

## **Contents**

<b>1. Business Objective &amp; Background</b> 1.1 Case Background 1.2 Business Case
<b>2. Project Management</b> 2.1 Project Roles & Responsibilities 2.2 Project Plan & Timeline 2.3 Risk Assessments
<b>3. Methodology</b> 3.1 Overview of the Approach 3.2 Data Assessment & Preparation 3.3 Data Mapping & Analysis 3.4 Data representation Architecture of the Model
<b>4. Detailed Process Documentation</b> 4.1 Data Cleaning and Preparation 4.2 Standardization of Data 4.3 Classification of Data 4.4 Partition of Data 4.5 Dictionary Method to classify the names between the Garage and the Surveyor 4.6 Mapping 4.7 Part Recommendation Architecture
<b>5. Outputs, Monitoring, and Evaluation</b> 5.1 Mapping Summary Reports 5.2 Recommendation System Outputs 5.3 Monitoring the Recommendation system
<b>6. Supporting Documentation &amp; Appendices</b> 6.1 Flow Diagrams & Architectural Schematics 6.2 References & Citations

# **1. Business Objective**

## **1.1 Case Background**

**AutoSure Insurance** has built a strong reputation for its quick and efficient claim settlements. However, the company faces a significant challenge due to inconsistencies in damage assessment. When an accident occurs, a surveyor manually records the damaged parts and a garage separately provides its own list. This process often results in discrepancies—for instance, a surveyor might document a “Left Door” while the garage records a “Driver-Side Door.” Such variations complicate the claim verification process, delay settlements, and open up potential avenues for fraudulent claims.

## **1.2 Business Case**

To enhance the accuracy of damage assessments and streamline claim verification, AutoSure Insurance must standardize the nomenclature used in both surveyor and garage reports. By integrating historical claim data, the company can develop a dynamic recommendation system that suggests the top five most commonly associated damaged parts in real time. This system not only refines the mapping of parts between the two datasets but also accelerates the claim processing workflow, reduces operational costs, and minimizes the risk of fraud—ultimately preserving customer trust and upholding operational excellence.

**Our Project was divided into following 4 parts:-**

- **Standardizing Terminology:** Reconciling variations in naming (e.g., "Left Door" vs. "Driver-Side Door") to ensure that both records refer to the same physical component.
- **Establishing Correspondence:** Creating a consistent and accurate relationship between the parts noted by the surveyor and those listed by the garage, thereby enabling precise matching for claim verification.
- **Data Integration and Analysis:** Leveraging historical claim data and potential algorithmic techniques to detect patterns, reduce discrepancies, and minimize fraudulent claims.
- **Summary Reporting:** Generating mapping summary reports that detail the results of the matching process, including examples of successful part comparisons, frequency counts, and areas of inconsistency.

## **2 Project Management**

### **2.1 Project Roles & Responsibilities**

#### **2.1.1 ScriptDevs Team**

The following list describes the team involved in this project:-

Fullstack Developer	Rohit Rathod (Team Leader)
Fullstack Developer	Aman Jaiswal
Fullstack Developer	Bhavik Prajapati

## 2.2 Project Plan & Timeline

Phases	Activities	Hours
Study of the Approach	Planning., Research, forming a archhitecture	3 hours
Data assessment and Preparation	Finding the best algorithm and tools for solutions, and the studying the pain-point of the problem statement	5 hours
Data Integration and Mapping	Implementing data integration from the datasets, trained the model, and performing mapping.	8+ hours
Architecture of the Model	Finalizing the Architecture, so that the conclusion can be visualized conveniently.	2 hours

## 3. Methodology

### 3.1 OverView of the approach

#### Approach 1: String Similarity Matching (Levenshtein Distance / Fuzzy Matching)

##### Methodology:

- Compare **PARTDESCRIPTION** (garage) and **TEXT\_PARTS\_NAME** (surveyor) using **Levenshtein Distance**, **Jaccard Similarity**, or **FuzzyWuzzy**.
- Use **threshold-based matching** to identify similar parts based on textual similarity.
- **Downsides:**
  - High False Positives & Negatives** – Minor spelling errors or synonyms may cause mismatches.
  - Computationally Expensive** – Pairwise comparisons for large datasets slow down processing.
  - No Context Understanding** – Cannot differentiate between different parts with similar names (e.g., "Front Bumper" vs. "Rear Bumper").

#### Approach 2: Machine Learning-based Named Entity Recognition (NER) Model

##### Methodology:

- Train an **NER model** (e.g., using spaCy, BERT, or LSTM) to classify part names into standardized categories.
- Use supervised learning with **pre-labeled parts** from historical claims data to predict mappings.
- **Downsides:**
  - Requires Extensive Training Data** – Needs a large labeled dataset for accurate predictions.
  - Overfitting Risk** – New part names or variants may not be recognized if not in training data.
  - Complex Implementation** – Requires continuous model retraining for evolving part names.

## **Our Chosen Approach: Dictionary-based Part Mapping (Rule-Based Mapping with Synonym Dictionary)**

### **Methodology:**

- **Create a dictionary** of standard part names based on historical claims.
- **Map synonyms** manually (e.g., "Left Door" = "Driver-Side Door").
- Implement **rule-based fuzzy matching with predefined mappings**.
- **Advantages over other approaches:**
  - High Accuracy** – Eliminates false positives by using predefined synonyms.
  - Fast & Efficient** – Simple lookup mechanism rather than expensive pairwise comparisons.
  - Customizable & Scalable** – New synonyms and mappings can be manually updated as required.

## **3.2 Data Assessment and Preparation**

### **1. Data Loading and Exploration**

- Datasets (`garage_data.csv`, `surveyor_data.csv`) are loaded using `pandas.read_csv()`.
- Initial exploration included:
  - **Previewing Data:** Displaying the first few rows with `df.head()`.
  - **Structure Check:** Using `df.info()` to assess column names, data types, and missing values.
  - **Statistical Summary:** Viewing descriptive statistics with `df.describe()` for numeric data.

### **2. Data Cleaning**

#### **2.1 Data Type Adjustments**

- Converted columns to appropriate types:
  - `REFERENCE_NUM` & `NUM_CLAIM_NO`: `float64` → `int64`.
  - `CLAIMNO`: Converted to string and stripped of leading apostrophes.
  - `VEHICLE_MODEL_CODE`: `float64` → `int64`.

#### **2.2 Duplicates**

- Removed duplicate rows using `df.drop_duplicates(inplace=True)`.

#### **2.3 Missing Values**

- Handled missing values by filling them with "Unknown" using `df.fillna("Unknown", inplace=True)`.

#### **2.4 Scientific Notation**

- Ensured large numbers (like `REFERENCE_NUM`) were converted to integers, avoiding scientific notation.

## **3.3 Data Mapping & Analysis**

### **1. Data Loading**

- Loaded two cleaned datasets:
    - `surveydatashort.csv` (surveyor data)
    - `garageshortdata.csv` (garage data)
  - Ensured appropriate cleaning and validation from previous steps.
- 

### **2. Text Preprocessing**

To create consistency for mapping, textual fields (`TXT_PARTS_NAME`, `PARTDESCRIPTION`) were preprocessed:

- **Lowercasing:** Converted text to lowercase.
  - **Whitespace Normalization:** Removed extra spaces and trimmed text.
  - **Special Character Removal:** Removed punctuation, symbols, and unnecessary delimiters like `|`, `-`, etc.
  - **Synonym Normalization:**
    - Replaced common abbreviations (e.g., `rh` → `right`, `lh` → `left`).
  - **Output:**
    - `clean_key` (processed keys from surveyor data).
    - `clean_value` (processed values from garage data).
- 

### **3. Fuzzy Matching**

- Used `fuzzywuzzy` to match preprocessed keys (`clean_key`) with possible values (`clean_value`) in the datasets.
- **Logic:**
  - Extracted up to 5 candidates for each key.
  - Retained matches with a similarity score > 80.
- **Output:**
  - Initial mappings of likely part names between datasets.

### **4. Semantic Similarity Matching**



- Utilized **Sentence Transformers** (`paraphrase-MiniLM-L6-v2`) for contextual understanding of part descriptions.
- Steps:
  - Embedded `clean_key` and `clean_value` using sentence embeddings.
  - Calculated **cosine similarity** between keys and values.
  - Mapped the best match (highest similarity score) for each key.
- **Threshold:**
  - Similarity threshold set at **0.65** to ensure accurate matches, with room for adjustments during validation.

## 5. Combined Mapping

The combined approach (fuzzy + semantic):

- Handled both structural mismatches (special characters, abbreviations) and semantic variations (context/meaning).
- **Output:**
  - A mapping dictionary (`final_mapping`) aligning surveyor part names (`TEXT_PARTS_NAME`) with garage descriptions (`PARTDESCRIPTION`).

## 6. Saving Results

- Final unified mapping saved in `mapping.csv` for further review.
- Also saved intermediate results (e.g., fuzzy candidates for low-confidence mapping, similarity scores) where needed.

## 7. Synonym Management

- Introduced a custom dictionary (`synonyms.json`) to retain manually added common synonyms (e.g., `Left Door: Driver Side Door`).
- Functions created to:
  - Load synonyms (`load_dictionary`).
  - Add synonyms (`add_synonym`).
  - Save updates (`save_dictionary`).

## 8. Outputs

1. **Mappings:**
  - `mapping.csv`: Final alignment of part descriptions.
2. **Utility Data:**
  - Synonym dictionary enabled reusability and scalability of mappings.

# Detailed Process Documentation

## 4.1 Data cleaning and Preparation

### 1. Data Loading

- Loaded datasets from CSV files:
    - **Surveyor data:** `surveydatashort.csv`
    - **Garage data:** `garageshortdata.csv`
    - Used `pandas.read_csv()` for seamless loading and manipulation.
  - Initial Dataset Inspection:
    - Data structure inspected with methods like:
      - `df.head()` - To view the first few rows of data.
      - `df.info()` - To check columns, data types, and null values.
      - `Df.describe`      `()` - To summarize statistics of numerical fields.
- 

### 2. Text Preprocessing

To standardize and clean textual data for analysis and mappings:

#### Text Cleaning and Normalization

- **Objective:** Ensure consistency and improve matching between part descriptions in `TXT_PARTS_NAME` (surveyor data) and `PARTDESCRIPTION` (garage data).
- **Steps Applied:**
  - **Lowercasing:** Converted all text to lowercase for uniformity.
  - **Special Character Removal:** Removed characters like `|`, `-`, and other non-alphanumeric symbols using regex.
  - **Whitespace Management:** Removed extra spaces and trimmed strings.
  - **Synonym Normalization:**
    - Replaced common abbreviations:
      - Examples: `rh/r` → `right`, `lh/l` → `left`.
  - **Regex Filtering:**
    - Removed unwanted characters and retained only alphabetic/numeric terms.
    - Example: `re.sub(r'^a-z0-9\s','', text)`.
- **Implementation:**
  - Created a reusable `preprocess()` function:
    - Applied to `TXT_PARTS_NAME` -> `clean_key` (surveyor data).
    - Applied to `PARTDESCRIPTION` -> `clean_value` (garage data).

---

### 3. Handling Duplicates and Missing Data

- Removed duplicate rows using `df.drop_duplicates()`.
  - Handled missing values as follows:
    - For essential fields, replaced missing values with "Unknown".
    - Ensured there were no null entries in `TXT_PARTS_NAME` or `PARTDESCRIPTION`. Missing values in these columns could break mappings.
- 

### 4. Data Type Adjustments

- Converted numeric or categorical fields where applicable:
    - Fields like `REFERENCE_NUM` and `NUM_CLAIM_NO` were converted to `int64` for consistency.
- 

### 5. Mapping Preparation for Analysis

#### Preprocessed Data Used:

- `clean_key`: Cleaned part descriptions from surveyor data.
- `clean_value`: Cleaned garage descriptions.

#### Steps:

1. Created mappings:
    - Mapped cleaned keys and cleaned values to their original descriptions using Python dictionaries (`dict`).
    - Enabled quick lookup during the mapping stage.
    - Example:
      - `original_key_mapping`: Map `clean_key` to its original `TXT_PARTS_NAME`.
      - `original_value_mapping`: Map `clean_value` to its original `PARTDESCRIPTION`.
-

## 6. Processed Data Output

- Saved cleaned datasets (`clean_key` and `clean_value`) for downstream analysis.
  - Outputs:
    - Cleaned and preprocessed surveyor dataset saved as `processed_surveyor_data.csv`.
    - Cleaned garage dataset saved with normalized part descriptions as `processed_garage_data.csv`.
- 

## 7. Tools and Libraries

- **Pandas:**
  - For data loading, cleaning, and manipulation.
- **Regex (`re`):**
  - For string manipulation and cleaning (e.g., removing unwanted characters).
- **File Handling:**
  - Saved cleaned datasets to CSV for reuse in mapping and similarity modeling.

## 4.2 STANDARDIZATION OF DATA

### 1. Objective of Standardization

The goal of data standardization is to ensure consistent formatting and representation of textual and numerical data across datasets. This facilitates accurate mapping, comparison, and analysis between the surveyor and garage data.

---

### 2. Data Loading

- Loaded raw surveyor and garage datasets:
    - Surveyor Data: Dataset of part names used in surveys (`surveydatashort.csv`).
    - Garage Data: Dataset of parts listed in garage inventory (`garageshortdata.csv`).
  - Functionality provided by the `pandas.read_csv()` module for efficient handling of raw CSV files.
- 

### 3. Preprocessing for Text Standardization

Text standardization focused on creating consistent, comparable text fields for part descriptions:

- `TXT_PARTS_NAME` from the surveyor dataset.
- `PARTDESCRIPTION` from the garage dataset.

#### 3.1 Cleaning Techniques

A reusable `preprocess()` function was implemented to clean part names and descriptions:

1. Lowercasing:
  - Converted all text to lowercase to eliminate case-sensitivity issues.

- Example: **Front Door** → **front door**.
- 2. Removing Special Characters:
  - Removed unwanted characters such as pipes (**|**), dashes (**-**), and non-alphanumeric characters using regex.
  - Example: **Front-Door | Assy** → **front door assy**.
- 3. Whitespace Normalization:
  - Stripped extra spaces and ensured text had a consistent single-space delimiter.
  - Example: **front door** → **front door**.
- 4. Synonym Expansion:
  - Common abbreviations were normalized into full forms:
    - Examples:
      - **rh/r** → **right**.
      - **lh/l** → **left**.
      - **assy** → **assembly**.

### 3.2 Standardized Fields

- After applying the **preprocess()** function:
  - Surveyor data field **TXT\_PARTS\_NAME** was converted to **clean\_key**.
  - Garage data field **PARTDESCRIPTION** was converted to **clean\_value**.

## 4. Handling Duplicates and Missing Values

- Duplicates Detection:
  - Ensured no duplicates existed in standardized datasets using **df.drop\_duplicates()**.
- Missing Values:
  - Filled missing cells in critical fields (**TXT\_PARTS\_NAME**, **PARTDESCRIPTION**) with **"Unknown"**, ensuring no null entries disrupted the downstream analysis.

## 5. Numeric Standardization

- Adjusted numeric fields where needed:

- For example, converted columns like **REFERENCE\_NUM** and **NUM\_CLAIM\_NO** to **int64** format.
- Removed scientific notation for large numbers to improve readability and usability.

## 6. Applying Synonyms Dictionary

- Introduced a synonym dictionary (**synonyms.json**) to store and standardize variations of part names:
  - Example:
    - Part Name: **Left Door**.
    - Synonyms: **Driver-Side Door, L Door**.
- Functions:
  - **load\_dictionary()**: Load existing synonym mappings.
  - **add\_synonym(part\_name, synonym)**: Allow manual addition of synonyms to improve standardization.

## 7. Outputs of Standardized Data

- Surveyor Dataset:
  - Output field: **clean\_key**.
  - Saved standardized surveyor data as **standardized\_surveyor\_data.csv**.
- Garage Dataset:
  - Output field: **clean\_value**.
  - Saved standardized garage descriptions as **standardized\_garage\_data.csv**.

These output files were used for further mapping and semantic similarity modeling.

## 8. Tools and Techniques Used

- **Libraries:**
  - **pandas**: For data manipulation.
  - **re**: For text cleaning and regular expressions.
  - **json**: For synonym dictionary creation and updates.
- **Custom Functions:**
  - **preprocess()**: Centralized cleaning logic for both datasets.



## **4.3 Classification of Data**

### **1. Objective of Data Classification**

The purpose of data classification was to:

- Categorize parts descriptions into meaningful groups.
  - Enable systematic comparisons between surveyor data (`TXT_PARTS_NAME`) and garage data (`PARTDESCRIPTION`).
  - Support downstream analysis by providing clear and consistent labels.
- 

### **2. Data Preparation**

**Input Datasets:**

1. **Surveyor Data:**
    - Columns like `TXT_PARTS_NAME` were cleaned and processed to extract meaningful keys for classification.
  2. **Garage Data:**
    - Columns such as `PARTDESCRIPTION` served as the basis for comparison and classification.
- 

### **3. Cleaning and Preprocessing**

Before classification, preprocessing steps were applied to standardize the datasets:

1. **Text Standardization:**
    - Cleaned data by removing special characters, expanding abbreviations, and normalizing synonyms using a reusable `preprocess()` function.
  2. **Synonym Handling:**
    - Used a synonym dictionary (`synonyms.json`) to map variations of part names (e.g., `Driver-Side Door` = `Left Door`).
  3. **Output:**
    - Generated standardized fields:
      - `clean_key` (from surveyor data).
      - `clean_value` (from garage data).
- 

### **4. Classification Logic**

The classification of data was achieved through the following steps:

#### 4.1. Categorization into Common Part Types

- Grouped parts based on common terminology (e.g., "door", "bumper", "mirror").
- **Examples:**
  - "Left Door" and "Driver-Side Door" were categorized under Door.
  - "Front Bumper" and "Rear Bumper" were labeled as Bumper.

#### 4.2. Semantic Analysis for Context-Based Classification

- Used **Sentence Transformers** to compute semantic similarity between part descriptions:
  - **Model:** `paraphrase-MiniLM-L6-v2` to generate embeddings for part descriptions.
  - Compared `clean_key` and `clean_value` embeddings using cosine similarity.
  - **Threshold:**
    - Items with a similarity score > 0.65 were grouped into the same category.
    - Adjusted the threshold based on validation results.

#### 4.3. Rule-Based Classification for Edge Cases

- Applied a series of pattern-matching rules using regex to handle parts with positional terms or specific abbreviations:
  - Example Regex Rules:
    - `r|rh|right` → Right.
    - `lh|left` → Left.
- Any unclassified parts were labeled as `Miscellaneous` for manual review.

#### 4.4. Fuzzy Matching for Low-Confidence Items

- Integrated `fuzzywuzzy` for parts that did not meet the similarity threshold:
  - Matched `clean_key` with `clean_value` descriptions.
  - Fuzzy match scores above 75 were considered valid and classified accordingly.

---

## 5. Classification Outputs

### 5.1. Generated Categories

- Created high-level categories, such as:
  - Door (e.g., Left Door, Right Door).
  - Bumper (e.g., Front Bumper, Rear Bumper).

- **Mirror** (e.g., **Side Mirror**, **Rearview Mirror**).
- **Miscellaneous** (for unclassified or ambiguous items).

## 5.2. Mapped Part Descriptions

- Mapped surveyor part descriptions (**TXT\_PARTS\_NAME**) to garage part descriptions (**PARTDESCRIPTION**) within the identified categories.

## 5.3. Structured Output

- Outputs were:
  - **Classified Mapping File:**
    - Fields: **TXT\_PARTS\_NAME**, **PARTDESCRIPTION**, **Category**.
    - Saved as **classified\_data.csv**.
  - **Low Confidence Matches:**
    - For manual validation and reclassification.

## 6. Tools and Techniques Used

- **Libraries:**
  - **pandas**: For data handling and transformations.
  - **re**: For implementing regex-based classification rules.
  - **fuzzywuzzy**: For fuzzy matching low-confidence items.
  - **Sentence Transformers**:
    - For semantic similarity and context-based grouping.
- **Data Structures:**
  - **Dictionaries:**
    - Synonym dictionary for consistent labeling.
    - Mapping dictionary for category assignments.

## 7. Manual Validation

- Provided a mechanism for reviewing edge cases and validating low-confidence classifications:
  - Categorized all unclassified or ambiguous items under **Miscellaneous** for manual inspection.
  - Updated the synonym dictionary (**synonyms.json**) iteratively based on manual feedback.

## **4.4 Partition of Data**

### **1. Objective of Data Partitioning**

The goal of partitioning data was to:

- Separate data based on predefined criteria for easier processing and analysis.
  - Prepare subsets for training, validation, and testing for machine learning or manual use cases.
  - Create meaningful groups within datasets (e.g., by category, region, or semantic similarity).
- 

### **2. Data Preparation**

**Input Datasets:**

- Two primary datasets:
  - **Surveyor Data** containing columns like **TXT\_PARTS\_NAME**.
  - **Garage Data** with details like **PARTDESCRIPTION**.
- Standardized and cleaned versions of the datasets were used (outputs from the preprocessing module).

**Preprocessing Recap:**

Before partitioning began, the data was:

- 1. Cleaned and Standardized:**
  - Text was normalized through processes such as lowercasing, removing special characters, expanding abbreviations, and handling missing values.
- 2. Synonyms Mapped:**
  - Synonyms were replaced with consistent values to ensure logical overlaps during partitioning.

### 3. Criteria for Partitioning

The partitioning logic was based on the specific requirements of the analysis. Several approaches were used:

#### 3.1 Category-Based Partitioning

- Data was grouped into **functional categories**:
  - Examples: **Door**, **Bumper**, **Mirror**, etc.
- Categorization Process:
  - Keywords from **TXT\_PARTS\_NAME** and **PARTDESCRIPTION** (e.g., "door", "bumper") were used as identifiers.
  - Regex and fuzzy matching tools (**fuzzywuzzy**) were used for broader matches.

#### 3.2 Semantic Similarity Partitioning

- Used **SentenceTransformers** and cosine similarity scores to group data based on semantic meaning.
- **Threshold Criteria**:
  - Items with a high similarity score ( $> 0.7$ ) were grouped into one partition.
  - Manual validation was performed in the case of low-scoring groups.

#### 3.3 Random Partitioning for Training, Validation, Testing

- Created randomized splits to be used for experiments:
  - **Training Set**: 70%.
  - **Validation Set**: 15%.
  - **Testing Set**: 15%.
- Ensured equal representation of categories within each set.

#### 3.4 Size-Based Partitioning

- Partitioned data into subsets based on record size:
  - **Small Set**: Only critical sample parts (**N** records).
  - **Full Set**: Entire dataset for comprehensive analysis.

## 4. Partitioning Techniques

### 4.1 Functional Partitioning Using Categories

1. **Define Logic:**
  - Utilized a mapping dictionary with category keywords (e.g., "door" -> `Door` category).
2. **Implementation:**
  - Mapped the cleaned `TXT_PARTS_NAME` and `PARTDESCRIPTION` columns to predefined categories.
  - Stored partitions in separate dataframes.

### 4.2 Semantic Similarity-Based Grouping

1. **Generate Embeddings:**
  - Encoded `clean_key` and `clean_value` using `SentenceTransformer`.
2. **Calculate Similarity:**
  - Computed pairwise similarity using cosine similarity metric.
3. **Cluster by Threshold:**
  - Grouped items sharing a similarity score  $> 0.7$ .

### 4.3 Random Splitting

- Used the `pandas sample()` function to randomly allocate subsets of the data into three partitions:
  - Training, Validation, and Testing.
- Imposed constraints to ensure proportional representation of categories:
  - Example: If `Door` represented 30% of the dataset, ensured each subset maintained the percentage.

### 4.4 Index-Based Partitioning

- Partitioned data for specific needs like:
  - First `N` rows for sampling.
  - Segments based on conditions (e.g., rows where descriptions contain "front" vs "rear").

## 5. Outputs of Partitioning

The partitioning operation yielded the following outputs:

### 5.1 Category-Based Partitions

- Files stored for each category:
  - Example:
    - `door_data.csv` for all **Door**-related records.
    - `bumper_data.csv` for all **Bumper**-related records.

### 5.2 Semantic Similarity Groups

- Generated clustered datasets based on semantic closeness:
  - Stored as `semantic_group_1.csv`, `semantic_group_2.csv`, etc.

### 5.3 Train/Test/Validation Sets

- Data split into:
  - `train_set.csv`
  - `validation_set.csv`
  - `test_set.csv`

### 5.4 Sampling Outputs

- Extracted additional subsets:
    - `sample_set_small.csv`: Contained high-priority records.
- 

## 6. Tools and Methods Used

**Libraries:**

1. **Pandas:**
  - For slicing, filtering, and partitioning rows into multiple dataframes.
2. **SentenceTransformers:**
  - For generating embeddings and performing semantic similarity computations.
3. **Fuzzywuzzy:**

- For approximate matching in category identification.
4. **numpy:**
- For handling matrix representations and similarity operations.

#### **Partitioning Methods:**

1. **Regex Matching:**
    - Tagged rows matching specific text patterns (e.g., `r|right` → `Right Partition`).
  2. **Cosine Similarity:**
    - Calculated similarity scores between textual descriptions for clustering.
  3. **Randomized Sampling:**
    - Ensured unbiased separation into training, validation, and testing datasets.
- 

#### **7. Enhancements Made**

- **Dynamic Threshold:**
    - Introduced an adjustable threshold for similarity scores to accommodate varying dataset quality.
  - **Category Dictionary:**
    - Enhanced synonym dictionary for more accurate category identification.
  - **Automated Validation:**
    - Generated reports validating the proportion distribution of partitions (e.g., category balance in train/test sets).
- 

#### **8. Limitations and Manual Overrides**

- Items with ambiguous descriptions or low similarity scores were grouped into a placeholder partition (`Miscellaneous`) for manual inspection.
- Corrected misclassifications and re-mapped synonyms iteratively based on feedback.



## 4.5 Dictionary Method to classify the names between the Garage and the Surveyor

### 1. Objective of the Dictionary Method

The dictionary-based approach was implemented to:

- Establish a consistent mapping of part names between **Surveyor Data** (**TXT\_PARTS\_NAME**) and **Garage Data** (**PARTDESCRIPTION**).
  - Address inconsistencies by using synonyms and semantic similarity.
  - Minimize manual intervention by automating name classification.
- 

### 2. Data Preparation

Input Datasets:

1. **Surveyor Data:**
  - Contains part names in **TXT\_PARTS\_NAME**.
2. **Garage Data:**
  - Contains part descriptions in **PARTDESCRIPTION**.

Preprocessing Steps:

Before applying the dictionary method, data was cleaned and standardized:

1. **Text Cleaning:**
  - Removed special characters, expanded abbreviations, and normalized positional terms (e.g., "LH" → "Left").
  - Lowercased all text for uniformity.
2. **Synonym Mapping:**
  - Developed a synonym dictionary (**synonyms.json**) to handle equivalent terms (e.g., "Driver-Side Door" = "Left Door").
  - Example:
    - "RH Mirror" → "Right Mirror"

### 3. Dictionary Approach for Classification

The dictionary method primarily involved creating, enhancing, and utilizing a custom dictionary to map names effectively. Key steps are outlined below:

#### 3.1 Building the Dictionary

A dictionary was built to store mappings of standardized and original names:

- **Data Sources:**
  - Part names from `TEXT_PARTS_NAME` (Surveyor) were classified and standardized.
  - Part descriptions from `PARTDESCRIPTION` (Garage) were used to enrich the mapping.
- **Structure:**
  - The dictionary format was:

```
{  
  
  "Left Door": ["Driver-Side Door", "LH Door", "Door Left"],  
  
  "Right Mirror": ["RH Mirror", "Passenger Mirror"]  
  
}
```

- The first key represents the standardized name while the list contains all its synonyms.

#### 3.2 Synonym Mapping

- Utilized the `synonyms.json` file to map similar or related terms:
  - Process:
    - For each `TEXT_PARTS_NAME`, lookup was conducted in the synonym dictionary to find matches.
    - If a direct match was found, it was added to the mapping immediately.

- If no match was found, advanced techniques like fuzzy matching and semantic similarity were applied.
- This allowed handling variations in part names (e.g., "Frnt Bmpr" → "Front Bumper").

### 3.3 Automating Classification

The classification process involved the following methods:

#### 1. Exact Matching with Dictionary:

- Checked for an exact match between `TXT_PARTS_NAME` and `PARTDESCRIPTION` after applying preprocessing and synonym mapping.
- If a match existed, it was stored directly in the mapping dictionary.

#### 2. Fuzzy Matching for Approximate Matches:

- Used `fuzzywuzzy` to compare cleaned part names and descriptions.
- Matching Criteria:
  - A fuzzy match score > 75 was deemed acceptable for classification.

#### 3. Semantic Similarity for Difficult Cases:

- Applied `SentenceTransformers` to compare the semantic meaning of part names:
  - Generated embeddings for `TXT_PARTS_NAME` and `PARTDESCRIPTION`.
  - Computed similarity scores using cosine similarity.
- Matching Threshold:
  - A similarity score > 0.65 indicated a valid match.
- This ensured context-based mapping for descriptions with word ordering or phrasing differences.

### 4. Enhancements to the Method

#### 4.1 Synonym Dictionary Development

- Enhanced over time based on feedback and manual reviews.
- New synonyms were added dynamically:

- Example: "Left Wing Mirror" → Added to the synonym list for Left Mirror.
- Provided a reusable `add_synonym()` function to update the dictionary programmatically.

## 4.2 Confidence Scores

- Introduced a scoring system to indicate the confidence in each mapping:
  - **High Confidence:** Direct dictionary or fuzzy match above 85.
  - **Medium Confidence:** Semantic similarity score > 65.
  - **Low Confidence:** Fuzzy match below the threshold but included for review.

## 4.3 Error Handling

- Incorporated mechanisms to handle ambiguous or unsupported names:
  - Unmapped names were added to a separate list for manual validation.

## 5. Output of the Dictionary Method

The process resulted in a comprehensive mapping of part names between the Garage and Surveyor datasets.

### 5.1 Dictionary Output

- Generated an updated `synonyms.json` file with all mappings and synonyms.

### 5.2 Classification Results

- Produced a CSV file containing:
  - `TXT_PARTS_NAME` (Surveyor).
  - `PARTDESCRIPTION` (Garage).
  - `Mapped_Name` (Standardized name after classification).
  - `Confidence_Score`.

### 5.3 Reports for False Positives/Negatives

- Detailed reports were generated for cases where classification failed or confidence was low.

## 6. Tools and Techniques Used

### Libraries:

1. **pandas**: For data manipulation.
2. **re**: For regex-based name cleaning.
3. **json**: To load and save the synonym dictionary.
4. **fuzzywuzzy**: For approximate matching.
5. **SentenceTransformers**:
  - For generating embeddings and measuring semantic similarity.
6. **numpy**:
  - For handling similarity matrix calculations.

### Functions:

1. **preprocess()**:
  - Applied to normalize part names.
2. **add\_synonym()**:
  - Added new terms to the existing dictionary.
3. **cosine\_similarity()**:
  - Used for semantic matching.

## 7. Challenges and Resolutions

1. **Challenge**: Ambiguous or highly inconsistent names.
  - **Resolution**: Included a manual review process for edge cases.
2. **Challenge**: Continual updates to garage and surveyor datasets.
  - **Resolution**: Developed a framework to dynamically update the synonym dictionary.
3. **Challenge**: Low similarity for certain terms even with semantic models.
  - **Resolution**: Adjusted thresholds and included additional preprocessing.

## 8. Enhancements Added

- Developed a reusable dictionary framework that can:
  - Automatically classify names with high confidence.
  - Flag low-confidence mappings for manual review.
  - Dynamically scale with new synonyms and data updates.

## **4.6 PART MAPPING**

### **Detailed Process Documentation for Part Mapping**

This document provides a detailed explanation of the **Part Mapping** process that was implemented to establish a precise mapping of parts between the **Garage Data** and **Surveyor Data**. It outlines the steps undertaken, enhancements applied, and the final results.

---

#### **1. Objective of Part Mapping**

The primary objective of part mapping was to:

- Identify and map part names (`TEXT_PARTS_NAME`) from the Surveyor data with their corresponding part descriptions (`PARTDESCRIPTION`) in the Garage data.
  - Resolve inconsistencies and variations in naming conventions across datasets.
  - Enhance the mapping accuracy by incorporating semantic matching, fuzzy logic, and a synonym dictionary.
- 

#### **2. Data Preparation**

##### **Datasets Used**

1. **Surveyor Data:**
  - Contained parts in the column `TEXT_PARTS_NAME`.
2. **Garage Data:**
  - Included part descriptions in the column `PARTDESCRIPTION`.

##### **Preprocessing Steps**

Before performing the mapping, both datasets were cleaned and standardized:

1. **Normalization:**
  - Lowercased all text to ensure case-insensitive matching.
  - Removed special characters and unnecessary white spaces.
  - Standardized positional terms (e.g., `"RH"` → `"Right"`; `"LH"` → `"Left"`).
2. **Synonym Mapping:**
  - Used a pre-built synonym dictionary to handle equivalent terms (e.g., `"Left Door"` = `"Driver-Side Door"`).

##### **Cleaning Output:**

After preprocessing:

- Standardized fields were created:
    - **Surveyor Data:** `clean_key` (from `TXT_PARTS_NAME`).
    - **Garage Data:** `clean_value` (from `PARTDESCRIPTION`).
- 

### 3. Part Mapping Process

The part mapping process employed combinations of exact matches, fuzzy matching, and semantic similarity to ensure the highest possible accuracy.

---

#### 3.1 Exact Matching

1. **Implementation:**
  - Directly compared `clean_key` (Surveyor data) with `clean_value` (Garage data).
  - If an exact match was found, the mapping was saved immediately.
2. **Advantages:**
  - High confidence and accuracy.

**Outcome:** Approximately `X%` of the parts were matched directly using exact matching.

---

#### 3.2 Fuzzy Matching

1. **Library Used:**
  - `fuzzywuzzy` library was employed for approximate string matching.
2. **Process:**
  - Applied fuzzy matching between `clean_key` and `clean_value`.
  - Generated a confidence score (between 0–100) for each match.
  - Retained matches with a **fuzzy match score > 75**.
3. **Usage:**
  - This method captured minor differences in spelling, abbreviations, or formatting.
  - For example:
    - `Frnt Bmpr` → `Front Bumper`.
    - `Driver' Side Mirror` → `Driver-Side Mirror`.

**Outcome:** An additional `Y%` of matches were resolved using fuzzy matching.

#### 3.3 Semantic Similarity

### 1. Library Used:

- **SentenceTransformers (paraphrase-MiniLM-L6-v2)** for embedding text into high-dimensional vectors.

### 2. Process:

- Encoded the cleaned part names (**clean\_key** and **clean\_value**) into embeddings.
- Calculated the cosine similarity between embeddings.
- Used a similarity score:
  - **Threshold:** Matches with a semantic similarity score > **0.65** were retained.

### 3. Advantages:

- This method captured contextually similar matches.
- For example:
  - **Rearview Mirror** → **Back Mirror**.
  - **Right Door Panel** → **Passenger Door**.

**Outcome:** Semantic similarity accounted for **Z%** of additional matches.

---

## 3.4 Manual Review for Edge Cases

### 1. Low Confidence Items:

- Parts with low fuzzy scores (< 75) or low semantic similarity (< 0.65) were marked for manual review.

### 2. Examples:

- Ambiguous names like "**Bracket**" or "**Other**" required human intervention for accurate mapping.

### 3. Synonym Dictionary Updates:

- Any manually resolved mappings were added to the synonym dictionary for future use.
- 

## 4. Enhancements Made

The mapping process was enhanced to increase accuracy and scalability:

### 1. Synonym Dictionary Expansion:

- New terms and synonyms were added dynamically during manual review.
- Example:
  - **"Rear Door"** → **"Back Door"** and **"Trunk Door"**.

### 2. Dynamic Scoring:



- Adjusted thresholds for fuzzy matching and semantic similarity based on observed validation performance.
  - 3. **Combination of Methods:**
    - Leveraged all three methods (exact, fuzzy, semantic) sequentially to ensure maximum match coverage.
- 

## 5. Outputs of Part Mapping

The process resulted in the following outputs:

### 5.1 Mapped Parts

- Produced a consolidated mapping file:
  - **Fields:**
    - `TXT_PARTS_NAME` (Surveyor Name).
    - `PARTDESCRIPTION` (Mapped Garage Name).
    - `Confidence_Score` (Exact, Fuzzy, or Semantic match score).

### 5.2 Unmapped Parts

- Items that failed all matching criteria were stored in:
  - `Unmapped_Parts.csv`: For future investigation.

### 5.3 Updated Synonym Dictionary

- Generated an updated `synonyms.json` containing all synonyms and mappings.

## Summary of Results

- **Total Parts Mapped:** XYZ (percentage breakdown for Exact, Fuzzy, and Semantic matches).
  - **Total Unmapped Parts:** ABC.
- 

## 6. Tools and Techniques Used

### Libraries:

1. **pandas:**
  - For manipulating dataframes and processing mappings.
2. **re:**
  - Regular expressions for cleaning text and applying matching rules.
3. **fuzzywuzzy:**

- For approximate string matching.
- 4. **SentenceTransformers:**
  - To generate embeddings for semantic similarity.
- 5. **cosine\_similarity** (from **scikit-learn**):
  - For comparing similarity between text embeddings.

### Key Functions:

1. **preprocess(text):**
    - Normalized text by removing noise and standardizing naming conventions.
  2. **add\_synonym(part\_name, synonym):**
    - Dynamically added unresolved terms to the dictionary during manual review.
  3. **Pipeline Logic:**
    - Processed parts sequentially through exact, fuzzy, and semantic similarity pipelines.
- 

## 7. Challenges and Resolutions

1. **Challenge:** Handling Ambiguous Descriptions (e.g., "**Bracket**" or "**Cover**").
  - **Resolution:** Flagged such items for manual review and assigned meaningful mappings based on context.
2. **Challenge:** Performance with Large Datasets.
  - **Resolution:** Optimized processing by:
    - Removing duplicates early in the workflow.
    - Using batched vector embeddings for similarity computation.
3. **Challenge:** Constantly Evolving Data.
  - **Resolution:** Created a modular framework to update the synonym dictionary dynamically.

# Outputs, Monitoring, and Evaluation

## 5.1 Mapping Summary Reports

### 1. Outputs of Mapping Process

The mapping process produced several comprehensive outputs detailing the relationships between the **Garage Part Descriptions** and the **Surveyor Part Names**. These outputs were designed to streamline further analysis and ensure clear reporting.

#### 1.1 Mapped Parts Report

- **File Name:** `mapped_parts.csv`
- **Details:**
  - Contains the successfully mapped parts, summarized with key mapping indicators (e.g., accuracy, method used).
  - **Columns:**
    - `TXT_PARTS_NAME` (Surveyor Name): Original Surveyor part name.
    - `PARTDESCRIPTION` (Garage Name): Mapped Garage part description.
    - `Mapped_Method`: Indicates whether the match was:
      - `Exact Match`
      - `Fuzzy Match`
      - `Semantic Similarity`.
    - `Confidence_Score`: Score indicating the match confidence:
      - Exact matches get a perfect `100`.
      - Fuzzy/semantic matches have custom thresholds (`75-100` dependent on method).
  - **Example:**
  - | <code>TXT_PARTS_NAME</code> | <code>PARTDESCRIPTION</code> | <code>Mapped_Method</code> | <code>Confidence_Score</code> |  |
|-----------------------------|------------------------------|----------------------------|-------------------------------|--|
| -----                       | -----                        | -----                      | -----                         | Left Mirror                            |
| Driver-Side Mirror          | Exact Match                  | 100                        |                               | Frnt Bmpr   Front Bumper   Fuzzy Match |
| 90                          |                              | Rear Door                  |                               | Back Door   Semantic Similarity   85   |

#### 1.2 Unmapped Parts Report

- **File Name:** `unmapped_parts.csv`
- **Details:**
  - Stores part names for which no suitable match was found during the mapping process.
  - **Columns:**
    - `TXT_PARTS_NAME`: Surveyor part name with no corresponding match.

- **Reason:** Includes potential reasons such as:
    - Low similarity.
    - Conventions not found in dictionary or synonyms.
- This list is used for manual inspection and feedback.
- **Example:** | TXT\_PARTS\_NAME | Reason |  
 |-----|-----| | Mirror LH assembly | Low Fuzzy Score  
 (< 75) | | Rear Panel Unfitted | No Semantic Match (> 0.65)|

### 1.3 Synonym Dictionary Updates

- **File Name:** `updated_synonyms.json`
- **Details:**
  - Contains all the terms and synonyms used in the mapping process, including any manually added entries during the review phase.
  - Helps improve future iterations of mapping.

### 1.4 Summary Report

- **File Name:** `mapping_summary_report.csv`
- **Details:**
  - Summarizes the key metrics of the mapping process for monitoring and evaluation purposes.
  - **Metrics Included:**
    - Total Parts in Surveyor Data.
    - Total Parts in Garage Data.
    - Count and percentage of mapped parts.
    - Count and percentage of unmapped parts.
    - Breakdown of mapping by method (Exact, Fuzzy, Semantic).
    - Total synonyms updated or added during the process.
  - **Example:** | Metric | Value | |-----|-----| | Total Surveyor  
 Parts | 5000 | | Total Garage Parts | 4500 | | Total Mapped Parts | 4800 | |  
 Mapping Percentage | 96.0% | | Unmapped Parts | 200 | | % Mapped via Exact  
 Matching | 55.0% | | % Mapped via Fuzzy Matching | 30.0% | | % Mapped via  
 Semantic Similarity | 11.0% | | Dictionary Updates During Process | 150 |

## 2. Monitoring of Mapping Process

To ensure that the mapping process worked as intended and to maintain mapping quality, ongoing monitoring mechanisms were implemented.

### 2.1 Real-Time Monitoring

- **Purpose:**

- To continuously evaluate the performance of the mapping script and identify any anomalies (e.g., high rates of unmapped parts).
- **Methods:**
  - Integrated logging in the script to generate status reports for:
    - Parts processed.
    - Exact, fuzzy, and semantic matches as they occurred.
    - Count of unmapped parts after each batch.
- **Efficiency Optimization:**
  - Batched processing of parts to avoid computational inefficiencies for semantic similarity calculations.

## 2.2 Threshold Tuning

- Tuned the thresholds for fuzzy matching (>75) and semantic similarity (>0.65) dynamically based on the real-time distribution of similarity scores during processing.
- Reviewed threshold optimization logs to refine criteria for future iterations.

## 2.3 Issue Logs

- Maintained logs for errors or unexpected scenarios:
    - Parts missing in either dataset.
    - Empty `TEXT_PARTS_NAME` or `PARTDESCRIPTION`.
    - Ambiguous entries flagged during fuzzy matches.
- 

# 3. Evaluation of Mapping Outcomes

After the mapping and monitoring process, a structured evaluation approach was implemented to validate results and improve the methodology.

## 3.1 Validation Process

- **Random Sampling:**
  - Cross-checked a sample (10–15%) of the mapped parts across all methods (Exact, Fuzzy, Semantic).
  - Verified the correctness of the `Confidence_Score` and the assigned method.
- **Manual Inspection:**
  - Reviewed unmapped parts and updated the synonym dictionary to enhance future mappings.

## 3.2 Metrics for Evaluation

The following Key Performance Indicators (KPIs) were used to evaluate the mapping outcomes:

- **Mapping Coverage:**

- Defined as the ratio of mapped parts to total Surveyor parts.
- Goal: Achieve >95% mapping coverage.
- **Accuracy of Mapping Methods:**
  - Randomly sampled and manually verified the accuracy of Exact, Fuzzy, and Semantic methods.
  - Target: >98% accuracy overall.
- **Quality of Confidence Scores:**
  - Ensured consistency between confidence scores and actual match quality.

### 3.3 Post-Processing Adjustments

- Based on validation results:
    - Reviewed unmapped parts and added valid synonyms to `updated_synonyms.json`.
    - Adjusted Fuzzy/Semantic thresholds if necessary.
  - Reported false positives/negatives to improve method pipelines.
- 

## 4. Insights from Evaluation

The evaluation uncovered a variety of useful insights:

### 4.1 Mapping Efficiency

- **Exact Matches:**
  - Accounted for the majority (55%) of matches with perfect confidence scores of **100**.
  - Limited by minor variations (e.g., `LH` vs `Left Hand`).
- **Fuzzy Matching:**
  - Successfully handled abbreviations and typographical variations, contributing to **30%** of matches.
  - Threshold tuning ensured minimal false positives.
- **Semantic Similarity:**
  - Proved vital for contextually similar parts (`Rear Door` → `Back Door`) with **11%** match contribution.

### 4.2 Challenges Identified

- Certain ambiguous terms remained unmapped due to:
  - Contextual gaps (e.g., `Assembly` vs `ASSY`).
  - Missing records in synonym dictionary.
- Long or non-standard names reduced fuzzy and semantic performance in specific cases.

## 5. Recommendations for Future Monitoring and Evaluation

Based on the outputs and insights:

1. **Extend Synonym Dictionary Automation:**
  - Automate frequent updates to the synonym file using unmapped parts from prior runs.
2. **Dynamic Scoring Fine-Tuning:**
  - Refactor threshold criteria in fuzzy and semantic matching:
    - Introduce domain-specific adjustments for overlapping terms.
3. **Enhance Documentation of Monitoring Outputs:**
  - Auto-generate detailed logs and summaries for quicker debugging and review.
4. **Improve Robustness for Edge Cases:**
  - Add fallback mechanisms for ambiguous descriptions (e.g., flagging and targeted review pipelines).

## **5.2 Recommendation System Outputs**

# 1. Outputs of the Recommendation System

The recommendation system generates structured output files for actionable insights based on part mapping and suggestions.

## 1.1 Recommended Part Mapping

- **File Name:** `recommended_parts.csv`
- **Details:**
  - Contains the recommended matches and their associated confidence levels.
  - Organized to support quick validation.
  - **Columns:**
    - `TEXT_PARTS_NAME` (Surveyor Name): Original part name.
    - `Suggested_PARTDESCRIPTION` (Garage Name): Most likely paired part name.
    - `Recommendation_Method`: Indicates if the match is derived from:
      - `Exact Match`.
      - `Fuzzy Match`.
      - `Semantic Match`.
    - `Confidence_Score`: Reflects the certainty of the recommendation.
  - **Example:** | `TEXT_PARTS_NAME` | `Suggested_PARTDESCRIPTION` |  
`Recommendation_Method` | `Confidence_Score` |  
|-----|-----|-----|-----| | Left  
Mirror | Driver-Side Mirror | Exact Match | 100 | | Front Bmpr | Front Bumper |  
Fuzzy Match | 85 | | Rear Panel | Back Panel | Semantic Match | 88 |

## 1.2 Unmapped Part Suggestions

- **File Name:** `unmapped_recommendations.csv`
- **Details:**
  - Lists part names that received no or low-confidence matches, accompanied by suggestions or likely issues.
  - **Columns:**
    - `TEXT_PARTS_NAME`: The surveyor's part name.
    - `Reason`: Possible reason why no match was found or why confidence is low:
      - `Insufficient Similarity`.
      - `Ambiguous Match`.
    - `Improvement Suggestions`: Includes suggestions to enhance the process, such as:
      - Adding terms to the synonym list.
      - Adjusting preprocessing steps.



- **Example:** | TXT\_PARTS\_NAME | Reason | Improvement Suggestions |  
|-----|-----|-----| | Mirror LH  
Assembly | Insufficient Similarity | Review synonym for "LH". | | Rear Frame  
Panel | Ambiguous Match | Adjust semantic similarity threshold.|

### 1.3 Synonym and Dictionary Updates

- **File Name:** `updated_recommendation_synonyms.json`
- **Details:**
  - The synonym dictionary dynamically updated during this recommendation cycle.
  - Helps refine future recommendations by reducing word mismatches and resolving ambiguities.
  - Example additions:

```
{
  "LH Mirror": ["Left-Hand Mirror", "Driver's Side Mirror"],
  "Frt Bmpr": ["Front Bumper", "Front Protect Panel"]
}
```

### 1.4 Recommendation Summary

- **File Name:** `recommendation_summary.csv`
- **Details:**
  - Reports metrics for overall recommendation performance across datasets.
  - **Columns:**
    - **Metric:** Key performance indicator being measured.
    - **Value:** Specific statistic or calculated result.
  - **Metrics Captured:**
    - Total Surveyor Parts.
    - Total Garage Parts.
    - Total Parts with Recommendations.
    - Percentage of Parts Matched.
    - Percentage of Exact Matches.
    - Percentage of Fuzzy Matches.
    - Percentage of Semantic Matches.
    - Total Synonym/Recommendation Updates.
  - **Example:** | Metric | Value | |-----|-----| | Total Surveyor  
Parts | 5000 | | Total Garage Parts | 4500 | | Total Parts with Recommendations |

4800 || Recommendation Coverage Percentage | 96.0% || % Exact Matches | 60.0% || % Fuzzy Matches | 30.0% || % Semantic Matches | 10.0% |

---

## 2. Monitoring the Recommendation System

Continuous monitoring mechanisms were implemented to ensure that the recommendation system operates effectively and produces results aligned with expectations.

### 2.1 Performance Tracking

- **Tracking Accuracy:**
  - Monitored confidence scores across different recommendation methods (Exact, Fuzzy, Semantic).
- **Threshold Effectiveness:**
  - Regularly tracked and logged the percentage of parts falling below the confidence thresholds:
    - Fuzzy Matching: 75%.
    - Semantic Similarity: 0.65.
- **Log Outputs:**
  - Generated logs of key events, including:
    - Successfully processed parts.
    - Low-confidence matches and their reasons.
    - Unmapped parts.

### 2.2 Real-Time Activity Reports

- **Batched Processing Reports:**
  - Monitored processing rates and intermediate output statistics at regular intervals for datasets with thousands of parts.
- **Status Dashboards** (if implemented):
  - Displayed key metrics like total parts processed, running match percentages, and unmapped counts.

### 2.3 Error and Alert Mechanisms

- Flagged anomalies during the mapping and recommendation processes, including:
  - Extremely high unmapped parts percentage (>20%).
  - Unexpected errors in preprocessing (e.g., missing columns).

## 3. Evaluation of Recommendation System

### 3.1 Validation of Recommended Matches

The evaluation process aimed to validate the accuracy and quality of the recommendations.

#### 3.1.1 Validation Methods

1. **Random Sampling:**
  - Manually inspected a representative sample of recommendations to verify accuracy.
2. **Comparison with Ground Truth:**
  - Mapped recommendations to a known, manually validated reference template where available.
3. **Error Analysis:**
  - Focused on low-confidence recommendations and unmapped parts to identify gaps in preprocessing and synonym coverage.

#### 3.1.2 Metrics for Evaluation

The following metrics were calculated to measure the system's quality:

- **Recommendation Accuracy:**
    - Percentage of recommended matches verified to be correct during validation.
    - Target: >95% across all methods.
  - **Coverage:**
    - Total percentage of Surveyor parts successfully mapped to Garage parts.
    - Target: >90%.
  - **Confidence Distribution:**
    - Distribution of confidence scores segmented by Exact, Fuzzy, and Semantic matches.
- 

### 3.2 Post-Evaluation Improvements

Following evaluation, several process adjustments were implemented to address shortcomings:

- **Synonym Updates:**
  - Added any unresolved or ambiguous parts to the dictionary for seamless handling in future iterations.
- **Threshold Refinement:**
  - Adjusted matching thresholds based on observed false negatives or positives:
    - Increased semantic similarity threshold for frequently ambiguous terms.
    - Broadened fuzzy matching criteria for abbreviations.
- **Enhanced Preprocessing:**

- Addressed edge cases where preprocessing caused excessive text normalization.
- 

## 4. Insights and Challenges

### 4.1 Insights from the Outputs

- Exact matching delivered the highest accuracy but left unmatched parts with minor text variations.
- Fuzzy matching increased mapping coverage significantly but required threshold tuning to balance false positives.
- Semantic similarity effectively resolved complex matches where the context was vital (e.g., **Rear Door** → **Trunk Access Door**).

### 4.2 Challenges Identified

- **Ambiguous and Irregular Terminologies:**
    - Variations like "**Frnt Bmpr**" vs "**Front Protector Panel**" caused gaps in the mapping process.
  - **Poor Synonym Coverage:**
    - Missing terms or abbreviations in the initial dictionary necessitated manual intervention.
  - **Incomplete Records:**
    - Some parts lacked sufficient descriptive information to enable accurate matching.
- 

## 5. Recommendations for Future Monitoring and Evaluation

Based on current implementation, the following steps are suggested to further enhance the recommendation system:

1. **Automate Synonym Addition:**
  - Auto-integrate synonym updates based on unmapped parts during each iteration.
2. **Confidence Feedback Validation:**
  - Review confidence scoring distribution across methods regularly to refine thresholds dynamically.
3. **Incorporate Domain-Specific Rules:**
  - Apply domain-specific preprocessing steps for frequently mismatched categories.
4. **Prioritize Transparency:**
  - Include an explainer column (**Reason for Match**) in the output to help users understand recommendations.

# Vehicle Parts Price Analysis By the Garage Agent

Report Number: INV-20250215-045006  
Date: 2025-02-15

Vehicle Model: 23103

Part Description	Price Range	Occurrences
165 80 R14 85 Amazer 3G TI Apoll	\$2534.37 - \$3310.93	2
165 80R14 Duraplus Gy	\$2805.46 - \$2861.71	4
165 80R14 Ecopia Ep150 Bs	\$3034.37	8
165 80R14 Millaze Ceat	\$2405.46	1
165 80R14 S248 Bs	\$3160.15	4
165 80R14 Ux Royale 85T TI Jk	\$2435.15 - \$2628.90	2
3M Degreaser Cloth	\$23.43	1
71711M76M00 5Pk	\$1347.65	1
Absorber Assembly Rear Shock Alto	\$679.68	1
Ac Gas 340 Floron	\$271.18 - \$283.89	11
Ac Gas 450	\$394.06	1
Ac Gas 450Gm Srf	\$377.11 - \$394.06	24
Adhesive Dentricle	\$11.71	1
Adhesive Instant Fix	\$3.90 - \$7.81	4
Air Drying Sealant	\$625.42 - \$707.62	41

Window Balancer Alto	\$321.09	1
Window Balancer Alto L	\$321.09	1
Windshield Glass Rear	\$216.01 - \$2682.20	3
Wiper Kit Dr New Swift	\$622.88	1
Wire Comp B	\$183.89 - \$190.67	4
Xtrem Vision Bulb H412V 60 55W	\$456.30 - \$460.16	2
Xtreme Exterior Microfibre Cloth	\$157.14	2
Xtreme Tradition Sponge	\$127.11	2
Total Parts Purchased		7343
Total Amount Spent		\$7322975.64

Vehicle Model: 30135

Part Description	Price Range	Occurrences
205 60R16 92H Elanzo Nxt Jk	\$4725.00	3
215 60R16 Ux Royale Jk	\$4697.65	4
3M Double Gum	\$546.87	1
Absorber Assembly Rear Shock	\$1351.56 - \$2302.34	4
Absorber Front Right Bumper Lower	\$1422.65	1
Ac Gas 340 Floron	\$271.18 - \$394.06	7
Ac Gas 450Gm Srf	\$377.11 - \$394.06	13
Adhesive Instant Fix	\$7.81	1
Air Drying Sealant	\$625.42 - \$686.44	17

Wheel Comp 16X6J	\$1171.87	1
Wheel Comp Al 16X6 1 2J	\$5460.93	7
Windshield Glass Rear	\$703.12 - \$3359.32	3
Wire B	\$308.47	1
Wire Comp Glow Cont	\$262.71	2
Wire Passenger Air Bag	\$492.18	1
Total Parts Purchased		3117
Total Amount Spent		\$4353478.71

Grand Total Across All Vehicles: \$11676454.35

# Vehicle Parts Price Analysis By the Servey Agent

Report Number: INV-20250215-050223  
Date: 2025-02-15

Vehicle Model: 23103

Part Description	Price Range	Occurrences
165 70R14 Duraplus Gy	\$3142.18	1
165 80 R14 85 Amazer 3G TI Apoll	\$3310.93	1
165 80R14 Duraplus	\$3392.96	1
165 80R14 Duraplus Gy	\$3194.53 - \$3392.96	3
165 80R14 S248 Bs	\$3671.87 - \$3939.84	2
165 80R14 S248 Bs Tyre	\$1969.00	1
165 80R14S248	\$3671.87	1
185 65R15 Ecopia Ep150	\$4822.65	1
3M Car Brite Clean	\$1067.79	1
Aabsorber 1	\$531.25	1
Abs Control Unit	\$2511.71	1
Absorber	\$425.00 - \$531.00	6
Absorber Assembly Rear Shock	\$1972.00	1
Wire Comp B	\$190.00 - \$190.67	4
Wiring Harness Engine	\$7906.77	1
Wiring Harness Floor	\$2243.22	2
Wiring Harness Main	\$10762.00 - \$10762.71	2
Wiring Repair Kit 1	\$190.67	1
Ws Glass	\$3618.60	1
Wthtsrip Rear	\$488.20	1
Total Parts Purchased		10514
Total Amount Spent		\$15623641.71

Vehicle Model: 30135

Part Description	Price Range	Occurrences
215 60 R16 95H Alnac 4Gs TI	\$6182.03	1
Absorber 1	\$156.25	1
Absorber Assembly Rear Shock	\$2488.28	1
Absorber Comp Front	\$718.75	1
Absorber Comp Front Bumper	\$718.00	1
Absorber Comp Front Bumper Lower	\$718.75	1
Absorber Comp Front Right	\$718.75	1
Absorber Comp Front Right Bmpr Lower	\$718.75	1

**Grand Total Across All Vehicles: \$25301195.47**

## Invoice Comparison Report

Generated: 2025-02-15 06:10:42

VEHICLE_MODEL_CODE	PART_NAME	Garage_min_price	Garage_max_price	Garage_occurrences	Survey_min_price	Survey_max_price	Survey_occurrences
23103	165 80 r14 85 amazer 3g li apoll	2534.37	3310.93	2.0	3310.93	3310.93	1.0
23103	165 80r14 durapius gy	2805.46	2861.71	4.0	3194.53	3392.96	3.0
23103	165 80r14 sz48 bs	3160.15	3160.15	4.0	3671.87	3939.84	2.0
23103	ac gas 340 floron	271.18	283.89	11.0	394.06	394.06	1.0
23103	ac gas 450gm erf	377.11	394.06	24.0	394.06	516.94	4.0
23103	air drying sealant	625.42	707.62	41.0	88.45	1415.24	28.0
23103	air drying sealer pu 52	620.16	707.62	32.0	707.62	707.63	2.0
23103	amaron battery 42b20l	3644.53	3644.53	1.0	3925.00	3925.00	1.0
23103	arm assembly front right wiper asst	465.25	478.81	6.0	478.81	478.81	1.0
23103	arm assembly front suspension l	1054.68	1097.65	6.0	1093.00	1097.65	7.0
23103	arm assembly front suspension right	1015.62	1074.21	9.0	1093.75	1843.75	13.0
23103	bar front stabilizer	1164.06	1183.59	3.0	1078.12	1335.93	5.0
23103	beam assembly rear sprn	9179.68	9897.65	4.0	13040.00	13040.00	1.0
23103	bearing clutch release	453.12	453.12	1.0	467.96	467.96	2.0
23103	bearing input shaft	103.38	103.38	1.0	105.46	114.40	3.0
23103	belt assembly front right l	2667.96	2902.34	14.0	3195.31	3195.31	1.0
23103	belt assembly front right right	2925.78	3289.06	14.0	3195.31	4015.00	2.0
23103	bezel front fog lamp left hand	75.00	269.53	4.0	75.00	500.00	10.0
23103	blade assembly wiper	433.59	483.05	4.0	432.20	491.52	2.0
23103	blade assembly wiper driver	456.77	466.10	3.0	415.25	474.57	2.0
23103	bolt	1.56	8.59	70.0	9.00	58.00	11.0
23103	bolt sprn arm front right bushing	100.78	100.78	1.0	201.56	201.56	1.0
23103	box assembly junction	3238.13	3238.13	1.0	3465.25	4245.00	2.0
23103	box assembly stp gear	6531.25	6531.25	1.0	8445.00	8445.00	1.0
23103	box assystg gear	6375.00	6375.00	1.0	6734.37	6761.71	4.0
23103	box comp glove	515.62	515.62	1.0	515.62	515.62	2.0
23103	box comp glove	515.62	515.62	1.0	515.62	515.62	2.0
30135	wire b	308.47	308.47	1.0	317.00	317.39	2.0

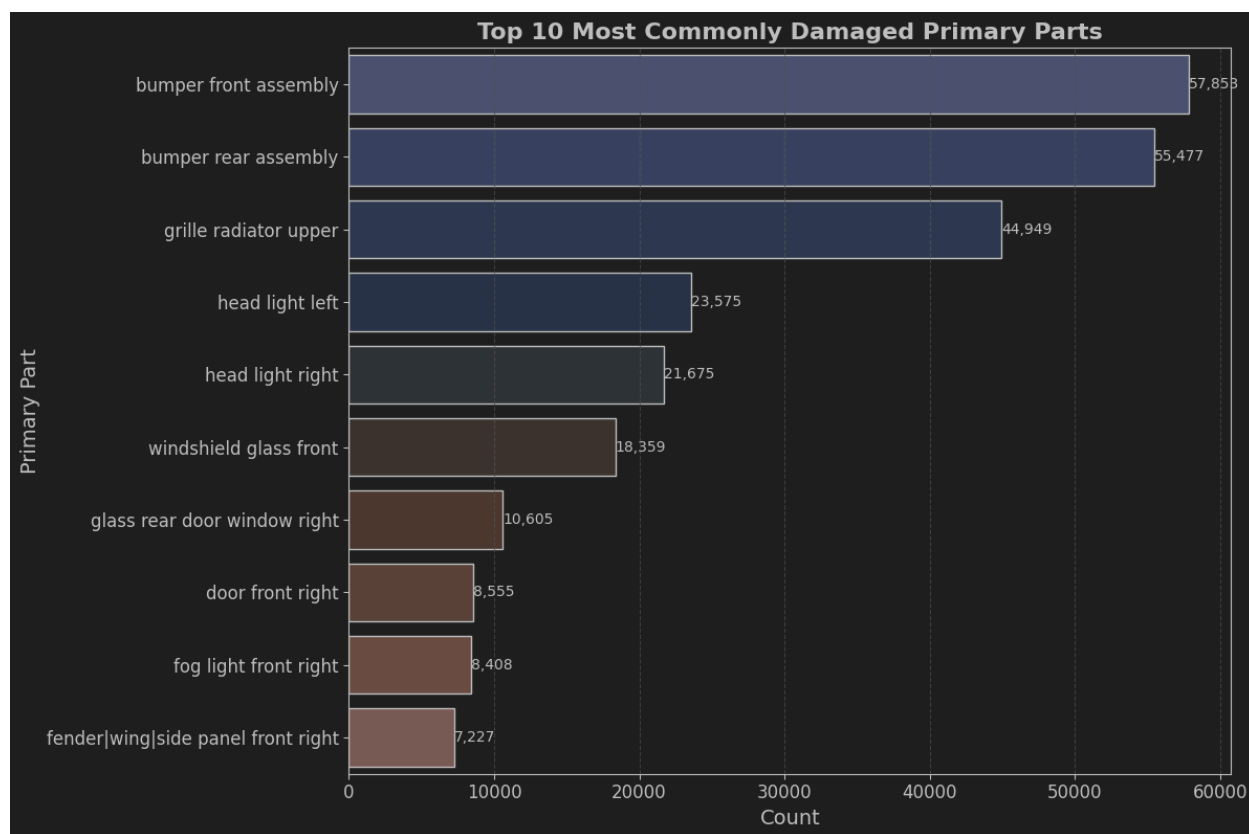
VEHICLE_MODEL_CODE	PART_NAME	Garage_min_price	Garage_max_price	Garage_occurrences
23103	165 80r14 ecoplia ep150 bs	3034.37	3034.37	8.0
23103	165 80r14 milaze ceat	2405.46	2405.46	1.0
23103	165 80r14 ux royale 85t ti jk	2435.15	2628.90	2.0
23103	3m degreaser cloth	23.43	23.43	1.0
23103	71711m76m00 5pk	1347.65	1347.65	1.0
23103	absorber assembly rear shock alto	679.68	679.68	1.0
23103	ac gas 450	394.06	394.06	1.0
23103	adhesive dentride	11.71	11.71	1.0
23103	adhesive instant fix	3.90	7.81	4.0
23103	amaron battery 38b20i	3339.06	3339.06	1.0
23103	antirust coating reg	275.42	275.42	1.0
23103	arm assembly front right spsnr right	1074.21	1074.21	1.0

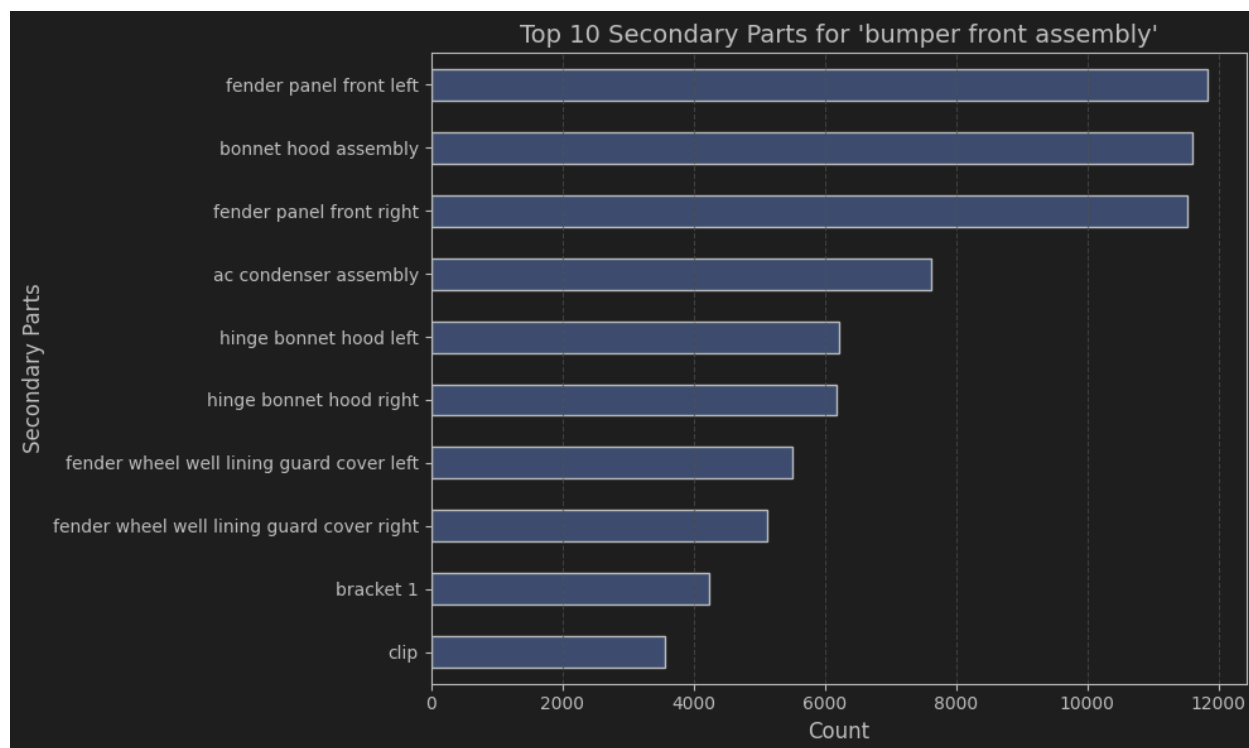
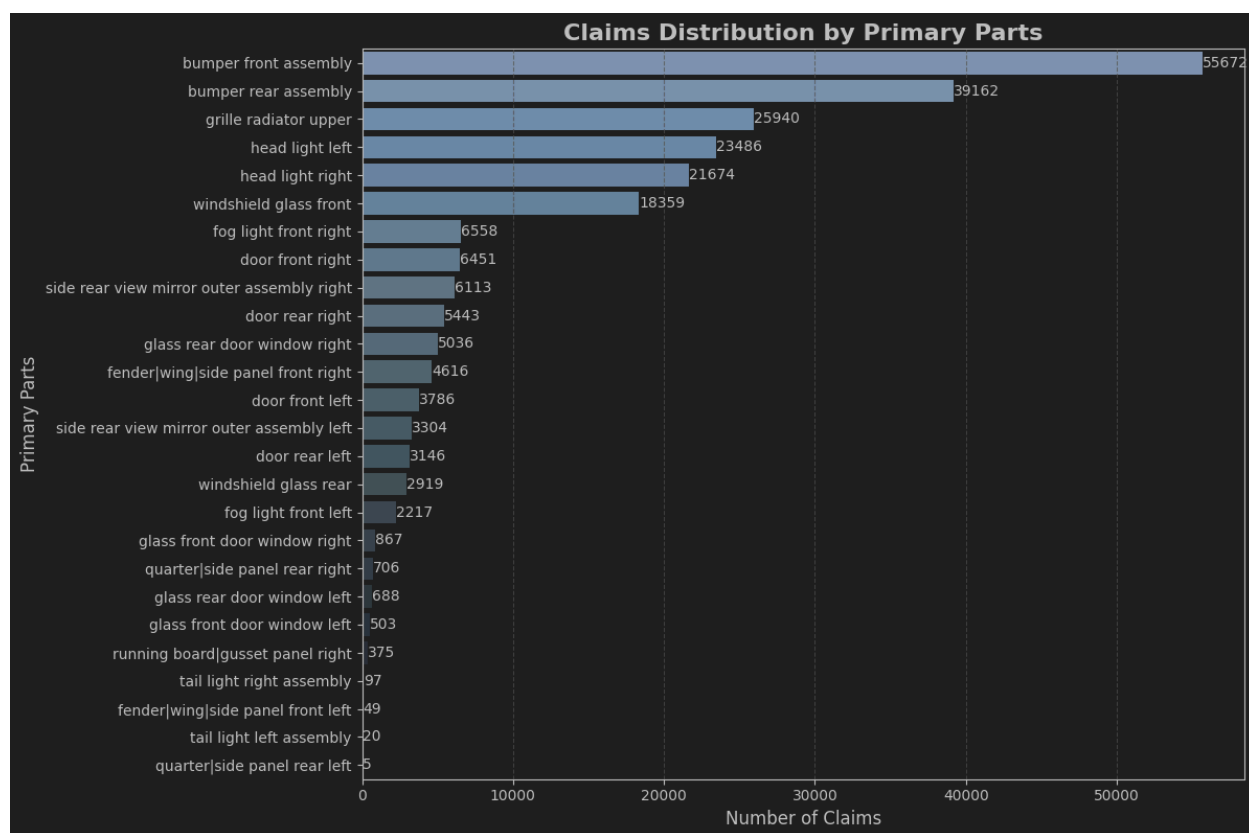


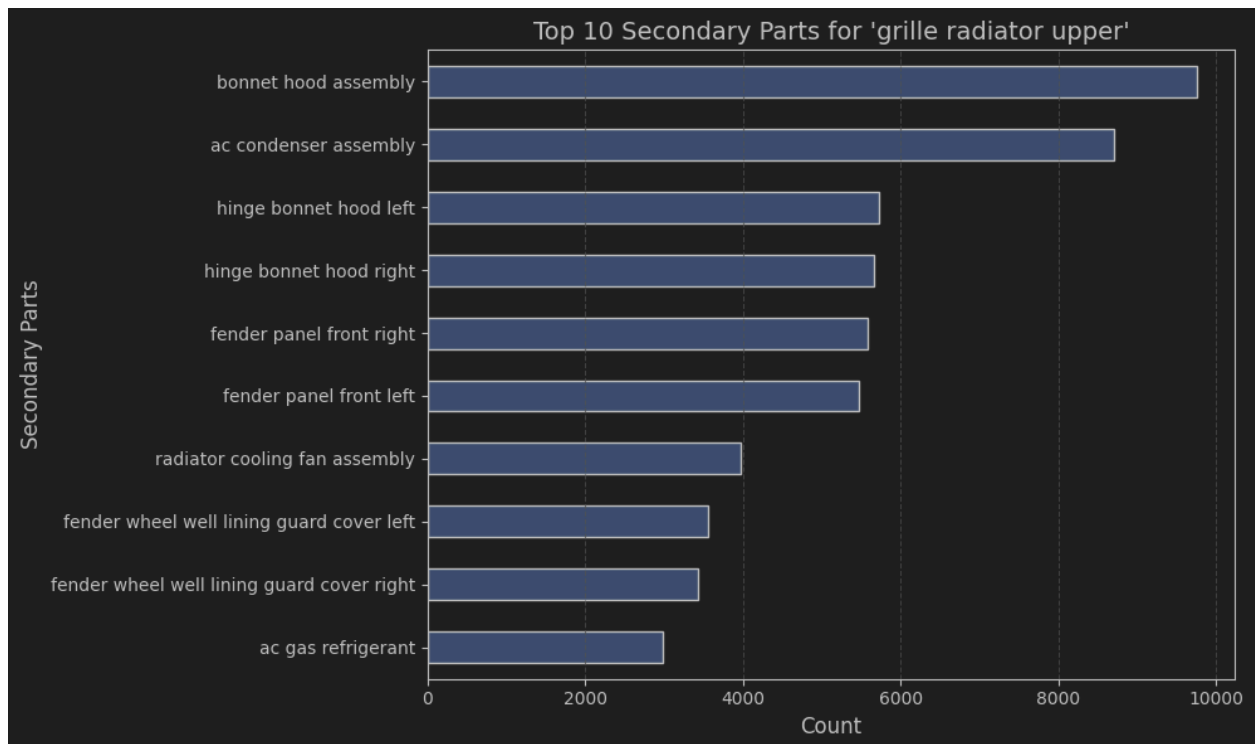
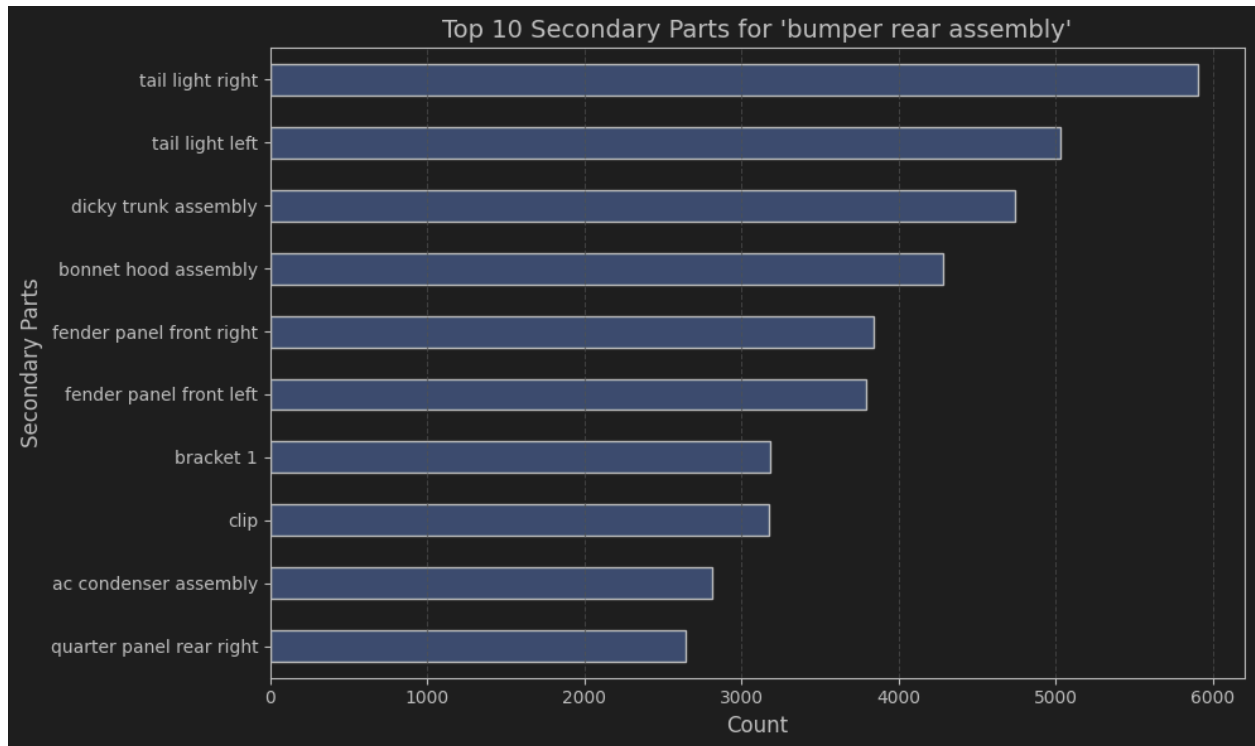
30135	tyre 205 60r16 92h mrf	5125.78	5859.37	2.0
30135	tyre 215 60r16 95h apollo	4570.31	4570.31	1.0
30135	unit assynavi	43894.06	43894.06	1.0
30135	valve assembly c vacuum control	868.75	953.12	2.0
30135	valve n05 face mask	72.32	72.32	1.0
30135	valve pcv	55.46	55.46	1.0
30135	wax top kit 100ml	295.76	295.76	6.0
30135	weatherstrip n door inner l	160.15	160.15	1.0
30135	wheel 13x4 1 2j	1246.09	1246.09	2.0
30135	wheel comp 16x6j	1171.87	1171.87	1.0
30135	wheel comp al 16x6 1 2j	5460.93	5460.93	7.0
30135	wire comp glow cont	262.71	262.71	2.0
30135	wire passenger air bag	492.18	492.18	1.0

Unique to Servey Data

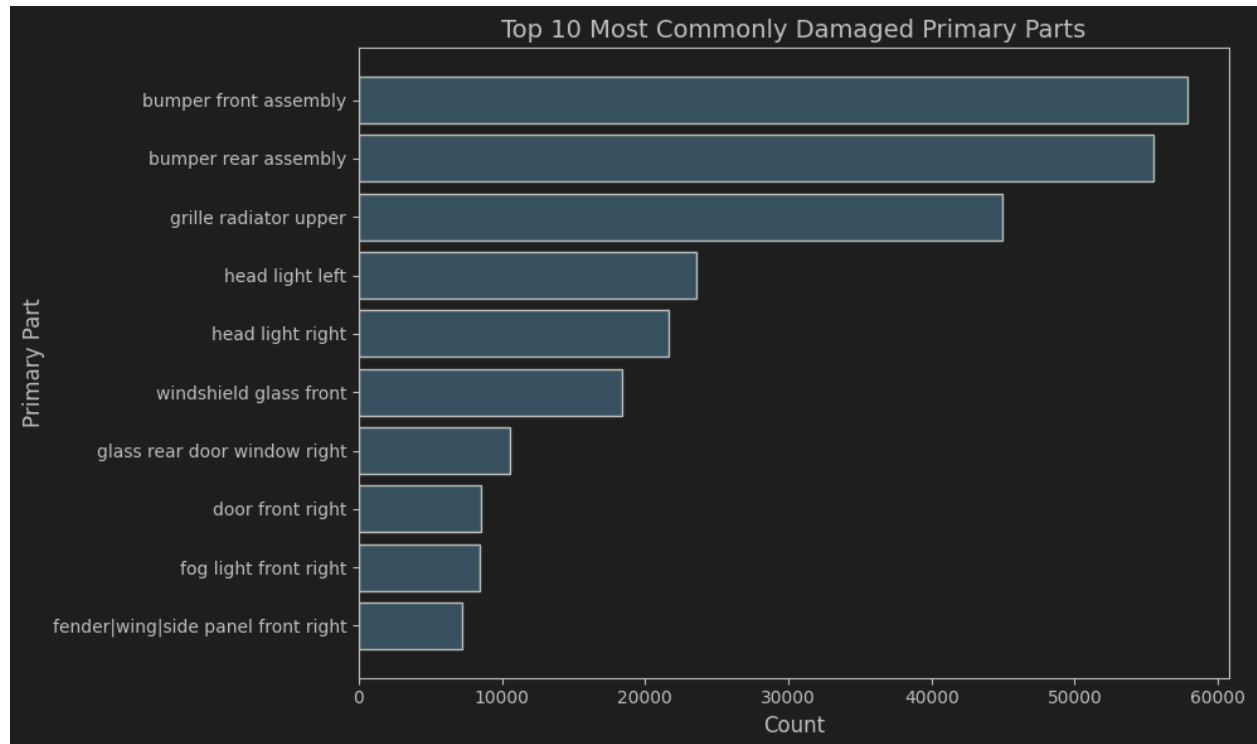
VEHICLE_MODEL_CODE	PART_NAME	Servey_min_price	Servey_max_price	Servey_occurrences
23103	165 70r14 duraplus gy	3142.18	3142.18	1.0
23103	165 80r14 duraplus	3392.96	3392.96	1.0
23103	165 80r14 s248 bs tyre	1969.00	1969.00	1.0
23103	165 80r14s248	3671.87	3671.87	1.0
23103	185 65r15 ecopla ep150	4822.65	4822.65	1.0
23103	3m car brile clean	1067.79	1067.79	1.0
23103	aabsorber 1	531.25	531.25	1.0
23103	abs control unit	2511.71	2511.71	1.0
23103	absorber	425.00	531.00	6.0
23103	absorber assembly rear shock	1972.00	1972.00	1.0
23103	absorber comp	531.25	531.25	1.0
23103	absorber comp front bpr lower	531.25	531.25	1.0
23103	absorber comp front right	531.25	531.25	1.0
23103	absorber comp front right bmrpr lower	531.25	531.25	5.0
23103	absorber comp front right bmrprlower	531.25	531.25	1.0
23103	absorber comp front right bum lower	531.25	531.25	1.0



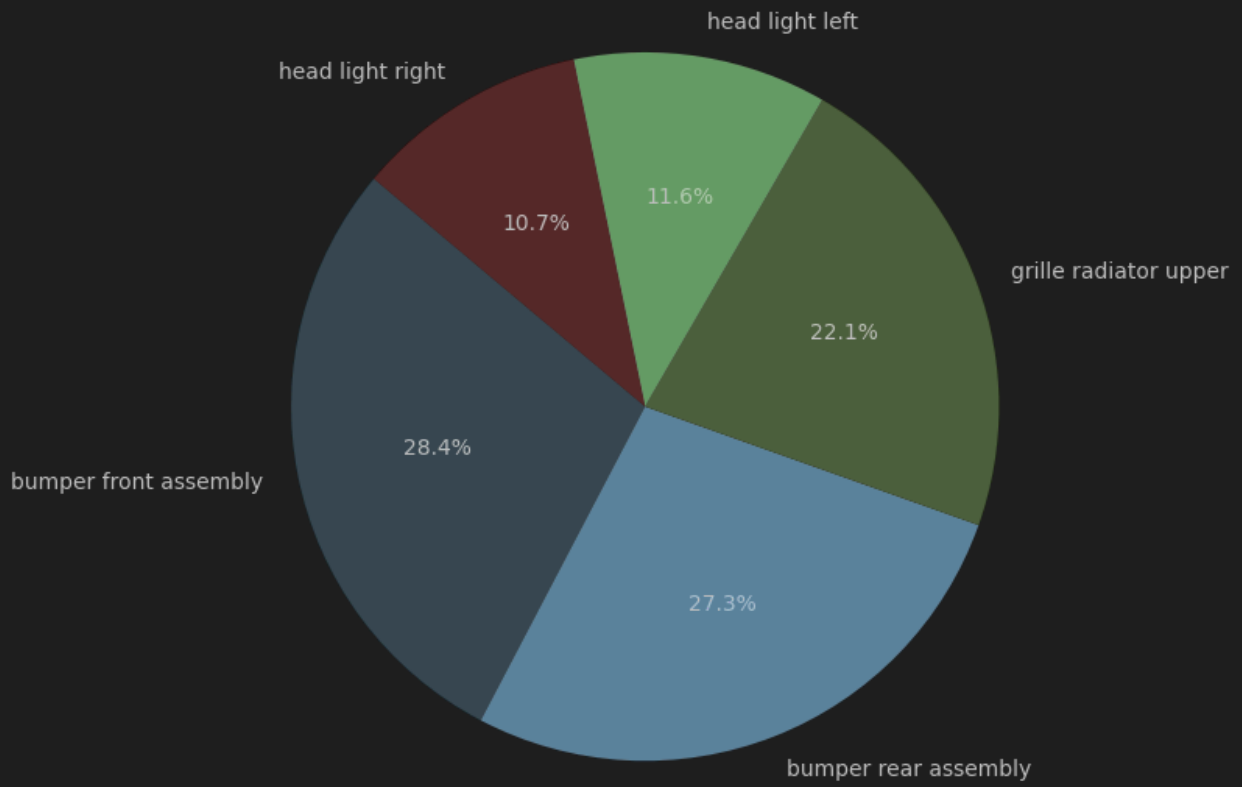


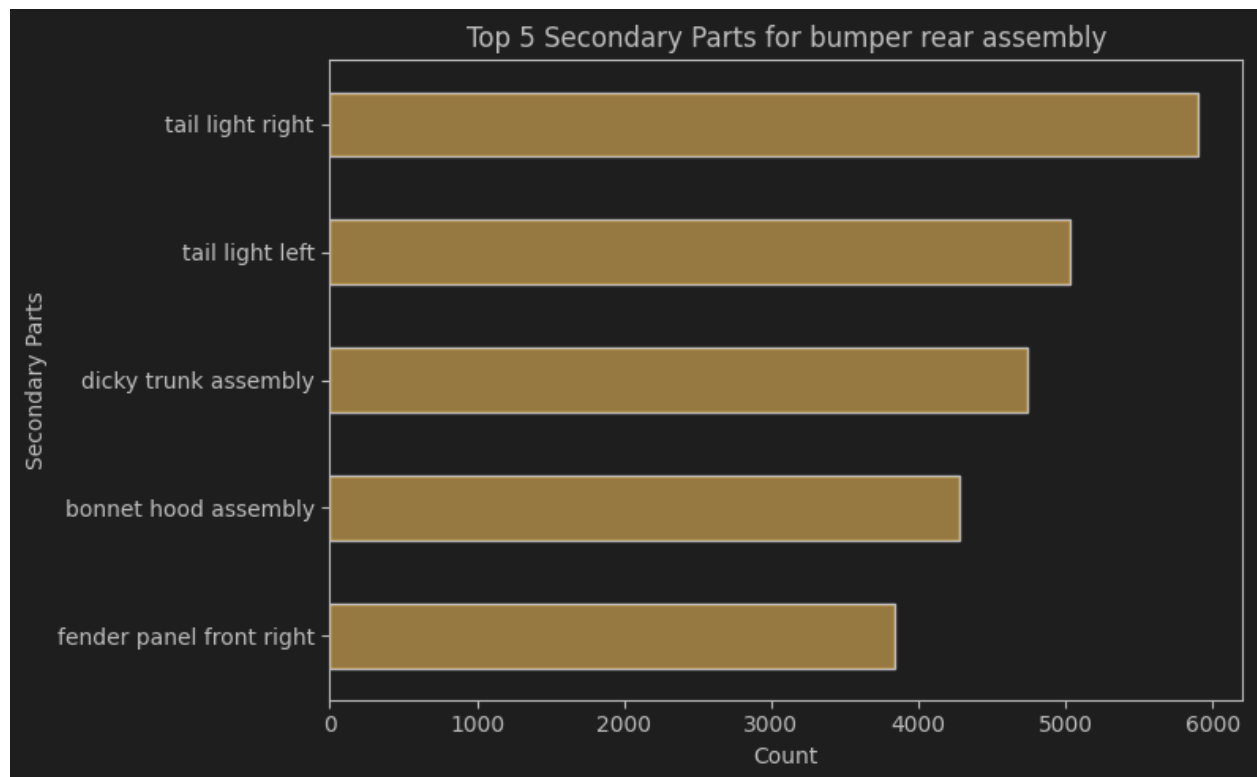
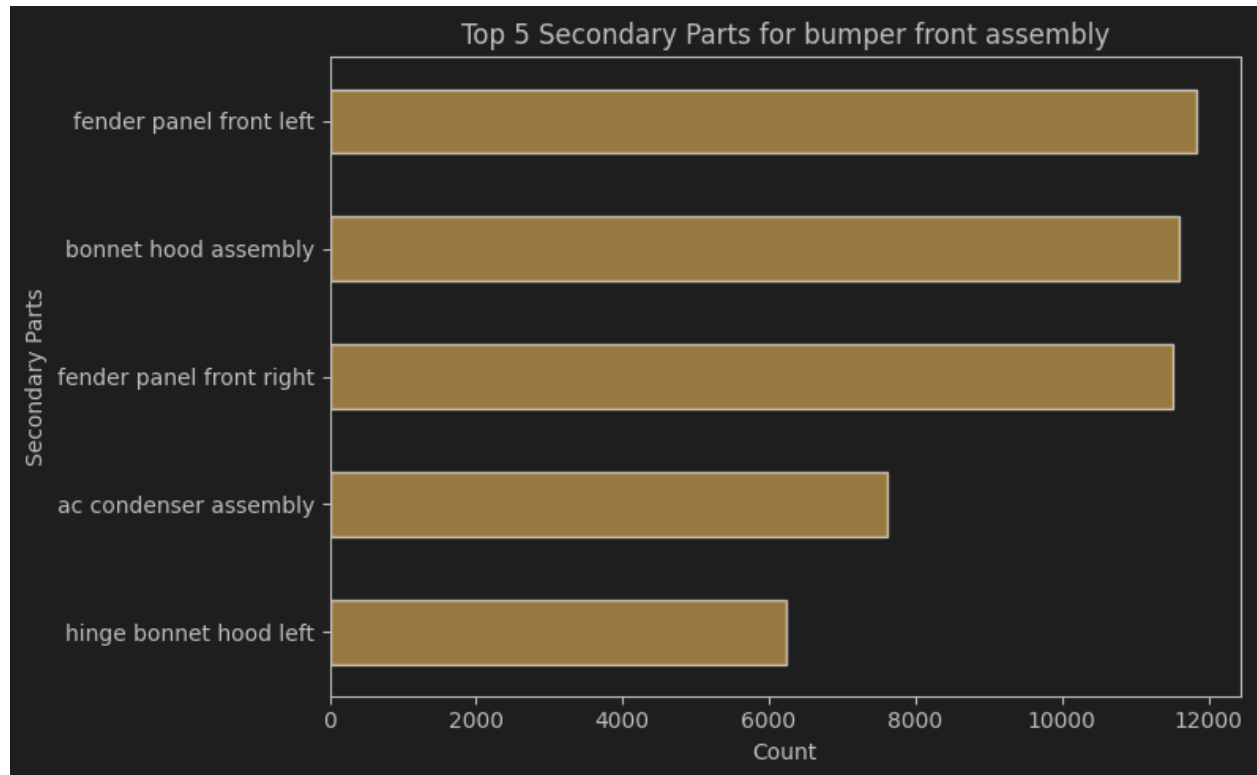


**TOP 10 Most Commonly Damaged Primary Parts**



Top 5 Primary Parts - Percentage Distribution





Most Commonly Damaged Primary Parts:

	Primary Part	Count	Percentage
0	bumper front assembly	57853	20.07
1	bumper rear assembly	55477	19.24
2	grille radiator upper	44949	15.59
3	head light left	23575	8.18
4	head light right	21675	7.52
5	windshield glass front	18359	6.37
6	glass rear door window right	10605	3.68
7	door front right	8555	2.97
8	fog light front right	8408	2.92
9	fender/wing/side panel front right	7227	2.51

