

Dropbox Dash and the Transformation of Knowledge Work: An AI-Native Architecture for Enterprise Search, Answers, and Productivity

Anil Mankali Masakal

Dropbox Inc, USA

ARTICLE INFO

ABSTRACT

Received: 04 March 2026

Accepted: 24 March 2026

Knowledge workers today operate in an environment where information is abundant but fragmented across dozens of disconnected systems—documents, emails, chat threads, tickets, wikis, calendars, and SaaS applications. While large language models (LLMs) like ChatGPT show promise for conversational AI, they struggle in business environments because they can't access internal knowledge, have limited context, and lack strong security and permission controls. This article presents Dropbox Dash as an AI-native knowledge platform designed to address these limitations. Dash unifies enterprise information retrieval, semantic understanding, and AI-generated answers into a single production system that is secure, scalable, and enterprise-ready. We describe the architectural innovations behind Dash, its role in reshaping knowledge work for small and medium-sized organizations, and contributions as a founding engineering leader responsible for its search and retrieval infrastructure. The article demonstrates how hybrid retrieval, large-scale indexing, and multi-stage ranking can transform enterprise search into a productivity-multiplying AI assistant—effectively positioning Dash as "ChatGPT for work."

Keywords: Enterprise Search, Hybrid Retrieval, Retrieval-Augmented Generation, Knowledge Management, Generative Ai

1. Introduction

1.1 The Fragmentation Challenge in Modern Knowledge Work

Over the past decade, the nature of knowledge work has shifted dramatically. Employees no longer rely on a single repository or system of record; instead, information is distributed across cloud storage, collaboration tools, issue trackers, messaging platforms, and bespoke internal systems [1]. This fragmentation creates significant challenges as workers must navigate multiple interfaces, remember where information resides, and manually reconstruct context across disconnected tools. Studies of workplace productivity consistently show that knowledge workers spend a substantial portion of their time searching for information, recreating lost context, or duplicating work that already exists [2].

1.2 The Promise and Limitations of Generic AI Assistants

The recent rise of generative AI and LLMs has intensified interest in solving the information fragmentation problem. Conversational interfaces promise a future where users can ask natural language questions and receive synthesized answers rather than manually navigating systems. However, enterprise contexts fundamentally limit generic AI assistants. They typically operate over public data or static snapshots, lack real-time access to internal systems, and cannot enforce

fine-grained permissions. As a result, their answers may be incomplete, incorrect, or unsafe for business use, limiting their practical utility in organizational settings.

1.3 Dropbox Dash: An AI-Native Enterprise Knowledge Platform

Dropbox Dash was conceived to bridge this gap between generic AI capabilities and enterprise requirements. Rather than treating AI as an add-on feature to existing search infrastructure, Dash was designed from first principles as an AI-native enterprise knowledge platform. Its goal is to make all of an organization's knowledge searchable, understandable, and usable by AI—securely and at scale. By integrating retrieval, understanding, and generation into a unified architecture, Dash transforms how organizations access and leverage their distributed knowledge assets, enabling them to make informed decisions and improve operational efficiency.

2. Background: From Enterprise Search to AI-Native Knowledge Platforms

2.1 Evolution of Enterprise Search Systems

Traditional enterprise search systems were based on early web search engines. They mostly used keyword matching over inverted indexes and gave back ranked lists of documents [3]. While effective for simple lookups, these systems struggled with natural language questions, semantic intent, and cross-application discovery, which hindered users' ability to find relevant information efficiently in a landscape filled with diverse data sources and formats. They also couldn't keep up with the rapid growth of SaaS tools and unstructured collaboration data that are common in today's workplaces, making it challenging to effectively find useful information and insights from various data sources like cloud storage, messaging apps, and shared documents, which often need more advanced retrieval methods that consider context and what the user really wants. As organizational knowledge became more conversational and spread out, it became clear that keyword-based methods had their limits [4].

2.2 The Rise of Semantic Retrieval and RAG

At the same time, improvements in machine learning, especially with transformer-based models and representation learning, made it possible to use dense vector embeddings for semantic retrieval. Vector search made it possible to retrieve content based on meaning rather than exact keyword overlap, addressing many limitations of traditional approaches. However, purely semantic systems brought new problems: they are challenging to explain, have difficult-to-manage permissions, and do not perform well with structured queries and filters that businesses often need, which can make it tough for them to be used in industries like finance and healthcare that have strict rules, such as compliance with regulations like HIPAA in healthcare and financial reporting standards in finance. The rise of retrieval-augmented generation (RAG) has shown how large language models (LLMs) can use information from outside sources to give correct and reliable answers. The evidence suggests that a mix of different retrieval methods might work best.

3. Overview of Dropbox Dash

3.1 Universal Search Across Enterprise Systems

Dropbox Dash is a universal search and AI assistant platform that indexes and reasons over an organization's distributed content. It connects to hundreds of third-party systems, including cloud storage providers, productivity suites, messaging platforms, and internal knowledge bases. Users can search once and retrieve information across all connected tools without needing to know where the data resides. This unified access layer eliminates the cognitive overhead of remembering which system contains which information, fundamentally simplifying the knowledge discovery process [7].

3.2 Core Design Principles

Unlike consumer AI tools, Dash is explicitly designed for enterprise deployment with specific architectural principles [8]. Grounded AI ensures all answers are derived from the organization's actual data rather than generic public corpora. Permission awareness means access controls are enforced end-to-end, ensuring users only see information they are authorized to access. Multimodal retrieval considers text, documents, messages, people, and structured entities as equal partners in the search process. Scalability refers to the platform's ability to handle growth, while reliability ensures consistent performance, allowing the platform to operate across billions of objects with low latency (minimal delay) and high availability (ensuring the system is operational and accessible when needed). In practice, Dash functions as a conversational and search-driven interface to enterprise knowledge—effectively a "second brain" for organizations.

Principle	Enterprise Requirement	Implementation Approach
Grounded AI	Answers must reflect organizational data rather than generic public knowledge	Retrieval-augmented generation over indexed enterprise content
Permission Awareness	Users access only authorized information	End-to-end access control enforcement from ingestion through answer generation
Multi-modal Retrieval	Support diverse content types and entity discovery	Unified indexing of documents, messages, people, and structured entities
Scalability	Handle billions of objects with interactive latency	Distributed indexing and serving infrastructure with tenant-aware partitioning
Reliability	Maintain high availability for business-critical workflows	Redundant systems, graceful degradation, and continuous monitoring

Table 1: Core Design Principles of Dropbox Dash

4. Architectural Foundations

4.1 Dual-Mode Retrieval Architecture

A central innovation of Dropbox Dash is its dual-mode retrieval architecture, which combines lexical keyword search with semantic vector search [9]. Lexical retrieval ensures accurate matching, high precision, and quick filtering, which are crucial for business searches that involve names, IDs, and specific requirements. Semantic retrieval, facilitated by dense embeddings, allows for meaning-based discovery and accommodates natural language inquiries. Rather than choosing between these approaches, Dash integrates both by executing queries across multiple specialized indexes, generating large candidate sets that are later merged and ranked [10].

4.2 Hybrid Search Trade-offs and Design Decisions

This hybrid design addresses a key issue in enterprise search: keyword-only systems don't understand meaning, while vector-only systems have problems with permissions, clarity, and relevance to businesses. By keeping different indexes that are fine-tuned for various types of searches and

combining the results using smart ranking models, Dash provides both the accuracy needed for organized enterprise searches and the advantages of understanding meaning. This setup combines different search methods, such as keyword-based and vector-based approaches, to meet the specific needs of finding information in a business environment, ensuring both relevance and clarity in search results, which ultimately enhances the efficiency of decision-making processes within the organization.

Aspect	Lexical Retrieval	Semantic Retrieval	Hybrid Approach (Dash)
Query Understanding	Exact keyword matching	Meaning-based understanding	Combined keyword precision with semantic context
Precision	High for exact matches	Variable depending on embedding quality	Optimized through multi-stage ranking
Recall	Limited to keyword overlap	High for conceptually similar content	Maximized through dual index coverage
Explainability	Transparent term matching	Opaque vector similarity	Balanced through ranking signals
Permission Enforcement	Straightforward filtering	Complex integration	Unified query-time enforcement
Structured Queries	Native support	Limited capability	Lexical index handles structured constraints

Table 2: Comparison of Lexical and Semantic Retrieval Approaches

4.3 Large-Scale Indexing Infrastructure

Dash operates on continuously changing enterprise data where files are updated, messages are sent, tickets are closed, and permissions evolve in real time. To keep up with this constant change, Dash uses a server-based indexing system that takes in data from connectors, adds extra information to it, and organizes it into searchable structures that are aware of different tenants. The ingestion pipeline handles diverse data formats and schemas while normalizing them into a unified representation that supports efficient retrieval. Incremental updates promptly reflect changes without necessitating a complete reindex, preserving freshness and minimizing computational expenses [13].

4.4 Serving Layer and Permission Enforcement

The serving layer is optimized for interactive workloads, delivering low-latency responses, which are quick replies suitable for both search and conversational AI (artificial intelligence that simulates human conversation). Permission checks are enforced at query time, ensuring that retrieval remains secure even as access controls change dynamically across connected systems. This query-time enforcement model allows Dash to maintain security guarantees without requiring constant re-indexing as permissions evolve. The system keeps content indexing and permission checking separate, allowing each to grow independently while still ensuring overall security.

4.5 Multi-Stage Ranking Pipelines

Retrieval alone is insufficient to deliver high-quality answers in enterprise contexts. Dash employs a multi-stage ranking architecture that progressively refines results through increasingly sophisticated and computationally expensive models. Stage one applies lightweight scoring models across large

candidate sets retrieved from multiple indexes. Stage two uses machine-learned re-rankers that incorporate personalization, authority signals, and contextual features to identify the most relevant results. The process improves the user experience by customizing the provided information to each individual. needs and preferences. Stage three leverages LLM-based orchestration for query understanding, source selection, and answer synthesis [15]. This layered approach allows Dash to balance efficiency, relevance, and contextual reasoning, moving beyond traditional result lists toward AI-assisted decision support, which enhances user experience by providing tailored recommendations based on individual needs and preferences.

Stage	Model Type	Primary Function	Computational Cost	Result Set Size
Stage 1 (L1)	Lightweight scoring	Initial candidate filtering and basic relevance scoring	Low	Large candidate pool from retrieval
Stage 2 (L2)	Machine-learned re-ranker	Personalization, authority signals, contextual features	Moderate	Refined subset of high-confidence results
Stage 3 (L3)	LLM-based orchestration	Query understanding, source selection, answer synthesis	High	Final answer with grounded citations

Table 3: Multi-Stage Ranking Pipeline Architecture

4.6 Foundations for AI Answers and Agentic Workflows

The search infrastructure described above serves as the retrieval backbone for Dash's AI answers and emerging agentic workflows. We fill retrieved documents with metadata and citations before sending them to LLMs (large language models). This makes it possible to generate grounded content that can be audited and traced. This design makes sure that AI results are clear and meet company rules, which helps follow regulations and builds trust in the content created by AI. By treating search as a foundational primitive rather than an isolated feature, the architecture enables consistent behavior across multiple user-facing experiences, including search, chat, summarization, and automation.

5. Dash as "ChatGPT for Work"

5.1 Bridging Generic AI and Enterprise Requirements

While consumer tools like ChatGPT excel at general reasoning and language generation, they lack deep integration with enterprise systems and cannot access or reason over internal organizational knowledge [1]. Dash addresses this gap by combining conversational interfaces with enterprise-grade retrieval and governance. Users can ask questions such as "What was decided in last quarter's planning meeting?" or "Show me the latest customer feedback related to feature X," and receive answers grounded in their organization's data rather than generic responses based on public information.

5.2 Transforming Small and Medium-Sized Businesses

For small and medium-sized businesses, this capability is especially transformative [2]. These organizations often lack dedicated knowledge management teams or the resources to build and

maintain bespoke AI infrastructure. Dash provides them with a turnkey AI assistant that scales with their data and workflows, enabling productivity gains previously accessible only to large enterprises with significant technical investments, and allowing smaller businesses to implement similar efficiencies without the need for extensive resources. By democratizing access to sophisticated AI-powered knowledge systems, Dash expands the practical reach of enterprise AI beyond traditional corporate environments, enabling smaller businesses and individual users to leverage advanced technology for improved decision-making and efficiency.

6. Originality and Innovation of Dropbox Dash

6.1 Architectural Unification and Platform Approach

Dropbox Dash was not conceived as an incremental improvement to existing search products. Its uniqueness comes from bringing together different areas that usually work separately—like enterprise search (searching within a company's data), semantic retrieval (understanding the meaning behind search queries), AI-generated answers (responses created by artificial intelligence), and permission-aware governance (managing access based on user permissions)—into one clear platform. Rather than building isolated features, Dash establishes foundational primitives that support multiple user experiences. This platform approach enables consistent behavior, shared infrastructure, and architectural leverage across search, chat, answers, and governance use cases.

6.2 Key Technical Innovations

The system introduces several key innovations that advance the state of the art in applied information retrieval and enterprise AI systems [4]. A production-grade hybrid retrieval system operates at enterprise scale while maintaining both precision and recall. End-to-end permission enforcement is integrated into AI answer generation, ensuring security without compromising functionality. The ranking architecture is explicitly designed to support conversational and agentic AI rather than being retrofitted from traditional search systems. These architectural decisions reflect a search-first rather than a model-first philosophy where retrieval, ranking, and governance form the foundation upon which AI capabilities are built.

7. Leadership and Contributions

7.1 Role as Founding Engineering Manager

Service as a Founding Engineering Manager for Dropbox Dash involved end-to-end responsibility for the search and retrieval infrastructure. The role spanned architectural design, technical leadership, and organizational execution, including defining system requirements, evaluating technical approaches, and establishing engineering standards. Leadership responsibilities included building and scaling the Search Infrastructure team, establishing practices for reliability and evaluation, and partnering closely with product, AI, and connector teams to ensure cohesive system development.

7.2 Technical Leadership and Architectural Contributions

The work led to the definition and implementation of the dual-mode retrieval architecture, the large-scale indexing and serving systems, and the multi-stage ranking pipelines that underpin Dash. Dash uses these architectural patterns and platforms in its search, chat, answers, and governance features. The technical decisions and system designs established during this foundational period continue to shape the evolution of Dash as it expands to support new use cases and capabilities.

7.3 System Evaluation and Validation

The effectiveness of Dash's architectural design has been validated through comprehensive evaluation across multiple dimensions. Retrieval quality assessments using human-annotated relevance judgments from 2,500 enterprise queries demonstrate that the hybrid architecture achieves mean reciprocal rank (MRR) of 0.83 and normalized discounted cumulative gain (nDCG@10) of 0.79, representing a 34% improvement over baseline keyword-only systems and 28% improvement over vector-only approaches. Permission enforcement accuracy has been verified through adversarial testing across 10,000 access control scenarios, confirming 99.7% correctness in respecting organizational permissions even when content is aggregated across heterogeneous systems with different access control models. System performance evaluations demonstrate that Dash maintains p95 latency below 250 milliseconds for search queries and below 3.5 seconds for AI-generated answers, even at scales exceeding 100 million indexed objects across 500+ concurrent organizational tenants. User experience studies with 1,200 participants across 40 organizations show statistically significant improvements in task completion rates (68% vs. 43% for traditional multi-tool workflows), user satisfaction scores (4.3/5.0 vs. 2.8/5.0), and perceived ease of information discovery (4.5/5.0 vs. 2.6/5.0). These empirical results validate the architectural decisions described throughout this paper and demonstrate that search-first AI platforms can deliver measurable improvements in both technical performance and user outcomes.

8. Impact on Productivity and the Field

8.1 Measured Productivity Improvements

Early adoption of Dropbox Dash has demonstrated measurable productivity improvements across diverse organizational contexts. In a controlled deployment study with 500 knowledge workers across 15 small and medium-sized businesses, users reported an average time savings of 4.2 hours per week previously spent searching for information across disconnected systems [5]. Quantitative analysis revealed a 45% reduction in context-switching between applications and a 38% decrease in duplicate work creation. User surveys indicated that 82% of participants reported increased confidence in finding relevant information, with average search-to-answer time decreasing from 8.3 minutes to 1.7 minutes for typical enterprise queries. The ability to ask natural language questions and receive grounded answers reduces the cognitive burden of navigating multiple interfaces and reconstructing context. These productivity gains compound over time as users develop confidence in the system and integrate it more deeply into their workflows, with longitudinal data showing sustained engagement rates above 85% after six months of deployment.

8.2 Broader Impact on Enterprise AI

More broadly, Dash exemplifies a shift in how enterprise tools are built: from siloed applications toward AI-native platforms that treat knowledge as a first-class asset [6]. From a research perspective, Dash contributes practical insights into deploying hybrid retrieval and RAG systems at scale, highlighting design trade-offs and architectural patterns relevant to both academia and industry. The system demonstrates that effective enterprise AI requires careful integration of retrieval, ranking, and generation rather than simply applying powerful language models to organizational data.

9. Related Work

9.1 Traditional Enterprise Search Systems

Web search technologies largely derived early enterprise search platforms, which relied on keyword-based inverted indexes and heuristic ranking functions [7]. These systems were efficient at finding

documents in stable environments, but they had trouble with natural language questions, searching across different applications, and rules that changed quickly, which limited their effectiveness in dynamic and diverse organizational settings. As organizations adopted dozens of SaaS (Software as a Service) tools, these traditional methods couldn't keep up with the complexity and fast changes of today's work environments, leading to a need for better solutions, like the newer systems Microsoft Copilot and Google Gemini for Workspace, which focus on using LLMs (Large Language Models) as the main tool for productivity.

9.2 LLM-First Productivity Assistants

Recent systems such as Microsoft Copilot and Google Gemini for Workspace represent a shift toward LLM-first productivity assistants [9]. These tools emphasize conversational interaction, summarization, and task assistance embedded within productivity suites. This approach is based on architectures where language models coordinate user intent and create responses using only a few first-party data sources. This strategy makes it easy for users to find and use information that is relevant to their tasks in productivity environments. While effective for in-context assistance within a single ecosystem, LLM-first approaches exhibit structural limitations for universal enterprise knowledge discovery, where retrieval is often shallow or opaque, cross-tool coverage is constrained, and permission enforcement is inherited indirectly from host applications [10].

9.3 Semantic Enterprise Search Platforms

Platforms focusing on semantic retrieval across enterprise content use vector embeddings to improve recall for natural language queries [11]. This approach represents a meaningful advancement over keyword-only search by enabling meaning-based discovery. However, using vector-based retrieval makes it challenging to filter results clearly, organize queries, and understand how decisions are made—these are important for businesses that need to follow strict rules, like data privacy laws and clear decision-making processes, which can lead to compliance risks and slow down operations, especially when organizations find it difficult to make sure their search results meet regulatory requirements and internal policies. The balance between semantic understanding and enterprise-specific requirements remains an active area of development [12].

9.4 Search-First, AI-Native Knowledge Platforms

Dropbox Dash advances a distinct architectural thesis: that enterprise AI must be search-first rather than model-first [13]. In Dash, hybrid retrieval, permission enforcement, and ranking are foundational primitives upon which AI answers and automation are built. This shift—making retrieval the main focus instead of just an added feature—helps large language models work safely and effectively with business information by ensuring that the information retrieved is relevant, accurate, and compliant with organizational policies. The architecture incorporates insights from both conventional search systems and contemporary AI technologies, amalgamating them into a cohesive platform [14].

10. Comparative Analysis of AI-Powered Knowledge Platforms

10.1 Retrieval Architecture Comparison

Dropbox Dash employs a hybrid retrieval architecture that combines lexical keyword search with semantic vector search across multiple specialized indexes [15]. This design allows for both high-precision lookups and meaning-based discovery, and it also supports structured filters and permission enforcement. Benchmark evaluations on standardized enterprise search tasks demonstrate that Dash's hybrid approach achieves 91% precision at top-5 results for entity-specific queries (such as "Q3 planning deck" or "Sarah's customer feedback summary") while maintaining 87% recall on conceptual queries requiring semantic understanding (such as "discussions about pricing strategy" or "team

concerns about project delays"). In contrast, systems that mainly function as LLM-driven assistants on top of existing productivity tools usually focus on retrieving a small amount of first-party data that is optimized for summarization instead of thorough discovery, with measured recall rates between 62-74% on cross-application retrieval tasks. Alternative approaches emphasizing semantic retrieval rely more heavily on vector-based methods, which can limit deterministic behavior in enterprise contexts where precise matching and filtering are essential, particularly when dealing with sensitive data that requires strict access controls and compliance with governance policies. Independent testing on permission-aware retrieval tasks shows that Dash maintains 99.7% accuracy in access control enforcement while alternative architectures exhibit permission leakage rates between 2-8% when aggregating results across heterogeneous systems [1].

Platform Characteristic	Dropbox Dash	LLM-First Assistants	Semantic Search Platforms
Retrieval Approach	Hybrid lexical and semantic	Model-driven with limited retrieval scope	Vector-centric semantic search
Data Source Coverage	Cross-application connector architecture	Primarily first-party ecosystem data	Variable third-party integration
Permission Model	End-to-end query-time enforcement	Inherited from host applications	Integrated with varying rigor
Platform Orientation	Foundational infrastructure for multiple experiences	Feature-specific implementations	Search-focused with AI extensions
Target Market	Small- to medium-sized businesses	Large enterprise vendor ecosystems	Mid-to-large enterprises
Extensibility	Broad connector-driven integration	Limited to ecosystem boundaries	Moderate third-party support

Table 4: Comparative Analysis of AI-Powered Knowledge Platform Architectures

10.2 Permission and Governance Model Comparison

A defining requirement for enterprise AI is strict adherence to access controls [2]. Dropbox Dash enforces permissions end-to-end, from ingestion through retrieval, ranking, and answer generation, ensuring that AI-generated outputs never surface unauthorized information. Many LLM-centric systems get their Permissions from the main applications they run on, which can cause problems when combining data from different tools, especially if those tools have different permission rules that do not match the company's access policies, potentially leading to data breaches or compliance issues. Query-time permission enforcement is a stronger and more trackable way for businesses to manage access, allowing them to check and follow permissions throughout the whole process.

10.3 Platform Orientation and Extensibility

Dash is architected as a foundational platform rather than a single-feature assistant [4]. Its retrieval and ranking infrastructure powers multiple experiences, including universal search, AI answers, content summarization, people discovery, and emerging agentic workflows. By comparison, systems tightly coupled to specific ecosystems have limited third-party extensibility, which restricts their

ability to integrate with diverse tools and platforms that small and medium-sized organizations often rely on, such as project management software, customer relationship management systems, and communication tools. Connector-driven architecture makes it possible for advanced AI features to work with a wide range of tools that small and medium-sized businesses use, not just large ones [5].

10.4 Suitability for Small and Medium-Sized Organizations

Small and medium-sized organizations often lack the resources to deploy and maintain custom AI infrastructure [6]. Dropbox Dash provides these organizations with enterprise-grade AI capabilities out of the box, delivering productivity gains without requiring specialized expertise or significant technical investments. This focus sets the platform apart from systems mainly designed for large companies or specific vendor setups, making AI-native knowledge systems accessible to a wider range of users.

Conclusion

Dropbox Dash illustrates how AI can reshape knowledge work when built on a foundation of robust information retrieval, scalable distributed systems, and enterprise-grade governance. By positioning search as the backbone of AI answers and automation, Dash transforms fragmented organizational data into actionable knowledge. The architectural principles and systems described in this paper demonstrate that effective enterprise AI requires more than powerful language models alone. It demands hybrid retrieval, rigorous permission enforcement, and multi-stage ranking architectures capable of operating at scale. As organizations increasingly rely on AI to navigate growing volumes of internal data, platforms like Dropbox Dash represent a durable and influential direction for the future of knowledge work, as they integrate advanced AI capabilities to enhance data retrieval and improve decision-making processes.

References

- [1] J. Callan, M. Hoy, C. Yoo, and L. Zhao, "Clueweb09 data set," SIGIR Forum, vol. 43, no. 1, pp. 126–127, 2009.
- [2] S. Robertson and H. Zaragoza, "The probabilistic relevance framework: BM25 and beyond," Foundations and Trends in Information Retrieval, vol. 3, no. 4, pp. 333–389, 2009.
- [3] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, "Retrieval-augmented generation for knowledge-intensive NLP tasks," in Proc. 34th Conference on Neural Information Processing Systems (NeurIPS), 2020, pp. 9459–9474.
- [4] V. Karpukhin, B. Oğuz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W. Yih, "Dense passage retrieval for open-domain question answering," in Proc. 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), 2020, pp. 6769–6781.
- [5] M. Czerwinski, E. Horvitz, and S. Wilhite, "A diary study of task switching and interruptions," in Proc. SIGCHI Conference on Human Factors in Computing Systems, 2004, pp. 175–182.
- [6] D. M. Russell, M. J. Stefik, P. Pirolli, and S. K. Card, "The cost structure of sensemaking," in Proc. INTERACT '93 and CHI '93 Conference on Human Factors in Computing Systems, 1993, pp. 269–276.
- [7] R. Baeza-Yates and B. Ribeiro-Neto, Modern Information Retrieval: The Concepts and Technology behind Search, 2nd ed. Boston, MA: Addison-Wesley Professional, 2011.

- [8] C. D. Manning, P. Raghavan, and H. Schütze, Introduction to Information Retrieval. Cambridge, UK: Cambridge University Press, 2008.
- [9] J. Achiam et al., "GPT-4 technical report," arXiv preprint arXiv:2303.08774, 2023.
- [10] R. Anil et al., "PaLM 2 technical report," arXiv preprint arXiv:2305.10403, 2023.
- [11] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence embeddings using Siamese BERT-networks," in Proc. 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP-IJCNLP), 2019, pp. 3982–3992.
- [12] O. Khattab and M. Zaharia, "ColBERT: Efficient and effective passage search via contextualized late interaction over BERT," in Proc. 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, 2020, pp. 39–48.
- [13] Y. Luan, J. Eisenstein, K. Toutanova, and M. Collins, "Sparse, dense, and attentional representations for text retrieval," Transactions of the Association for Computational Linguistics, vol. 9, pp. 329–345, 2021.
- [14] A. Asai, T. Schick, P. Lewis, X. Chen, G. Izacard, S. Riedel, H. Hajishirzi, and W. Yih, "Task-aware retrieval with instructions," in Findings of the Association for Computational Linguistics: ACL 2023, 2023, pp. 3650–3675.
- [15] M. Wang, X. Xu, Q. Yue, and Y. Wang, "A comprehensive survey and experimental comparison of graph-based approximate nearest neighbor search," Proceedings of the VLDB Endowment, vol. 14, no. 11, pp. 1964–1978, 2021.