



Building Intelligent Search Systems: Advances in AI-Based Information Retrieval

Oleksii Segeda

Senior Data Engineer, Mapbox Washington, D.C., USA

OPEN ACCESS

SUBMITTED 11 April 2025

ACCEPTED 26 May 2025

PUBLISHED 04 June 2025

VOLUME Vol.07 Issue 06 2025

CITATION

Oleksii Segeda. (2025). Building Intelligent Search Systems: Advances in AI-Based Information Retrieval. The American Journal of Applied Sciences, 7(06), 06–11.
<https://doi.org/10.37547/tajas/Volume07Issue06-02>

COPYRIGHT

© 2025 Original content from this work may be used under the terms of the creative commons attributes 4.0 License.

Abstract: The exponential growth of digital content has driven the need for more intelligent, context-aware information retrieval systems. While traditional keyword-based search engines remain foundational, they often fall short of capturing deeper semantic meaning. This article explores the evolution, methodologies, and recent developments in intelligent information retrieval systems powered by artificial intelligence. Special attention is given to the use of machine learning, natural language processing (NLP), and neural networks to improve relevance, personalization, and contextual understanding, including the application of learning-to-rank techniques. The paper contrasts the strengths and limitations of conventional search technologies with those of AI-driven models. A critical part of the study focuses on potential risks associated with AI-based search engines, including environmental concerns linked to the heavy water consumption of data centers relying on water-based cooling systems. The research concludes that a holistic approach is needed in the design and implementation of AI-powered search systems—one that integrates ethical, cognitive, and environmental considerations. This article will be of interest to professionals in media and information technology, researchers, and developers engaged in building intelligent search infrastructures.

Keywords: information, artificial intelligence, search system, environmental risks.

Introduction: The ability to locate and make sense of information has always sat at the heart of scholarship and innovation. Over the past decade—particularly since the early 2020s—technical progress has lowered the barriers to retrieval while vastly enlarging what can be found. The most disruptive change is the infusion of artificial intelligence (AI) into every stage of the search pipeline. Contemporary engines do far more than list documents: they infer intent, distill arguments, and in many cases weave together new knowledge.

Classic retrieval frameworks—Boolean logic and the vector-space model chief among them—excel at matching strings but falter when a query is ambiguous, nuanced, or purely conceptual. Their limitations have spurred a turn toward “intelligent” search, grounded in machine learning and natural-language processing (NLP). By embedding statistical, linguistic, and behavioural signals, these systems evolve from static indexes into adaptive, user-aware advisers.

Deep learning methods drive much of this shift. Neural networks trained on massive corpora detect latent patterns and preferences that older heuristics overlook. Simultaneously, advances in NLP allow queries to be parsed at both syntactic and semantic levels, aligning system understanding more closely with human intent. Transformer architectures—BERT, GPT, and their many descendants—anchor today’s state-of-the-art: they capture context, gauge relevance, and generate personalised responses at scale.

Yet these benefits carry costs. Widespread deployment raises questions about privacy, transparency, energy consumption, and even the cognitive impact of outsourcing judgment to opaque models. The present article surveys the technical foundations of modern Interactive Information Retrieval (IIR), traces their historical trajectory, illustrates applications in the wild, and reflects on the broader societal and environmental stakes.

METHODS AND MATERIALS

To address our research questions we combined several complementary strategies:

- Comparative analysis and systematisation of prevailing retrieval models, highlighting convergences and divergences.

- Case-study review, juxtaposing theoretical constructs with field deployments.

- Synthesis of findings from academic journals, industrial white papers, and practitioner reports to generate a multidimensional perspective.

Although the scholarly literature on AI-enhanced search is still nascent, its relevance is undeniable. We therefore mapped key contributions across domains. Allan et al. [1] chart the prospects of generative AI for retrieval, spotlighting transformers and their integration into search platforms. Hambarde and Proença [2] trace the evolution from term-based ranking through semantic methods to neural approaches. Garlough-Shah [3] probes how AI reshapes user behaviour and search advertising, while White [4] examines agent-mediated interaction and the new tasks such agents enable. Hersh’s monograph [5] reminds us that classical techniques retain value amid AI expansion.

On the modeling front, Trabelsi et al. [6] review neural ranking architectures and outline future research avenues. Looking ahead, Zhu et al. [7] survey the incorporation of large language models (LLMs) into retrieval workflows, and Hersh [5] analyses the academic implications of generative AI. Zhang et al. [8] introduce *Agentic Information Retrieval*, in which LLM-driven agents enrich traditional pipelines with context-aware dialogue. Finally, Siddiqui [9] offers a granular, practice-oriented overview of AI adoption in libraries and information centres.

Together, these sources furnish both the empirical material and the conceptual scaffolding for the present study, enabling us to situate our analysis within the evolving landscape of AI-powered information retrieval.

RESULTS AND DISCUSSION

The exponential growth of digital content has intensified the need for retrieval tools that grasp meaning rather than merely match strings. Classical approaches—Boolean filters or vector-space scoring—anchor their judgments in exact keywords and therefore misread intent or overlook latent semantics [6]. By contrast, the current generation of search systems relies on artificial-intelligence techniques, most notably machine- and deep-learning, to migrate from surface-level matching to genuine semantic interpretation [3]. Below, we

review the principal models that now define intelligent information retrieval (IIR) [1].

Machine learning (ML) enables a search application to observe user behaviour, incorporate feedback, and update its ranking logic without hand-coded rules. In essence, ML algorithms infer statistical regularities from data and generalise them to unseen inputs. Deep learning—a rapidly advancing ML subfield—exploits multilayer neural networks whose expressive power eclipses that of earlier classifiers and regressors. These networks now dominate tasks ranging from document categorisation to query expansion and synthetic data generation.

Because of their versatility, ML methods permeate countless domains: natural-language processing, computer vision, speech-to-text transcription, spam detection, medical decision support, precision agriculture, and industrial robotics [6]. Predictive analytics, where organisations model customer churn, anticipate market swings, or quantify operational risk, likewise leans on the statistical inference and optimisation principles that ground ML [5]. What was once a niche research pursuit has, by 2025, matured into a foundational layer for digital infrastructure—from web search and recommender engines to bioinformatics and financial modelling [4].

Within the supervised-learning family, learning-to-rank (LTR) algorithms remain indispensable. They train on query–document pairs labelled for relevance and learn scoring functions that order previously unseen lists in a way that better mirrors human judgment [7]. LTR’s utility extends well beyond web search: recommendation engines, e-commerce catalogues, conversational agents, and social-media feeds all rely on it to surface the most pertinent items. Personalised services and rising user expectations ensure that LTR

continues to attract research and industrial attention in 2025 [4].

Natural-language processing (NLP) gives search engines the ability to parse synonyms, homonyms, grammatical nuance, and discourse context. Firms adopt NLP to automate customer service, power chatbots, and extract insight from text at scale [6]. The same techniques allow voice assistants to emulate natural dialogue, boosting capacity while trimming costs.

The arrival of transformer architectures—BERT, GPT, T5, and their many task-specific offshoots—has lifted retrieval quality markedly. By encoding entire sequences, transformers recover the full contextual meaning of both query and document, uncover implicit cues, resolve ambiguity, and model intricate linguistic dependencies [3]. Specialised variants for search, such as ColBERT or DistilBERT-QA, outperform earlier pipelines in question answering and fact extraction [5]. Trained through masked-token prediction and sentence-pair objectives, transformer models empower systems to:

- infer user intent instead of merely tallying keywords;
- carry out context-aware search over complex phrasing;
- enable multimodal retrieval that spans text, images, and spoken input;
- tailor results through fine-grained analysis of individual interests and histories.

Thus, modern intelligent search systems draw on a variety of information retrieval models, each offering distinct theoretical and practical solutions for data extraction and ranking. These models differ in their design and operational strategies, reflecting diverse approaches to search optimization.

Figure 1 illustrates the key models that underpin the construction of such systems.

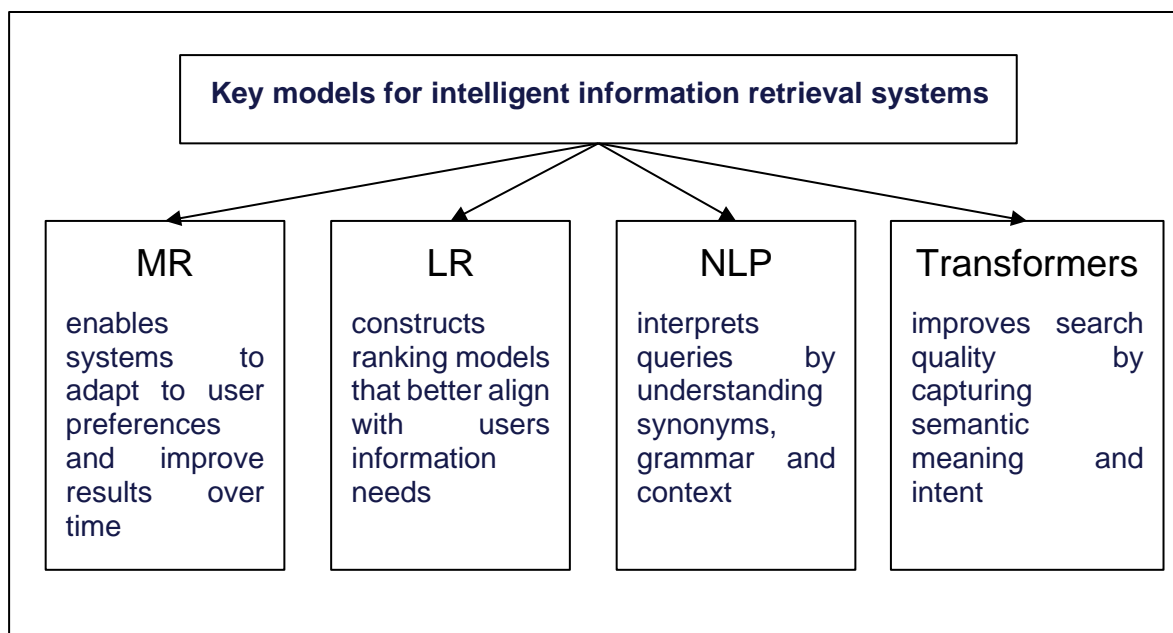


Figure 1 — Core models of information retrieval (compiled by the author based on original research)

Figure 1 summarises the principal models that underpin today's IIR platforms. Collectively, they transform search from a passive listing service into an interactive aide capable of conversation, anticipation, and autonomous knowledge extraction. Personalised ranking blends collaborative filtering, matrix factorisation, and neural embeddings, markedly improving recommendations on media and commerce sites [7]. Yet heightened personalisation also raises persistent concerns: safeguarding privacy, preventing manipulation, and clarifying how opaque models reach their decisions [4]. In short, the trajectory of intelligent retrieval is unmistakably towards greater adaptivity and user-centrism, but responsible deployment demands equal attention to transparency and trust.

AI-driven search remains among the fastest-moving frontiers in contemporary information technology. Whereas earlier engines concentrated on literal word overlap, today's semantic systems tailor results to the contextual meaning of a query and to an individual's profile—search history, location, and other behavioural signals—thus supporting large-scale analytics and highly contextualised retrieval [3]. Breakthroughs in machine learning, deep learning, and natural-language processing have opened a new era of search characterised by richer, context-aware answers and genuinely interactive machine–human exchanges [6].

links, the system delivers concise summaries, explanations, or analyses and then invites follow-up

Crucially, these systems grow more accurate with every session: the larger the stream of queries and feedback, the sharper their inferences become [5].

Since the early 2020s, information access has pivoted toward transformer-based conversational models such as OpenAI's GPT series [2]. Unlike conventional engines—Google, for instance—which index pages, tokenize terms, match postings, invoke ranking functions (BM25, PageRank), and finally present hyperlinks, transformer systems ingest a prompt, interpret and reorganise pertinent knowledge, and return an answer that approximates human reasoning [4]. Classic engines work well when confronted with a precise, unambiguous query that aligns with their indexing logic; transformers cope with vague wording, idioms, and cross-domain requests that once confounded search technology [5]. Artificial intelligence therefore streamlines retrieval, sparing users the time-consuming task of sifting through multiple sources [9].

Modern transformer platforms such as ChatGPT combine large-scale pre-training with sophisticated attention mechanisms, yielding a deep semantic grasp of language, nuanced context modelling, and natural conversational flow [7]. A typical interaction starts with free-form text input, proceeds through contextual interpretation and intent recognition, and culminates in a tailored response. Instead of offering a ranked list of dialogue [8]. Search is no longer a passive lookup operation but an active, conversational partnership.

This shift makes information discovery more intuitive, especially for users with limited digital literacy. In domains where rapid comprehension of dense material is essential—education, research, law, healthcare—the benefit is immediate and substantial. Yet the technology also introduces new trade-offs, summarised in Table

Table 1. Comparison of the advantages and disadvantages of AI-based intelligent search systems (compiled by the author based on original research)

Advantages	Disadvantages
Semantic query understanding — interprets the intent, not just keywords	Potential hallucinations — may generate plausible but inaccurate or false information
Personalization — adapts to user behavior, preferences, and context	Privacy concerns — users rely on answers without verifying sources
Speed and convenience — delivers structured, ready-to-use responses	Lack of transparency — models often do not cite sources or explain reasoning clearly
Natural language support — accepts queries in conversational form	Possible algorithmic bias — may inherit social, cultural, or political biases from training data. Natural language carries ambiguity in certain situations
Multimodal capabilities — integrates text, images, and audio inputs	Dependency on specific platforms — users may be locked into particular AI ecosystems
Conversational interface — supports dynamic dialogue and follow-up	High computational demands — requires significant processing power and may lack offline access
Summarization and synthesis — condenses and contextualizes large volumes of data	Ethical and legal concerns — including authorship, licensing, data privacy, and content accuracy

Despite the clear gains—streamlined access and higher processing efficiency—significant downsides persist, extending even to environmental impact. Key open problems include model interpretability, robustness against adversarial noise, and the continual need for balanced, high-quality training data [6]. Transparency is a prominent concern: unless a system is wrapped in a retrieval-augmented generation (RAG) pipeline, it rarely discloses its sources, undermining traceability [8]. Unlike legacy engines that visibly rank and link documents, generative models synthesise answers without explicit references [5]. Verifying such content becomes difficult,

trust can erode, and users may gradually relinquish the habit of cross-checking facts [5].

Over time this convenience risks dulling critical-thinking skills, independent inquiry, and comparative reasoning [9]. In educational contexts—where information literacy and cognitive autonomy are foundational—the danger is especially acute [7]. Concentration of influence is another worry: dominance by a handful of conversational systems (ChatGPT, Anthropic Claude, and the like) may consolidate control over information flows and introduce subtle ideological bias [3],

potentially suppressing or skewing knowledge representation [2].

CONCLUSION

Intelligent search technologies grounded in contemporary AI have begun to redefine the entire experience of information seeking. By combining large-scale pattern recognition with personalised modelling, they lighten cognitive effort, accelerate discovery, and adapt results to each user's context, thereby turning digital environments into far more responsive partners. Recent leaps in machine learning have translated directly into higher retrieval accuracy across medicine, commerce, and education, where timely, relevant insight carries concrete social value. Yet these same advances expose a counter-trend: the richer the automation, the thinner the role left for human judgement. When a single prompt elicits a fully formed answer, the skills of source evaluation, cross-reference, and critical reflection risk atrophy. Such dependency also amplifies exposure to misinformed or intentionally distorted content, while displacing whole categories of professional expertise.

REFERENCES

- Allan, J., Choi, E., Lopresti, D. P., & Zamani, H. (2024). Future of Information Retrieval Research in the Age of Generative AI. arXiv preprint arXiv:2402.12345. Retrieved April 1, 2025, from <https://arxiv.org/abs/2412.02043>
- Hambarde, K. A., & Proença, H. (2023). Information Retrieval: Recent Advances and Beyond. Universidade da Beira Interior. Retrieved April 27, 2025. Retrieved April 3, 2025, from <https://arxiv.org/abs/2301.08801>
- Garlough-Shah, Gabriel. The Rise of AI-powered Search Engines: Implications for Online Search Behavior and Search Advertising. MS thesis. University of Minnesota, 2024. Retrieved April 5, 2025
- White R. W. Advancing the Search Frontier with AI Agents //Communications of the ACM. – 2024. – T. 67. – №. 9. – C. 54-65. Retrieved April 7, 2025, from <https://arxiv.org/abs/2311.01235>
- Hersh W. Search Still Matters: Information Retrieval in the Era of Generative AI //Journal of the American Medical Informatics Association. – 2024. – T. 31. – №. 9. – C. 2159-2161.
- Trabelsi, M., Chen, Z., Davison, B. D., & Heflin, J. (2021). Neural Ranking Models for Document Retrieval. Information Retrieval Journal, 24(6), 400-444.
- Zhu, Y., Yuan, H., Wang, S., Liu, J., Liu, W., Deng, C., Chen, H., Liu, Z., Dou, Z., & Wen, J. (2023). Large Language Models for Information Retrieval: A Survey. arXiv preprint arXiv:2308.07107. Retrieved April 11, 2025, from <https://arxiv.org/abs/2308.07107>
- Zhang, W., Liao, J., Li, N., Du, K., & Lin, J. (2024). Agentic Information Retrieval. arXiv preprint arXiv:2410.09713. Retrieved April 11, 2025, from <https://arxiv.org/abs/2410.09713>
- Siddiqui, S. (2024). Artificial Intelligence in Information Retrieval: AI-based Techniques for Improving Search and Information Retrieval Systems in Both Libraries and Other Knowledge Hubs. Retrieved April 10, 2025, from <https://www.researchgate.net/publication/384805881>