

---

# FinAgentBench: A Benchmark Dataset for Agentic Retrieval in Financial Question Answering

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 information retrieval (IR) is critical in the financial domain, where investors must  
2 identify relevant information from large collections of documents. Traditional IR  
3 methods—whether sparse or dense—often fall short in retrieval accuracy, as it  
4 requires not only capturing semantic similarity but also performing fine-grained  
5 reasoning over document structure and domain-specific knowledge. Recent ad-  
6 vances in large language models (LLMs) have opened up new opportunities for  
7 *retrieval with multi-step reasoning*, where the model ranks passages through itera-  
8 tive reasoning about which information is most relevant to a given query. However,  
9 there exists no benchmark to evaluate such capabilities in the financial domain. To  
10 address this gap, we introduce FINAGENTBENCH, the first large-scale benchmark  
11 for evaluating retrieval with multi-step reasoning in finance – a setting we term  
12 *agentic retrieval*. The benchmark consists of 3,429 expert-annotated examples on  
13 S&P-100 listed firms and assesses whether LLM agents can (1) identify the most  
14 relevant document type among candidates, and (2) pinpoint the key passage within  
15 the selected document. Our evaluation framework explicitly separates these two  
16 reasoning steps to address context limitations. This design enables to provide a  
17 quantitative basis for understanding retrieval-centric LLM behavior in finance. We  
18 evaluate a suite of state-of-the-art models and further demonstrated how targeted  
19 fine-tuning can significantly improve agentic retrieval performance.

## 20 1 Introduction

21 Information Retrieval (IR) is a foundational research field that studies how to effectively search for and  
22 retrieve relevant information from large-scale text collections [21, 23]. Its practical importance in real-  
23 world applications has made it one of the central and long-standing areas in computer science since  
24 the early days of computing [15]. IR has evolved from sparse, term-frequency-based methods [24],  
25 which rely on exact keyword matches, to dense neural retrieval models that embed text into continuous  
26 latent spaces to capture deeper semantic [10]. In finance, accurate retrieval is critical, as investors  
27 depend on precise access to vast filings and reports to make high-stakes, time-sensitive decisions.  
28 To support this need, finance-specific IR benchmarks have been developed [8, 6], enabling rigorous  
29 evaluation of both sparse and dense retrievers in this complex, data-rich domain.

30 However, recent studies reveal persistent accuracy ceilings for both sparse and dense retrieval  
31 methods [12, 2], especially in domains that demand fine-grained understanding and structured  
32 reasoning over complex documents. To overcome these limitations, recent research has turned to  
33 large language models (LLMs), which bring strong language understanding and the ability to process  
34 long contexts [7, 13, 29]. One promising direction is *generative retrieval*, where an LLM generates  
35 the index of the most relevant document given a query and document collection [16]. These results  
36 show impressive gains in performance, indicating a shift toward retrieval systems that embed deeper  
37 reasoning capabilities [1, 18].

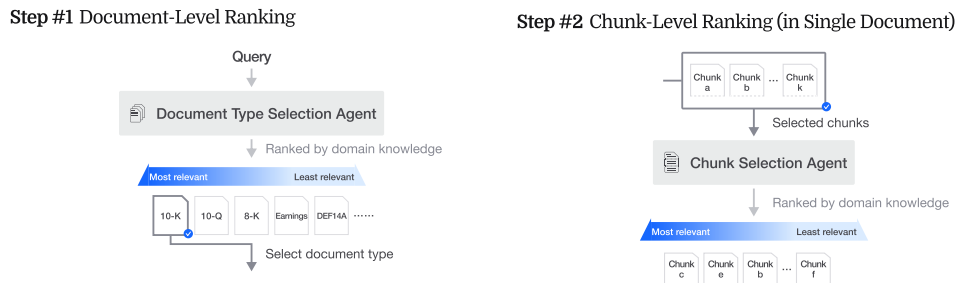


Figure 1: Agentic retrieval pipeline in FINAGENTBENCH. The process consists of two stages. (Top) Given a natural-language query, the **Document Type Selection Agent** ranks five SEC filing types (10-K, 10-Q, 8-K, earnings transcripts, and DEF-14A) and selects the most relevant type. (Bottom) The corresponding document is segmented into paragraph-level chunks, from which the **Chunk Selection Agent** identifies the top- $k$  passages.

38 In the financial domain, this shift is particularly relevant. Financial documents are typically long  
 39 and dense, often requiring multiple steps of reasoning: first to identify which document type (e.g.,  
 40 10-K, earnings transcript) best matches the information need, and then to locate the specific evidence  
 41 within it. As such, simple one-shot retrieval is often insufficient—accurate performance hinges on  
 42 multi-step reasoning that reflects how professionals actually search for information. Despite the  
 43 practical importance of such multi-step reasoning, no benchmark currently exists to evaluate LLMs  
 44 in this setting. This leaves open critical questions about whether LLMs can serve as effective retrieval  
 45 agents in high-stakes domains like finance—where precision, explainability, and structured navigation  
 46 are essential.

47 To address this gap, we introduce the first large-scale benchmark, FINAGENTBENCH, for evaluating  
 48 generative retrieval systems in the financial domain, focusing on a setting we term *agentic retrieval*.  
 49 Unlike prior benchmarks that assess retrieval in a single stage, our proposed benchmark evaluates the  
 50 ability of LLM agents to reason and retrieve the relevant information through a two-stage pipeline: (1)  
 51 identifying which document to retrieve, evaluated via document-level ranking, and (2) determining  
 52 which part of the document to focus on, evaluated via chunk-level ranking. The benchmark contains  
 53 3,429 samples and is constructed from real-world financial documents paired with expert-written  
 54 queries and annotations, reflecting authentic use cases faced by professional investors. By capturing  
 55 both retrieval accuracy and reasoning depth, it provides a foundation for systematically analyzing the  
 56 strengths and limitations of LLM-based generative retrieval in high-stakes domains like finance.

57 The main contribution of this work is three-fold:

- 58 • We propose FINAGENTBENCH, to the best of our knowledge the first large-scale benchmark  
 59 for evaluating agentic retrieval in finance, featuring 3,429 expert-annotated samples across  
 60 document- and chunk-level ranking tasks.
- 61 • We evaluate a range of state-of-the-art LLMs on our benchmark, revealing their performance  
 62 in accuracy and reasoning when applied to real-world financial retrieval scenarios.
- 63 • We further investigate the impact of fine-tuning the LLMs on agentic retrieval tasks, demon-  
 64 strating that targeted supervision can substantially improve both document selection and  
 65 chunk-level reasoning performance.

## 66 2 FINAGENTBENCH Dataset

67 In this section, we describe the concept of FINAGENTBENCH, a large-scale benchmark for evaluating  
 68 agentic retrieval in finance. To retrieve accurate information, financial retrieval requires multi-step  
 69 reasoning due to both the volume of data and the regularity of financial disclosures [5]. When  
 70 leveraging the reasoning capabilities of large language models to improve accuracy in retrieval, the  
 71 extensive length and redundancy of financial documents—where even a single 10-K can exceed  
 72 hundreds of pages—pose a significant challenge, making it inefficient to process all content without  
 73 any filtering or prioritization [19]. Therefore, to maintain efficiency, a system should first select  
 74 the document type most likely to contain the answer—feasible information because filings follow  
 75 predictable conventions, with different information consistently organized by document type (e.g.,

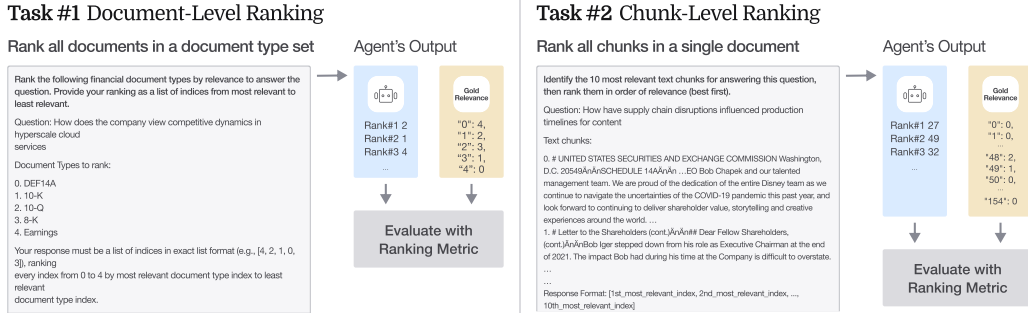


Figure 2: Examples of the two retrieval tasks in FINAGENTBENCH. The benchmark comprises: (Left) Document-level ranking, where the agent ranks five SEC document types based on their relevance to the input query. All document types are ordered by relevance, and the full ranking is used as the gold label during evaluation. (Right) Chunk-level ranking, where the agent selects and orders the top- $k$  most relevant passages from the selected document. Each chunk is annotated with a ground-truth relevance label—0 (irrelevant), 1 (partially relevant), or 2 (directly relevant)—and predictions are evaluated using MRR, MAP, and nDCG.

Table 1: Evaluation of reasoning LLMs on the *Document Ranking* and *Chunk Ranking* tasks.

Model	Document Ranking			Chunk Ranking		
	nDCG@5	MAP@5	MRR@5	nDCG@5	MAP@5	MRR@5
GPT-o3	0.770	0.829	0.875	0.351	0.257	0.538
Claude-Opus-4	0.773	0.840	0.875	0.418	<b>0.307</b>	<b>0.568</b>
Claude-Sonnet-4	<b>0.783</b>	<b>0.849</b>	<b>0.892</b>	<b>0.419</b>	0.296	0.567

76 risk factors in 10-Ks, strategic commentary in earnings calls). It should then identify the relevant  
77 chunk or passage within the selected document.

78 This motivates a two-stage retrieval process—document selection followed by passage selec-  
79 tion—which we term *agentic retrieval*, as it reflects the sequential reasoning steps taken by experts.  
80 Figure 1 provides an overview of the *agentic retrieval* workflow we benchmark. At test time, the  
81 system is provided with a natural-language query  $q$ , typically issued by a professional investor about  
82 a specific firm, along with a document collection  $\mathcal{D}$  containing over 35K U.S. corporate disclosures  
83 spanning from 2010 to 2024. In FINAGENTBENCH, we include 10-K, 10-Q, 8-K, earnings call  
84 transcripts, and DEF 14A proxy statements, some of the most commonly used public filings that are  
85 crucial for financial retrieval. The retrieval task follows a two-stage reasoning pipeline:

86 **Stage 1: Document-Level Ranking.** Rather than searching the entire corpus, the agent first  
87 identifies the document type most likely to contain the answer. Given the type set

$$\mathcal{T} = \{10\text{-K}, 10\text{-Q}, 8\text{-K}, \text{Earnings}, \text{DEF14A}\} \quad (1)$$

88 , the model produces a ranking over  $\mathcal{T}$ . This stage evaluates the model’s understanding of finance-  
89 specific reporting conventions—for instance, risk factors typically appear in 10-Ks, while shareholder  
90 proposals are found in DEF-14A filings.

91 **Stage 2: Chunk-Level Ranking.** The selected document ( $d_{t^*} \in \mathcal{D}$ ) is split into non-overlapping  
92 passages  $\mathcal{C}(d_{t^*}) = \{c_1, \dots, c_M\}$ . The agent ranks these chunks and returns the top  $k$  passages  
93 indexed as  $\langle c_{(1)}, \dots, c_{(k)} \rangle$ .

94 **Annotations.** Each query in FINAGENTBENCH is generated and annotated by domain experts. At  
95 the document-type level, a gold label  $t^G$  in  $\mathcal{T}$  is provided with an gold ranking . At the chunk level,  
96 every chunk element  $\mathcal{C}_{(i)}$  in  $\mathcal{C}(d_{t^*})$  is annotated as  $\mathcal{C}_{(i)}^G$  with a gold relevance score. For annotating  
97 those relevance scores, we followed TREC Eval [17]: 0 (irrelevant), 1 (partially relevant), and 2  
98 (directly relevant).

Table 2: Impact of reinforcement fine-tuning on GPT-o4-mini across both retrieval tasks.

Task	w/o fine-tuning			w/ fine-tuning		
	nDCG@5	MAP@5	MRR@5	nDCG@5	MAP@5	MRR@5
Document Ranking	0.758	0.826	0.872	<b>0.808</b>	<b>0.865</b>	<b>0.933</b>
Chunk Ranking	0.345	0.256	0.526	<b>0.371</b>	<b>0.274</b>	<b>0.587</b>

### 99 3 Experiments

100 We evaluate a range of reasoning-capable LLMs on FINAGENTBENCH to assess their ability to  
 101 perform agentic retrieval in finance. Our experiments focus on two subtasks: document-type ranking  
 102 and chunk-level passage selection. We also study the effect of domain-specific fine-tuning on retrieval  
 103 performance.

#### 104 3.1 Experimental Setup

105 We benchmark three commercial LLMs—GPT-o3, Claude-Opus-4, and  
 106 Claude-Sonnet-4—using zero-shot prompting. For each query, models are prompted to  
 107 complete two tasks: (1) rank document types from a set of five SEC filing categories, and (2) rank  
 108 paragraph-level chunks from the selected document. We split both the document-type ranking and  
 109 chunk-level relevance tasks into training and evaluation sets using an 80/20 split. To assess the  
 110 impact of domain adaptation, we additionally evaluate GPT-o4-mini before and after reinforcement  
 111 fine-tuning provided by OpenAI<sup>1</sup> on held-out randomly sampled 10% of training splits from  
 112 FINAGENTBENCH.

113 Performance is measured using standard ranking metrics measure by top-5 results: normalized  
 114 Discounted Cumulative Gain (nDCG), Mean Average Precision (MAP), and Mean Reciprocal Rank  
 115 (MRR). For chunk-level evaluation, we report metrics over the top ranked passages, evaluated against  
 116 expert-annotated graded relevance scores. All models operate in a retrieval-only setting without  
 117 access to external tools or retrieval augmentation.

#### 118 3.2 Results

119 **Task Performance.** Table 1 compares reasoning-capable LLMs on document- and chunk-level  
 120 ranking tasks. At the document level, all models perform substantially above random, indicating  
 121 that LLMs possess strong priors over financial reporting structure. The best-performing model  
 122 reaches near-perfect retrieval, while smaller gaps across models suggest that architectural or scale  
 123 differences still matter for capturing subtle type-level cues. In contrast, chunk-level ranking proves  
 124 more challenging: overall scores are lower, reflecting the complexity of fine-grained retrieval within  
 125 long disclosures. Here, the two Claude variants perform comparably, with one slightly stronger in  
 126 overall ranking quality and the other leading in precision-oriented metrics. This highlights that while  
 127 LLMs can reliably handle document-type discrimination, chunk-level reasoning remains a demanding  
 128 setting that exposes model differences more sharply.

129 **Impact of Fine-Tuning.** Table 2 evaluates the impact of reinforcement fine-tuning on GPT-o4-mini  
 130 across both retrieval stages. Fine-tuning yields substantial gains: for document-type ranking, nDCG  
 131 improves from 0.758 to 0.808, and MRR increases from 0.872 to 0.933. Similar improvements are  
 132 observed in chunk ranking, with MRR rising from 0.526 to 0.587. These results demonstrate that  
 133 domain-specific supervision significantly enhances retrieval accuracy, both in selecting the correct  
 134 document type and in pinpointing relevant information within the document. This highlights the  
 135 importance of adapting LLMs to financial reasoning tasks through task-aligned training signals.

### 136 4 Conclusion

137 We introduced FINAGENTBENCH, a large-scale benchmark designed to evaluate agentic retrieval  
 138 capabilities of large language models in the high-stakes domain of finance. The benchmark simulates  
 139 realistic investor queries over a diverse set of corporate filings, requiring models to reason over both  
 140 document types and intra-document content.

141 FINAGENTBENCH provides a new foundation for studying end-to-end retrieval behaviors in complex  
 142 domains. Future work may explore enhancing agentic retrieval performance and joint modeling of  
 143 retrieval and generation for investment decision support.

<sup>1</sup>OpenAI Reinforcement Fine-Tuning

144 **References**

145 [1] Hongru Cai, Yongqi Li, Ruifeng Yuan, Wenjie Wang, Zhen Zhang, Wenjie Li, and Tat-Seng  
146 Chua. Exploring training and inference scaling laws in generative retrieval. *arXiv preprint*  
147 *arXiv:2503.18941*, 2025.

148 [2] Brian J Chan, Chao-Ting Chen, Jui-Hung Cheng, and Hen-Hsen Huang. Don't do rag: When  
149 cache-augmented generation is all you need for knowledge tasks. In *Companion Proceedings of*  
150 *the ACM on Web Conference 2025*, pages 893–897, 2025.

151 [3] Sanchit Chanana, Hyung Won Chung, Nan Du, Jeffrey Zhao, et al. Don't do RAG: When  
152 cache-augmented generation is all you need. *arXiv preprint arXiv:2412.15605*, 2024.

153 [4] Zhiyu Chen, Wenhui Chen, Sameena Shah, et al. Finqa: A dataset of numerical reasoning over  
154 financial data. In *EMNLP*, 2021.

155 [5] Jaeyoung Choe, Jihoon Kim, and Woohwan Jung. Hierarchical retrieval with evidence curation  
156 for open-domain financial question answering on standardized documents. *arXiv preprint*  
157 *arXiv:2505.20368*, 2025.

158 [6] Chanyeol Choi, Jihoon Kwon, Jaeseon Ha, Hojun Choi, Chaewoon Kim, Yongjae Lee, Jy yong  
159 Sohn, and Alejandro Lopez-Lira. Finder: Financial dataset for question answering and evaluat-  
160 ing retrieval-augmented generation. *arXiv preprint arXiv:2504.15800*, 2025.

161 [7] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun,  
162 Haofen Wang, and Haofen Wang. Retrieval-augmented generation for large language models:  
163 A survey. *arXiv preprint arXiv:2312.10997*, 2(1), 2023.

164 [8] Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie  
165 Vidgen. Financebench: A new benchmark for financial question answering. *arXiv preprint*  
166 *arXiv:2311.11944*, 2023.

167 [9] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov,  
168 Danqi Chen, and Wen tau Yih. Dense passage retrieval for open-domain question answering. In  
169 *EMNLP*, 2020.

170 [10] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick SH Lewis, Ledell Wu, Sergey Edunov,  
171 Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering.  
172 In *EMNLP (1)*, pages 6769–6781, 2020.

173 [11] X Li, J Jin, Y Zhou, Y Zhang, P Zhang, Y Zhu, and Z Dou. From matching to generation: A sur-  
174 vey on generative information retrieval (2024). URL [https://api.semanticscholar.org/CorpusID/](https://api.semanticscholar.org/CorpusID/269303210)  
175 [269303210](https://api.semanticscholar.org/CorpusID/269303210).

176 [12] Yu-An Liu, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng.  
177 Robust neural information retrieval: An adversarial and out-of-distribution perspective. *arXiv*  
178 *preprint arXiv:2407.06992*, 2024.

179 [13] Kun Luo, Minghao Qin, Zheng Liu, Shitao Xiao, Jun Zhao, and Kang Liu. Large language  
180 models as foundations for next-gen dense retrieval: A comprehensive empirical assessment.  
181 *arXiv preprint arXiv:2408.12194*, 2024.

182 [14] Maurício Sousa Maia, Alberto H. F. Laender, P. Gnther, and Maria da Graça Campos Pimentel.  
183 Fiqa: A collection for aspect-based opinion mining and question answering. In *ECIR*, 2018.

184 [15] Mandar Mitra and BB Chaudhuri. Information retrieval from documents: A survey. *Information*  
185 *retrieval*, 2:141–163, 2000.

186 [16] Thong Nguyen and Andrew Yates. Generative retrieval as dense retrieval. *arXiv preprint*  
187 *arXiv:2306.11397*, 2023.

188 [17] Joao Palotti, Harris Scells, and Guido Zuccon. Trectools: an open-source python library for  
189 information retrieval practitioners involved in trec-like campaigns. SIGIR'19. ACM, 2019.

- 190 [18] Ronak Pradeep, Kai Hui, Jai Gupta, Adam D Lelkes, Honglei Zhuang, Jimmy Lin, Donald  
191 Metzler, and Vinh Q Tran. How does generative retrieval scale to millions of passages? *arXiv*  
192 *preprint arXiv:2305.11841*, 2023.
- 193 [19] Varshini Reddy, Rik Koncel-Kedziorski, Viet Dac Lai, Michael Krumdick, Charles Lovering,  
194 and Chris Tanner. Docfinqa: A long-context financial reasoning dataset. *arXiv preprint*  
195 *arXiv:2401.06915*, 2024.
- 196 [20] Stephen E. Robertson and Steve Walker. Some simple effective approximations to the 2-poisson  
197 model for probabilistic weighted retrieval. In *SIGIR*, pages 232–241, 1994.
- 198 [21] Gerard Salton. *Automatic information organization and retrieval*. McGraw Hill Text, 1968.
- 199 [22] Gerard Salton and Christopher Buckley. Term-weighting approaches in automatic text retrieval.  
200 *Information Processing & Management*, 11(5):513–523, 1975.
- 201 [23] Hinrich Schütze, Christopher D Manning, and Prabhakar Raghavan. *Introduction to information*  
202 *retrieval*, volume 39. Cambridge University Press Cambridge, 2008.
- 203 [24] Karen Sparck Jones. A statistical interpretation of term specificity and its application in retrieval.  
204 *Journal of documentation*, 28(1):11–21, 1972.
- 205 [25] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei.  
206 Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368*,  
207 2024.
- 208 [26] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.  
209 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,  
210 2023.
- 211 [27] Xiaoju Ye, Zhichun Wang, and Jingyuan Wang. Infinite retrieval: Attention-enhanced llms in  
212 long-context processing. *arXiv preprint arXiv:2502.12962*, 2025.
- 213 [28] Lei Zhu, Fandong Meng, Zhen Wu, Haiyang Song, and Jie Zhou. TAT-QA: A question  
214 answering benchmark on a hybrid of tabular and textual content in finance. In *ACL*, 2021.
- 215 [29] Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan  
216 Chen, Zheng Liu, Zhicheng Dou, and Ji-Rong Wen. Large language models for information  
217 retrieval: A survey. *arXiv preprint arXiv:2308.07107*, 2023.

## 218 A Related Work

219 **Retrieval Systems.** IR systems has progressed from early sparse methods like TF-IDF and  
220 BM25 [22, 20], which rely on exact lexical overlap, to dense dual encoders that map queries  
221 and documents into a shared semantic space. Notable advances such as Dense Passage Retrieval  
222 (DPR) [9] and E5-Mistral [25] have improved performance on open-domain QA by capturing deeper  
223 semantics. Despite these gains, both sparse and dense retrievers struggle with multi-hop reasoning  
224 and long-context queries—highlighting inherent limitations in fixed retrieval pipelines.

225 To address this, *generative retrieval* reframes the task as sequence generation, enabling models to  
226 directly produce relevant document identifiers or content [11]. Recent approaches such as CAG [3],  
227 and Infinite Retrieval [27] leverage LLMs’ language modeling and attention capabilities by caching the  
228 documents to bypass traditional indexing. In parallel, agentic reasoning frameworks like ReAct [26]  
229 incorporate tool use and step-by-step reasoning, suggesting the potential of the LLMs to iteratively  
230 decide what to retrieve and why—bridging retrieval with planning and decision-making.

231 **Finance-Retrieval Benchmarks.** Retrieval in the financial domain presents unique challenges:  
232 documents are long, queries are complex, and precision is critical. Existing datasets such as FinQA [4],  
233 TAT-QA [28], and FiQA [14] address specific financial reasoning needs, while more recent corpora  
234 like FinanceBench [8] and FinDER [6] support open-domain and retrieval-augmented generation  
235 tasks. However, these benchmarks still assume retrieval as a fixed subroutine, typically relying on

236 vector search or keyword search—and do not evaluate whether LLMs can reason about what to  
237 retrieve or where to look within long documents.

238 Our proposed FINAGENTBENCH fills this gap by directly evaluating agentic retrieval in finance. It  
239 assesses whether an LLM agent can (1) generate the correct document identifier and (2) rank the most  
240 relevant chunks within that document, providing a rigorous testbed for end-to-end reasoning-driven  
241 retrieval in high-stakes settings.