
1801		
1802		
1803		
1804	Stephen E Robertson and K Sparck Jones. 1976. Relevance weighting of search terms. <i>Journal of the American Society for Information science</i> , 27(3):129–146.	
1805		
1806		
1807		
1808	Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. 1995. Okapi at trec-3. <i>Nist Special Publication Sp</i> , 109:109.	
1809		
1810		
1811		
1812	Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works . <i>Transactions of the Association for Computational Linguistics</i> , 8:842–866.	
1813		
1814		
1815		
1816	Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. Improving passage retrieval with zero-shot question generation . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3781–3797, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1819
1817		1820
1818		1821
		1822
		1823
	Gerard Salton, Edward A Fox, and Harry Wu. 1983. Extended boolean information retrieval. <i>Communications of the ACM</i> , 26(11):1022–1036.	1824
		1825
		1826
	Gerard Salton, Anita Wong, and Chung-Shu Yang. 1975. A vector space model for automatic indexing. <i>Communications of the ACM</i> , 18(11):613–620.	1827
		1828
		1829
	V Sanh. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> .	1830
		1831
		1832
	Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, et al. 2022. Multitask prompted training enables zero-shot task generalization. In <i>International Conference on Learning Representations</i> .	1833
		1834
		1835
		1836
		1837
		1838
	Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2022. ColBERTv2: Effective and efficient retrieval via lightweight late interaction . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3715–3734, Seattle, United States. Association for Computational Linguistics.	1839
		1840
		1841
		1842
		1843
		1844
		1845
		1846
		1847
	Naomi Saphra and Sarah Wiegrefe. 2024. Mechanistic? <i>arXiv preprint arXiv:2410.09087</i> .	1848
		1849
	Hinrich Schütze, Christopher D Manning, and Prabhakar Raghavan. 2008. <i>Introduction to information retrieval</i> , volume 39. Cambridge University Press Cambridge.	1850
		1851
		1852
		1853
	Rulin Shao, Jacqueline He, Akari Asai, Weijia Shi, Tim Dettmers, Sewon Min, Luke Zettlemoyer, and Pang Wei Koh. 2024. Scaling retrieval-based language models with a trillion-token datastore. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	1854
		1855
		1856
		1857
		1858
		1859
	Rulin Shao, Rui Qiao, Varsha Kishore, Niklas Muenighoff, Xi Victoria Lin, Daniela Rus, Bryan Kian Hsiang Low, Sewon Min, Wen-tau Yih, Pang Wei Koh, et al. 2025. Reasonir: Training retrievers for reasoning tasks. <i>arXiv preprint arXiv:2504.20595</i> .	1860
		1861
		1862
		1863
		1864
		1865
	Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. 2014. Learning semantic representations using convolutional neural networks for web search . In <i>Proceedings of the 23rd International Conference on World Wide Web, WWW ’14 Companion</i> , page 373–374, New York, NY, USA. Association for Computing Machinery.	1866
		1867
		1868
		1869
		1870
		1871
		1872

1873	Noah Shinn, Federico Cassano, Ashwin Gopinath,	Yi Tay, Vinh Q. Tran, Mostafa Dehghani, Jianmo Ni,	1929
1874	Karthik Narasimhan, and Shunyu Yao. 2023. Re-	Dara Bahri, Harsh Mehta, Zhen Qin, Kai Hui, Zhe	1930
1875	flexion: Language agents with verbal reinforcement	Zhao, Jai Gupta, Tal Schuster, William W. Cohen,	1931
1876	learning. <i>Advances in Neural Information Process-</i>	and Donald Metzler. 2022. Transformer memory as	1932
1877	<i>ing Systems</i> , 36:8634–8652.	a differentiable search index . In <i>Advances in Neural</i>	1933
		<i>Information Processing Systems</i> .	1934
1878	Ashudeep Singh and Thorsten Joachims. 2018. Fairness	Michael Taylor, John Guiver, Stephen Robertson, and	1935
1879	of exposure in rankings. In <i>Proceedings of the 24th</i>	Tom Minka. 2008. Softrank: optimizing non-smooth	1936
1880	<i>ACM SIGKDD international conference on knowl-</i>	rank metrics. In <i>Proceedings of the 2008 Interna-</i>	1937
1881	<i>edge discovery & data mining</i> , pages 2219–2228.	<i>tional Conference on Web Search and Data Mining</i> ,	1938
		pages 77–86.	1939
1882	Chan Hee Song, Jiaman Wu, Clayton Washington,	Wilson L. Taylor. 1953. “cloze procedure” : A new	1940
1883	Brian M Sadler, Wei-Lun Chao, and Yu Su. 2023.	tool for measuring readability . <i>Journalism & Mass</i>	1941
1884	Llm-planner: Few-shot grounded planning for em-	<i>Communication Quarterly</i> , 30:415 – 433.	1942
1885	odied agents with large language models. In <i>Pro-</i>		
1886	<i>ceedings of the IEEE/CVF International Conference</i>	Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019.	1943
1887	<i>on Computer Vision</i> , pages 2998–3009.	BERT rediscovers the classical NLP pipeline . In	1944
		<i>Proceedings of the 57th Annual Meeting of the Asso-</i>	1945
1888	Fei Song and W Bruce Croft. 1999. A general language	<i>ciation for Computational Linguistics</i> , pages 4593–	1946
1889	model for information retrieval. In <i>Proceedings of</i>	4601, Florence, Italy. Association for Computational	1947
1890	<i>the eighth international conference on Information</i>	<i>Linguistics</i> .	1948
1891	<i>and knowledge management</i> , pages 316–321.		
		Nandan Thakur, Nils Reimers, Andreas Rücklé, Ab-	1949
1892	Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen,	hishek Srivastava, and Iryna Gurevych. 2021. Beir:	1950
1893	Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-	A heterogeneous benchmark for zero-shot evaluation	1951
1894	Rong Wen. 2025. R1-searcher: Incentivizing the	of information retrieval models. In <i>Thirty-fifth Con-</i>	1952
1895	search capability in llms via reinforcement learning.	<i>ference on Neural Information Processing Systems</i>	1953
1896	<i>arXiv preprint arXiv:2503.05592</i> .	<i>Datasets and Benchmarks Track (Round 2)</i> .	1954
1897	Heydar Soudani, Evangelos Kanoulas, and Faegheh Ha-	Yu Tian, Xiao Yang, Jingyuan Zhang, Yinpeng Dong,	1955
1898	sibi. 2024. Fine tuning vs. retrieval augmented gener-	and Hang Su. 2023. Evil geniuses: Delving into	1956
1899	ation for less popular knowledge . In <i>Proceedings of</i>	the safety of llm-based agents. <i>arXiv preprint</i>	1957
1900	<i>the 2024 Annual International ACM SIGIR Confer-</i>	<i>arXiv:2311.11855</i> .	1958
1901	<i>ence on Research and Development in Information</i>		
1902	<i>Retrieval in the Asia Pacific Region, SIGIR-AP 2024</i> ,	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	1959
1903	page 12–22, New York, NY, USA. Association for	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	1960
1904	Computing Machinery.	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	1961
		Bhosale, et al. 2023. Llama 2: Open founda-	1962
1905	Karen Sparck Jones. 1972. A statistical interpretation	tion and fine-tuned chat models. <i>arXiv preprint</i>	1963
1906	of term specificity and its application in retrieval.	<i>arXiv:2307.09288</i> .	1964
1907	<i>Journal of documentation</i> , 28(1):11–21.		
		Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	1965
1908	Hongjin Su, Howard Yen, Mengzhou Xia, Weijia	and Ashish Sabharwal. 2023. Interleaving retrieval	1966
1909	Shi, Niklas Muennighoff, Han-yu Wang, Haisu Liu,	with chain-of-thought reasoning for knowledge-	1967
1910	Quan Shi, Zachary S Siegel, Michael Tang, et al.	intensive multi-step questions . In <i>Proceedings of</i>	1968
1911	2024. Bright: A realistic and challenging bench-	<i>the 61st Annual Meeting of the Association for Com-</i>	1969
1912	mark for reasoning-intensive retrieval. <i>arXiv preprint</i>	<i>putational Linguistics (Volume 1: Long Papers)</i> ,	1970
1913	<i>arXiv:2407.12883</i> .	pages 10014–10037, Toronto, Canada. Association	1971
		for Computational Linguistics.	1972
1914	Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang	C Van Rijsbergen. 1979. Information retrieval: the-	1973
1915	Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and	ory and practice. In <i>Proceedings of the joint</i>	1974
1916	Zhaochun Ren. 2023. Is chatgpt good at search?	<i>IBM/University of Newcastle upon tyne seminar on</i>	1975
1917	investigating large language models as re-ranking	<i>data base systems</i> , volume 79, pages 1–14.	1976
1918	agents. In <i>Proceedings of the 2023 Conference on</i>		
1919	<i>Empirical Methods in Natural Language Processing</i> ,	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	1977
1920	pages 14918–14937.	Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz	1978
		Kaiser, and Illia Polosukhin. 2017. Attention is all	1979
1921	Jiejun Tan, Zhicheng Dou, Wen Wang, Mang Wang,	you need . In <i>Advances in Neural Information Pro-</i>	1980
1922	Weipeng Chen, and Ji-Rong Wen. 2024. Html-	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	1981
1923	rag: Html is better than plain text for modeling re-		
1924	trieved knowledge in rag systems. <i>arXiv preprint</i>	Jean-Philippe Vert, Koji Tsuda, and Bernhard Schölkopf.	1982
1925	<i>arXiv:2411.02959</i> .	2004. A primer on kernel methods.	1983
1926	Jiabin Tang, Tianyu Fan, and Chao Huang. 2025. Au-		
1927	toagent: A fully-automated and zero-code framework		
1928	for llm agents. <i>arXiv e-prints</i> , pages arXiv–2502.		

1984	Shengxian Wan, Yanyan Lan, Jiafeng Guo, Jun Xu,	A modern bidirectional encoder for fast, memory	2040
1985	Liang Pang, and Xueqi Cheng. 2016. A deep archi-	efficient, and long context finetuning and inference.	2041
1986	architecture for semantic matching with multiple posi-	<i>arXiv preprint arXiv:2412.13663</i> .	2042
1987	tional sentence representations. In <i>Proceedings of the</i>		
1988	<i>Thirtieth AAAI Conference on Artificial Intelligence</i> ,	Weaviate. 2024. What is agentic rag .	2043
1989	AAAI'16, page 2835–2841. AAAI Press.		
1990	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	Cong Wei, Yang Chen, Haonan Chen, Hexiang Hu,	2044
1991	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	Ge Zhang, Jie Fu, Alan Ritter, and Wenhui Chen.	2045
1992	Xu Chen, Yankai Lin, et al. 2024a. A survey on large	2025. Uniir: Training and benchmarking univer-	2046
1993	language model based autonomous agents. <i>Frontiers</i>	salmultimodal information retrievers. In <i>European</i>	2047
1994	<i>of Computer Science</i> , 18(6):186345.	<i>Conference on Computer Vision</i> , pages 387–404.	2048
		Springer.	2049
1995	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	2050
1996	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	2051
1997	Xu Chen, Yankai Lin, et al. 2024b. A survey on large	et al. 2022. Chain-of-thought prompting elicits rea-	2052
1998	language model based autonomous agents. <i>Frontiers</i>	soning in large language models. <i>Advances in neural</i>	2053
1999	<i>of Computer Science</i> , 18(6):186345.	<i>information processing systems</i> , 35:24824–24837.	2054
2000	Liang Wang, Haonan Chen, Nan Yang, Xiaolong	Orion Weller, Benjamin Chang, Sean MacAvaney, Kyle	2055
2001	Huang, Zhicheng Dou, and Furu Wei. 2025. Chain-	Lo, Arman Cohan, Benjamin Van Durme, Dawn	2056
2002	of-retrieval augmented generation. <i>arXiv preprint</i>	Lawrie, and Luca Soldaini. 2024a. Followir: Evaluat-	2057
2003	<i>arXiv:2501.14342</i> .	ing and teaching information retrieval models to fol-	2058
		low instructions. <i>arXiv preprint arXiv:2403.15246</i> .	2059
2004	Liang Wang, Nan Yang, Xiaolong Huang, Binxing	Orion Weller, Benjamin Van Durme, Dawn Lawrie,	2060
2005	Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder,	Ashwin Paranjape, Yuhao Zhang, and Jack Hessel.	2061
2006	and Furu Wei. 2022a. Text embeddings by weakly-	2024b. Promptriever: Instruction-trained retriev-	2062
2007	supervised contrastive pre-training. <i>arXiv preprint</i>	ers can be prompted like language models. <i>arXiv</i>	2063
2008	<i>arXiv:2212.03533</i> .	<i>preprint arXiv:2409.11136</i> .	2064
2009	Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang,	Thomas D Wilson. 2000. Human information behavior.	2065
2010	Rangan Majumder, and Furu Wei. 2023. Improving	<i>Informing science</i> , 3:49.	2066
2011	text embeddings with large language models. <i>arXiv</i>		
2012	<i>preprint arXiv:2401.00368</i> .	SK Michael Wong and YY Yao. 1989. A probability	2067
2013	Xuanhui Wang, Michael Bendersky, Donald Metzler,	distribution model for information retrieval. <i>Informa-</i>	2068
2014	and Marc Najork. 2016. Learning to rank with se-	<i>tion Processing & Management</i> , 25(1):39–53.	2069
2015	lection bias in personal search. In <i>Proceedings of</i>		
2016	<i>the 39th International ACM SIGIR conference on Re-</i>	Junde Wu, Jiayuan Zhu, and Yuyuan Liu. 2025. Agentic	2070
2017	<i>search and Development in Information Retrieval</i> ,	reasoning: Reasoning llms with tools for the deep	2071
2018	pages 115–124.	research. <i>arXiv preprint arXiv:2502.04644</i> .	2072
2019	Xuanhui Wang, Nadav Golbandi, Michael Bendersky,	Qiang Wu, Christopher JC Burges, Krysta M Svore,	2073
2020	Donald Metzler, and Marc Najork. 2018. Position	and Jianfeng Gao. 2010. Adapting boosting for in-	2074
2021	bias estimation for unbiased learning to rank in per-	formation retrieval measures. <i>Information Retrieval</i> ,	2075
2022	sonal search. In <i>Proceedings of the eleventh ACM</i>	13:254–270.	2076
2023	<i>international conference on web search and data min-</i>		
2024	<i>ing</i> , pages 610–618.	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu,	2077
2025	Yujing Wang, Yingyan Hou, Haonan Wang, Ziming	Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang,	2078
2026	Miao, Shibin Wu, Hao Sun, Qi Chen, Yuqing Xia,	Xiaoyun Zhang, and Chi Wang. 2023. Auto-	2079
2027	Chengmin Chi, Guoshuai Zhao, Zheng Liu, Xing	gen: Enabling next-gen llm applications via multi-	2080
2028	Xie, Hao Sun, Weiwei Deng, Qi Zhang, and Mao	agent conversation framework. <i>arXiv preprint</i>	2081
2029	Yang. 2022b. A neural corpus indexer for document	<i>arXiv:2308.08155</i> .	2082
2030	retrieval . In <i>Advances in Neural Information Pro-</i>		
2031	<i>cessing Systems</i> .	Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and	2083
2032	Zora Zhiruo Wang, Akari Asai, Xinyan Velocity Yu,	Hang Li. 2008. Listwise approach to learning to	2084
2033	Frank F Xu, Yiqing Xie, Graham Neubig, and Daniel	rank: theory and algorithm. In <i>Proceedings of the</i>	2085
2034	Fried. 2024c. Coderag-bench: Can retrieval augment	<i>25th international conference on Machine learning</i> ,	2086
2035	code generation? <i>arXiv preprint arXiv:2406.14497</i> .	pages 1192–1199.	2087
2036	Benjamin Warner, Antoine Chaffin, Benjamin Clavié,	Chong Xiang, Tong Wu, Zexuan Zhong, David Wagner,	2088
2037	Orion Weller, Oskar Hallström, Said Taghadouini,	Danqi Chen, and Prateek Mittal. 2024a. Certifiably	2089
2038	Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom	robust rag against retrieval corruption. <i>arXiv preprint</i>	2090
2039	Aarsen, et al. 2024. Smarter, better, faster, longer:	<i>arXiv:2405.15556</i> .	2091

2092	Zhen Xiang, Linzhi Zheng, Yanjie Li, Junyuan Hong,	Zhichao Xu, Hansi Zeng, Juntao Tan, Zuohui Fu,	2147
2093	Qinbin Li, Han Xie, Jiawei Zhang, Zidi Xiong,	Yongfeng Zhang, and Qingyao Ai. 2023. A reusable	2148
2094	Chulin Xie, Carl Yang, et al. 2024b. Guardagent:	model-agnostic framework for faithfully explain-	2149
2095	Safeguard llm agents by a guard agent via knowledge-	able recommendation and system scrutability. <i>ACM</i>	2150
2096	enabled reasoning. <i>arXiv preprint arXiv:2406.09187</i> .	<i>Transactions on Information Systems</i> , 42(1):1–29.	2151
2097	Quanting Xie, So Yeon Min, Tianyi Zhang, Kedi Xu,	Le Yan, Zhen Qin, Honglei Zhuang, Rolf Jagerman,	2152
2098	Aarav Bajaj, Russ Salakhutdinov, Matthew Johnson-	Xuanhui Wang, Michael Bendersky, and Harrie Oost-	2153
2099	Roberson, and Yonatan Bisk. 2024. Embodied-RAG:	erhuis. 2024. Consolidating ranking and relevance	2154
2100	General non-parametric embodied memory for	predictions of large language models through post-	2155
2101	retrieval and generation . In <i>Language Gamification -</i>	processing . In <i>Proceedings of the 2024 Conference</i>	2156
2102	<i>NeurIPS 2024 Workshop</i> .	<i>on Empirical Methods in Natural Language Process-</i>	2157
2103	Chenyan Xiong, Zhuyun Dai, Jamie Callan, Zhiyuan	<i>ing</i> , pages 410–423, Miami, Florida, USA. Associa-	2158
2104	Liu, and Russell Power. 2017. End-to-end neural	tion for Computational Linguistics.	2159
2105	ad-hoc ranking with kernel pooling . In <i>Proceedings</i>	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	2160
2106	<i>of the 40th International ACM SIGIR Conference on</i>	Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,	2161
2107	<i>Research and Development in Information Retrieval</i> ,	Chengen Huang, Chenxu Lv, et al. 2025. Qwen3	2162
2108	SIGIR '17, page 55–64, New York, NY, USA. Asso-	technical report. <i>arXiv preprint arXiv:2505.09388</i> .	2163
2109	ciation for Computing Machinery.		
2110	Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang,	Songlin Yang, Jan Kautz, and Ali Hatamizadeh. 2024.	2164
2111	Jialin Liu, Paul N. Bennett, Junaid Ahmed, and	Gated delta networks: Improving mamba2 with delta	2165
2112	Arnold Overwijk. 2020. Approximate nearest neigh-	rule. <i>arXiv preprint arXiv:2412.06464</i> .	2166
2113	bor negative contrastive learning for dense text re-		
2114	trieval . <i>ArXiv</i> , abs/2007.00808.	Tao Yang, Zhichao Xu, and Qingyao Ai. 2023a. Vertical	2167
2115	Jun Xu, Xiangnan He, and Hang Li. 2018. Deep learn-	allocation-based fair exposure amortizing in ranking.	2168
2116	ing for matching in search and recommendation. In	In <i>Proceedings of the Annual International ACM SI-</i>	2169
2117	<i>The 41st International ACM SIGIR Conference on</i>	<i>GIR Conference on Research and Development in In-</i>	2170
2118	<i>Research & Development in Information Retrieval</i> ,	<i>formation Retrieval in the Asia Pacific Region</i> , pages	2171
2119	pages 1365–1368.	234–244.	2172
2120	Shicheng Xu, Danyang Hou, Liang Pang, Jingcheng	Tao Yang, Zhichao Xu, Zhenduo Wang, and Qingyao	2173
2121	Deng, Jun Xu, Huawei Shen, and Xueqi Cheng.	Ai. 2023b. Fara: Future-aware ranking algorithm for	2174
2122	2024a. Invisible relevance bias: Text-image retrieval	fairness optimization. In <i>Proceedings of the 32nd</i>	2175
2123	models prefer ai-generated images . In <i>Proceedings</i>	<i>ACM International Conference on Information and</i>	2176
2124	<i>of the 47th International ACM SIGIR Conference on</i>	<i>Knowledge Management</i> , pages 2906–2916.	2177
2125	<i>Research and Development in Information Retrieval</i> ,	Tao Yang, Zhichao Xu, Zhenduo Wang, Anh Tran, and	2178
2126	SIGIR '24, page 208–217, New York, NY, USA. As-	Qingyao Ai. 2023c. Marginal-certainty-aware fair	2179
2127	sociation for Computing Machinery.	ranking algorithm. In <i>Proceedings of the Sixteenth</i>	2180
2128	Zhichao Xu. 2024. Rankmamba, benchmarking	<i>ACM International Conference on Web Search and</i>	2181
2129	mamba’s document ranking performance in the era	<i>Data Mining</i> , pages 24–32.	2182
2130	of transformers. <i>arXiv preprint arXiv:2403.18276</i> .	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	2183
2131	Zhichao Xu, Aosong Feng, Yijun Tian, Haibo Ding, and	Shafraan, Karthik R Narasimhan, and Yuan Cao. 2023.	2184
2132	Lin Lee Cheong. 2025a. Csplade: Learned sparse	React: Synergizing reasoning and acting in language	2185
2133	retrieval with causal language models. <i>arXiv preprint</i>	models . In <i>The Eleventh International Conference</i>	2186
2134	<i>arXiv:2504.10816</i> .	<i>on Learning Representations</i> .	2187
2135	Zhichao Xu, Hemank Lamba, Qingyao Ai, Joel	Zonghai Yao, Aditya Parashar, Huixue Zhou, Won Seok	2188
2136	Tetreault, and Alex Jaimes. 2024b. Cfe2: Counter-	Jang, Feiyun Ouyang, Zhichao Yang, and Hong Yu.	2189
2137	factual editing for search result explanation. In <i>Pro-</i>	2024. Mcqg-srefine: Multiple choice question gener-	2190
2138	<i>ceedings of the 2024 ACM SIGIR International Con-</i>	ation and evaluation with iterative self-critique, cor-	2191
2139	<i>ference on Theory of Information Retrieval</i> , pages	rection, and comparison feedback. <i>arXiv preprint</i>	2192
2140	145–155.	<i>arXiv:2410.13191</i> .	2193
2141	Zhichao Xu, Jinghua Yan, Ashim Gupta, and Vivek	Zeynep Akkalyoncu Yilmaz, Wei Yang, Haotian Zhang,	2194
2142	Srikumar. 2025b. State space models are strong text	and Jimmy J. Lin. 2019. Cross-domain modeling of	2195
2143	rerankers . In <i>Proceedings of the 10th Workshop on</i>	sentence-level evidence for document retrieval . In	2196
2144	<i>Representation Learning for NLP (RepL4NLP-2025)</i> ,	<i>Conference on Empirical Methods in Natural Lan-</i>	2197
2145	pages 152–169, Albuquerque, NM. Association for	<i>guage Processing</i> .	2198
2146	Computational Linguistics.	Soyoung Yoon, Eunbi Choi, Jiyeon Kim, Hyeongu Yun,	2199
		Yireun Kim, and Seung-won Hwang. 2024. ListT5:	2200
		Listwise reranking with fusion-in-decoder improves	2201

2202	zero-shot retrieval. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2287–2308, Bangkok, Thailand. Association for Computational Linguistics.	2259	Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Min Zhang, and Shaoping Ma. 2020. Repbert: Contextualized text embeddings for first-stage retrieval. <i>arXiv preprint arXiv:2006.15498</i> .	2260
2203		2261		2262
2204		2263	Crystina Zhang, Minghan Li, and Jimmy Lin. 2024a. CELI: Simple yet effective approach to enhance out-of-domain generalization of cross-encoders. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)</i> , pages 188–196, Mexico City, Mexico. Association for Computational Linguistics.	2264
2205		2265		2266
2206		2267		2268
2207	Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2024. Making retrieval-augmented language models robust to irrelevant context. In <i>ICLR 2024 Workshop on Large Language Model (LLM) Agents</i> .	2269		2270
2208		2271	Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and Weizhu Chen. 2023a. RepoCoder: Repository-level code completion through iterative retrieval and generation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2471–2484, Singapore. Association for Computational Linguistics.	2272
2209		2273		2274
2210		2275		2276
2211	Puxuan Yu and James Allan. 2020. A study of neural matching models for cross-lingual ir. In <i>Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 1637–1640.	2277		2278
2212		2279	Hanqi Zhang, Chong Chen, Lang Mei, Qi Liu, and Jiaxin Mao. 2024b. Mamba retriever: Utilizing mamba for effective and efficient dense retrieval. In <i>Proceedings of the 33rd ACM International Conference on Information and Knowledge Management</i> , pages 4268–4272.	2280
2213		2281		2282
2214		2283		2284
2215		2285	Xinyu Zhang, Sebastian Hofstätter, Patrick Lewis, Raphael Tang, and Jimmy Lin. 2023b. Rank-without-gpt: Building gpt-independent listwise rerankers on open-source large language models. <i>arXiv preprint arXiv:2312.02969</i> .	2286
2216	Puxuan Yu, Antonio Mallia, and Matthias Petri. 2024a. Improved learned sparse retrieval with corpus-specific vocabularies. In <i>European Conference on Information Retrieval</i> , pages 181–194. Springer.	2287		2288
2217		2289		2290
2218		2291	Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamalloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. 2023c. Miracl: A multilingual retrieval dataset covering 18 diverse languages. <i>Transactions of the Association for Computational Linguistics</i> , 11:1114–1131.	2292
2219		2293		2294
2220	Puxuan Yu, Luke Merrick, Gaurav Nuti, and Daniel Campos. 2024b. Arctic-embed 2.0: Multilingual retrieval without compromise. <i>arXiv preprint arXiv:2412.04506</i> .	2295		2296
2221		2297	Xinyu Zhang, Andrew Yates, and Jimmy Lin. 2021. Comparing score aggregation approaches for document retrieval with pretrained transformers. In <i>Advances in Information Retrieval</i> , pages 150–163, Cham. Springer International Publishing.	2298
2222		2299		2300
2223		2301		2302
2224		2303	Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Livia Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, et al. 2024. Sglang: Efficient execution of structured language model programs. <i>Advances in Neural Information Processing Systems</i> , 37:62557–62583.	2304
2225	Puxuan Yu, Razieh Rahimi, and James Allan. 2022. Towards explainable search results: a listwise explanation generator. In <i>Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 669–680.	2305		2306
2226		2307		2308
2227		2309	Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. 2025. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. <i>arXiv preprint arXiv:2504.03160</i> .	2310
2228		2311		2312
2229		2313		2314
2230	Wenhao Yu, Hongming Zhang, Xiaoman Pan, Peixin Cao, Kaixin Ma, Jian Li, Hongwei Wang, and Dong Yu. 2024c. Chain-of-note: Enhancing robustness in retrieval-augmented language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 14672–14685, Miami, Florida, USA. Association for Computational Linguistics.	2315	Zexuan Zhong, Ziqing Huang, Alexander Wettig, and Danqi Chen. 2023. Poisoning retrieval corpora by	
2231				
2232				
2233				
2234				
2235				
2236				
2237				
2238	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 <i>Proceedings of the 27th ACM international conference on information and knowledge management</i> , pages 497–506.			
2239				
2240				
2241				
2242				
2243				
2244				
2245	Hansi Zeng, Chen Luo, Bowen Jin, Sheikh Muhammad Sarwar, Tianxin Wei, and Hamed Zamani. 2024. Scalable and effective generative information retrieval. In <i>Proceedings of the ACM Web Conference 2024, WWW '24</i> , page 1441–1452, New York, NY, USA. Association for Computing Machinery.			
2246				
2247				
2248				
2249				
2250				
2251	Chengxiang Zhai and John Lafferty. 2004. A study of smoothing methods for language models applied to information retrieval. <i>ACM Transactions on Information Systems (TOIS)</i> , 22(2):179–214.			
2252				
2253				
2254				
2255	ChengXiang Zhai et al. 2008. Statistical language models for information retrieval a critical review. <i>Foundations and Trends® in Information Retrieval</i> , 2(3):137–213.			
2256				
2257				
2258				

injecting adversarial passages. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13764–13775, Singapore. Association for Computational Linguistics.

Yujia Zhou, Yan Liu, Xiaoxi Li, Jiajie Jin, Hongjin Qian, Zheng Liu, Chaozhuo Li, Zhicheng Dou, Tsung-Yi Ho, and Philip S Yu. 2024. Trustworthiness in retrieval-augmented generation systems: A survey. *arXiv preprint arXiv:2409.10102*.

Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan Chen, Zheng Liu, Zhicheng Dou, and Ji-Rong Wen. 2023. Large language models for information retrieval: A survey. *arXiv preprint arXiv:2308.07107*.

Honglei Zhuang, Zhen Qin, Kai Hui, Junru Wu, Le Yan, Xuanhui Wang, and Michael Bendersky. 2024a. [Beyond yes and no: Improving zero-shot LLM rankers via scoring fine-grained relevance labels](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 358–370, Mexico City, Mexico. Association for Computational Linguistics.

Honglei Zhuang, Zhen Qin, Rolf Jagerman, Kai Hui, Ji Ma, Jing Lu, Jianmo Ni, Xuanhui Wang, and Michael Bendersky. 2023a. [Rankt5: Fine-tuning t5 for text ranking with ranking losses](#). In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '23*, page 2308–2313, New York, NY, USA. Association for Computing Machinery.

Shengyao Zhuang, Hang Li, and Guido Zuccon. 2021. Deep query likelihood model for information retrieval. In *Advances in Information Retrieval: 43rd European Conference on IR Research, ECIR 2021, Virtual Event, March 28–April 1, 2021, Proceedings, Part II* 43, pages 463–470. Springer.

Shengyao Zhuang, Bing Liu, Bevan Koopman, and Guido Zuccon. 2023b. [Open-source large language models are strong zero-shot query likelihood models for document ranking](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8807–8817, Singapore. Association for Computational Linguistics.

Shengyao Zhuang, Honglei Zhuang, Bevan Koopman, and Guido Zuccon. 2024b. [A setwise approach for effective and highly efficient zero-shot ranking with large language models](#). In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '24*, page 38–47, New York, NY, USA. Association for Computing Machinery.

A Supplement Materials on Traditional IR Models

Boolean Models The most basic Boolean model can be extended by incorporating term weighting, allowing both queries and documents to be represented as sets of weighted terms. Then, the logical implication $\mathcal{D} \rightarrow \mathcal{Q}$ is also weighted. The commonly used weighted approaches for the logical implication $\mathcal{D} \rightarrow \mathcal{Q}$ include using a fuzzy set extension of the Boolean logic (Radecki, 1979; Kraft and Buell, 1983) and p -norm (Salton et al., 1983).

Vector Space Model The weights q_i or d_i can be represented by other sophisticated schema, such as TF-IDF (Sparck Jones, 1972) and BM25 (Robertson et al., 1995). The extracted abundant features can improve the capacity and accuracy of the vector space models. Besides, given the vector representations of query \mathcal{Q} and document \mathcal{D} , the most commonly used is cosine similarity, defined as

$$\text{sim}(\mathcal{Q}, \mathcal{D}) = \frac{\mathcal{Q} \cdot \mathcal{D}}{|\mathcal{Q}| \times |\mathcal{D}|},$$

where $\mathcal{Q} \cdot \mathcal{D}$ is the dot product and $|\mathcal{Q}|, |\mathcal{D}|$ denotes the length of the vector.

Probabilistic Model In Probabilistic Model, the relevance score of a document \mathcal{D} to a query \mathcal{Q} depends on a set of events $\{x_i\}_1^n$ representing the occurrence of term t_i in this document. The simplest probabilistic model is the binary independence retrieval model (Robertson and Jones, 1976), which assumes terms are independent so only $x_i = 1$ and $x_i = 0$ exist in the representation. Given a set of sample documents whose relevance is judged, the estimation of the relevance score can be derived as $\text{Score}(\mathcal{Q}, \mathcal{D}) \propto \sum_{(x_i=1) \in \mathcal{D}} \log \frac{r_i(T-n_i-R+r_i)}{(R-r_i)(n_i-r_i)}$, where T and R are the total number of sampled judged documents and relevant samples, and n_i and r_i denote the number of samples and relevant samples containing t_i , respectively. The smooth mechanisms (Baeza-Yates et al., 1999) are necessary to deal with zero occurrences of the t_i .

Except for the binary independence retrieval model, more sophisticated probabilistic models are proposed in the literature (Wong and Yao, 1989; Fuhr, 1992), such as the inter-dependency between terms (Van Rijsbergen, 1979).

Statistical Language Model The general idea of a statistical language model is to estimate the relevance score of a document \mathcal{D} to a query \mathcal{Q} via $\mathcal{P}(\mathcal{D}|\mathcal{Q})$ (Ponte and Croft, 1998). Based on Bayes Rule, $\mathcal{P}(\mathcal{D}|\mathcal{Q})$ can be derived as directly proportional to $\mathcal{P}(\mathcal{Q}|\mathcal{D})\mathcal{P}(\mathcal{D})$. For simplification, most studies assume a uniform distribution for $\mathcal{P}(\mathcal{D})$. The main focus is on modeling $\mathcal{P}(\mathcal{Q}|\mathcal{D})$ as a ranking function by treating the query as a set of independent terms as $\mathcal{Q} = \{t_i\}_{i=1}^n$, thus $\mathcal{P}(\mathcal{Q}|\mathcal{D}) = \prod_{t_i \in \mathcal{Q}} \mathcal{P}(t_i|\mathcal{D})$. The probability $\mathcal{P}(t_i|\mathcal{D})$ is determined using a statistical language model θ_D that represents the document, then the relevance is estimated by log-likelihood as $\text{Score}(\mathcal{Q}, \mathcal{D}) = \log \mathcal{P}(\mathcal{Q}|\theta_D) = \sum_{t_i \in \mathcal{Q}} \log \mathcal{P}(t_i|\theta_D)$, where the estimation of the language model θ_D is usually achieved by maximum likelihood.

The statistic language models for IR (Miller et al., 1999; Berger and Lafferty, 1999; Song and Croft, 1999) also encounter the problem of the zero occurrences of the query term t_i , i.e., the probability $\mathcal{P}(\mathcal{Q}|\theta_D)$ becomes zero, if a query term t_i does not appear in the document. This is too restrictive for IR, as a document can still be relevant even if it contains only some of the query terms. To address this zero-probability issue, smoothing techniques are applied, assigning small probabilities to terms that do not appear in the document. The principle behind smoothing is that any text used to model a language captures only a limited subset of its linguistic patterns (or terms, in this case). The commonly used smoothing methods (Zhai and Lafferty, 2004; Zhai et al., 2008) include Jelinek-Mercer smoothing, Dirichlet smoothing, etc.

B Supplement Materials on Learning-to-Rank Architecture and Training Strategy

We present a list of learning-to-rank works and their backbone architectures in Table 1. A significant portion of the literature focuses on loss functions and feature transformers (Qin et al., 2021; Bruch et al.,

Name	Model	Backbone Architecture	Loss Function
MART (Friedman, 2001)	ML	Boosting	Pointwise
RANKBOOST (Freund et al., 2003)	ML	Boosting	Pairwise
RANKNET (Burges et al., 2005)	Neural Nets	DNN	Pairwise
RANKSVM (Joachims, 2006)	ML	SVM	Pairwise
LAMBDARANK (Burges et al., 2006)	Neural Nets	DNN	Pairwise
LISTNET (Cao et al., 2007)	Neural Nets	DNN	Listwise
SOFT-RANK (Taylor et al., 2008)	Neural Nets	DNN	Listwise
LISTMLE (Xia et al., 2008)	ML	Linear	Listwise
LAMBDAMART (Burges, 2010)	ML	GBDT	Listwise
APPROXNDCG (Qin et al., 2010)	ML	Linear	Listwise
DLCM (Ai et al., 2018a)	Neural Nets	DNN	Listwise
GSF (Ai et al., 2019)	Neural Nets	DNN	Listwise
APPROXNDCG (Bruch et al., 2019)	Neural Nets	DNN	Listwise
SETRANK (Pang et al., 2020)	Neural Nets	Self Attention Blocks	Listwise

Table 1: A list of learning-to-rank works and their model architectures.

2019; Burges, 2010). Additionally, some studies focus on unbiased relevance estimation using biased feedback (Wang et al., 2016; Joachims et al., 2017b,a; Ai et al., 2018c,b; Wang et al., 2018; Hu et al., 2019; Ren et al., 2022) while other focus on jointly optimizing effectiveness and fairness of the ranking systems (Singh and Joachims, 2018; Biega et al., 2018; Morik et al., 2020; Patro et al., 2020; Oosterhuis, 2021; Yang et al., 2023a,c,b). We omit detailed discussions here and refer readers to the original papers.

C Supplement Materials on Neural Ranking Models

Representation-based Models Representation-based neural ranking models can be regarded as extensions of vector space models (§ 3), which independently encode queries and documents into a latent vector space. The relevance ranking of a document is determined by computing the similarity (e.g., cosine similarity) between the query and document embeddings.

The Deep Structured Semantic Model (DSSM) (Huang et al., 2013) is an early example of a representation-based neural ranking model. It utilizes word hashing and multilayer perceptrons (MLPs) to independently encode term vectors of queries and documents into a shared semantic space, enabling the computation of ranking scores based on the cosine similarity of their embeddings. Research has focused on enhancing DSSM by modifying its encoder network to improve the model’s ability to capture richer semantic and contextual information. For instance, Convolutional DSSM (Shen et al., 2014) leverages a CNN architecture to project vectors within a context window to a local contextual feature vector. These local features are then aggregated using a max-pooling layer to produce a representation of the entire query or document. Another variation of DSSM replaces MLPs with a Long Short-Term Memory (LSTM) network (Palangi et al., 2016; Wan et al., 2016). By leveraging LSTM’s memory mechanism, such models can capture both local and global context information without the pooling layer, thus better suited for handling longer documents.

Interaction-based Models Interaction-based models process queries and documents jointly through neural networks. The model’s output is typically a score that measures the relevance of the input query-document pair. Various network architectures have been proposed to jointly encode queries and documents. For instance, MATCHPYRAMID (Pang et al., 2016) employs CNN over the interaction matrix between query and document terms. This approach treats the interaction matrix as an image, allowing the CNN to capture local matching patterns. The model then aggregates these patterns through convolution and pooling operations to produce a relevance score, effectively modeling the hierarchical matching structures between queries and documents (Hu et al., 2014). Building upon the concept of interaction-focused models, Guo et al. (2016a) highlighted the importance of exact term matches in neural ranking models and proposed the Deep Relevance Matching Model (DRMM). The model constructs matching histograms for each query term to capture the distribution of matching signals across

Name	Architecture	Backbone	Embeddings
DSSM (Huang et al., 2013)	Representation-based	MLP	Word Hashing
CDSSM (Shen et al., 2014)	Representation-based	CNN	Word Hashing
ARC-I (Hu et al., 2014)	Representation-based	CNN	Word2Vec
ARC-II (Hu et al., 2014)	Interaction-based	CNN	Word2Vec
MATCHPYRAMID (Pang et al., 2016)	Interaction-based	CNN	Randomly Initialized
LSTM-RNN (Palangi et al., 2016)	Representation-based	LSTM	Randomly Initialized
MV-LSTM (Wan et al., 2016)	Representation-based	Bi-LSTM	Word2Vec
DRMM (Guo et al., 2016a)	Interaction-based	MLP	Word2Vec
DESM (Nalisnick et al., 2016)	Interaction-based	MLP	Word2Vec
K-NRM (Xiong et al., 2017)	Interaction-based	MLP + RBF kernels	Word2Vec
CONV-KNRM (Dai et al., 2018)	Interaction-based	CNN	Word2Vec
TK (Hofstätter et al., 2020c)	Interaction-based	Transformer + Kernel	GloVe
TKL (Hofstätter et al., 2020a)	Interaction-based	Transformer + Kernel	GloVe
NDRM (Mitra et al., 2021)	Interaction-based	Conformer + Kernel	BERT

Table 2: A list of neural ranking models and their model architectures.

document terms. These histograms are then processed through a feed-forward neural network to learn hierarchical matching patterns. Xiong et al. (2017) introduced the Kernel-Based Neural Ranking Model (K-NRM), which further advanced interaction-based approaches. K-NRM employs a translation matrix to compute interactions between query and document terms based on their embeddings. It then applies Radial Basis Function (RBF) kernels to transform these word-level interactions into informative ranking features. Later, they extended the RBF kernel approach to a convolutional neural network (Dai et al., 2018).

Word Embeddings In addition to advancements in network architecture, pre-trained textual representations have also contributed to neural ranking models’ performance (Guo et al., 2016b). GloVe (Pennington et al., 2014) and Word2Vec (Mikolov, 2013) learn dense vector representations for each vocabulary term from large-scale text corpora. Pre-trained embeddings provide semantic-based term representations to enable neural ranking models to focus on learning relevance matching patterns. Both representation-based and interaction-based models adopt pre-trained word embeddings as input representations to their networks, facilitating training convergence and improved performance (Levy et al., 2015). Interaction-based models with cross-lingual word embeddings (Joulin et al., 2018) for cross-lingual reranking have also been explored (Yu and Allan, 2020).

Table 2 shows a list of neural ranking models and backbone architectures. Researchers have explored different backbone neural network architectures in this era, including Convolutional Neural Network (CNN, LeCun et al., 1989, 1998), Long Short Term Memory (LSTM, Hochreiter and Schmidhuber, 1997) and kernel methods (Vert et al., 2004; Chang et al., 2010; Xiong et al., 2017).

Notably, a line of research explores integrating kernel methods with the TRANSFORMER architecture (Vaswani et al., 2017). The main distinction between this line of research and the models discussed in § 6 is that the transformer modules here are not pre-trained on large-scale corpora like Wikipedia and C4 (Devlin et al., 2019; Raffel et al., 2020). We consider this line of research as an intersection between neural ranking models (§ 5) and retrieval with pre-trained transformers (§ 6). TK (Hofstätter et al., 2020c) uses a shallow transformer neural network (up to 3 layers) to encode the query Q and document D separately. After encoding, the contextualized representations are input to one single interaction match matrix, similar to model architecture shown in Fig. 2b. The entire model is trained end-to-end and is able to achieve better performance-efficiency trade-off compared to BERT-based reranker (Nogueira et al., 2019). The main bottleneck of applying transformer architectures to long document reranking is $O(n^2)$ time complexity, where n denotes the document length. TKL (Hofstätter et al., 2020b) further improves upon TK with a local attention mechanism and leads to performance improvement on long document ranking.

D Supplement Materials on Pre-trained Language Models for Information Retrieval

We show a list of models and their corresponding architectures in Table 3, including reranking models, learned dense retrieval, multi-vector representations and learned sparse retrieval. A majority of the models use BERT (Devlin et al., 2019) as the backbone language models, with a few exceptions using DISTILBERT (Sanh, 2019), RoBERTa (Liu, 2019) and encoder part of T5 family models (Raffel et al., 2020; Sanh et al., 2022; Mo et al., 2023; Chung et al., 2024).

One line of work aims to combine the benefits of learned dense retrieval and sparse retrieval. (Gao et al., 2021b,a; Ma et al., 2021; Lin and Ma, 2021; Cormack et al., 2009). Ranklist fusion techniques (e.g., Reciprocal Rank Fusion, Cormack et al., 2009) directly fuses ranklists from different retrievers and has been shown to improve retrieval performance. COIL (Gao et al., 2021a) proposes to enhance traditional bag-of-words retrieval method with semantic embeddings from BERT encoder. UNICOIL (Lin and Ma, 2021) further simplifies reduces the dimension of semantic embeddings to 1 — equivalent to learned term weight in learned sparse retrieval models like SPLADE (Formal et al., 2021b,a).

A few works fall into the intersection of learned sparse retrieval and multi-vector representations. For example, SLIM (Li et al., 2023c) first maps each contextualized token vector to a sparse, high-dimensional lexical space before performing late interaction between these sparse token embeddings. SPLATE (Formal et al., 2024) take an alternative approach to first encodes contextualized token vectors, then maps these token vectors to a sparse vocabulary space with a partially learned SPLADE module. Both models achieve performance improvement compared to learned sparse retrieval baselines such as SPLADE (Formal et al., 2021b,a).

Instead of improving retrieval performance from the modeling perspective, a separate line of works aim to enhance the backbone language models via domain adaptation or continued pre-training, which has been proven successful by prior works in NLP (Howard and Ruder, 2018; Gururangan et al., 2020). Lee et al. (2019) propose to pre-train BERT model with Inverse-Cloze Task (Taylor, 1953) for better text representations. CONDENSER (Gao and Callan, 2021) propose to “condense” text representations into [CLS] token via a dedicated pre-training architecture and corresponding training objective. COCO-DR further extends upon CONDENSER via a technique named implicit Distributionally Robust Optimization to mitigate distribution shift problem in dense retrieval. We refer readers to original papers for details.

As we noted in § 9, one desiderata of future IR models is interpretability and truthfulness. A few works have attempted to interpret transformer-based neural retrieval models’ representations, i.e., mechanistic interpretability (Elhage et al., 2021; Saphra and Wiegrefe, 2024). For example, MacAvaney et al. (2022) showed that neural retrieval models rely less on exact match signals and instead encodes rich semantic information. Ram et al. (2023) project dense retrievers’ intermediate representations to vocabulary space and show the connection of dense retrieval and traditional bag-of-words sparse retrieval methods. Instead of providing model-intrinsic explanations, a few works design IR systems to provide model-agnostic explanations (Rahimi et al., 2021; Yu et al., 2022; Xu et al., 2024b) in order to meet certain desiderata such as faithfulness (Jacovi and Goldberg, 2020; Xu et al., 2023). As IR systems become an integral part of other applied ML domains, we believe it is important to study and design interpretable, truthful and trustworthiness IR models.

E Supplement Materials on LLM for IR

We summarize a list of works that study LLM for retrieval (Table 4) and reranking (Table 5). For generative retrieval, we point to a dedicated survey (Li et al., 2025c). Another comprehensive survey (Mo et al., 2024b) could be referred to for conversational information retrieval. Modern IR systems require extensive labeled data to achieve good performance. One line of work studies the proposal of using LLMs for synthesizing training data (Bonifacio et al., 2022; Boytsov et al., 2024; Dai et al., 2023; Lee et al., 2024; Mo et al., 2024a,c). From the evaluation perspective, LLMs’ superior natural language understanding capability also raise the question of whether they can be used for relevance judgments. A separate line of work tackle the relevance judgments problem (Faggioli et al., 2023, 2024; Clarke and Dietz, 2024). As our focus of this survey is on model architectures, we skip the discussion and point to original papers for further details.

Name	Model	Architecture	Backbone LM	Training strategy
MONOBERT (Nogueira et al., 2019)	Reranking	Cross-encoder	BERT	Classification
CEDR (MacAvaney et al., 2019)	Reranking	Cross-encoder	BERT	Contrastive Learning
BERT-MAXP (Dai and Callan, 2019b)	Reranking	Cross-encoder	BERT	Pairwise Loss
Gao et al. (2020)	Reranking	Cross-encoder	BERT	Distillation
TART-FULL (Asai et al., 2023)	Reranking	Cross-encoder	FLAN-T5-ENC	Instruction Tuning
DPR (Karpukhin et al., 2020)	LDR	Bi-encoder	BERT	Contrastive Learning
ANCE (Xiong et al., 2020)	LDR	Bi-encoder	ROBERTA	Contrastive Learning
REPBERT (Zhan et al., 2020)	LDR	Bi-encoder	BERT	In-batch negatives
MARGIN-MSE (Hofstätter et al., 2020a)	LDR	Bi-encoder	DISTILBERT	Distillation
TAS-B (Hofstätter et al., 2021)	LDR	Bi-encoder	BERT	Distillation
ROCKETQA (Qu et al., 2020)	LDR	Bi-encoder	ERNIE	Contrastive Learning
ROCKETQA-v2 (Ren et al., 2021)	LDR	Bi-encoder	ERNIE	Distillation
GTR (Ni et al., 2022b)	LDR	Bi-encoder	ENCT5	Contrastive Learning
TART-DUAL (Asai et al., 2023)	LDR	Bi-encoder	CONTRIEVER	Instruction Tuning
E5 (Wang et al., 2022a)	LDR	Bi-encoder	BERT	Contrastive Learning
GTE (Li et al., 2023e)	LDR	Bi-encoder	BERT	Contrastive Learning
POLY-ENCODER (Humeau et al., 2020)	Multi-vector	Misc	BERT	In-batch Negatives
ME-BERT (Luan et al., 2021)	Multi-vector	Bi-encoder	BERT	Contrastive Learning
COLBERT (Khattab and Zaharia, 2020)	Multi-vector	Bi-encoder	BERT	Pairwise Loss
COIL (Gao et al., 2021a)	Multi-vector	Bi-encoder	BERT	Contrastive Learning
COLBERT-v2 (Santhanam et al., 2022)	Multi-vector	Bi-encoder	BERT	Distillation
COLBERTER (Hofstätter et al., 2022)	Multi-vector	Bi-encoder	BERT	Distillation
DEEPCCT (Dai and Callan, 2019a)	LSR	Bi-encoder	BERT	Unsupervised
SPARTERM (Bai et al., 2020)	LSR	Bi-encoder	BERT	Contrastive Learning
SPLADE (Formal et al., 2021b)	LSR	Bi-encoder	BERT	Contrastive Learning
SPLADE-v2 (Formal et al., 2021a)	LSR	Bi-encoder	BERT	Distillation
DEEPImpact (Mallia et al., 2021)	LSR	Bi-encoder	BERT	Contrastive Learning
SPARSEMBED (Kong et al., 2023)	LSR	Bi-encoder	BERT	Contrastive Learning
SLIM (Li et al., 2023c)	LSR + Multi-vector	Bi-encoder	BERT	Contrastive Learning
SLIM++ (Li et al., 2023c)	LSR + Multi-vector	Bi-encoder	BERT	Distillation
SPLATE (Formal et al., 2024)	LSR + Multi-vector	Bi-encoder	BERT	Distillation

Table 3: Summary of IR model architecture for passage retrieval and passage ranking based on pre-trained transformers. LDR and LSR denote learned dense retrieval and learned sparse retrieval, respectively. DEEPCCT (Dai and Callan, 2019a) is trained without labeled training set. The "late interaction" mechanism introduced in (Khattab and Zaharia, 2020; Santhanam et al., 2022) can be considered a special case of multi-vector architecture. Contrastive Learning and in-batch negatives means listwise loss function is used.

Name	Architecture	Backbone LM	Training strategy
CPT-TEXT (Neelakantan et al., 2022)	LLM Encoder	GPT-3	Listwise Loss
SGPT-BE (Muennighoff, 2022)	LLM Encoder	GPT-J & GPT-NEOX	Listwise Loss
GTR (Ni et al., 2022b)	LLM Encoder	T5	Listwise Loss
REPLAMA (Ma et al., 2024b)	LLM Encoder	LLAMA	Listwise Loss
E5-MISTRAL (Wang et al., 2023)	LLM Encoder	MISTRAL	Synthetic Data + Listwise Loss
LLARA (Li et al., 2023a)	LLM Encoder	LLAMA	Adaptation + Contrastive Training
MAMBARETRIEVER (Zhang et al., 2024b)	LLM Encoder	MAMBA	Listwise Loss
LLM2VEC (Behnamghader et al., 2024)	LLM Encoder	LLAMA & MISTRAL	Adaptation + Contrastive Pre-training
GRIT-LM (Muennighoff et al., 2025)	LLM	MISTRAL & MIXTRAL 8x7B	Generative/Embedding Joint Training
NVEMBED (Lee et al., 2025)	LLM Encoder	MISTRAL	Adaptation + Synthetic Data + Listwise Loss
GTE-QWEN2-INSTRUCT (Li et al., 2023e)	LLM Encoder	QWEN	Adaptation + Synthetic Data + Listwise Loss

Table 4: Summary of IR model architecture utilizing large language models as retrieval backbone.

F Extended Discussions on Challenges and New Directions

F.1 Autonomous Search Agents

We discuss recent progress on developing autonomous agents for search and information seeking purposes. While these works do not focus on improving IR models *per se*, we believe it is important for IR researchers to adapt to these new use cases of search/retrieval.

Prior works have studied methods to augment language models with search/retrieval to improve generation quality, which we term as retrieval-augmented generation. Early practices include KNN-

Name	Architecture	Backbone LM	Training / Prompting Strategy
<i>Fine-tune LLM for Reranking</i>			
MONOT5 (Nogueira et al., 2020)	LM	T5	Classification
Nogueira dos Santos et al. (2020)	LM	BART	Unlikelihood
QLM-T5 (Zhuang et al., 2021)	LM	T5	Language Modeling
DUOT5 (Pradeep et al., 2021)	LM	T5	Pairwise Loss
RANKT5 (Zhuang et al., 2023a)	LLM Encoder + Prediction Layer	T5	Listwise Loss
LISTT5 (Yoon et al., 2024)	Fusion-in-decoder	T5	Listwise Loss
SGPT-CE (Muennighoff, 2022)	LM	GPT-J & GPT-NEO	Listwise Loss
RANKLLAMA (Ma et al., 2024b)	LLM Encoder + Prediction Layer	LLAMA	Listwise Loss
RANKMAMBA (Xu, 2024)	LLM Encoder + Prediction Layer	MAMBA	Listwise Loss
RANKVICUNA (Pradeep et al., 2023b)	LM	VICUNA	Listwise
RANKZEPHYR (Pradeep et al., 2023c)	LM	ZEPHYR	Listwise
Zhang et al. (2023b)	LM	CODE-LLAMA-INSTRUCT	Listwise
Liu et al. (2024b)	Embedding + LM	MISTRAL	Listwise
<i>Prompt LLM for Reranking</i>			
Zhuang et al. (2023b)	LM	Multiple	Pointwise Prompting
Zhuang et al. (2024a)	LM	FLAN-PALM-S	Pointwise Prompting
UPR (Sachan et al., 2022)	LM	T5 & GPT-NEO	Pointwise Prompting
PRP (Qin et al., 2024a)	LM	FLAN-UL2	Pairwise Prompting
Yan et al. (2024)	LM	FLAN-UL2	Pairwise Prompting
Zhuang et al. (2024b)	LM	FLAN-T5	Pairwise & Setwise Prompting
LRL (Ma et al., 2023)	LM	GPT-3	Listwise Prompting
RANKGPT-3.5 (Sun et al., 2023)	LM	GPT-3.5	Listwise Prompting
RANKGPT-4 (Sun et al., 2023)	LM	GPT-4	Listwise Prompting

Table 5: Summary of IR model architecture utilizing large language models for reranking. Nogueira dos Santos et al. (2020) and Zhuang et al. (2021) revisit the statistic language model problem with modern transformer-based models, including BART (Lewis et al., 2020a) T5 (Raffel et al., 2020) and GPT-2 (Radford et al., 2019).

LM (Khandelwal et al., 2020b), REALM (Gua et al., 2020), RAG (Lewis et al., 2020b), *inter alia*. With more powerful models such as ChatGPT (OpenAI, 2022), researchers begin to design systems to handle daily tasks autonomously with LLM backbones. We refer to such systems as LLM-based agents (Wang et al., 2024b; Guo et al., 2024). WEBGPT (Nakano et al., 2021) leverages reinforcement learning to train GPT-3-based language models for web browsing, which is one of the earliest works along this direction. Due to the instrumental role of retrieval in tasks solving, popular agent frameworks (Wu et al., 2023; Li et al., 2023b) have supported built-in retrieval functionality (commonly referred to as *agentic memory* in agent literature). Most of existing general purpose agent frameworks treat retrieval as one of the available tools, and use LLMs to plan and orchestrate workflows accordingly for task completion, with techniques such as self-refine (Madaan et al., 2023), reflexion (Shinn et al., 2023) and critique (Gou et al., 2023).

Earlier works by the IR and NLP community – such as FLARE (Jiang et al., 2023) and SELF-RAG (Asai et al., 2024a) – have proposed methods to build autonomous search systems, i.e., to enable the system to know when, where and how to search. FLARE (Jiang et al., 2023) explored using prompting method while SELF-RAG (Asai et al., 2024a) focuses on data synthesis and supervised fine-tuning. Popularized by large reasoning models such as GPT-4o (Hurst et al., 2024) and DEEPSEEK-R1 (Guo et al., 2025), there is a surge of recent works aiming to incorporate retrieval to augment LLM reasoning, or to train large reasoning models to use search tool for better performance with reinforcement learning techniques (Li et al., 2025a; Jin et al., 2025; Li et al., 2025b; Chen et al., 2025; Zheng et al., 2025; Guan et al., 2025; Song et al., 2025; Wu et al., 2025; Hu et al., 2025; Wang et al., 2025; Gao et al., 2025). As these works mainly focus on optimizing the generator component of RAG systems, we refer the readers to these individual works for further details.

Existing agentic RAG methods, including single agent RAG systems or multi-agent systems (Nachimovsky et al., 2025; Chang et al., 2024; Weaviate, 2024) treat retriever as a static component of the system, and focus on improving the generator via prompting or model optimization. While these methods do not directly propose new IR modeling or training strategies, we believe it is critical for IR researchers to contextualize common agentic use cases and propose new IR model architectures better suited for these application scenarios.

F.2 Deployment of Modern IR Systems

Efficiency and Effectiveness Tradeoff Traditional retrieval systems face significant challenges when scaling to web-scale document corpus, and to deploy such systems requires a blend of science and engineering expertises (Dean et al., 2009; Huang et al., 2020; Li et al., 2021). In recent years, retrieval-augmented generation, conversational search and agentic systems with memory have been widely adopted for information access (Guu et al., 2020; Lewis et al., 2020b; Google, 2019; OpenAI, 2024; Google, 2024, *inter alia*). These applications often require multiple rounds of retrieval and dynamic corpus, urging for efficient and effective retrieval. Mainstream inference optimization frameworks such as vLLM (Kwon et al., 2023) and SGLang (Zheng et al., 2024) have provided support for embedding models. From the modeling perspective, an open question is to design and pre-train models for retrieval purposes (Warner et al., 2024; Nussbaum et al., 2025; Günther et al., 2023).

Robustness in Noisy Environment We discuss a few challenges in IR models’ deployment in noisy environment, especially when used in retrieval-augmented generation systems. We should note that while these challenges have been studied by prior works, it remains an open question on how to mitigate these challenges from the perspective of IR modeling and architectures.

- **Robustness to AI generated content.** With the advent of LLMs, the amount of AI-generated content is also increasing. Dai et al. (2024) show that neural retrievers are biased towards AI-generated documents. Xu et al. (2024a) show that similar problems persist in text-image retrieval models. Future IR modeling research should also consider the robustness of models to AI-generated content.
- **Robustness to adversarial attacks.** Recent works on RAG LLM safety have discussed the threat of corpus poisoning where injected harmful documents lead to unsafe outputs (Zhong et al., 2023; Xiang et al., 2024a; Deng et al., 2024, *inter alia*). This topic is also relevant to the safety of LLM agents using tools (Deng et al., 2025; Tian et al., 2023; Xiang et al., 2024b), noting the importance of IR models being robust to adversarial attacks for downstream applications.
- **Robustness to bias and toxicity.** As noted by a recent work (An et al., 2025), documents that contains biases and toxic materials can potentially jailbreak aligned LLMs. This observation highlights the importance for IR models to be robust to bias and toxic contents.
- **Robustness to imperfect retrieval results.** Different works have pointed out that existing RAG systems show performance degradation when the retrieval results contain irrelevant documents (Yoran et al., 2024; Chang et al., 2024; Yu et al., 2024c, *inter alia*). Therefore, the RAG paradigm demands more precise results from the retrieval models.
- **Robustness to out-of-distribution input.** Given the fact that modern neural retrieval models are trained with data-driven approaches, perhaps it is not surprising to find their performance may vary with different linguistic properties of the queries and documents, i.e., out-of-distribution input from the training data. Cao et al. (2025) conduct a rigorous benchmarking, and find formality, readability, politeness and grammatical correctness – fundamental aspects of real-world user-LLM queries – can lead to significant performance variances of retrievers and RAG systems. This observation highlights the importance of retrieval models’ robustness to OOD input (Gupta et al., 2024).

We refer readers for more detailed discussions on IR models’ robustness to dedicated surveys (Asai et al., 2024b; Liu et al., 2025; Zhou et al., 2024).