# FTEC 5660 HW1 Report

Tian Zixiao 1155244606

January 30, 2026

## 1 Analysis

**Multimodal Understanding and OCR Challenges:** The core of this task lies in processing dual inputs: "images + text." The system must not only recognize text via OCR but also possess sophisticated semantic understanding. Given the vast differences in receipt layouts across various supermarkets, the Agent must precisely locate specific key fields (e.g., "Total Payable," "Amount Deducted") rather than simply capturing all numerical data.

**Intent Routing and Task Logic:** The system must accurately dispatch logic based on the user's query:

- **Query 1 (Key Data Extraction):** The objective is to directly identify and extract the explicitly labeled "Final Total Amount" from the bill image. This requires high visual grounding precision to distinguish between confusing items like "Subtotal," "Total," and "Change."

- **Query 2 (Reverse Logical Calculation):** More complex than simple extraction, this task requires the Agent to identify all negative values or discount items and logically sum them with the paid amount to reconstruct the "Original Price before discount."

- **Irrelevant Queries (Rejection Mechanism):** The Agent must maintain clear semantic boundaries to recognize and reject irrelevant or invalid instructions, preventing the model from generating hallucinations.

**Guaranteeing Numerical Accuracy:** While Query 1 focuses on extraction, Query 2 involves calculation, and aggregating data from multiple bills still requires numerical processing. To ensure 100% accuracy, the system should employ structured extraction (e.g., using Pydantic) combined with deterministic calculation logic to avoid the probabilistic errors inherent in LLMs when handling numbers.

## 2 Design

### 2.1 Flowchart

To visualize the execution logic and data flow of the Agent, the overall system workflow is illustrated below. This flowchart depicts the end-to-end pipeline, spanning from multimodal input acquisition and intent-based routing to structured data extraction and final deterministic calculation.
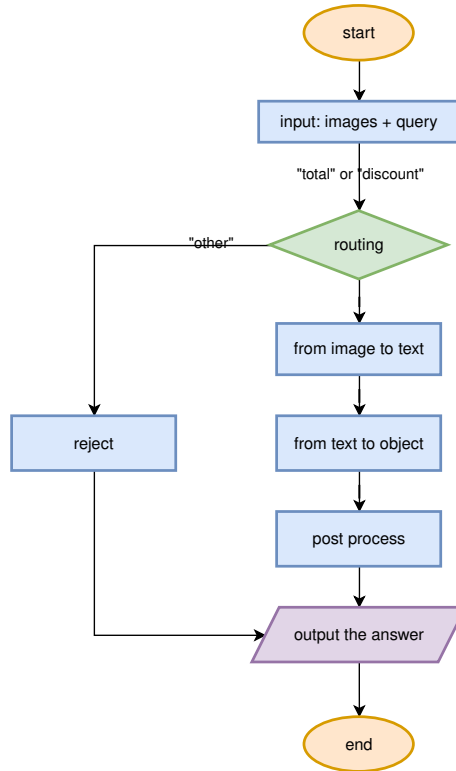
Figure 1: End-to-End System Workflow

## 2.2 Prompts

**Prompt for Routing:**

```
system_prompt = """
You are a query classifier. Focus only on the user's text and
    ignore any image/base64 data.
Assign a 'question_type' as an integer:
- 'total': if the query is similar to "How much money did I spend
     in total for these bills?".
- 'discount': if the query is similar to "How much would I have
    had to pay without the discount?".
- 'other': any other irrelevant queries.
ONLY output one word: 'total', 'discount', or 'other'.
"""
```

Listing 1: System Prompt for Intent Classification

**Prompt for Extracting Information from Images:**

```
transcription_prompt = """
Analyze the provided receipt image(s).
Transcribe all text from the receipt(s) verbatim, line by line.
Pay special attention to numbers, and symbols like '-' or '
    Discount'.
Do not summarize. Just give me the raw text found in the images.
"""
```

Listing 2: Prompt for Image Transcription

## 2.3 Data Structure

To facilitate structured parsing and efficient processing of bill data, I designed a stream-lined data model. This model is engineered to encapsulate core information from receipts—with maximum precision and zero redundancy. By ensuring the integrity of the data schema, the Agent can maintain high-quality inputs for subsequent logical reasoning.
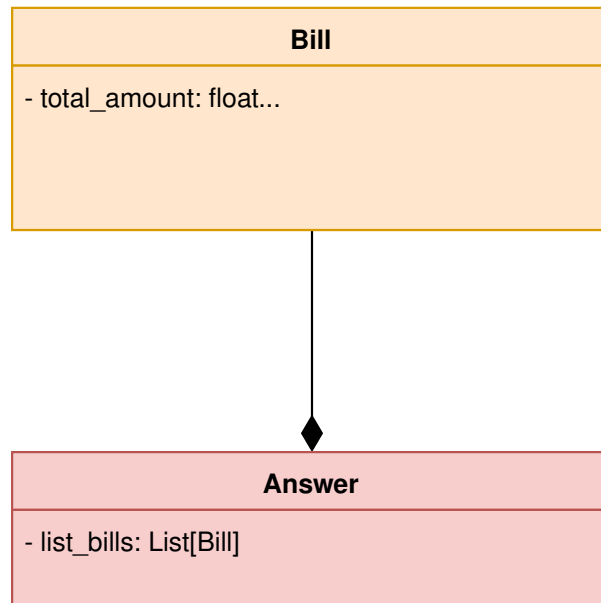


Figure 2: UML Class Diagram for Bill and Answer Classes

As shown in Figure 2, the **Answer** class maintains a composition relationship with the **Bill** class, reflecting a one-to-many mapping for multi-image processing scenarios.

# 3 Outcome

## 3.1 Intermediate Module Validation

### 3.1.1 Router Chain

This module performs semantic classification on user queries to dispatch tasks to the appropriate logic branch.

### 3.1.2 Transcription Chain

This component leverages multimodal capabilities to convert raw receipt images into machine-readable text. The output of the transcription chain serves as the direct input for the extraction chain, ensuring a seamless data flow.

### 3.1.3 Extraction Chain

This module parses unstructured text into structured `Answer` objects based on the pre-defined Pydantic schema.

```python
# List 1: Queries about total amount spent (Question Type 1)
total_spent_queries = [
    "What is the grand total amount paid for these receipts?",
    "Could you calculate the final sum of all these bills?",
    "What was my total expenditure across all these transactions?",
    "How much did everything cost in total after all deductions?"
]

# List 2: Queries about total before discounts (Question Type 2)
pre_discount_queries = [
    "What was the total original price before any discounts were applied?",
    "How much would I have paid if there were no coupons or savings?",
    "What is the gross total of these items before any price reductions?",
    "Can you tell me the full price of these items without the discounts?"
]
# List 3
other_queries = [
    "What is the current weather forecast for Hong Kong?",
    "Can you help me write a Python script to scrape a website?",
    "Who won the Academy Award for Best Picture last year?"
]

for query in total_spent_queries+pre_discount_queries+other_queries:
  result = router_chain.invoke({"question": query})
  print(f"{query:<70} | {result:<15}")
```

```
What is the grand total amount paid for these receipts?                | total
```

Figure 3: Execution output of the routing module

```python
# Test transcription_chain

image_paths = ["/content/receipt1.jpg","/content/receipt2.jpg"]

image_data_urls = [get_image_data_url(path) for path in image_paths]

vision_chain_resp = transcription_chain.invoke({"image_data_urls": image_data_urls})
print(vision_chain_resp)
```

```
**Image 1:**
www.moneyback.com.hk
084213 韭菜豬肉雲吞20粒裝
包裝變形
084213 韭菜豬肉雲吞20粒裝
包裝變形
084213 韭菜豬肉雲吞20粒裝
包裝變形
395092 IF100% COCONUT WATER
數量: 2
Buy 2 Save $12.8
044228 玉芒
數量: 3
Buy 3 Save $10.8
490948 Fresh綜合餐汁
數量: 2
Buy 2 Save $3.9
126894 雀巢脫脂高鈣牛奶飲品
055421 蔬菜先生包裝蕃茄(中
047827 蔬菜先生翠玉瓜1磅(
089288 鹹蛋蒸肉餅
```

Figure 4: Execution output of the vision module (Part 1)

```python
# Test extraction_chain

extraction_chain_resp = extraction_chain.invoke(vision_chain_resp)
print(extraction_chain_resp)
```

```
list_bills=[Bill(total_amount=394.7, list_discount=[-12.4, -12.4, -12.4, -12.8, -10.8, -3.9, -20.78, -0.02]), Bil
```

Figure 5: Execution output of the vision module (Part 2)

## 3.2  End-to-End System Testing

This section evaluates the entire pipeline's performance in real-world scenarios, from image input to final numerical output.

```python
# List 1: Queries about total amount spent (Question Type 1)
total_spent_queries = [
    "What is the grand total amount paid for these receipts?",
    "Could you calculate the final sum of all these bills?",
    "What was my total expenditure across all these transactions?"
]

# List 2: Queries about total before discounts (Question Type 2)
pre_discount_queries = [
    "What was the total original price before any discounts were applied?",
    "How much would I have paid if there were no coupons or savings?",
    "What is the gross total of these items before any price reductions?"
]

question = total_spent_queries[0]
image_paths = ["/content/receipt1.jpg","/content/receipt2.jpg"]
image_data_urls = [get_image_data_url(path) for path in image_paths]

request={
    "question": question,
    "image_data_urls": image_data_urls
}

resp = coordinator_agent.invoke(request)
print(resp)
type(resp)
```

```
710.8
float
```

Figure 6: Execution output of the complete agent

# 4  Conclusion

This project successfully developed a bill-analysis Agent based on a modular architecture. By decoupling the workflow into **intent routing**, **structured extraction**, and **deterministic calculation**, the system effectively addresses the inherent unreliability of LLMs in multi-step reasoning and precise numerical computation. Experimental results demonstrate that integrating Python-based post-processing ensures 100% accuracy for financial data, while the intent recognition mechanism provides robust rejection of irrelevant queries.

**Future Work:** While the system performs exceptionally on standard receipts, the robustness of OCR under extreme blurriness or handwritten conditions could be further improved. Future iterations may explore ensemble modeling to enhance the stability of data extraction in edge cases.