# Milestone 3 (Week 9): Core Model Development & Integration

Group 5: Aliyyah Jackhan, Mohammed Aadil and Jonathan Chacko
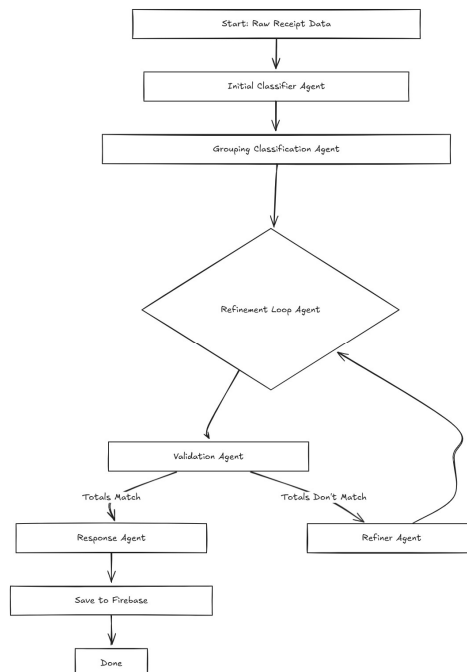
## 1. Core Machine Learning Model(s) / AI Components

- **Vertex Document AI (Pre-Trained Model):**

    o Vertex has Document AI which Provides Us with Option of Training or Selecting a Pre-Trained Model to Use.

    o Since Expense Bills are available under the Pre-Trained Model we selected it and have the Option to Fine Tune it in the Future.

- **Classification**:

    o Implemented a highly structured ADK (Agent Development Kit) model pipeline for receipt item classification.

    o Pipeline uses strict, well-defined input and output formats for reliable and explainable classification results.

    o Switched from the previous Ollama-based LLM approach (Milestone 2) to the new multi-agent ADK pipeline for better control, structure, and validation.



- **Regression Model (Upcoming)**:

    o A regression model for spend prediction is planned and will be integrated in the next stage (along with the Dashboard).

## 2. System Integration

- **Integration Demo**:

  - The gcp_adk_classification.py module demonstrates seamless integration between the main application and the ADK classification pipeline.

  - End-to-end flow: Discord bot → Flask API (main_api.py) → GCP Document AI OCR → ADK agent pipeline → Firebase storage.

- **Upcoming Dashboard**:

  - A user dashboard to visualize classified and predicted data is planned for the next milestone.
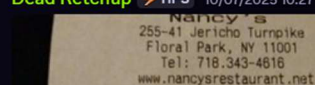
- **Screenshots from the Demo of the Deployed system so far:**

## Agent Development Kit

receipt_classifier ▾

Trace | Events | State | Artifacts | Sessions | Eval

**Conversations**    ✓ Events | Trace

0 text:{ "classified": [ { "item": "gl Imp White", "quan...

1 text:{ "grouped": [ { "category": "Others", "items": [...

2 functionCall:calculate_final_total

3 functionResponse:calculate_final_total

4 functionCall:exit_function

5 functionResponse:exit_function

6 functionCall:save_to_firebase

7 functionResponse:save_to_firebase

8 text:```json { "result": "success", "message": "Data s...

SESSION ID 484083a9-cfd9-41b2-871e-bb459a54997c    Token Streaming  + New Session  🗑  ⬇

⚡ save_to_firebase

✓ save_to_firebase

```
{
  "result": "success",
  "message": "Data saved to Firebase.",
  "data": [
    {
      "category": "Others",
      "items": [
        "gl Imp White"
      ],
      "total_price": "7.75"
    },
    {
      "category": "Fast Food",
      "items": [
        "Blue Moon Tap",
        "Mozzarella&Tomato",
        "Pork Quesadilla",
        "Fren Onion Soup",
        "Pork Chop",
        "Hanger Sizzle"
      ],
      "total_price": "78.75"
    }
  ]
}
```

Type a Message...

Spendify ▾ | Cloud Firestore

Data | Rules | Indexes | Disaster Recovery | Usage | 🧩 Extensions

💲 Avoid surprises on your bill by creating a budget to monitor incurred charges   Create budget   ✕

Panel view | Query builder

/DATA/SUMMARISED_DATA/2025-07-10/af621247-622a-48e1-b18b-85635c798271

**SUMMARISED_DATA**    **2025-07-10**    **af621247-622a-48e1-b18b-85635c798271**

+ Start collection    + Add document    + Start collection
  2025-06-05           af621247-622a-48e1-b18b-85635c798271  >    + Add field
  **2025-07-10**   >                                         ▾ categories
  2025-07-13                                                    ▾ 0
                                                                  category: "Others"
+ Add field                                                     ▾ items
                                                                    0  "gl Imp White"
                                                                    1  "Blue Moon Tap"
                                                                  total_price: "13.75"
                                                                ▾ 1
                                                                  category: "Fast Food"
                                                                ▾ items
                                                                    0  "Mozzarella&Tomato"
This document does not exist, it will not appear in queries or snapshots.   1  "Pork Quesadilla"
Learn more                                                          2  "Fren Onion Soup"
                                                                    3  "Pork Chop"
                                                                    4  "Hanger Sizzle"
                                                                  total_price: "72.75"
                                                              final_total: 86.5
                                                              notes: "Computed total matches the net amount (excluding tax)."
                                                              status: 1
                                                              summary: "Final receipt classification summary."

📍 Database location: nam5

- The full pipeline was tested end-to-end:

  - Discord Bot uploads a receipt image

  - API receives and processes the image, runs OCR

  - Data sent to ADK classification pipeline (Initial Classifier → Grouping Agent → Validation/Refinement Loop → Final Response Agent)

  - Classification output is validated and stored in Firebase

- **Result:** API endpoint returns predictions successfully.

## 3. Initial Deployment Steps

Deployment guides and step-by-step instructions are provided in the following repository links. For each major component, use the commands below for deployment:

a. **ADK Classification Agent**

- Guide: deploy-adk.md

- **Cloud Run Deployment Command:**

```
adk deploy cloud_run \
  --project=$GOOGLE_CLOUD_PROJECT \
  --region=$GOOGLE_CLOUD_LOCATION \
  --service_name=$SERVICE_NAME \
  --app_name=$SERVICE_NAME \
  --port=8080 \
  --log_level=info \
  --with_ui \
  ./receipt_classifier
```

b. **Flask API**

- Guide: deploy-api.md

- **Build & Deploy to Cloud Run:**

```
gcloud builds submit --tag gcr.io/$GOOGLE_CLOUD_PROJECT/$SERVICE_NAME .
gcloud run deploy $SERVICE_NAME \
  --image gcr.io/$GOOGLE_CLOUD_PROJECT/$SERVICE_NAME \
  --platform managed \
  --region $GOOGLE_CLOUD_LOCATION \
  --allow-unauthenticated
```

   c. **Discord Bot**

- Guide: [deploy-bot.md](deploy-bot.md)

- **VM Deployment Command:**

```
ssh -i ~/.ssh/gcp_key user@VM_IP
tmux new -s discordbot
python3 bot.py
```

Each document contains Dockerfile usage, environment variable setup, GCP Cloud Run deployment, and troubleshooting tips specific to each core component.

## 4. Progress, Planning & Next Steps

**From Milestone 2 → Milestone 3**

- Major change: Replaced Ollama/LLM-based classification with the robust ADK model pipeline.

- Improved reliability, structure, and validation by enforcing strict input/output schemas and a multi-agent review/refine process.

**Detailed Task Breakdown (Next 3 Weeks, by Role)**

- **Regression model implementation & integration (Aliyyah):**

    o Predict user spend (time series/regression)

    o Integrate into existing backend pipeline

- **Dashboard development (Aadil):**

    o Build initial UI for receipt and prediction visualization

    o Connect dashboard to Firebase/API

- **Deployment & system integration (Jonathan):**

    o Deployment of Regression Model and Dashboard

    o Full system end-to-end testing

## 5. Challenges Faced to Achieve Milestone 3

- Ollama proved inaccurate and unreliable for classification, which led us to adopt GCP ADK as the primary classification engine.

- Swapping from Ollama to ADK allows us to access the larger computational power of GCP and benefit from a more scalable, production-grade environment.

- Integrating the highly structured ADK model pipeline with strict input/output, which required reworking the entire classification workflow.

- GCP ADK is an amazing tool for building agent pipelines, but it lacks comprehensive documentation and real-world examples, which significantly slowed down development.

- Deployment was difficult for GCP ADK for the same reasons as above, with a lot of trial and error required to achieve a working cloud setup.

- Ensuring compatibility and smooth data handoff between Discord Bot, Flask API, GCP Document AI, and ADK pipeline.

- Debugging Firestore data serialization and managing Firestore schema changes.

- Handling OCR inconsistencies in receipt formats from GCP Document AI.

- Adapting deployment strategies (switching from Ollama to ADK), requiring updates to containerization and cloud setup.

## 6. Remaining Challenges

- **Regression Model:** Need to design and validate regression predictions; may need additional user data.

- **Dashboard:** Ensuring robust connection between dashboard and backend, and handling real-time updates.

- **Validation:** Ensuring ground-truth data for accurate model validation (classification & regression).

- **Performance:** Scaling API/bot for multiple concurrent users; latency testing on cloud.