# AI Expense Categorizer - Technical Documentation

## 1. Categorization Logic

The categorization system follows a hybrid approach combining rule-based mapping and LLM-based classification to ensure both consistency and flexibility.

### Rule-Based Mapping

Known merchant keywords are directly mapped to categories (e.g., SWIGGY → Meals, AWS → Software).

This ensures deterministic results, faster processing, and reduced LLM calls.

Rule-based categorization also improves explainability and auditability.

### LLM-Based Classification

If no rule matches, the LLM classifies the transaction using its description context.

The LLM output is validated against a strict JSON schema and enforced category list.

Fallback to 'Other' occurs if the LLM output is invalid.

## 2. Prompt Structure

The LLM prompt enforces structured output and deterministic behavior.

Temperature is set to 0 to ensure consistent categorization.

The prompt instructs the LLM to return JSON with category, confidence score, and reason.

Category validation ensures outputs remain within predefined categories.

## 3. Anomaly Detection Approach

### Statistical High Amount Detection

Uses Median Absolute Deviation (MAD) to detect unusually high expenses.

Transactions above threshold (median + 4*MAD) are flagged as statistical anomalies.

MAD provides robustness against skewed financial distributions.

**Duplicate Detection**

Transactions with identical date, amount, and normalized description are flagged as possible duplicates.

**Category Outlier Detection**

MAD is applied within each category to identify abnormal spending patterns.

Example: unusually high meal expense flagged within Meals category.

## 4. Limitations

Small datasets may produce false anomaly flags due to limited statistical context.

Currency normalization assumes single-currency datasets.

Anomaly detection is dataset-relative without historical baseline.

Rule-based mapping requires manual updates.

PDF reporting focuses on summary-level insights.