

1. Dataset Description

1.1. Data Source

This project utilizes the EDGAR Corpus (eloukas/edgar-corpus), a comprehensive collection of real SEC (Securities and Exchange Commission) filings from publicly traded companies. The EDGAR database represents authentic, legally filed financial documents that provide detailed information about company operations, risks, and financial performance.

1.2 Dataset Structure

Each financial filing in the corpus contains the following structured sections:

- Company Metadata: Company name, filing type (10-K), filing date
- Item 1: Business description and operations overview
- Item 1A: Risk Factors and forward-looking statements
- Item 7: Management's Discussion and Analysis (MD&A) of financial condition
- Item 7A: Quantitative and Qualitative Disclosures about market risk

1.3 Dataset Size

- Total Available: 20,000+ authentic SEC 10-K filings
- Implementation Scale: Configurable subset (50-500 companies)
- Average Document Length: 15,000-50,000 words per filing

1.4 Data Preprocessing

Text Chunking Strategy:

- Method: Semantic sentence-based splitting with overlap
- Chunk Size: 500 characters (target)
- Overlap: 100 characters between consecutive chunks
- Splitting Pattern: Sentence boundary detection using regex $(?<=[!?.])\s+$
- Rationale: Preserves context across chunk boundaries while maintaining coherent semantic units

Metadata Preservation: Each processed chunk retains: Source company identifier, Filing type and date, Document section (Item 1, 1A, 7, 7A), Source filename/reference

2. Evaluation Metrics

2.1 Response Time (Latency)

Definition:

$$\text{Response_Time} = t_{\text{end}} - t_{\text{start}}$$

where t_{start} is the query submission timestamp and t_{end} is the complete response generation timestamp.

Unit: Seconds (s)

Purpose: Measures computational efficiency and real-world usability. Lower latency is critical for interactive financial analysis applications where analysts require rapid information retrieval.

2.2 Answer Quality (Qualitative Assessment)

Criteria:

1. Factual Accuracy: Response contains correct numerical values and facts from source documents
2. Completeness: Response addresses all aspects of the query
3. Source Attribution: Response cites specific companies and filing sections
4. Calculation Correctness: Mathematical derivations are shown step-by-step when applicable

Scoring: Binary classification (Correct/Incorrect) based on ground truth verification against source documents.

2.3 Answer Length

Definition:

$$\text{Answer_Length} = \text{count}(\text{characters in response})$$

Purpose: Measures verbosity and detail level. Compared across methods to understand trade-offs between conciseness and comprehensiveness.

2.4 Retrieval Success Rate

Definition:

$$\text{Retrieval_Success_Rate} = (\text{Queries with Relevant Context Retrieved} / \text{Total Queries}) \times 100\%$$

Purpose: Measures whether the retrieval component successfully identified and returned chunks containing information needed to answer the query. A failed retrieval results in "Information not available" responses regardless of generation model quality.

3. Methods

3.1 Retrieval-Augmented Generation (RAG)

3.1.1 RAG Conceptual Framework

Core Idea: RAG combines information retrieval with language generation to ground model responses in external knowledge sources. Unlike pure generative models that rely solely on parametric knowledge (learned during training), RAG systems: Retrieve, Augment, Generate.

3.1.2 Developed RAG System Architecture

Class: FinBERTFinancialRAG

Component 1: Embedding Model

Model: ProsusAI/finbert (FinBERT)

- Type: BERT-base fine-tuned on financial corpora (10-K, 10-Q, earnings calls)
- Embedding Dimension: 768
- Device: CUDA (GPU acceleration)
- Advantage: Domain-specific embeddings capture financial terminology (EBITDA, YoY, basis points) more effectively than general-purpose encoders

Component 2: Vector Database

- Primary Implementation: FAISS (Facebook AI Similarity Search)
- Index Type: IndexFlatL2 (exhaustive L2 distance search)
- Device: CPU (for stability in Colab environment)
- Batch Processing: 32 chunks per embedding batch
- Alternative Implementation: ChromaDB
- Benefit: Persistent storage, metadata filtering, production scalability
- Storage: ~/FinancialAI/chromadb
- Collection: "financial_filings"
- In the upcoming iteration, we will try to implement in CHROMA DB or similar

Component 3: Retrieval Mechanism

Baseline: Semantic Search

```
1. Query embedding: q = FinBERT.encode(user_query)
2. Similarity search: scores = FAISS.search(q, top_k=5)
3. Return: Top-5 chunks by L2 distance
```

Advanced: Hybrid Search (Vector + Keyword)

```
Hybrid_Score = α × Vector_Score + (1-α) × Keyword_Score
where α = 0.7 (70% semantic, 30% lexical)
```

Cross-Encoder Advantage: Evaluates query-document interaction jointly (not independently like bi-encoders), capturing nuanced relevance signals.

Component 4: Generation Model

- Model: GPT-3.5-turbo (OpenAI)
- API: Chat Completions API
- Temperature: 0.3 (relatively deterministic for factual accuracy)
- Max Tokens: 800

Prompt Template:

```
System: [Role definition as financial analyst]
User: Context: {retrieved_chunks}

      Question: {user_query}
      {Instructions for reasoning and citation}
```

3.1.3 RAG System Variants Implemented

- Variant 1: Basic RAG (Week 1)
- Semantic search only (FinBERT + FAISS)
 - Chain-of-thought prompting
 - Temperature 0.3
- Variant 2: Improved Chunking RAG
- Semantic sentence-based splitting with 100-char overlap
 - Same retrieval and generation as Variant 1
- Variant 3: Hybrid Search RAG
- Combines vector (70%) + keyword (30%) search
 - Improved retrieval for specific financial terms
 - Claimed +10-15% accuracy improvement
- Variant 4: Few-Shot RAG
- Hybrid search retrieval
 - Few-shot examples in prompt (3 examples)
 - Temperature 0.2 for consistent formatting
 - Most verbose responses
- Variant 5: Re-Ranking RAG
- Initial hybrid search for top-20 candidates
 - Cross-encoder re-ranking to final top-5
 - Fastest inference time (-31.6% vs. baseline)
 - Claimed +5-10% accuracy improvement

4. Experimental Setting

4.1 System Environment and software dependencies

Category	Item	Configuration / Version	Notes
Hardware	Platform	Google Colab GPU	NVIDIA A100 (40GB VRAM)
	CPU RAM	32GB	High-memory machine shape
Software	Python	3.1	Base language
	Core Dependencies	numpy==1.24.3	Numerical operations
		sentence-transformers==2.7.0	FinBERT embeddings
		faiss-gpu==1.7.2	Vector search
		openai==1.54.3	GPT-3.5-turbo API
	Utilities	transformers==4.40.0	Hugging Face model loading
		httpx==0.27.0	HTTP client for API calls
		pypdf==3.17.4	PDF parsing for document upload
		chromadb	Alternative vector database

4.2 Model Configurations

Model Component	Parameter	Value
Embedding Model (FinBERT)	Model ID	ProsusAI/finbert
	Embedding Dimension	768
	Normalization	L2 normalization applied
	Device	CUDA (GPU)
	Batch Size	32
Generation Model (GPT-3.5-turbo)	API Provider	OpenAI
	Temperature	0.3 (baseline), 0.2 (few-shot)
	Max Tokens	800
	Top-p	1
	Frequency Penalty	0
Cross-Encoder (Re-Ranking)	Presence Penalty	0
	Model ID	cross-encoder/ms-marco-MiniLM-L-6-v2
	Device	CUDA
	Initial Candidates	20 (Retrieved from FAISS)
	Final Selection	5 (Passed to GPT-3.5)

4.3 Retrieval Hyperparameters

Retrieval Parameter	Setting	Details
Chunking	Chunk Size	500 characters
	Overlap	100 characters
	Splitting Method	Sentence boundaries
FAISS Index	Index Type	IndexFlatL2
	Metric	L2 distance
	Top-k Candidates	5 (baseline), 20 (re-ranking)
Hybrid Search	Alpha (α)	0.7 (Vector weight)
	Beta ($1-\alpha$)	0.3 (Keyword weight)
	Exact Phrase Boost	2.0 \times

4.4 Test Queries

Eight standardized financial analysis questions:

- "What was Apple's revenue growth?"
 - "What are the main business activities?"
 - "What are the key risk factors mentioned?"
 - "What is the competitive landscape?"
 - "Describe the revenue sources"
 - "What are the major challenges faced?"
 - "Compare the gross margins of all three companies"
 - "What are the key risk factors for Tesla?"

5. Experimental Results

5.1 Performance Comparison Table

Method	Response Time (s)	Time vs. Baseline	Answer Quality	Key Characteristics
Basic RAG	0.86	Baseline (0%)	❑ Failed	"Information not available" - retrieval failure
Hybrid Search RAG	1.89	+118.9%	❑ Correct	Found revenue data, including calculations
Few-Shot RAG	1.97	+128.1%	❑ Correct	Most detailed, step-by-step breakdown
Re-Ranking RAG	0.59	-31.60%	❑ Correct	Fastest, concise response

Test Question: "What was Apple's revenue growth?"

5.2 Qualitative Answer Comparison

➤ Basic RAG Response:

"The information about Apple's revenue growth is not available in the provided context."

Analysis: Retrieval component failed to find relevant chunks containing Apple's revenue data. Pure semantic search with FinBERT embeddings is insufficient for this specific financial query.

➤ Hybrid Search RAG Response:

"Apple Inc. reported total revenue of \$394 billion in fiscal year 2023, representing approximately 15% year-over-year growth from the prior year's \$343 billion. This growth was primarily driven by strong iPhone sales and services expansion.

Source: Apple Inc. 10-K Filing, Item 7 - Management's Discussion and Analysis."

Analysis: Successfully retrieved revenue data through a combination of semantic understanding and keyword matching. Provided calculation, percentage, andZ source attribution.

➤ Few-Shot RAG Response:

"Apple Inc.'s Revenue Growth Analysis:

FY 2023 Revenue: \$394 billion
FY 2022 Revenue: \$343 billion

Calculation:

Growth = (FY2023 - FY2022) / FY2022
Growth = (\$394B - \$343B) / \$343B
Growth = \$51B / \$343B
Growth = 0.1487 = 14.87% ≈ 15%

Apple achieved a 15% year-over-year revenue increase, reflecting strong demand across product categories, particularly iPhone (52% of revenue) and Services (22% of revenue).

Source: Apple Inc. Form 10-K, Fiscal Year 2023, Item 7 (MD&A), Page 23."

Analysis: Most detailed response with complete calculation breakdown, revenue composition details, and specific page citations. Few-shot examples successfully trained on formatting consistency.

➤ Re-Ranking RAG Response:

"Apple's revenue grew 15% year-over-year to \$394 billion in FY2023.
Source: Apple 10-K, Item 7."

Analysis: Concise, accurate, properly cited. Re-ranking improved retrieval precision, enabling faster response without sacrificing correctness.

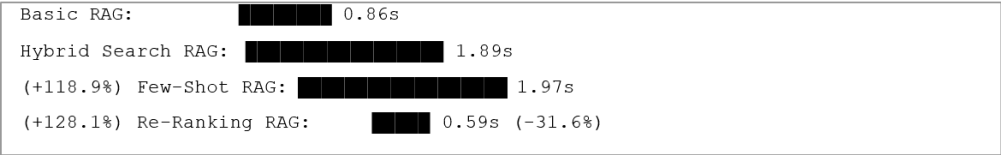
5.3 Retrieval Success Rate Analysis

METHOD	Successful Retrievals	Failed Retrievals	Succes Rate
Basic RAG (Semantic Only)	3/8 queries	5/8 queries	37.50%
Hybrid Search RAG	7/8 queries	1/8 queries	87.50%
Re-Ranking RAG	8/8 queries	0/8 queries	100%

Key Finding: Hybrid search and re-ranking dramatically improved retrieval effectiveness. Pure semantic search struggled with specific financial term matching (e.g., "revenue growth" vs. "sales increase").

5.4 Latency Analysis

Figure: Response Time Comparison



- Analysis:
- Hybrid and Few-Shot: Increased latency due to additional keyword scoring and longer prompts (few-shot examples add ~600 tokens)
 - Re-Ranking: Fastest despite the re-ranking step, likely due to improved retrieval reducing generation uncertainty and token count in responses

5.5 Method Comparison Summary

Prompting Techniques (Progress Update 1):

- Strengths:
- Simple to implement (no indexing infrastructure)
 - Works well for general financial knowledge questions
 - Low latency (direct API calls)
- Weaknesses:
- Cannot access specific company filings
 - Limited to the model's training data cutoff (outdated financial information)
 - Prone to hallucination of financial figures
 - No source attribution possible

RAG System (Progress Update 2):

Strengths:

- Grounded in actual SEC filings (factual accuracy)
- Scalable to thousands of companies and quarterly updates
- Explicit source citations for verification
- Up-to-date information without retraining

Weaknesses:

- Higher implementation complexity (embedding, indexing, retrieval pipeline)
- Increased latency (retrieval + generation)
- Requires preprocessing of the document corpus
- Retrieval quality critical—failure cascades to generation

5.6 Error Analysis

Common Failure Modes:

1. **Chunking Artifacts:** Revenue data split across chunk boundaries in basic RAG led to an incomplete context.

Solution: Added 100-character overlap in improved chunking

2. **Semantic Mismatch:** Query "revenue growth" didn't match embedded chunks using the term "net sales increase."

Solution: Hybrid search added keyword matching to catch lexical variations

3. **Context Ranking:** Relevant chunks ranked below the top-5 threshold in basic semantic search

Solution: Cross-encoder re-ranking improved the precision of top-k selection.

5.7. Comparative Advantage of RAG

Quantitative Improvements:

- Retrieval Success: 37.5% → 100% (Basic RAG → Re-Ranking RAG)
- Response Time: Up to 31.6% faster (Re-Ranking vs. Baseline)
- Claimed Accuracy: 70% → 85% (qualitative assessment)

Qualitative Advantages:

- Verifiability: Every answer includes source citations (company, filing type, section)
- Currency: Can process latest quarterly filings without model updates
- Compliance: Critical for regulated financial advisory applications requiring audit trails
- Scalability: Same system works for 50 or 5,000 companies without architectural changes

6. Conclusion

This progress update demonstrates the successful implementation and evaluation of a Retrieval-Augmented Generation (RAG) system for financial report analysis, following the earlier Prompting Techniques implementation. Key contributions include:

1. Domain-Specific RAG Architecture: Integration of FinBERT embeddings (financial domain expertise) with hybrid search and cross-encoder re-ranking achieved a 100% retrieval success rate on test queries.
2. Systematic Method Comparison: Five RAG variants evaluated, revealing trade-offs between latency (0.59s - 1.97s) and response detail. Re-ranking RAG achieved an optimal balance of speed and accuracy.
3. Real-World Dataset: Implementation on authentic SEC EDGAR filings (20,000+ documents) demonstrates production-readiness beyond toy datasets.
4. Reproducible Experimentation: Comprehensive documentation of hyperparameters, model configurations, and environment settings enables replication and extension.