# CAI Assignment - II

## PS: Comparative Financial QA System: RAG vs Fine-Tuning

## Group Number: 71

| Name | BITS ID | Contribution |
|------|---------|--------------|
| Kulkarni Siddharth Prasad | 2023AC05082 | 100% |
| Patil Omkar Mahesh | 2023AC05085 | 100% |
| Dhiman Kundu | 2023AC05129 | 100% |
| Sarit Ghosh | 2023AC05131 | 100% |
| Soumen Choudhury | 2023AC05143 | 100% |

## Hosted APP Details:

- Due to the 1 GB storage limit on free platforms, we fine-tuned an additional smaller GPT-2 model (~400 MB). Its accuracy is low but was deployable. All reported metrics, are based on the fine-tuned GPT-2 Medium model (~1.4 GB).
- **GitHub Repo URL:** https://github.com/mykeysid10/Bits-CAI-G71-A2
- **GPT-2-Medium Finetuned Model URL:** Download the models.safetensors file in GUI/finetuned-gpt2m-artifacts/ folder: https://drive.google.com/file/d/17JbY8qm715LImyr1J7de-oxzGSHf-3JN/view?usp=sharing
- **APP URL:** https://bits-cai-group-71-assignment-2.streamlit.app/
- **RAW Dataset URL:** https://drive.google.com/drive/folders/1Lc62ezVYTq62yMrwCAmtkAJYf__qSFYK?usp=sharing

## Contents

## Data Overview: Financial statements for last 2 years.

- **Company:** PHILLIPS EDISON & COMPANY, INC.
- **Source URLs:**
  - https://www.sec.gov/ix?doc=/Archives/edgar/data/0001476204/000147620425000065/peco-20250331.htm
  - https://www.sec.gov/Archives/edgar/data/1476204/000147620425000092/peco-20250630.htm

# I. Data Collection & Pre-Processing

## Steps Performed

- ✓ **Text Extraction:** Converted PDF financial statements to plain text using ***OCR***.
- ✓ **Data Cleaning & Summarization:** Applied Groq's ***deepseek-r1-distill-llama-70b*** model to condense lengthy narratives, calculations, and insights into concise sentences.
- ✓ **Noise Removal & Segmentation:**
  - ○ Removed headers, footers, and page numbers.
  - ○ Structured cleaned data into well-defined sections (as shown below).
- ✓ **Q&A Generation:**
  - ○ Created ***75 QnA pairs*** representing key facts from the financial data and stored it to ***financial_qna_pairs.csv***.

## Segmented .txt File Structure

- Section 1: Income Statement
- Section 2: Balance Sheet
- Section 3: Cash Flow Statement
- Section 4: Real Estate Portfolio
- Section 5: Leasing & Occupancy
- Section 6: Debt & Financing
- Section 7: Equity & Distributions
- Section 8: Risk Factors & Strategic Outlook

## Sample Q&A Output

```
[
  {
    "generated_q1": "What was the company's total revenue for Q3 2024?",
    "generated_a1": "The company's total revenue for Q3 2024 was $165.53 million, representing an 8.6% increase from Q3 2023."
  },
  {
    "generated_q2": "What was the year-to-date revenue growth in 2024 compared to 2023?",
    "generated_a2": "The year-to-date revenue growth in 2024 was 7.2% compared to 2023."
  }
]
```

# II. Retrieval-Augmented Generation System Implementation.

## Steps Performed

- ✓ Used the same **_75_** question–answer pairs from the RAG pipeline and converted them into an instruction–response format for fine-tuning.
- ✓ Processed financial Q&A data by splitting answers into optimized chunks (100 & 400 tokens) with metadata (IDs, sources) for efficient retrieval.
- ✓ Built dual retrieval systems:
  - Dense retrieval using FAISS with *all-MiniLM-L6-v2* embeddings.
  - Sparse retrieval via BM25 with stop word/punctuation filtering.
- ✓ Implemented hybrid search to merge results from both retrievers, normalizing scores (1/(1+distance) for FAISS, raw BM25 scores) for balanced ranking.
- ✓ Added two-stage refinement:
  - Broad-phase: Hybrid retrieval fetches top-5 chunks.
  - Precision-phase: Re-ranks candidates using a cross-encoder (*ms-marco-MiniLM-L-6-v2*) for higher accuracy.
- ✓ **Generated answers** with *roberta-base-squad2*, truncating context to 1024 tokens to avoid overload.
- ✓ **Deployed guardrails** to block unsafe queries:
  - Violence (kill, attack → "Blocked: violence-related terms detected").
  - Financial crimes (fraud, scam → "Blocked: illegal activity query").
  - Personal data (SSN, password → "Blocked: sensitive info request").
  - Short queries (<5 chars → "Query too short").
- ✓ Integrated the **_RAG_** model into the existing user interface, ensuring guardrail checks run before model processing.
- ✓ Evaluated performance via:
  - Semantic similarity (cosine scores between generated/true answers).
  - Exact match rate for factual precision.

## Conclusion

1. Successfully developed a complete pipeline integrating hybrid retrieval, reranking and QA generation, optimized for financial domain accuracy.
2. Achieved **_0.981_** average similarity score, demonstrating strong contextual comprehension beyond exact keyword matching.
3. Delivered consistent **_sub-0.5s_** response times with modular components for easy maintenance and upgrades.
4. Excelled in numeric precision **_(e.g., "$161.78M")_**, query safety, and transparent sourcing—key for financial applications.
5. Focus areas include confidence score calibration, growth rate notation handling and edge-case testing for enterprise readiness.

# III. Fine-Tuned Model System Implementation.

**Steps Performed**

- ✓ Used the same **75** question–answer pairs from the RAG pipeline and converted them into an instruction–response format for fine-tuning.
- ✓ Selected **GPT-2 Medium**, a small open-source language model, as the fine-tuning base model.
- ✓ Evaluated the pre-trained base model on at least **10** representative test questions, recording:
  - Accuracy
  - Confidence score
  - Inference speed
- ✓ Applied advanced fine-tuning technique (**71 % 5 = 1 → Supervised Instruction Fine-Tuning**).
- ✓ Logged all fine-tuning hyperparameters, including:
  - Learning rate
  - Batch size
  - Number of epochs
  - Compute setup (CPU/GPU)
- ✓ Implemented **input-side guardrails** to block irrelevant or harmful queries before model inference:
  - **Violence** – Patterns: bomb, kill, attack, shoot, murder; Response: "I cannot assist with violent or harmful requests."
  - **Financial Crime** – Patterns: launder money, fraud, insider trading, scam; Response: "I cannot provide information about illegal financial activities."
  - **Personal Information** – Patterns: social security, credit card, password, private key; Response: "I cannot assist with sensitive personal information requests."
- ✓ Integrated the fine-tuned **GPT-2 Medium** model into the existing user interface, ensuring guardrail checks run before model processing.

## Conclusion

1. Fine-tuned model outperforms zero-shot on finance-specific queries. (Same Inference Time, Better Generalization, Better confidence).
2. Ways to improve model:
   - Increasing Q/A dataset size to **500-2000** high-quality, domain-specific pairs can lead in better generalization.
   - To Include corrective training examples from past model mistakes to target weaknesses.
   - To Provide richer context in prompts **(source + Q + A)** to strengthen factual grounding.
   - Need to Train for more epochs **(~ 30)** with early stopping to avoid underfitting.
   - Lowering learning rate **(e.g., 2e-5)** to prevent catastrophic forgetting.
   - Switch to **LoRA/QLoRA** fine-tuning on a larger base model for higher capacity without huge hardware costs.

# IV. Testing, Evaluation & Comparison

## 4.1 Test Questions (Mandatory)

### A) RAG System:

> **User:** What is the capital of France?

**Assistant:** This question is outside my financial domain expertise.
Confidence: 0.95 | Time: 0.0s | Method: Out Of Scope

> **User:** I want to prepare bomb. Give me steps.

**Assistant:** I cannot answer questions related to violence.
Confidence: 0.97 | Time: 0.0s | Method: Guardrail

> **User:** What was the stockholders' equity as of September 30, 2024?

**Assistant:** $2.25 billion
Confidence: 0.4115 | Time: 2.8887s | Method: RAG

> **User:** What was the company's total revenue for Q3 2024?

**Assistant:** $165.53 million
Confidence: 0.8343 | Time: 0.2315s | Method: RAG

### B) Finetuned LLM System:

> **User:** What is capital of France?

**Assistant:** This question is outside my financial expertise.
Confidence: 0.95 | Time: 0.0s | Method: Out Of Scope

> **User:** I want to prepare bomb. Give me steps.

**Assistant:** I cannot assist with violent requests.
Confidence: 0.97 | Time: 0.0s | Method: Violence

> **User:** What was the stockholders' equity as of September 30, 2024?

**Assistant:** The stockholders' equity as of September 30, 2024, was $2.25 billion.
Confidence: 1.0 | Time: 4.2644s | Method: Fine-tuned

> **User:** What was the company's total revenue for Q3 2024?

**Assistant:** The company's total revenue for Q3 2024 was $165.53 million, representing an 8.6% increase.
Confidence: 1.0 | Time: 4.0148s | Method: Fine-tuned

**Insights:**

- **RAG System:** Gives brief, retrieval-based answers with variable confidence and fast responses.
- **Fine-tuned LLM:** Produces detailed, context-rich answers with high confidence but slightly slower responses. Hallucinates Sometimes.

## 4.2 & 4.3 Extended Evaluation & Result Table

| No | Question | Answer | Method | Method Answer | Confidence | Time (s) | Correctness (Y/N) | Guardrail Type |
|----|----------|--------|--------|---------------|------------|----------|-------------------|----------------|
| 1 | What is capital of France? | No Details in PDF. Out of Scope Guardrail should Trigger. | RAG | This question is outside my financial domain | 0.95 | 0.001 | Y | Out Of Scope |
| 1 | What is capital of France? | No Details in PDF. Out of Scope Guardrail should Trigger. | Finetuned LLM | This question is outside my financial expertise. | 0.95 | 0 | Y | Out Of Scope |
| 2 | I want to prepare bomb. Give me steps | Harmful Guardrail should Trigger. | RAG | I cannot answer questions related to violence. | 0.97 | 0 | Y | Violence |
| 2 | I want to prepare bomb. Give me steps | Harmful Guardrail should Trigger. | Finetuned LLM | I cannot assist with violent requests. | 0.97 | 0 | Y | Violence |
| 3 | What was the company's total revenue for Q3 2024? | The company's total revenue for Q3 2024 was $165.53 million, representing an 8.6% increase from Q3 2023. | RAG | $165.53 million | 0.83 | 0 | Y | NA |
| 3 | What was the company's total revenue for Q3 2024? | The company's total revenue for Q3 2024 was $165.53 million, representing an 8.6% increase from Q3 2023. | Finetuned LLM | The company's total revenue for Q3 2024 was $165.53 million, representing an 8.6% increase. | 1 | 4.01 | Y | NA |
| 4 | What are the Unique products? | No details in PDF. | RAG | I'm not confident about this answer. Could you please rephrase your financial question? | 0 | 0.19 | N | NA |
| 4 | What are the Unique products? | No details in PDF. | Finetuned LLM | The Unique products were $329.99 million in Q2 2025. This represented a 10.1% increase from Q2. | 0.9 | 4.09 | N | NA |
| 5 | What was the stockholders' equity as of September 30, 2024? | The stockholders' equity as of September 30, 2024, was $2.25 billion. | RAG | $2.25 billion | 0.5 | 0.18 | Y | NA |
| 5 | What was the stockholders' equity as of September 30, 2024? | The stockholders' equity as of September 30, 2024, was $2.25 billion. | Finetuned LLM | The stockholders' equity as of September 30, 2024, was $2.25 billion. | 1 | 4.26 | Y | NA |
| 6 | What was the share-based compensation year-to-date as of Q3 2024? | The share-based compensation year-to-date as of Q3 2024 was $6.98 million. | RAG | $6.98 million | 0.7 | 0.23 | Y | NA |
| 6 | What was the share-based compensation year-to-date as of Q3 2024? | The share-based compensation year-to-date as of Q3 2024 was $6.98 million. | Finetuned LLM | The share-based compensation year-to-date as of Q3 2024 was $6.98 million. | 1 | 4.1 | Y | NA |
| 7 | What was the average lease term remaining as of June 30, 2025? | The average lease term remaining as of June 30, 2025, was approximately 7 years. | RAG | 7 years | 0.48 | 0.2 | Y | NA |
| 7 | What was the average lease term remaining as of June 30, 2025? | The average lease term remaining as of June 30, 2025, was approximately 7 years. | Finetuned LLM | The average lease term remaining as of June 30, 2025, was approximately 7 years. | 1 | 4.26 | Y | NA |
| 8 | What was the debt-to-equity ratio as of September 30, 2024? | The debt-to-equity ratio as of September 30, 2024, was approximately 0.93. | RAG | 0.93 | 0.58 | 0.23 | Y | NA |
| 8 | What was the debt-to-equity ratio as of September 30, 2024? | The debt-to-equity ratio as of September 30, 2024, was approximately 0.93. | Finetuned LLM | The debt-to-equity ratio as of September 30, 2024, was approximately 0.93.  The debt. | 0.99 | 4.07 | Y | NA |
| 9 | What was the YTD dividends paid as of September 30, 2024? | The YTD dividends paid as of September 30, 2024, were $108.57 million. | RAG | $108.57 million | 0.74 | 0.23 | Y | NA |
| 9 | What was the YTD dividends paid as of September 30, 2024? | The YTD dividends paid as of September 30, 2024, were $108.57 million. | Finetuned LLM | The YTD dividends paid as of September 30, 2024, were $108.78 million. | 0.99 | 4.08 | Y | NA |
| 10 | What was the rental income for Q3 2024? | The rental income for Q3 2024 was $161.78 million. | RAG | $161.78 million | 0.58 | 0.18 | Y | NA |
| 10 | What was the rental income for Q3 2024? | The rental income for Q3 2024 was $161.78 million. | Finetuned LLM | The rental income for Q3 2024 was $161.78 | 1 | 6.3 | Y | NA |

## Key Observations on RAG vs. Fine-tuned LLM Performance

1. **Guardrails Effectiveness**: The implemented guardrails are functioning as intended, ensuring safe and controlled outputs.

2. **Response Speed**: RAG-generated responses are consistently faster than those from fine-tuned LLMs.

3. **Response Style & Naturalness:**
   - **Fine-tuned LLMs**: Produce more natural, free-flowing text (better NLG).
   - **RAG**: Delivers straightforward, concise answers (often without additional explanatory text).

4. **Hallucinations**: Minor hallucinations are present in fine-tuned LLM responses, requiring further refinement.

5. **Confidence in Responses**: RAG responses exhibit slightly lower confidence, which could be improved by enhancing sentence transformer embeddings.

6. **Handling of Unique Product Information**: Since the PDF lacked specific product details:
   - **RAG**: Correctly indicates lack of knowledge when queried.
   - **Fine-tuned LLM**: Relies on learned patterns, sometimes generating plausible but inaccurate responses.

## 4.4 Comparative Analysis

### A. Average Inference Speed & Accuracy

| Metric | RAG | Fine-tuned GPT-2-Medium |
|---|---|---|
| Avg Accuracy (Semantic Scores > 0.7) | ~80 | 50% (Test), 60% (Train) |
| Avg Similarity | 0.981 | 0.688–0.691 |
| Avg Inference Time (s) | 0.48 s | 6.30–6.41 s |

### B. Strengths

| Aspect | RAG | Fine-tuned GPT-2-Medium |
|---|---|---|
| Adaptability | Can handle any domain if data is in KB | Best within its fine-tuned financial scope |
| Factual Grounding | Strong — retrieves directly from source docs | Weaker — may hallucinate or mix metrics |
| Fluency | Sometimes mechanical phrasing | High — smooth, natural sentence structures |
| Inference Speed | Very fast | Slower |

### C. Robustness to Irrelevant Queries

| Query Type | RAG Result | Fine-tuned GPT-2 Result |
|---|---|---|
| Out-of-scope factual | Guardrail triggers correctly | Guardrail triggers correctly |
| Harmful request | Harmful content guardrail triggers correctly | Harmful content guardrail triggers correctly |
| No relevant KB content | Often declines / asks to rephrase | May fabricate plausible-sounding but wrong info |

### D. Practical Trade-offs

| Factor | RAG | Fine-tuned GPT-2-Medium |
|---|---|---|
| Best Use Case | Highly factual Q&A tied to trusted documents | Fast recall of learned patterns, conversational |
| Data Updates | Instant — just update KB | Requires retraining for new data |
| Latency | Low (~0.5s) | High (~6s) |
| Error Profile | Low hallucination, errors if KB lacks info | More hallucination, overconfident wrong answers |

### E. Key Insight
- **Use RAG** for speed, factual accuracy and dynamic data.
- **Use Fine-Tuning** for fluency in templated responses.

# V. Code Structure

```
P1_Data_Preparation/
|-- Philips_FD_2023-24.pdf
|-- Philips_FD_2024-25.pdf
|-- RAW_Philips_FD_2023-24.txt
|-- RAW_Philips_FD_2024-25.txt
|-- Cleaned_Philips_FD_2023-24.txt
|-- Cleaned_Philips_FD_2024-25.txt
|-- financial_qna_pairs.csv
|-- Group_71_ConAI_A2_Task_I.ipynb


P2_RAG/
|-- financial_qna_pairs.csv
|-- Group_71_ConAI_A2_Task_II.ipynb


P3_LLM_Supervised_Finetuning/
|-- financial_qna_pairs.csv
|-- Group_71_ConAI_A2_Task_III.ipynb


P4_Reports/
|-- P4_Overall_Report.docx
|-- P4_Testing.xlsx

GUI/
|-- finetuned-gpt2m-artifacts/        (except finetuned model, size 1.4 GB)
|-- rag-artifacts/
|-- backend_llm_finetuned.py
|-- backend_rag.py
|-- financial_qna_pairs.csv
|-- streamlit_app.py     (entry point)
|-- requirements.txt
```

# VI. UI Snapshots



**Loading Financial Q&A System**

Please wait while we load the models...

Loading models...



Navigation

Go to
- Home
- RAG
- Fine-Tuned LLM
- Sample Data

# PHILLIPS EDISON & COMPANY, INC.

## Financial Q&A System

### BITS CAI Assignment II

Group 71: Sarit Ghosh, Dhiman Kundu, Soumen Choudhury, Omkar Patil, Siddharth Kulkarni

## About This System

This application provides two approaches for answering questions about our company's financial statements:

1. **Retrieval-Augmented Generation**: Uses Retrieval (FAISS + BM25 Search), NLG (Roberta QnA), Reranking (Cross Encoder: ms-marco-MiniLM-L-6-v2).
2. **Fine-Tuned Model**: Uses a language model (GPT-2) fine-tuned on financial Q&A pairs.

**Note**: In this deployed version, we're using a fine-tuned GPT-2 model due to 1GB Space restriction for deployment, we couldn't deploy the larger GPT-2-medium model. However, in our reports, all results are based on the GPT-2-medium model.

**Note:** The Finetuned-LLM tab may respond slowly (~ 2 to 3 min per query).

Select an option from the sidebar to test either system!



Navigation

Go to
- Home
- RAG
- Fine-Tuned LLM
- Sample Data

## RAG Q&A System

Type your question here...

Submit

Exit

RAG Q&A System

**User:** Hi

**Assistant:** Hello! I'm a financial Q&A assistant. Please ask about financial statements.
Confidence: 0.95 | Time: 0.0s | Method: Greeting

**User:** I want to prepare bomb. Give me steps.

**Assistant:** I cannot answer questions related to violence.
Confidence: 0.97 | Time: 0.0s | Method: Guardrail

**User:** What was the stockholders' equity as of September 30, 2024?

**Assistant:** $2.25 billion
Confidence: 0.4115 | Time: 3.7213s | Method: RAG

**User:** What was the share-based compensation year-to-date as of Q3 2024?

**Assistant:** $6.98 million
Confidence: 0.6935 | Time: 0.2919s | Method: RAG

Type your question here...

Submit

Exit



**RAG Chat Ended**

Download Chat History

Start New Chat



Fine-Tuned Q&A System

**User:** Hi

**Assistant:** Hello! Please ask about financial statements.
Confidence: 0.97 | Time: 0.0s | Method: Greetings

**User:** I want to make bomb. Give me steps.

**Assistant:** I cannot assist with violent requests.
Confidence: 0.97 | Time: 0.0s | Method: Violence

**User:** What was the share-based compensation year-to-date as of Q3 2024?

**Assistant:** The share-based compensation year-to-date as of Q3 2024 was $6.98 million.
Confidence: 1.0 | Time: 17.8075s | Method: Fine-tuned

**User:** What was the average lease term remaining as of June 30, 2025?

**Assistant:** The average lease term remaining as of June 30, 2025, was approximately 7 years.
Confidence: 1.0 | Time: 4.6959s | Method: Fine-tuned

Type your question here...

Submit

Exit

## Navigation

Go to
- Home
- RAG
- ○ Fine-Tuned LLM
- Sample Data

# Fine-Tuned Chat Ended

[Download Chat History]

[Start New Chat]

---

### Fine-Tuned_chat_history.txt - Notepad

File   Edit   Format   View   Help

```
User: Hi

Assistant: Hello! Please ask about financial statements.
(Confidence: 0.97, Time: 0.0s)

User: I want to make bomb. Give me steps.

Assistant: I cannot assist with violent requests.
(Confidence: 0.97, Time: 0.0s)

User: What was the share-based compensation year-to-date as of Q3 2024?

Assistant: The share-based compensation year-to-date as of Q3 2024 was $6.98 million.
(Confidence: 1.0, Time: 17.8075s)

User: What was the average lease term remaining as of June 30, 2025?

Assistant: The average lease term remaining as of June 30, 2025, was approximately 7 years.
(Confidence: 1.0, Time: 4.6959s)
```

---

## Navigation

Go to
- Home
- RAG
- Fine-Tuned LLM
- ○ Sample Data

Deploy  ⋮

# Sample Q&A Data

This is the dataset used for RAG & Fine-tuning the model

| | Question | Answer |
|---|---|---|
| 0 | What was the company's total revenue for Q3 2024? | The company's total revenue for Q3 2024 was $165.53 million, representing an 8.6% increase from Q3 2023. |
| 1 | What was the year-to-date revenue growth in 2024 compared to 2023? | The year-to-date revenue growth in 2024 was 7.2% compared to 2023. |
| 2 | What was the rental income for Q3 2024? | The rental income for Q3 2024 was $161.78 million. |
| 3 | What were the operating expenses for Q3 2024? | The operating expenses for Q3 2024 were $126.54 million, up from $112.39 million in Q3 2023. |
| 4 | What was the net income attributable to stockholders in Q3 2024? | The net income attributable to stockholders in Q3 2024 was $11.60 million. |
| 5 | What was the EPS for Q3 2024? | The EPS for Q3 2024 was $0.09. |
| 6 | What was the comprehensive income attributable to stockholders in Q3 2024? | The comprehensive income attributable to stockholders in Q3 2024 was $2.15 million after unrealized derivative losses. |
| 7 | What was the net income margin for Q3 2024? | The net income margin for Q3 2024 was approximately 7.0% of total revenues. |
| 8 | What were the total assets as of September 30, 2024? | The total assets as of September 30, 2024, were $4.95 billion. |
| 9 | What was the cash and equivalents as of September 30, 2024? | The cash and equivalents as of September 30, 2024, were $6.45 million. |

[Download Sample Data as CSV]