

## **Solution Description**

This homework implements a multimodal receipt understanding system based on Gemini, which is designed to answer two types of queries from supermarket receipt images:

Query 1: *How much money did I spend in total for these bills?*

Query 2: *How much would I have had to pay without the discount?*

The system takes a set of receipt images together with a natural language query as input, and returns either a numerical result or a rejection depending on the query type.

In the implementation, each receipt image is first processed by the Gemini. Through designed prompts, the model extracts structured numerical information from the image and returns the final amount actually paid after discount (`paid_total`) and the total amount before discount (`pre_discount_total`). To ensure reliability, the model output is constrained to pure numerical values. If discount information is not explicitly available on the receipt, the pre-discount total is set equal to the paid total. If the model cannot reliably identify the total amount, null values are returned to avoid misleading results. After extracting receipt-level information, the system determines how to process the user query by means of a query routing module. Queries are classified into three categories: queries asking for the total amount actually paid (Query 1), queries asking for the total amount without discount (Query 2), and all other queries, which are treated as out-of-domain. The routing process is also implemented using a language model, with outputs strictly constrained to a predefined label set to ensure deterministic behavior.

When a query is classified as Query 1 or Query 2, the system aggregates the corresponding numerical fields across all receipt images and returns a single floating-point value, as required by the automatic evaluation program. When a query is classified as out-of-domain, the system explicitly rejects the request and does not return any numerical output.

At the implementation level, caching and retry mechanisms are introduced to improve system robustness. Extracted results for each receipt image are cached to avoid redundant model calls, and an exponential backoff strategy is applied when temporary

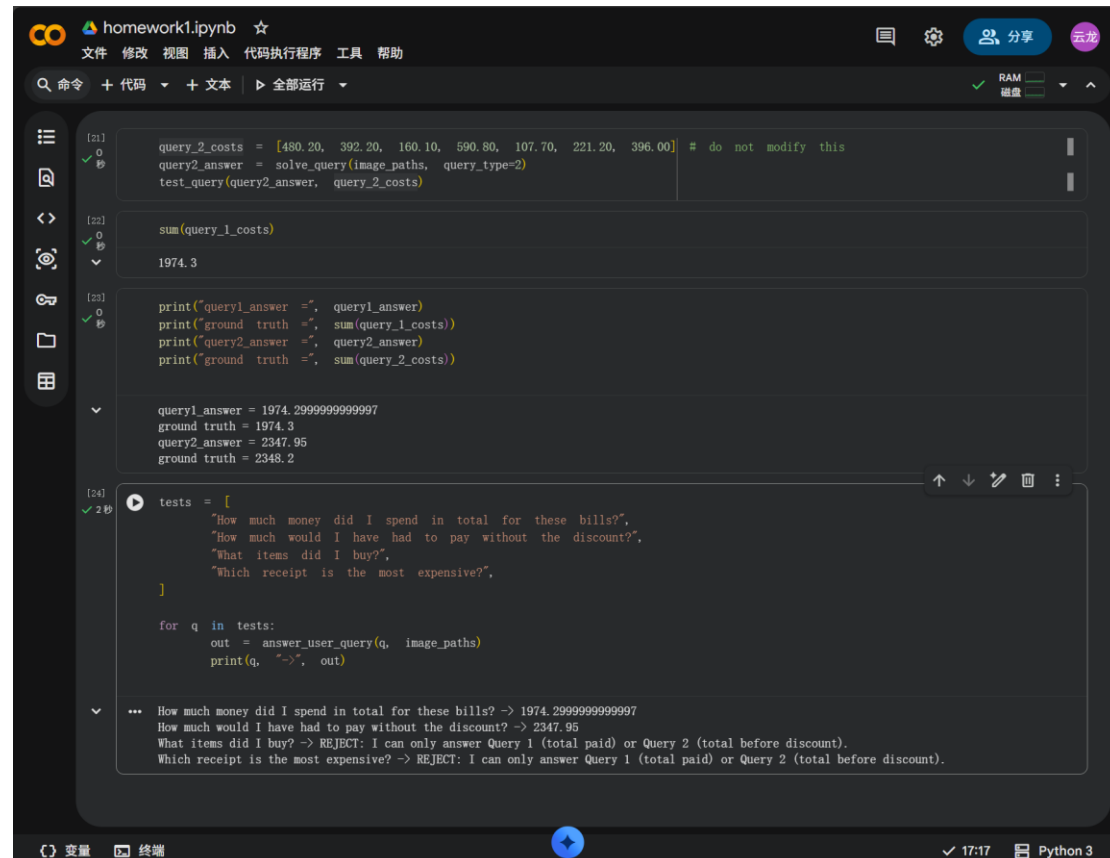
service unavailability occurs, reducing the risk of failure due to API overload.

## Testing Procedure

The system is tested exclusively through the unified interface

`answer_user_query(user_query, image_paths)`.

## Results



The screenshot shows a Jupyter Notebook titled 'homework1.ipynb'. The code is written in Python and includes the following:

```
[21] query_2_costs = [480.20, 392.20, 160.10, 590.80, 107.70, 221.20, 396.00] # do not modify this
      query2_answer = solve_query(image_paths, query_type=2)
      test_query(query2_answer, query_2_costs)

[22] sum(query_1_costs)

1974.3

[23] print("query1_answer =", query1_answer)
      print("ground truth =", sum(query_1_costs))
      print("query2_answer =", query2_answer)
      print("ground truth =", sum(query_2_costs))

query1_answer = 1974.2999999999997
ground truth = 1974.3
query2_answer = 2347.95
ground truth = 2348.2

[24] tests = [
      "How much money did I spend in total for these bills?",
      "How much would I have had to pay without the discount?",
      "What items did I buy?",
      "Which receipt is the most expensive?",
    ]

    for q in tests:
        out = answer_user_query(q, image_paths)
        print(q, "->", out)

... How much money did I spend in total for these bills? -> 1974.2999999999997
How much would I have had to pay without the discount? -> 2347.95
What items did I buy? -> REJECT: I can only answer Query 1 (total paid) or Query 2 (total before discount).
Which receipt is the most expensive? -> REJECT: I can only answer Query 1 (total paid) or Query 2 (total before discount).
```

The system was evaluated using multiple receipt images and both in-domain and out-of-domain queries.

For Query 1 (total amount actually paid), the system returned a value of 1974.30, which matches the ground truth within the allowed error tolerance.

For Query 2 (total amount before discount), the system returned 2347.95, which is very close to the ground truth value of 2348.20 and also satisfies the evaluation criteria.

Additional tests with out-of-domain queries show that the system consistently returns an explicit rejection rather than a numerical output.

Overall, the results demonstrate that the system can accurately handle both target queries and reliably reject non-relevant queries, satisfying the requirements of the assignment.