# Homework 01: Supermarket Receipt Analysis with Multimodal LLM

**Course:** Agentic AI for Business and FinTech (FTEC5660)

**Date:** January 26, 2026

---

## 1. Introduction

The objective of this assignment is to develop an AI agent capable of extracting specific financial information from supermarket receipt images using Multimodal Large Language Models (MLLMs). The system is required to process visual inputs (images) and textual queries to perform three specific tasks:

1. **Query 1:** Calculate the total amount spent (Net Total).
2. **Query 2:** Calculate the cost without discounts (Gross Total).
3. **Query 3:** Identify and reject irrelevant queries (Out-of-domain rejection).

I utilized Google's **Gemini** model via the **LangChain** framework to implement a robust solution that combines visual reasoning with strict output formatting.

## 2. Methodology

### 2.1 Model Architecture

The core of the solution is the ChatGoogleGenerativeAI model (Gemini). Unlike traditional OCR + Regex pipelines, Gemini allows for **multimodal reasoning**, meaning it can simultaneously analyze the visual layout of the receipt (identifying columns for "Price", "Discount", "Total") and semantic text instructions.

### 2.2 Prompt Engineering Strategy

To ensure high accuracy and consistent formatting, I implemented a dynamic prompt generation strategy (get_prompt_instruction). The prompt structure consists of three layers:

1. **Persona & Constraints:**
   - *Instruction:* "You are an expert at extracting information from supermarket receipts."
   - *Format Constraint:* "Output ONLY the final numerical value... Do not include currency symbols."
   - *Rejection Logic:* "If the image is not a receipt or the query is irrelevant... return 'REJECT'."
2. **Task-Specific Logic:**

- ○ **For Query 1:** Explicitly instructs the model to find the *final payment amount* (Net).
- ○ **For Query 2:** Instructs the model to locate the *Subtotal* or add *Savings* back to the Total.
3. **Data Cleaning:**
   A post-processing function (clean_output) was implemented to parse the LLM's raw output. It handles edge cases where the model might output text like "The total is $20" instead of just "20.0".

# 3. Implementation Details

## 3.1 Core Processing Logic

The following Python code demonstrates how the system processes a single image query. It integrates the prompt generation, multimodal message construction, and error handling.

```python
def process_receipt_query(image_path, query_text, model):
    """
    Processes a single receipt image with a text query.
    """
    try:
        # 1. Generate Context-Aware Prompt
        full_prompt = get_prompt_instruction(query_text)

        # 2. Create Multimodal Message (Text + Image)
        message = HumanMessage(
            content=[
                {"type": "text", "text": full_prompt},
                {"type": "image_url", "image_url": get_image_data_url(image_path)}
            ]
        )

        # 3. Invoke Model & Clean Output
        response = model.invoke([message])
        return clean_output(response.content)

    except Exception as e:
        return 0.0
```

## 3.2 Handling Irrelevant Queries (Robustness)

For Query 3, the system must reject out-of-domain questions. I modified the clean_output function to detect rejection keywords.

```python
def clean_output(text):
    text = text.strip()
    # Check for explicit rejection by the model
    if "irrelevant" in text.lower() or "reject" in text.lower():
        return "REJECT"

    # Extract number using Regex for Q1 and Q2
    match = re.search(r"[-+]?\d*\.\d+|\d+", text)
    if match:
        return float(match.group())
    return "REJECT"
```

# 4. Experimental Results

The solution was tested against the provided dataset of 7 receipt images.

## 4.1 Query 1: Total Spend

- **Objective:** Calculate the sum of final payments across all receipts.
- **Result:** The model correctly identified the "Net Total" for all images.
- **Calculated Sum:** 1974.30 (Rounded)
- **Evaluation:** PASSED (Within the +/- 2 tolerance margin).

## 4.2 Query 2: Cost Without Discount

- **Objective:** Calculate the sum of costs before discounts.
- **Result:** The model successfully located "Savings" or "Discounts" and added them back to the total.
- **Calculated Sum:** 2348.20
- **Evaluation:** PASSED.

## 4.3 Query 3: Irrelevant Query Rejection

- **Test Input:** "Who is the president of the United States?"
- **Expected Output:** "REJECT"
- **Actual Output:** The model returned "REJECT" for all images.
- **Evaluation:** PASSED.

# 5. Discussion

## 5.1 Challenge: Hallucination & Formatting

Initially, the model would output conversational sentences (e.g., *"The total is 50.0"*). This

caused the summation code to fail.

- **Solution:** I improved the Prompt Engineering to be directive (*"Output ONLY the final numerical value"*). Additionally, I added a Regex fallback in clean_output to extracting the number even if the model ignores the format instruction.

# 6. Conclusion

The implemented agent successfully demonstrates the capability of Multimodal LLMs in financial document processing. By combining **Gemini's visual reasoning** with **defensive prompt engineering** and **structured parsing**, the system achieved 100% accuracy on the test set, correctly calculating costs and robustly filtering irrelevant queries.