

FINCARDS: Card-Based Analyst Reranking for Financial Document Question Answering

Anonymous ACL submission

Abstract

Financial question answering (QA) over long corporate filings requires evidence to satisfy strict constraints on entities, financial metrics, fiscal periods, and numeric values. However, existing LLM-based rerankers primarily optimize semantic relevance, leading to unstable rankings and opaque decisions on long documents. We propose FINCARDS, a structured reranking framework that reframes financial evidence selection as *constraint satisfaction* under a finance-aware schema. FINCARDS represents filing chunks and questions using aligned schema fields (entities, metrics, periods, and numeric spans), enabling deterministic field-level matching. Evidence is selected via a multi-stage tournament reranking protocol with stability-aware aggregation, producing auditable decision traces. Across two corporate filing QA benchmarks, FINCARDS substantially improves early-rank retrieval quality over both lexical and LLM-based reranking baselines, while significantly reducing ranking variance, without requiring model fine-tuning or unpredictable inference budgets. The code is available at <https://anonymous.4open.science/r/Fincards-0414>.

1 Introduction

Financial question answering (QA) over corporate filings is often framed as a retrieval problem, but in practice it is a reranking problem under strict financial constraints. Correct evidence must simultaneously match the queried metric, fiscal period, and entity, and often include an explicit numerical value. These signals are sparsely distributed within filings that span hundreds of pages and are heavily interleaved with boilerplate disclosures and recurring statements. As a result, the problem reduces to reliably reranking candidate passages within a single document (Figure 1) so that the top results satisfy all required financial conditions (Chen et al., 2021, 2022; Zhu et al., 2021).

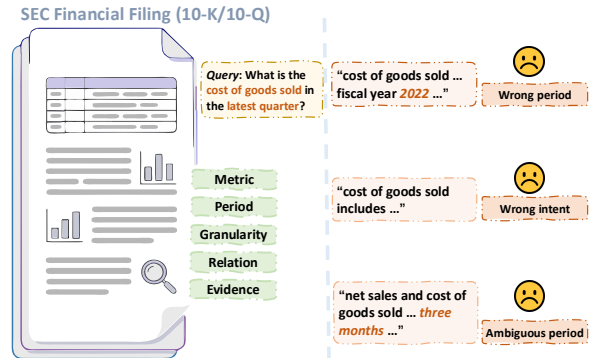


Figure 1: **Key challenge in financial QA.** Reranking must satisfy the correct metric and fiscal period (often numeric), not just semantic relevance. The illustration is based on U.S. SEC corporate financial filings and shows typical failure modes.

A common agent-style approach is to feed large batches of text chunks into a large language model (LLM) and ask it to rank or select evidence. In long corporate filings, this approach breaks down for two practical reasons. First, *scale*: multi-hundred-page reports quickly exceed LLM context budgets, and expanding the input window leads to prohibitive token and latency costs. Second, *opacity*: monolithic, prompt-driven rankings provide little insight into why a particular passage is selected, making the decisions difficult to inspect or audit, which is unacceptable in regulated financial analysis. This leads to systematic errors such as selecting evidence from the wrong fiscal period, misaligning the metric intent of the query, or returning temporally ambiguous passages (Sun et al., 2023; Ma et al., 2023; Choi et al., 2025b). As a result, generic LLM rerankers tend to optimize surface-level semantic relevance rather than explicit constraint satisfaction, making their decisions brittle in financial settings.

Unlike generic LLM rerankers that focus on semantic relevance, FINCARDS aligns evidence using explicit financial fields such as metrics and

066 periods, producing auditable decision traces. Our
067 design is motivated by how financial analysts reason
068 over long filings: evidence is not evaluated all
069 at once, but progressively narrowed, ordered, and
070 adjudicated under explicit criteria. A junior analyst
071 first screens likely evidence, a senior analyst im-
072 poses a coherent global ordering, and a committee
073 resolves close calls when distinctions are subtle.

074 We operationalize this workflow as a machine-
075 usable alignment contract. Document chunks are
076 abstracted into structured *cards* under a shared
077 schema, questions are mapped to explicit in-
078 tent specifications over the same fields, and
079 a lightweight tournament-style review performs
080 screening, global ordering, and adjudication to pro-
081 duce a ranked set of evidence passages.

082 This staged formulation improves numeric and
083 temporal grounding, supports auditability, and
084 keeps computation within predictable cost budgets.
085 Section 3 details the implementation of each stage.

086 Our focus is the *intra-document* ranking setting:
087 given a single, pre-selected filing, surface the most
088 relevant chunks for a question. This isolates the
089 core retrieval bottleneck in financial QA: locating
090 grounded evidence within long documents before
091 generation (Zhu et al., 2021; Chen et al., 2021).

092 We make the following contributions:

- 093 • We reformulate financial question answer-
094 ing over long corporate filings as an *intra-*
095 *document evidence reranking* problem under
096 strict numeric, temporal, and entity con-
097 straints, shifting the modeling focus away
098 from monolithic long-context reasoning.
- 099 • We propose FINCARDS, a structured repre-
100 sentation that abstracts document chunks into
101 auditable evidence cards and maps questions
102 into explicit intent specifications, enabling de-
103 terministic and interpretable alignment.
- 104 • We introduce a zero-shot, tournament-style
105 reranking pipeline that produces stable ranked
106 evidence sets under finance-aware criteria,
107 and demonstrate consistent improvements in
108 early precision over strong baselines.

109 2 Related Work

110 **Financial QA benchmarks and evidence struc-**
111 **ture.** Early financial QA emphasized numerical
112 reasoning grounded in text and tables. Chen et al.
113 (2021) introduced FinQA with program annota-
114 tions to make multi-step arithmetic explicit, while

115 Chen et al. (2022) extended this to conversational
116 settings with chained reasoning. Zhu et al. (2021)
117 (TAT-QA) highlighted hybrid text-table reasoning
118 drawn from real reports, making clear that answer
119 correctness depends on metric, period, and entity
120 alignment rather than surface similarity. Subse-
121 quent resources (e.g., FinanceBench) found that
122 even strong LLMs misfire in realistic enterprise-
123 style questions, underscoring retrieval bottlenecks
124 and hallucination risks (Islam et al., 2023). More
125 recently, FinDER (Choi et al., 2025a) stressed
126 retrieval-augmented generation in finance with ex-
127 pert triplets that reflect terse practitioner queries,
128 and FinAgentBench (Choi et al., 2025b) isolated
129 *chunk-level* ranking within a selected document,
130 reporting that the paragraph-level setting remains
131 challenging for LLMs and that reinforcement fine-
132 tuning of small models can improve MRR and
133 nDCG. The present work focuses on this intra-
134 document chunk-ranking setting, which remains
135 challenging for LLMs in zero-shot scenarios.

Retrieval and LLM reranking. Classical IR re-
136 lies on lexical methods such as BM25, while mod-
137 ern pipelines incorporate dense retrievers and cross-
138 encoders (Xiong et al., 2020). Recently, LLMs
139 have been used directly as zero-shot rerankers: list-
140 wise approaches like RankGPT and LRL show
141 that instruction-tuned LLMs can reorder candidates
142 competitively without task-specific training (Sun
143 et al., 2023; Ma et al., 2023). However, listwise
144 prompts can be input-order sensitive and context-
145 length constrained. Pairwise prompting (*A vs. B?*)
146 improves calibration and stability (Qin et al., 2024),
147 while setwise/tournament strategies mitigate order
148 sensitivity and scale better with long lists (Zhuang
149 et al., 2024; Chen et al., 2025). Simple rank fusion
150 such as reciprocal rank fusion (RRF) remains a
151 strong baseline to aggregate noisy rankings (Cor-
152 mack et al., 2009). Recent work suggests that list-
153 wise and pairwise comparisons play complemen-
154 tary roles in robust reranking, from global order-
155 ing to resolving close decisions. Our work builds
156 on these insights and adapts them to the financial
157 domain by structuring such comparisons within a
158 constraint-driven review process that enforces ex-
159 plicit numeric and temporal alignment. 160

Agentic reasoning and human review processes.
161 A parallel thread models how humans search and
162 reason. ReAct interleaves reasoning and actions
163 (Yao et al., 2023b), Tree-of-Thoughts explores mul-
164 tiple solution paths (Yao et al., 2023a), Reflexion
165

adds self-critique with episodic memory (Shinn et al., 2023), and self-consistency improves chain-of-thought reliability (Wang et al., 2023). Work on LLM-as-a-judge surveys reliability and variance when models evaluate content, which is directly relevant to reranking and motivates explicit numeric and temporal guardrails. In finance, documents are visually and numerically dense; layout-aware models (e.g., LayoutLMv3) benefit document QA (Huang et al., 2022), while chart/table QA benchmarks emphasize cross-modal reasoning (Masry et al., 2022; Zhu et al., 2021). Surveys of RAG for LLMs show that most failures are caused by retrieval quality rather than model capacity, highlighting the importance of domain-aligned retrieval (Gao et al., 2024). In this context, our work draws inspiration from agentic and human-aligned reasoning by structuring reranking as a staged review process that emphasizes grounding and explicit constraints, without requiring any model updates.

3 FINCARDS

3.1 Overview

We study *intra-document evidence reranking* for financial question answering: given a user question and all chunks from a single long SEC filing (e.g., 10-K/10-Q), the goal is to rank chunks so that the top- k results satisfy the required financial conditions (metric, fiscal period, entity, and often explicit numbers). Figure 2 summarizes our approach, FINCARDS, which decomposes this problem into three components. **(1) Card abstraction** converts each chunk into a compact, structured *card* that records finance-relevant fields (entities, metrics, periods, numbers, and section cues) for auditable matching. **(2) Query intent mapping** converts the question into a structured intent that specifies the demanded entities/metrics/periods and whether numeric evidence is required. **(3) Tournament reranking** performs staged, zero-shot reranking over cards, combining a screening step, a global listwise ordering step, and a targeted adjudication step, followed by lightweight fusion and post-hoc alignment to produce the final top- k evidence list.

In addition to the final ranking, the pipeline produces an explicit *audit trace* for each selected chunk, recording which Card fields were matched, how the chunk was retained or filtered at each stage, and how its final rank was determined.

3.2 Card Abstraction

Why Card Abstraction. The *card abstraction* module makes financial evidence explicit and auditable before any ranking decisions are made. It converts raw text chunks from long SEC filings into compact, structured records that explicitly encode entities, financial metrics, fiscal periods, and verbatim numeric spans. By operating on these schema fields rather than raw text, downstream stages compare candidates through field-level matching instead of free-form semantic similarity.

By operating on schema fields rather than unconstrained text, downstream stages can perform reliable comparisons under strict numeric and temporal constraints.

Motivation. This design directly addresses recurring failure modes in financial QA over long filings. Lexical retrieval is particularly vulnerable to numeric drift, temporal misalignment, and boilerplate repetition, where legally mandated disclosures dominate surface signals without conveying substantive evidence. Card abstraction mitigates these issues by enforcing explicit temporal normalization, verbatim numeric copying, and boilerplate awareness at the representation level.

Analyst analogy. Conceptually, card abstraction mirrors the preparatory work of a junior financial analyst: relevant numbers, periods, and contextual cues are recorded explicitly so that subsequent reviewers can assess relevance without repeatedly reinterpreting raw text. The resulting Card corpus defines the evidence space for the tournament-style reranking stages described in Section 3.4.

Formal definition. Let $\mathcal{X} = \{x_1, \dots, x_L\}$ denote the text chunks from a single SEC filing. Each chunk x_i is mapped to a structured *Chunk Card* via a schema-constrained extraction function with deterministic decoding:

$$f : \mathcal{X} \rightarrow \mathcal{C}, \quad c_i = f(x_i) = \text{ChunkCard}(x_i). \quad (1)$$

Each Chunk Card $c_i \in \mathcal{C}$ is a typed record with a full schema, from which we derive a compact *alignment core* for intent matching, while remaining auxiliary fields are used only for screening and stability control:

$$c_i = (c_i^{\text{core}}, c_i^{\text{aux}}), \quad c_i^{\text{core}} = (T_i, E_i, M_i, N_i, P_i, S_i, D_i, \Xi_i). \quad (2)$$

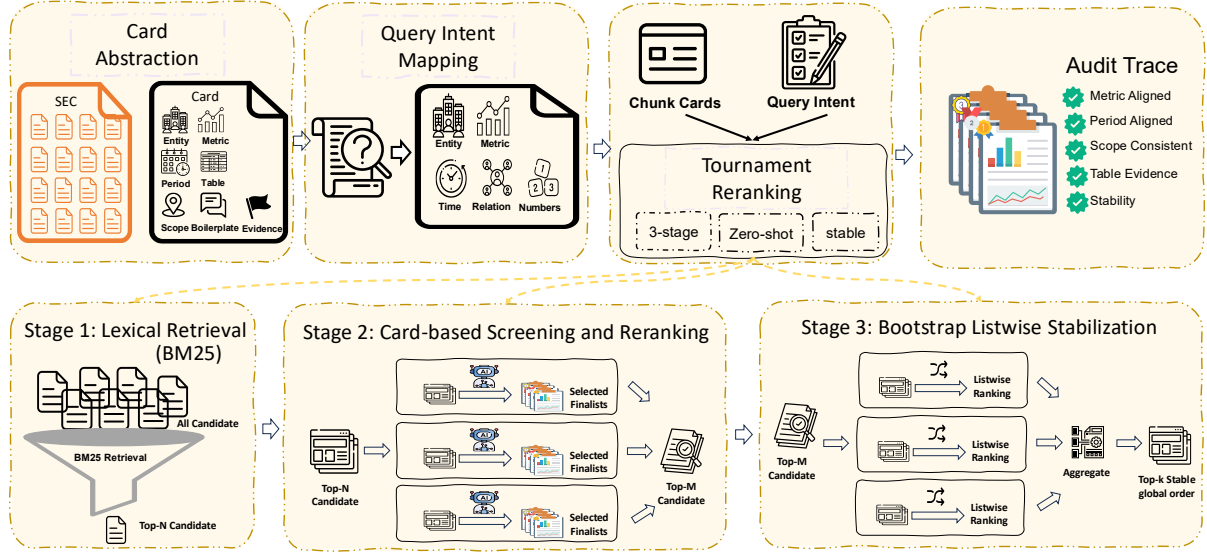


Figure 2: **Overview of the FINCARDS pipeline.** From an SEC filing and a user question, the system constructs structured Cards, a structured query intent, and a tournament reranking module that produces the final Top- k evidence chunks.

The core schema c_i^{core} is an *alignment core (projection)* of the full Chunk Card, used exclusively for intent alignment. Auxiliary fields c_i^{aux} are never used for alignment and are employed only for screening and stability control.

Core fields. $T_i \in \mathcal{T}$ is a topic label; $E_i \subseteq \mathcal{E}$ the set of explicitly mentioned entities; $M_i \subseteq \mathcal{M}$ financial metrics; N_i a (possibly empty) multiset of verbatim numeric spans; $P_i \in \mathcal{P}$ a normalized fiscal period or interval; $S_i \in \mathcal{S}$ the section identifier; $D_i \in \mathcal{D}$ a derived entity–metric–period triple; and Ξ_i evidence spans linking all fields to the original text (*audit trace*).

Auxiliary fields. c_i^{aux} includes derived cues such as scope descriptors, table signatures, and a boilerplate flag. These fields are derived via lightweight rules or parsers over x_i and are used only to guide Stage 2/3 screening and stability control.

3.3 Query Intent Mapping

On the query side, we map each natural-language question to a structured *intent* representation that explicitly encodes entities, metrics, temporal constraints, and numeric requirements. This representation enables direct field-level matching against Card schemas, rather than relying on unconstrained text similarity.

Financial questions are often underspecified in surface form. For example, the query “How did revenue change last quarter?” implicitly requires

numeric evidence, a specific fiscal period, and a comparison relation. Intent mapping resolves this ambiguity by decomposing questions into dimensions that can be directly aligned with Card fields.

Formally, each question q_j is mapped to an intent object

$$\text{Intent}(q_j) = (T_j, E_j, M_j, R_j, \Theta_j, \nu_j, K_j), \quad (3)$$

where T_j is a topic label, E_j entities, M_j metrics, R_j the relation type (e.g., comparison or trend), Θ_j temporal constraints, ν_j whether explicit numeric evidence is required, and K_j lexical keywords.

Together, Card abstraction and intent mapping expose a shared schema for field-level alignment in the tournament reranking stages (Section 3.4).

3.4 Tournament Reranking

We now describe the tournament-style reranking module used to select evidence within a single filing. The pipeline assumes (i) a Card corpus \mathcal{C} derived from filing chunks (Section 3.2), and (ii) a structured query intent $\text{Intent}(q)$ extracted from the question (Section 3.3).

Ranking proceeds in three stages: recall-oriented candidate generation, Card-based semantic filtering, and stability-aware listwise aggregation. This staged design reflects how financial analysts progressively narrow, order, and adjudicate evidence under strict numeric and temporal constraints.

3.4.1 Preliminaries

Given a question q and a filing, let $\mathcal{X} = \{x_1, \dots, x_L\}$ denote the set of text chunks and $\mathcal{C} = \{c_1, \dots, c_L\}$ the corresponding Chunk Cards. We extract a structured intent representation $\text{Intent}(q)$ as described in Section 3.3.

The goal of tournament reranking is to return an ordered list of chunk indices $\pi = (\pi_1, \dots, \pi_k)$ corresponding to the top- k evidence chunks in the filing, optimized for early-rank relevance.

3.4.2 Stage 1: Lexical Retrieval (BM25)

Stage 1 constructs a high-recall candidate set via lexical retrieval within the same filing. Concretely, we score all chunks using the BM25 ranking function (Robertson and Zaragoza, 2009) and keep the top- N candidates:

$$\mathcal{S}_1(q) = \text{TopN}(\text{BM25}(q, x_i)). \quad (4)$$

To make the candidate budget comparable across filings of different lengths, we use a length-adaptive cutoff

$$N = \text{clamp}(\lceil rL \rceil, N_{\min}, N_{\max}), \quad (5)$$

with $r = 0.5$, $N_{\min} = 60$, and $N_{\max} = 150$ (and $N = L$ if $L < N_{\min}$). Stage 1 serves as a recall-oriented starting point for downstream reranking and highlights cases where relevant evidence receives low lexical scores.

3.4.3 Stage 2: Card-based Screening and Reranking

Stage 2 reduces the Stage 1 candidate set using structured Card representations, without accessing raw chunk text. Starting from $\mathcal{S}_1(q)$, we partition candidates into groups via round-robin assignment so that each group contains a mix of high-, mid-, and low-ranked Stage 1 candidates. Let $\{\mathcal{G}_1, \dots, \mathcal{G}_m\}$ denote groups of size approximately g (default $g = 25$).

For each group \mathcal{G}_t , an LLM agent selects a small set of finalists:

$$\mathcal{F}_t = \text{SELECT}(\mathcal{G}_t, q), |\mathcal{F}_t| \in [k_{\min}, k_{\max}]. \quad (6)$$

The selection rubric emphasizes alignment between the query intent and Card fields, including (i) metric overlap, (ii) temporal compatibility, (iii) scope consistency (company-wide vs. segment/region/product), (iv) appropriate evidence type (e.g., table-centric cards for quantitative queries),

and (v) down-weighting boilerplate content unless it is uniquely relevant.

To avoid over-filtering, we enforce lightweight coverage constraints. For instance, trend queries must retain at least one temporally grounded, table-bearing card, while definition or policy queries must retain at least one explanatory card.

Group-level selections are merged by union and deduplication, retaining the highest relevance score when a candidate appears multiple times:

$$\mathcal{S}_2(q) = \text{Dedup}\left(\bigcup_{t=1}^m \mathcal{F}_t\right). \quad (7)$$

Ties are optionally broken using the original Stage 1 score. If gold retention relative to Stage 1 drops below a threshold, we reshuffle group assignments and repeat selection for a small number of retries.

3.4.4 Stage 3: Bootstrap Listwise Stabilization

While Stage 2 substantially reduces the candidate set using structured Card cues, the resulting ranking can still be unstable due to the sensitivity of single-pass LLM judgments to grouping and comparison context. Stage 3 addresses this issue by enforcing *ranking stability* through multi-round bootstrap listwise aggregation.

Given the Stage 2 candidate set $\mathcal{S}_2(q)$ of size M , we select a group size $g \in [15, 25]$ (adapted to M) and perform up to $R_{\max} = 5$ bootstrap rounds. In each round r , candidates are randomly shuffled with a fixed seed and partitioned into groups $\{\mathcal{H}_{r,1}, \dots, \mathcal{H}_{r,p_r}\}$. For each group, an LLM produces a complete listwise ranking. Random re-grouping exposes each candidate to multiple comparison contexts, mitigating bias introduced by any single partition.

To aggregate rankings across groups and rounds, we employ *normalized Borda scores*. This choice is motivated by two considerations: (i) Borda aggregation preserves fine-grained relative ordering information, rather than relying on hard selection or voting, and (ii) normalization ensures comparability across groups of different sizes, which naturally arise under bootstrap partitioning. For a group \mathcal{H} of size $|\mathcal{H}|$, an item ranked at position ρ receives a score

$$s(\rho; \mathcal{H}) = \frac{|\mathcal{H}| - \rho}{|\mathcal{H}| - 1} \in [0, 1]. \quad (8)$$

Scores are accumulated over all appearances of each candidate:

$$S(i | q) = \sum_r \sum_t s(\rho_{r,t}(i); \mathcal{H}_{r,t}). \quad (9)$$

407 Sorting candidates by $S(i | q)$ yields a stable global
408 ranking π , from which the top- k evidence chunks
409 are returned.

410 To control computational cost, we monitor con-
411 vergence of the current top- k set. If the Jaccard
412 similarity between top- k results from consecutive
413 rounds exceeds a threshold (0.9), the procedure
414 terminates early.

415 3.4.5 Reproducibility and Cost

416 All LLM interactions use deterministic decoding
417 (temperature = 0) and strict JSON schema val-
418 idation to ensure reproducibility. In addition to
419 the final Top- k ranking, the pipeline outputs an
420 explicit *audit trace* for each selected chunk, record-
421 ing stage-wise candidate lists, grouping decisions,
422 and per-round ranks in Stage 3. This trace exposes
423 which structured Card fields (e.g., metric, period,
424 scope, and table cues) were matched, how each
425 chunk was retained or filtered across stages, and
426 how its final rank was determined, enabling trans-
427 parent inspection and controlled multi-model com-
428 parisons under identical algorithms and prompts.

429 In terms of cost, let L be the number of chunks
430 in a filing, N the Stage 1 candidate cutoff, and
431 $M = |\mathcal{S}_2(q)|$ the Stage 2 candidate size. Stage 1
432 scores all chunks using BM25 in $O(L)$ time per
433 query. Stage 2 requires $m = \lceil N/g \rceil$ LLM calls
434 (one per group), and Stage 3 performs at most
435 R_{\max} bootstrap rounds with $\lceil M/g \rceil$ listwise calls
436 per round. In practice, the total number of LLM
437 calls is often substantially reduced by early stop-
438 ping once the Top- k ranking stabilizes.

439 4 Experiments

440 4.1 Experimental Setup

441 We evaluate our multi-stage intra-document re-
442 trieval and ranking system on the **FinAgentBench**
443 (Choi et al., 2025b) benchmark, which consists of
444 financial question answering tasks derived from
445 U.S. SEC filings (10-K and 10-Q). For each query,
446 the system is provided with a *single financial docu-*
447 *ment* and must identify and rank the most relevant
448 evidence chunks within that document. This setting
449 isolates the challenge of *intra-document retrieval*,
450 where relevant evidence is often sparse, temporally
451 constrained, and interleaved with boilerplate dis-
452 closures.

453 All experiments are conducted in Python 3.x
454 with large language models accessed via API calls.
455 We use deterministic decoding (temperature = 0)

and structured outputs throughout all experiments.

456 4.2 System Variants

457 We compare traditional lexical retrieval, zero-shot
458 LLM-based reranking, and several variants of
459 our multi-stage pipeline. All systems share the
460 same document chunking and evaluation proto-
461 col. Throughout this paper, **Stage 1** refers to
462 BM25-based lexical retrieval. For fairness, the
463 zero-shot LLM reranking baseline operates on the
464 same Stage 1 candidate set, with identical candi-
465 date budgets.

466 **Stage 1.** The lexical retrieval baseline that returns
467 the top-ranked chunks from Stage 1 without any
468 semantic reranking.

469 **Zero-shot LLM Reranking.** A standard LLM-
470 based reranking baseline that applies a single-pass,
471 zero-shot LLM to reorder the *Stage 1 candidate*
472 *set* using raw chunk text. Unlike our approach,
473 this baseline does not leverage structured Card rep-
474 resentations, candidate grouping, or multi-round
475 stabilization.

476 **Stage 1 + Stage 3.** An ablated variant that ap-
477 plies the Stage 3 bootstrap listwise ranking directly
478 on the Stage 1 candidate set, without Card-based
479 filtering. This setting isolates the effect of stability-
480 aware ranking independent of structured semantic
481 screening.

482 **Stage 1 + Stage 2.** A two-stage variant that ap-
483 plies Card-based semantic screening and reranking
484 on top of Stage 1 retrieval, but does not include the
485 bootstrap-based stability mechanism of Stage 3.

486 **Stage 1 + Stage 2 + Stage 3 (Full Pipeline).** Our
487 full system, which sequentially combines Stage 1
488 lexical retrieval, Stage 2 Card-based semantic fil-
489 tering, and Stage 3 bootstrap listwise stabilization.
490 This progressive design compresses the candidate
491 set while jointly improving ranking accuracy and
492 stability.

493 4.3 Evaluation Measures

494 We adopt standard information retrieval metrics at
495 rank 10. We report metrics at $k = 10$ because each
496 question in FinAgentBench is typically associated
497 with a small set of relevant evidence chunks (on
498 the order of ten), making early-rank quality the
499 primary evaluation focus. **nDCG@10** measures
500 graded relevance with emphasis on early ranks,
501 **MAP@10** captures precision across the top-ranked
502

System	nDCG@10	MAP@10	MRR@10	Cand. Size
<i>Traditional Baseline</i>				
Stage 1	44.26	60.01	69.88	25
<i>LLM-based Baseline</i>				
Zero-shot LLM Reranking	55.80	66.50	78.20	100 → 25
<i>Ablation of Our Method</i>				
Stage 1+Stage 3	58.23	68.42	77.56	100 → 25
Stage 1+Stage 2	63.66	73.09	82.95	100 → 40
<i>Our Full Pipeline</i>				
Stage 1+Stage 2+Stage 3	71.58	78.69	89.17	100 ⁴⁰ → 25

Table 1: **Main results on FinAgentBench under intra-document retrieval.** All retrieval metrics (nDCG@10, MAP@10, MRR@10) are reported as percentages. The proposed three-stage pipeline substantially improves early-rank accuracy while progressively reducing the candidate set size (Cand. Size) at each stage.

503 results, and **MRR@10** reflects how quickly the
504 first relevant evidence chunk appears. All mea-
505 sures are reported as averages over all evaluation
506 queries with scores multiplied by 100 and reported
507 as percentages.

508 4.4 Main Results

509 Table 1 summarizes the main results on FinAgent-
510 Bench under the intra-document retrieval setting.
511 Our three-stage pipeline achieves a large and con-
512 sistent improvement over both traditional and LLM-
513 based baselines across all metrics.

514 Compared to Stage 1, the full pipeline improves
515 nDCG@10 by over **+27 points** and MRR@10 by
516 nearly **+20 points**. Even relative to a strong zero-
517 shot LLM reranking baseline, our approach yields
518 substantial gains (+15.8 nDCG@10), demonstrat-
519 ing that naive LLM reranking is insufficient for
520 financial evidence selection.

521 Importantly, these accuracy gains are achieved
522 while *progressively reducing the candidate set size*
523 from roughly 100 chunks to fewer than 25, indicat-
524 ing that the proposed design improves both ranking
525 quality and retrieval efficiency.

526 4.5 Analysis

527 The results in Table 1 provide clear evidence that
528 each stage of the proposed pipeline contributes
529 meaningfully to retrieval effectiveness.

530 First, zero-shot LLM reranking improves over
531 BM25 but remains limited, highlighting that re-
532 placing lexical scores with unstructured LLM judg-
533 ments does not adequately resolve temporal mis-
534 match, scope ambiguity, or boilerplate interference
535 in financial filings.

Variant	nDCG@10	MAP@10	MRR@10
Full Card (Baseline)	63.66	73.09	82.95
w/o temporal_data	59.80	69.20	78.50
w/o financial_metrics	61.20	70.85	80.20
w/o tables	58.50	67.80	77.20
w/o scope	62.15	71.80	81.50
Only summary	54.20	63.50	72.80
Raw chunks (no Card)	49.50	60.80	70.15

Table 2: **Stage 2 ablation results.** Each variant removes one component from the Card representation.

536 Second, introducing Card-based filtering in
537 Stage 2 yields a large performance jump (+7.9
538 nDCG@10 over zero-shot reranking), confirming
539 that *structured intermediate representations are*
540 *crucial* for aligning query intent with financial evi-
541 dence.

542 Finally, Stage 3 further improves early-rank met-
543 rics by stabilizing rankings across multiple com-
544 parison contexts. This demonstrates that *ranking*
545 *stability*, rather than additional semantic filtering
546 alone, is essential for reliable early precision in
547 long, noisy financial documents.

548 Overall, the analysis validates the core design
549 principles of the proposed system: structured rea-
550 soning, progressive filtering, and stability-aware
551 aggregation, all achieved without task-specific fine-
552 tuning.

553 4.6 Ablation Study

554 4.6.1 Stage 2: Card Component Ablation

555 We ablate individual components of the Card rep-
556 resentation used in Stage 2, while keeping the
557 pipeline and model fixed.

558 Table 2 shows that structured Card fields are es-
559 sential for effective semantic filtering. Removing
560 temporal information causes the largest degrada-
561 tion, highlighting the importance of temporal align-
562 ment in financial QA. Ablating table indicators or
563 financial metrics also leads to substantial drops,
564 indicating that identifying quantitative evidence is
565 critical even without exposing raw numbers. Us-
566 ing only free-text summaries performs poorly, and
567 operating directly on raw chunks yields the worst
568 results. Overall, the Card abstraction is a neces-
569 sary intermediate representation rather than a mere
570 efficiency optimization.

571 4.6.2 Stage 3: Ranking Strategy Ablation

572 We further analyze Stage 3 by ablating key de-
573 sign choices in the bootstrap-based listwise ranking

Variant	nDCG@10	MAP@10	MRR@10	Rank Var.
Bootstrap (R=3-5)	71.58	78.69	89.17	0.0342
Single Round (R=1)	68.20	75.20	85.80	0.0856
Fixed Grouping	69.15	76.35	86.95	0.0621
Mean Rank Aggregation	69.80	77.05	87.50	0.0498
Voting Aggregation	67.35	74.80	85.10	0.0723
No Early Stopping	71.65	78.75	89.25	0.0318

Table 3: **Stage 3 ablation results.** Rank variance measures ranking stability across bootstrap rounds.

Stage / Model	nDCG@10	MAP@10	MRR@10
Stage 1 (BM25)	44.26	60.01	69.88
<i>Stage 2: Card-based Filtering</i>			
GPT-5 Mini	63.66	73.09	82.95
GPT-4 Mini	64.63	75.63	84.66
Claude-4.5-Opus	70.89	81.03	89.67
Claude-4.5-Sonnet	67.15	77.16	85.71
Gemini 2.5 Flash	63.28	73.32	82.46
Gemini 3 Pro (Preview)	67.06	76.53	85.98
<i>Stage 3: Bootstrap Stable Ranking</i>			
GPT-5 Mini	71.58	78.69	89.17
GPT-4 Mini	71.15	77.77	87.98
Claude-4.5-Opus	75.72	81.20	89.94
Claude-4.5-Sonnet	74.42	79.88	90.46
Gemini 2.5 Flash	74.87	79.46	89.96
Gemini 3 Pro (Preview)	76.52	81.39	91.76

Table 4: **Robustness across LLM backbones on FinAgentBench.**

procedure. In addition to ranking quality, we report *rank variance* as a stability metric, computed across bootstrap rounds.

Table 3 shows that bootstrap-based aggregation is crucial for both accuracy and stability. Single-round ranking exhibits much higher variance, confirming the instability of one-shot LLM judgments. Random regrouping consistently outperforms fixed grouping, indicating that exposure to diverse comparison contexts reduces bias. Normalized Borda aggregation outperforms simpler aggregation schemes, suggesting that fine-grained relative ordering is important. Disabling early stopping yields only marginal gains while increasing computation, as rankings typically converge within 3-4 rounds. Overall, Stage 3 acts as a stability control layer that improves early precision while controlling variance.

4.7 Robustness Across LLM Backbones

To test model dependence, we run Stage 2 and Stage 3 with six different LLM backbones. All other components of the pipeline remain unchanged.

Analysis. Table 4 shows that the proposed pipeline is robust across diverse LLM backbones.

First, all models exhibit a substantial improvement from Stage 1 to Stage 2, indicating that Card-based semantic filtering consistently improves evidence selection regardless of model capacity. This suggests that the gains are primarily driven by the structured intermediate representation rather than backbone-specific reasoning ability.

Second, Stage 3 further improves or stabilizes ranking quality for every model, with consistent gains in nDCG@10 and MRR@10. This confirms that bootstrap-based listwise aggregation effectively mitigates the variance of single-pass LLM rankings across architectures.

Finally, while stronger backbones achieve higher absolute scores, the relative improvements introduced by Stage 2 and Stage 3 remain similar across model families. This demonstrates that the proposed framework is model-agnostic and complements advances in base LLM capability, rather than relying on them.

5 Conclusions and Future Work

We introduced FINCARDS, a tournament-style, zero-shot intra-document reranking framework for financial QA over long SEC filings. The key idea is to replace monolithic relevance ranking with structured evidence selection, combining Card-based abstractions and a staged reranking protocol that enforces metric, temporal, and scope constraints.

Across extensive experiments on FinAgentBench, FINCARDS consistently outperforms both lexical baselines and strong zero-shot LLM rerankers, while progressively reducing the candidate set size. These gains are achieved without task-specific fine-tuning, demonstrating that reliability in financial QA can be substantially improved through structured intermediate representations and stability-aware decision procedures.

More broadly, our results indicate that reliable zero-shot reranking in financial QA depends on explicit intermediate structure rather than single-pass relevance judgments.

Future work will explore adaptive budget control across stages, cross-document evidence selection, and structured Card interfaces for downstream answer generation.

644 Limitations

645 Despite its strong empirical performance, the pro-
646 posed framework has several limitations.

647 First, the multi-stage pipeline incurs non-trivial
648 computational cost. In particular, the listwise and
649 bootstrap-based ranking stages require multiple
650 LLM calls per query, which may limit scalabil-
651 ity in large-scale or latency-sensitive deployments.
652 While candidate compression and early stopping
653 mitigate this cost in practice, efficiency remains an
654 important consideration.

655 Second, the current study focuses exclusively
656 on intra-document retrieval within a single SEC
657 filing. Many real-world financial analysis tasks
658 require reasoning across multiple documents or
659 heterogeneous sources, such as press releases and
660 earnings calls. The effectiveness of the tournament-
661 style design in such multi-document settings has
662 not yet been evaluated.

663 Third, although the approach avoids task-
664 specific fine-tuning, it may still be sensitive to
665 prompt design and schema choices. While we em-
666 ploy strict structured outputs and deterministic de-
667 coding to improve stability, further work is needed
668 to understand robustness under prompt variation
669 and evolving model behaviors.

670 Ethical Considerations

671 This work studies retrieval and reranking methods
672 for financial question answering over publicly avail-
673 able regulatory filings. The proposed framework
674 operates solely on textual disclosures released by
675 companies and does not involve personal data, user
676 profiling, or sensitive individual information.

677 The primary goal of this research is to improve
678 evidence selection and interpretability in financial
679 analysis and decision support. By emphasizing
680 structured intermediate representations and trans-
681 parent ranking procedures, the approach aims to
682 support more responsible use of large language
683 models in high-stakes financial settings.

684 Potential risks include over-reliance on auto-
685 mated systems and misinterpretation of retrieved
686 evidence if used without appropriate human over-
687 sight. Accordingly, the proposed system is in-
688 tended as an assistive tool for analysts, rather than
689 a replacement for professional judgment. Beyond
690 risks commonly associated with automated infor-
691 mation retrieval systems, we do not anticipate sig-
692 nificant negative societal impact.

Data License All experiments in this work are
conducted on publicly available datasets derived
from U.S. SEC filings. The underlying documents,
such as 10-K and 10-Q reports, are released un-
der public disclosure requirements and are freely
accessible for research purposes.

The FinAgentBench benchmark used in our ex-
periments follows the original data collection and
usage terms specified by its authors. This work
does not redistribute the original filings, nor does
it impose additional licensing constraints beyond
those associated with the source datasets.

References

- Yiqun Chen, Qi Liu, Yi Zhang, Weiwei Sun, Xinyu Ma,
Wei Yang, Daiting Shi, Jiabin Mao, and Dawei Yin.
2025. [Tourrank: Utilizing large language models
for documents ranking with a tournament-inspired
strategy](#). *Preprint*, arXiv:2406.11678.
- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena
Shah, Iana Borova, Dylan Langdon, Reema Moussa,
Matt Beane, Ting-Hao Huang, Bryan Routledge, and
William Yang Wang. 2021. [FinQA: A dataset of nu-
merical reasoning over financial data](#). In *Proceedings
of the 2021 Conference on Empirical Methods in Nat-
ural Language Processing*, pages 3697–3711, Online
and Punta Cana, Dominican Republic. Association
for Computational Linguistics.
- Zhiyu Chen, Shiyang Li, Charese Smiley, Zhiqiang Ma,
Sameena Shah, and William Yang Wang. 2022. [Con-
vFinQA: Exploring the chain of numerical reasoning
in conversational finance question answering](#). In *Pro-
ceedings of the 2022 Conference on Empirical Meth-
ods in Natural Language Processing*, pages 6279–
6292, Abu Dhabi, United Arab Emirates. Association
for Computational Linguistics.
- Chanyeol Choi, Jihoon Kwon, Jaeseon Ha, Hojun Choi,
Chaewoon Kim, Yongjae Lee, Jy-yong Sohn, and
Alejandro Lopez-Lira. 2025a. [FinDER: Financial
Dataset for Question Answering and Evaluating
Retrieval-Augmented Generation](#), page 638–646. As-
sociation for Computing Machinery, New York, NY,
USA.
- Chanyeol Choi, Jihoon Kwon, Alejandro Lopez-Lira,
Chaewoon Kim, Minjae Kim, Juneha Hwang, Jae-
seon Ha, Hojun Choi, Suyeol Yun, Yongjin Kim, and
Yongjae Lee. 2025b. [FinAgentBench: A Benchmark
Dataset for Agentic Retrieval in Financial Question
Answering](#), page 632–637. Association for Comput-
ing Machinery, New York, NY, USA.
- Gordon V. Cormack, Charles L A Clarke, and Stefan
Buettcher. 2009. [Reciprocal rank fusion outperforms
condorcet and individual rank learning methods](#). In
*Proceedings of the 32nd International ACM SIGIR
Conference on Research and Development in Infor-
mation Retrieval, SIGIR '09*, page 758–759, New

748	York, NY, USA. Association for Computing Machinery.	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models . In <i>The Eleventh International Conference on Learning Representations</i> .	803
749			804
750	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. Retrieval-augmented generation for large language models: A survey . <i>Preprint</i> , arXiv:2312.10997.		805
751			806
752			807
753			808
754			
755	Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. Layoutlmv3: Pre-training for document ai with unified text and image masking . In <i>Proceedings of the 30th ACM International Conference on Multimedia</i> , MM '22, page 4083–4091, New York, NY, USA. Association for Computing Machinery.	Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. Approximate nearest neighbor negative contrastive learning for dense text retrieval . <i>Preprint</i> , arXiv:2007.00808.	809
756			810
757			811
758			812
759			813
760		Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023a. Tree of thoughts: Deliberate problem solving with large language models . <i>Preprint</i> , arXiv:2305.10601.	814
761			815
762	Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie Vidgen. 2023. Financebench: A new benchmark for financial question answering . <i>Preprint</i> , arXiv:2311.11944. ArXiv preprint.	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023b. React: Synergizing reasoning and acting in language models . <i>Preprint</i> , arXiv:2210.03629.	816
763			817
764			818
765			
766			
767	Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. 2023. Zero-shot listwise document reranking with a large language model . <i>Preprint</i> , arXiv:2305.02156.	Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. 2021. TAT-QA: A question answering benchmark on a hybrid of tabular and textual content in finance . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3277–3287, Online. Association for Computational Linguistics.	823
768			824
769			825
770			826
771	Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. ChartQA: A benchmark for question answering about charts with visual and logical reasoning . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2263–2279, Dublin, Ireland. Association for Computational Linguistics.		827
772			828
773			829
774			830
775			831
776			832
777		Shengyao Zhuang, Honglei Zhuang, Bevan Koopman, and Guido Zuccon. 2024. A setwise approach for effective and highly efficient zero-shot ranking with large language models . In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '24, page 38–47, New York, NY, USA. Association for Computing Machinery.	833
778	Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, and Michael Bendersky. 2024. Large language models are effective text rankers with pairwise ranking prompting . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 1504–1518, Mexico City, Mexico. Association for Computational Linguistics.		834
779			835
780			836
781			837
782			838
783			839
784			840
785			
786			
787	Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: BM25 and beyond. <i>Foundations and Trends in Information Retrieval</i> , 3(4):333–389.		
788			
789			
790			
791	Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning . <i>Preprint</i> , arXiv:2303.11366.		
792			
793			
794			
795	Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. 2023. Is ChatGPT good at search? investigating large language models as re-ranking agents . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 14918–14937, Singapore. Association for Computational Linguistics.		
796			
797			
798			
799			
800			
801			
802			

A Candidate Set and Evaluation Fairness

All reranking-based systems operate on the same Stage 1 (BM25) candidate pool. Specifically, for each query, we retrieve the top- N chunks using BM25 ($N \in [60, 150]$ depending on document length). Zero-shot LLM reranking, Stage 2 filtering, and Stage 3 ranking all receive identical candidate sets at their respective entry points. This design ensures that all reported improvements arise from reranking quality rather than candidate recall advantages or budget discrepancies. In particular, no method is allowed to retrieve additional candidates beyond the Stage 1 pool.

B Prompts and Templates

This appendix reports the prompts used to instantiate the structured representations in FINCARDS. All prompts enforce strict JSON-only input and output schemas to ensure deterministic behavior and reproducibility. We stress that FINCARDS does not rely on prompt wording or prompt-specific heuristics; instead, the prompts serve solely to instantiate fixed interfaces defined by the method.

B.1 Card Abstraction Prompt

The card abstraction prompt converts each document chunk into a structured *Chunk Card*. The card explicitly encodes the chunk’s evidence role, temporal anchoring, scope, and verifiability signals, and serves as the primary alignment interface for downstream reranking stages. An excerpt of the prompt used for this abstraction is shown in Figure 3.

B.2 Query Intent Mapping Prompt

The query intent mapping prompt converts a natural-language financial question into a structured *Query Intent*. This representation explicitly encodes the topical focus, requested financial metrics, temporal constraints, and relational form of the question, and defines the demand-side requirements used for alignment with Chunk Cards during reranking. An excerpt of the corresponding prompt is shown in Figure 4.

B.3 Stage 2: Batch Selection Prompt

Stage 2 performs group-wise evidence selection within the candidate pool retrieved by Stage 1. Given a user question and a small group of candidate chunks with their corresponding Chunk

System: Financial document analysis expert.

Task: Given a document chunk extracted from a 10-K filing, generate a structured *Chunk Card* that specifies under what conditions the chunk can serve as evidence for reranking.

Output Format (JSON only):

- Identity: chunk_id, section_path, chunk_index
- claim_role (primary_evidence, supporting_context, definition, caveat, boilerplate, structural_heading)
- evidence_type (table_numeric, narrative_numeric, qualitative_explanation, policy_text, guidance, none)
- answerability_profile (single_fact, comparison, trend, aggregation, attribution)
- temporal_anchor (quality and normalized span)
- scope_signature (entity scope, geography, product)
- measurement_basis (GAAP, non-GAAP, adjusted, reported, unknown)
- verifiability (numeric claims, table presence, comparison cues)
- risk_signals (boilerplate likelihood, limitations)
- semantic_sketch (one-sentence claim summary and topic anchors)

Constraints:

- Output must be valid JSON.
- Structural headings have zero evidence capability.
- Temporal fields use YYYY-MM or null.
- Use unknown if scope or measurement basis is uncertain.

Figure 3: Excerpt of the prompt used to instantiate a **Chunk Card (full schema)**.

Cards, the model selects a bounded number of relevant chunks based solely on card-level information. This stage does not rely on external retrieval scores and serves to filter and structure the candidate set before stability-oriented reranking in Stage 3. An excerpt of the batch selection prompt is shown in Figure 5.

B.4 Stage 3: Listwise Ranking Prompt

Stage 3 performs listwise reranking over the filtered candidate sets produced by Stage 2. Given a user question and a group of candidate chunks represented only by their Chunk Cards, the model produces a complete relative ordering from most relevant to least relevant. This stage explicitly avoids numerical calculation and absolute scoring, and is designed to provide stable ordinal judgments that can be aggregated across multiple rounds. An excerpt of the listwise ranking prompt is shown in

System: Financial QA intent extractor for questions over SEC filings.

Task: Given a user question, produce a structured *Query Intent* that captures the information need of the question.

Output Format (JSON only):

- topic (e.g., Revenue, Costs/Expenses, Profitability, Liquidity, Guidance/Outlook, Risk)
- entities (explicitly mentioned companies or segments, if any)
- metrics (requested financial metrics)
- temporal_scope (type, normalized periods, granularity)
- requires_numeric_evidence (true / false)
- relation (lookup, trend, comparison, explanation, definition, policy)
- keywords (salient lexical cues)

Constraints:

- Output must be valid JSON.
- If uncertain, use conservative defaults (e.g., Other, none, or []).
- Temporal expressions should prefer fiscal normalization (e.g., FY2023, latest quarter).

Figure 4: Excerpt of the prompt used to instantiate Query Intents.

Figure 6.

C Error Analysis and Failure Modes

C.1 Methodology

To complement the aggregate retrieval metrics, we conduct a structured error study to analyze the behavior of the proposed pipeline under different failure conditions. Rather than relying on anecdotal examples, our analysis is grounded in *stage-wise retrieval traces* collected from the full system, enabling verification and reproducibility.

We first select a small but representative set of six queries, covering both successful and failed retrieval scenarios. Specifically, the selected cases include: (i) queries where BM25 fails to retrieve any gold evidence but subsequent stages recover relevant chunks; (ii) queries where Card-based alignment remains ineffective; and (iii) queries where bootstrap-based aggregation introduces performance degradation. This selection strategy ensures coverage of all major pipeline components.

For each selected query, we analyze retrieval outcomes at all three stages. At Stage 1, we inspect the full BM25 ranking over the document and record

gold chunk statistics, including the total number of gold chunks, their absolute BM25 ranks and scores, and whether they appear in the dynamic Top- N candidate set. This allows us to distinguish between recall failures and ranking failures at the lexical retrieval level.

At Stage 2, we examine the final candidate set produced by Card-based filtering and reranking. For top-ranked candidates as well as representative hard negatives, we extract structured Card attributes, including matched financial metrics, temporal information, table presence, and scope alignment. We additionally record the selection rationale generated by the Stage 2 agent, enabling direct attribution of ranking decisions to specific Card fields.

At Stage 3, we analyze the stability of bootstrap-based reranking for cases where the final ranking differs from Stage 2 or is used to demonstrate convergence. We log per-round top- K candidate sets, early stopping behavior, and aggregation statistics such as rank variance, top-5 frequency, and Borda score accumulation. These signals allow us to identify whether performance changes arise from instability across randomized grouping contexts or from systematic evidence reweighting.

Finally, each case is annotated with a small set of failure-type labels (e.g., lexical mismatch, temporal misalignment, schema coverage gap, or implicit reasoning) and a coarse query intent category (quantitative lookup, trend/comparison, or qualitative impact). This taxonomy enables cross-case comparison and clarifies which error patterns are addressed by the proposed design and which remain open challenges.

Together, this error study provides a fine-grained, auditable analysis of the pipeline, explaining not only *whether* the method succeeds or fails, but also *why* these outcomes occur at different stages.

C.2 Representative Cases

Table 5 presents six representative queries used to illustrate typical success and failure patterns across different stages of the pipeline.

C.3 Findings

Our error analysis yields several consistent and instructive findings about the behavior of multi-stage retrieval under financial QA settings.

First, lexical retrieval failures dominate early-stage errors. Across multiple cases, Stage 1 BM25 fails to retrieve any gold evidence within

978 the dynamic Top- N candidate set, resulting in zero
979 nDCG@10. These failures are primarily caused by
980 lexical and semantic mismatches, including abbreviations (e.g., “SG&A” vs. “selling, general and
981 administrative expenses”), paraphrased financial
982 concepts (e.g., “cash from operations” vs. “net cash
983 provided by operating activities”), and implicit temporal
984 constraints (e.g., “latest quarter”). This confirms that term-based retrieval alone is insufficient
985 for financial documents, where equivalent concepts
986 are frequently expressed using heterogeneous terminology.
987

990 **Second, Card-based alignment effectively recovers gold evidence when the query intent is structurally expressible.** For quantitative lookup
991 queries involving explicit financial metrics and
992 time scopes, Stage 2 substantially improves retrieval quality. In these cases, gold chunks are
993 consistently promoted due to matched financial
994 metrics, explicit temporal annotations, and the
995 presence of structured tables. Notably, these improvements occur even when Stage 1 recall is zero,
996 demonstrating that Card-based reasoning can compensate for lexical failures by leveraging schema-level
997 alignment rather than surface text overlap.
998

1000 **Third, Card-based methods degrade gracefully but predictably under implicit or explanatory queries.** For qualitative or impact-oriented
1001 questions (e.g., interest rate implications or strategic decision-making), both Stage 2 and Stage 3
1002 exhibit limited gains and, in some cases, performance degradation. Error traces indicate that such
1003 failures arise not from ranking instability but from schema coverage gaps: the Card representation
1004 lacks fields to encode implicit reasoning, causal effects, or cross-section narrative synthesis. As a
1005 result, the model over-selects superficially related but ultimately non-answering chunks, revealing a
1006 fundamental limitation of schema-driven filtering for abstract reasoning tasks.
1007

1008 **Fourth, bootstrap-based aggregation improves ranking stability but cannot correct systematic misalignment.** Stage 3 bootstrap reranking
1009 consistently reduces rank variance and yields highly stable Top- K sets when Stage 2 candidates
1010 are well-aligned with the query. However, when Stage 2 admits structurally mismatched candidates,
1011 bootstrap aggregation reinforces these errors rather than correcting them. This indicates that Stage 3
1012 primarily serves as a stabilizer rather than a semantic repair mechanism.
1013

1014 Overall, the error study demonstrates that the

1030 proposed pipeline is highly effective when query
1031 intent can be decomposed into explicit schema-aligned constraints (metric, time, scope, and evidence type). Conversely, failures predominantly
1032 arise from intent types that exceed the expressive capacity of the current Card schema, rather than
1033 from ranking noise or stochasticity.
1034
1035
1036

System: Financial document analysis expert.

Task: Given a user question and a group of candidate document chunks, select the most relevant evidence chunks to answer the question. Selection must be based exclusively on the provided Chunk Cards.

Inputs:

- Question text
- A group of candidate chunks with their Chunk Cards
- A required selection range (k_{\min} to k_{\max})
- Optional coverage quotas (hard constraints)

Selection Criteria:

- Metric matching between the question and chunk content
- Temporal alignment with the question requirements
- Scope consistency (entity, segment, or region)
- Content type suitability (table vs. narrative)
- Overall relevance inferred from card summaries and signals

Table Handling Warning: Table-based chunks require additional verification of structure, headers, temporal coverage, and metric relevance. Sequential or continuous tables must be interpreted cautiously and should not be selected solely due to the presence of tabular data.

Output Format (JSON only):

- selected_chunks: an ordered list of between k_{\min} and k_{\max} chunks
- For each chunk: chunk_id, selection_reasons, relevance_score (0–100)

Constraints:

- Selection must satisfy mandatory coverage quotas, if provided.
- Chunks are evaluated purely on card information.
- Results are ordered by descending relevance.

Figure 5: Excerpt of the Stage 2 batch selection prompt used for group-wise evidence filtering.

System: Financial document analysis expert.

Task: Given a user question and a group of candidate evidence chunks, rank all chunks from most relevant to least relevant for answering the question. Ranking must rely exclusively on the provided Chunk Cards.

Inputs:

- Question text
- A group of candidate chunks with their Chunk Cards

Ranking Criteria:

- Metric matching between the question and chunk content
- Temporal alignment with the question requirements
- Scope consistency (entity, segment, or region)
- Content type suitability (table vs. narrative)
- Overall relevance inferred from card summaries and signals

Critical Constraints:

- Do not access the original chunk text.
- Do not perform numerical calculations.
- Do not assign absolute relevance scores.
- Rank *all* chunks in the group using relative ordering only.

Output Format (JSON only):

- ranked_chunks: a complete ordered list of all chunks
- For each chunk: chunk_id, rank (1 = most relevant), and a brief reason

Figure 6: Excerpt of the Stage 3 listwise ranking prompt used for stability-oriented reranking.

Case ID	Query Type	Stage 1	Stage 2	Stage 3	Failure / Success Mode	Key Evidence from Trace
qdfaa5e169d37	Quantitative lookup (SG&A, latest quarter)	Fail	Recover	Stable	Lexical + temporal mismatch fixed by Card	Gold ranked >60 by BM25; Card matches financial_metrics=SG&A and quarterly temporal data
q2a3f6f1492e8	Quantitative lookup (GAAP operating expense)	Fail	Recover	Stable	Abbreviation mismatch fixed by Card	BM25 misses "GAAP OPEX"; Card aligns metric + income-statement table
q83bb50eb29cd	Qualitative risk disclosure (cybersecurity)	Fail	Strong	Degrade	Bootstrap aggregation instability	Stage 2 near-perfect ranking; Stage 3 shows high rank variance across bootstrap rounds
q4652caf2531b	Quantitative + geographic provenance	Fail	Fail	Fail	Schema coverage gap	Gold evidence requires geography; no Card field expresses origin outside U.S.
q7c5dd2e1ba56	Quantitative lookup (revenue, latest quarter)	Fail	Fail	Fail	Temporal aggregation gap	Query requires cross-quarter aggregation; Card lacks temporal trend encoding
q80a13d0df306	Qualitative strategy / impact	Partial	Partial	Fail	Implicit reasoning beyond schema	Relevant evidence dispersed across narrative sections; no localized Card alignment

Table 5: **Error study summary across representative cases.** Stage-level outcomes are annotated as **Recover/Stable**, **Partial/Degrade**, or **Fail**. The table highlights how Card-based alignment resolves lexical and temporal mismatches, while remaining failures concentrate in schema coverage gaps and implicit reasoning queries.