



SECP3133 HIGH PERFORMANCE DATA PROCESSING

SEMESTER 2 2024/2025

## **Project Report:**

# **Optimising High-Performance Data Processing for Carlist.my**

|               |  |
|---------------|--|
| GROUP         | 01   |
| SECTION       | 02   |
| GROUP MEMBERS | <ol style="list-style-type: none"><li>1. MARCUS JOEY SAYNER (A22EC0193)</li><li>2. MUHAMMAD LUQMAN HAKIM BIN MOHD RIZAUDIN (A22EC0086)</li><li>3. CAMILY TANG JIA LEI (A22EC0039)</li><li>4. GOH JING YANG (A22EC0052)</li></ol> |
| LECTURER      | DR. ARYATI BINTI BAKRI   |

## TABLE OF CONTENTS

|  |           |
|--|-----------|
| <b>1.0 Introduction.....</b>                     | <b>1</b>  |
| 1.1 Objectives.....                              | 1         |
| 1.2 Target Website and Data to be Extracted..... | 2         |
| <b>2.0 System Design and Architecture.....</b>   | <b>3</b>  |
| 2.1 Description of Architecture.....             | 3         |
| 2.2 Tools and Frameworks Used.....               | 4         |
| 2.3 Roles of Team Members.....                   | 4         |
| <b>3.0 Data Collection.....</b>                  | <b>5</b>  |
| 3.1 Crawling and Scraping Method.....            | 5         |
| 3.2 Number of Records Collected.....             | 5         |
| 3.3 Ethical Considerations.....                  | 5         |
| <b>4.0 Data Processing.....</b>                  | <b>7</b>  |
| 4.1 Cleaning Methods.....                        | 7         |
| 4.2 Data Structure (CSV/JSON/database).....      | 7         |
| 4.3 Transformation and Formatting.....           | 13        |
| 4.3.1 Combining Redundant Columns.....           | 13        |
| 4.3.2 Removing Noisy and Irrelevant Columns..... | 13        |
| 4.3.3 String Cleaning and Formatting.....        | 14        |
| 4.3.3.1 Installment Column.....                  | 14        |
| 4.3.3.2 Condition Column.....                    | 15        |
| 4.3.3.3 Sales Channel Column.....                | 15        |
| 4.3.3.4 Mileage Column.....                      | 15        |
| 4.3.4 Handling Missing Data.....                 | 16        |
| 4.3.5 Renaming Columns for Clarity.....          | 17        |
| 4.3.6 Removing Duplicate Rows.....               | 17        |
| <b>5.0 Optimisation Techniques.....</b>          | <b>18</b> |

|  |           |
|--|-----------|
| 5.1 Methods Used.....                          | 18        |
| 5.2 Code Overview.....                         | 18        |
| <b>6.0 Performance Evaluation.....</b>         | <b>19</b> |
| 6.1 Before vs After Optimisation.....          | 19        |
| 6.2 Time, Memory, CPU Usage, Throughput.....   | 19        |
| 6.3 Charts and Graphs.....                     | 21        |
| <b>7.0 Challenges and Limitations.....</b>     | <b>27</b> |
| 7.1 Reasons of Challenges and Limitations..... | 27        |
| 7.2 Solutions.....                             | 28        |
| <b>8.0 Conclusion and Future Work.....</b>     | <b>29</b> |
| 8.1 Summary of Findings.....                   | 29        |
| 8.2 Improvement for Future Works.....          | 30        |
| <b>9.0 References.....</b>                     | <b>32</b> |
| <b>10.0 Appendices.....</b>                    | <b>34</b> |
| 10.1 Appendix A.....                           | 34        |
| 10.2 Appendix B.....                           | 36        |
| 10.3 Appendix C.....                           | 39        |
| 10.4 Appendix D.....                           | 47        |
| 10.5 Appendix E.....                           | 51        |
| 10.6 Appendix F.....                           | 76        |
| 10.7 Appendix G.....                           | 79        |

## **1.0 Introduction**

Digital world-roaming societies today rely heavily on data to make their decisions. The modern world requires the ability to store and process data from massive online spaces in order to stay in touch with one's realistic goals and forecasting in data-centric decision-making. The goal of this project is to create and design a high-performance web crawler that will retrieve structured job-related data from Carlist.my.

An important source of information for automotive industry research is Carlist.my, a well-known website that offers real-time listings for new, used, and reconditioned cars in Malaysia. Its broad coverage of cars from private sellers as well as dealerships provides a thorough understanding of the market, which is helpful for research on consumer demand, pricing fluctuations, and market trends. However, several challenges may impact the efficiency of data collection from this site, including slow response times, crawl rate limitations, and the necessity to adhere to ethical scraping practices.

To address these challenges, High Performance Data Processing (HPDP) techniques—such as multithreading and parallel processing—will be applied to enhance system efficiency. The main objective is to optimise the data processing pipeline to achieve greater speed, scalability, and reliability while maintaining responsible and ethical standards in data scraping practices.

## **1.1 Objectives**

To address the challenges of large-scale job data extraction and processing, this project aims to accomplish the following primary goals:

- To create a web crawler that can retrieve at least 100,000 structured records from Carlist.my.
- To clean, process, and store the collected data in a database or other usable formats like CSV or JSON.
- To apply high-performance computing techniques, including multithreading and parallel processing, to optimise crawling and processing speed.

- To conduct a comparative analysis of the system's performance with and without optimisation, focusing on speed, CPU usage, memory usage, and throughput.

## 1.2 Target Website and Data to be Extracted

Carlist.my is a leading online platform for vehicle listings in Malaysia, offering an extensive database of new, used, and reconditioned vehicles across the country. The platform provides detailed information on various car models, including their specifications, pricing, seller information, and geographical locations. The types of data to be extracted for analysis are summarised in Appendix A (Table A1), which presents an overview of the key car attributes available on Carlist.my.

The data is collected and cleaned for uniformity by removing duplicates and irrelevant entries and stored in a structured format, MongoDB, for analysis and visualisation.

## 2.0 System Design and Architecture

This section presents the overall system architecture and the collaborative development process used in this project. It outlines the workflow from web scraping to data transformation and evaluation, detailing the architecture involved, the tools and frameworks used by each team member, and the specific roles assigned to ensure an effective division of tasks. The goal is to demonstrate how various technologies and teamwork were integrated to build a scalable and efficient data pipeline.

### 2.1 Description of Architecture

Figure 1 represents the workflow of a web scraping system followed by big data tools to process, clean, and evaluate scraped data.

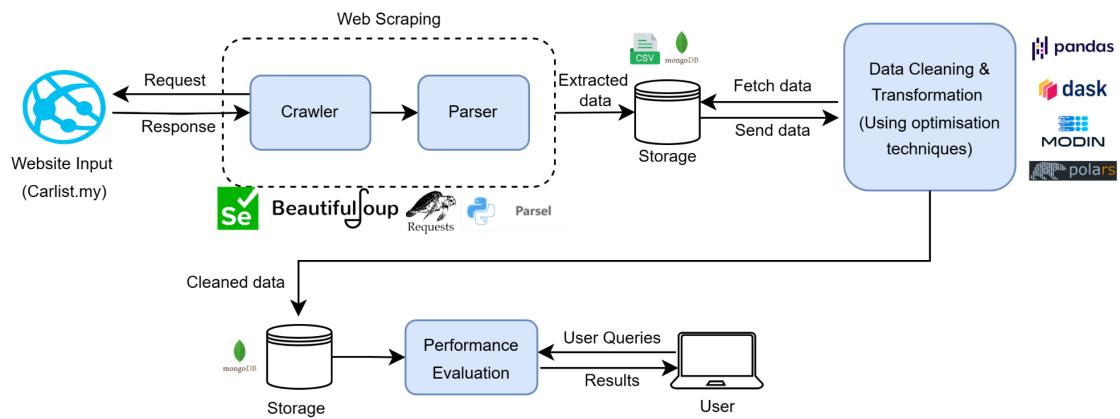


Figure 1: Web Scraping and Big Data Processing Architecture

The components involved in this architecture are described in Appendix B (Table B1), including details on the crawler, parser, data collection process, data transformation, optimisation techniques, and system evaluation metrics.

## 2.2 Tools and Frameworks Used

In this project, various libraries and frameworks were utilised to efficiently perform web scraping and big data processing tasks. The libraries used by each team member, categorised by their functionality, are summarised in Appendix B (Table B2).

## 2.3 Roles of Team Members

Each team member was assigned specific roles to ensure efficient division of work and comprehensive coverage of the scraping and data processing tasks, as shown in Table 1.

Table 1 Role Description of Each Team Member

| Team Member                             | Role                       | Description  |
|---|----------------------------|--|
| Marcus Joey Sayner                      | Group Leader and Developer | Coordinated the project and implemented crawling and big data processing using httpx, parsel and Polars. |
| Muhammad Luqman Hakim bin Mohd Rizaudin | Developer                  | Implemented crawling and scraping using urllib and Selectolax, and big data processing using Modin.      |
| Camily Tang Jia Lei                     | Developer                  | Implemented crawling and scraping using urllib and BeautifulSoup and big data processing using Polars.   |
| Goh Jing Yang                           | Developer                  | Implemented crawling and scraping using requests and BeautifulSoup, and big data processing using Dask.  |

## **3.0 Data Collection**

The data collection process involved scraping used car listings from the Carlist.my website. Each team member independently targeted specific pages on the website and extracted structured information such as car brand, model, manufacture year, mileage, location, price, transmission type, body type, fuel type, and other relevant attributes.

To ensure a wide coverage of data, different page ranges were assigned to each member. After the scraping tasks were completed, the individual datasets from all four members were consolidated into a single comprehensive dataset for further transformation and big data processing.

### **3.1 Crawling and Scraping Method**

The web data extraction process consisted of two primary stages: crawling and scraping. Crawling involved systematically navigating through multiple listings pages on the Carlist.my website to retrieve raw HTML content. Scraping referred to the parsing of this content to extract required data fields such as car name, model, and price. To avoid overloading the website's servers, each crawler was configured with page limits, custom user-agent headers, and delays between requests. The source code for the crawling and scraping implementations is provided in Appendix C (Figures C1 to C3).

### **3.2 Number of Records Collected**

Each team member was tasked with collecting approximately 43,000 records through their individual web scraping efforts. After combining the datasets from all members, a total of around 174,150 car listing records were successfully gathered. This large dataset forms the basis for the subsequent data transformation, optimisation, and analysis tasks.

### **3.3 Ethical Considerations**

Throughout the data collection process, ethical web scraping practices were observed:

- **Polite Scraping:** A user-agent header was used to mimic a real browser, and appropriate delays were implemented between requests to avoid overloading the server.
- **No Data Misuse:** Only publicly available data was collected. Any sensitive information, user personal details, or any hidden information was not targeted.
- **Transparency and Attribution:** The source of the data (Carlist.my) is acknowledged in the report, respecting intellectual property rights.

## 4.0 Data Processing

Data processing is considered a critical phase in any data-driven project, during which raw data is cleaned, transformed, and structured to meet the requirements of subsequent analysis. The reliability of the final results is shaped by the quality and consistency of the data, both of which are ensured during this stage. In this project, four libraries—Pandas, Polars, Dask, and Modin—were employed for data processing. Each library is recognised for its distinctive strengths in data handling and manipulation, offering scalable and high-performance solutions suited to datasets of varying sizes.

### 4.1 Cleaning Methods

Data cleaning is regarded as a critical step in ensuring the consistency, accuracy, and usability of a dataset. During the preprocessing stage, various operations were undertaken to resolve redundancies, address missing values, and standardise data formats. As presented in Appendix D (Table D1), key cleaning techniques were applied, including column merging, string manipulation, and the treatment of null entries, in order to prepare the dataset for subsequent analysis.

### 4.2 Data Structure (CSV/JSON/database)

In this project, the capabilities of four robust data processing libraries—Pandas, Polars, Modin, and Dask—were demonstrated in managing structured data from CSV and JSON files. The datasets were cleaned, transformed, and efficiently exported to a MongoDB database. Each library was selected for its optimisation features in data manipulation and processing, with distinct strengths observed:

- **Pandas** was recognised for its simplicity and flexibility in handling small to medium-sized datasets.
- **Polars**, implemented in Rust, was optimised for high performance and memory efficiency.
- **Modin** offered scalability by enabling parallel processing through Dask or Ray with minimal code changes.

- Dask supported adaptive parallel computation, allowing processing of datasets that exceeded memory capacity by dividing tasks into smaller partitions.

To maintain consistency, data in both CSV and JSON formats were processed and stored in MongoDB using a uniform collection schema across all libraries. This uniformity allows for a reliable comparison of performance and resource utilisation across different execution environments.

The entire process reflects the practical applicability of each library within real-world ETL (Extract, Transform, Load) pipelines, facilitating the identification of the most efficient tool for varying data workloads. Figure 2 presents the list of cleaned datasets stored in the MongoDB database.

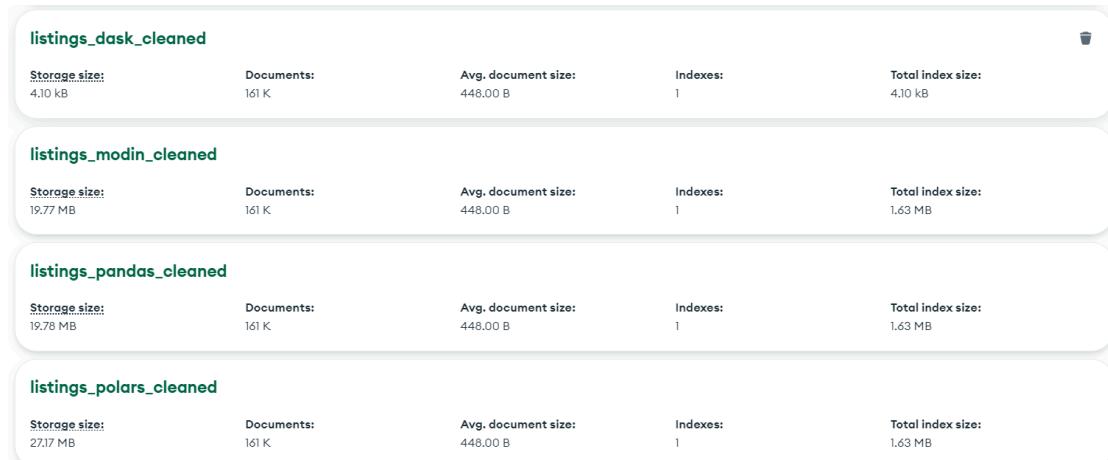


Figure 2: Cleaned Datasets Stored in MongoDB Database

The cleaned datasets from Pandas, Polars, Modin, and Dask were all processed through similar data cleaning workflows: handling missing values, removing irrelevant records, and reordering fields. Each dataset was then transformed into tabular and JSON formats and stored as collections in MongoDB for efficient querying. Figures 3 to 10 display the JSON and tabular format outputs generated using four data processing libraries: Pandas, Polars, Modin, and Dask.

### Pandas (listings\_pandas\_cleaned):

Pandas was used for its straightforward syntax and strong community support. Figures 3 and 4 show the cleaned data in tabular and JSON formats, respectively, stored in MongoDB.

|    | Car Name   | Car Brand       | Car Model       | Manufacture Year | Body Type   | Fuel Ty  |
|----|--|-----------------|-----------------|------------------|-------------|----------|
| 1  | "2023 Lexus RX350 2.4 F Sport SUV"   | "Lexus"         | "RX350"         | 2023             | "SUV"       | "Petrol  |
| 2  | "2010 Toyota Estima 2.4 Aeras MPV Hot Mpv Car In Market / Tip Top Condition" | "Toyota"        | "Estima"        | 2010             | "MPV"       | "Petrol  |
| 3  | "2020 Porsche Cayenne Coupe"   | "Porsche"       | "Cayenne"       | 2020             | "Coupe"     | "Petrol  |
| 4  | "2021 Honda City 1.5 V i-Matic"  | "Honda"         | "City"          | 2021             | "Hatchback" | "Petrol  |
| 5  | "2022 Toyota Corolla Cross"  | "Toyota"        | "Corolla Cross" | 2022             | "SUV"       | "Petrol  |
| 6  | "2018 Mazda CX-5 2.0 SKYActiv-D"   | "Mazda"         | "CX-5"          | 2018             | "SUV"       | "Petrol  |
| 7  | "2020 Porsche Cayenne Coupe"   | "Porsche"       | "Cayenne"       | 2020             | "Coupe"     | "Petrol  |
| 8  | "2020 Porsche Cayenne Coupe"   | "Porsche"       | "Cayenne"       | 2020             | "Coupe"     | "Petrol  |
| 9  | "2012 Proton Preve 1.6 CF"   | "Proton"        | "Preve"         | 2012             | "Sedan"     | "Petrol  |
| 10 | "2015 Nissan X-Trail 2.0."   | "Nissan"        | "X-Trail"       | 2015             | "SUV"       | "Petrol  |
| 11 | "2015 Lamborghini Huracan"   | "Lamborghini"   | "Huracan"       | 2015             | "Coupe"     | "Petrol  |
| 12 | "2020 Porsche Cayenne 4.0."  | "Porsche"       | "Cayenne"       | 2020             | "SUV"       | "Petrol  |
| 13 | "2024 Toyota Alphard 2.4."   | "Toyota"        | "Alphard"       | 2024             | "MPV"       | "Petrol  |
| 14 | "2020 Bentley Flying Spur V8"  | "Bentley"       | "Flying Spur"   | 2020             | "Sedan"     | "Petrol  |
| 15 | "2014 Honda Jazz 1.3 Hybrid"   | "Honda"         | "Jazz"          | 2014             | "Hatchback" | "Hybrid" |
| 16 | "2022 Porsche 911 3.0 Carerra S"   | "Porsche"       | "911"           | 2022             | "Coupe"     | "Petrol  |
| 17 | "2023 Mercedes-Benz C300 Wagon"  | "Mercedes-Benz" | "C300"          | 2023             | "Sedan"     | "Petrol  |
| 18 | "2022 Honda HR-V 1.5 V SUV"  | "Honda"         | "HR-V"          | 2022             | "SUV"       | "Petrol  |

Figure 3: Table-Format View of Pandas-Processed Dataset in MongoDB

```

{
  "_id": {},
  "Car Name": "2023 Lexus RX350 2.4 F Sport SUV",
  "Car Brand": "Lexus",
  "Car Model": "RX350",
  "Manufacture Year": 2023,
  "Body Type": "SUV",
  "Fuel Type": "Petrol - Unleaded (ULP)",
  "Mileage (K KM)": "ns - 10",
  "Transmission": "Automatic",
  "Color": "Black",
  "Price (RM)": 375000,
  "Installment (RM)": 4862,
  "Condition": "Refurbished",
  "Location": "Selangor, Klang",
  "Sales Channel": "Sales Agent",
  "Seat Capacity": 5
}

{
  "_id": {},
  "Car Name": "2010 Toyota Estima 2.4 Aeras MPV Hot Mpv Car In Market / Tip Top Condition",
  "Car Brand": "Toyota",
  "Car Model": "Estima",
  "Manufacture Year": 2010,
  "Body Type": "MPV",
  "Fuel Type": "Petrol - Unleaded (ULP)",
  "Mileage (K KM)": "115 - 120",
  "Transmission": "Automatic",
  "Color": "White",
  "Price (RM)": 100000,
  "Installment (RM)": 726,
  "Condition": "Used",
  "Location": "Selangor, Kajang",
  "Sales Channel": "Sales Agent"
}

```

Figure 4: JSON Output of Pandas-Processed Dataset in MongoDB

## Polars (listings\_polars\_cleaned) :

Polars, written in Rust, offered noticeably faster execution and lower memory usage during transformation. See Figures 5 and 6 for its output formats.

|    | Car Name String                                   | Car Brand String | Car Model String | Manufacture Year Int32 | Body Type String | Fuel T   |
|----|---|------------------|------------------|------------------------|------------------|----------|
| 1  | "2020 MAZDA CX-8 2.5 SKYACTIV-G 4WD"              | "Mazda"          | "CX-8"           | 2020                   | "SUV"            | "Petrol" |
| 2  | "2011 Toyota Land Cruiser Prado 4.0 V6 VVT-i 4WD" | "Toyota"         | "Land Cruiser"   | 2011                   | "SUV"            | "Petrol" |
| 3  | "2021 Toyota Innova Crysta 2.0 V CVT 4WD"         | "Toyota"         | "Innova"         | 2021                   | "MPV"            | "Petrol" |
| 4  | "2022 Honda N-Box Custom 1.2 V MT 4WD"            | "Honda"          | "N-Box Custom"   | 2022                   | "Hatchback"      | "Petrol" |
| 5  | "2017 Mercedes-Benz GLC250 4MATIC 2.0 V6 VGT 4WD" | "Mercedes-Benz"  | "GLC250"         | 2017                   | "SUV"            | "Petrol" |
| 6  | "2015 Mazda 2 1.5 SKYACTIV-G 4WD"                 | "Mazda"          | "#2"             | 2015                   | "Hatchback"      | "Petrol" |
| 7  | "2020 BMW M135i 2.0 xDrive 4WD"                   | "BMW"            | "M135i"          | 2020                   | "Hatchback"      | "Petrol" |
| 8  | "2019 Porsche Macan 2.0 PDK 4WD"                  | "Porsche"        | "Macan"          | 2019                   | "SUV"            | "Petrol" |
| 9  | "2016 Honda CR-V 2.0 i-VTEC 4WD"                  | "Honda"          | "CR-V"           | 2016                   | "SUV"            | "Petrol" |
| 10 | "2015 FORD FIESTA 1.5 FULL AUTOMATIC"             | "Ford"           | "Fiesta"         | 2015                   | "Hatchback"      | "Petrol" |
| 11 | "2015 Kia Sorento 2.4 SUV"                        | "Kia"            | "Sorento"        | 2015                   | "SUV"            | "Petrol" |
| 12 | "2021 Toyota Vellfire 3.5 Elegance 4WD"           | "Toyota"         | "Vellfire"       | 2021                   | "MPV"            | "Petrol" |
| 13 | "2023 Toyota Alphard 2.5 Elegance 4WD"            | "Toyota"         | "Alphard"        | 2023                   | "MPV"            | "Petrol" |
| 14 | "2018 Perodua Bezza 1.3 X 4WD"                    | "Perodua"        | "Bezza"          | 2018                   | "Sedan"          | "Petrol" |
| 15 | "2018 Toyota Land Cruiser Prado 4.0 V6 VGT 4WD"   | "Toyota"         | "Land Cruiser"   | 2018                   | "SUV"            | "Petrol" |
| 16 | "2016 Toyota Avanza 1.5 S 4WD"                    | "Toyota"         | "Avanza"         | 2016                   | "MPV"            | "Petrol" |
| 17 | "2018 Mercedes-Benz C350e 2.0 V6 VGT 4WD"         | "Mercedes-Benz"  | "C350 e"         | 2018                   | "Sedan"          | "Hybrid" |
| 18 | "2016 Toyota Vellfire 2.5 Elegance 4WD"           | "Toyota"         | "Vellfire"       | 2016                   | "MPV"            | "Petrol" |

Figure 5: Table-Format View of Polars-Processed Dataset in MongoDB

```

{
  "_id": ObjectId('6823e5a225dbf9874782dab0'),
  "Car Name": "2022 Toyota Corolla Cross 1.8 V SUV",
  "Car Brand": "Toyota",
  "Car Model": "Corolla Cross",
  "Manufacture Year": 2022,
  "Body Type": "SUV",
  "Fuel Type": "Petrol - Unleaded (ULP)",
  "Mileage (K KM)": "30 - 35",
  "Transmission": "Automatic",
  "Color": "White",
  "Price (RM)": 106800,
  "Installment (RM)": 1385,
  "Condition": "Used",
  "Location": "Johor, Johor Bahru",
  "Sale Channel": "Dealer",
  "Seat Capacity": 5
}

{
  "_id": ObjectId('6823e5a225dbf9874782dab0'),
  "Car Name": "2019 Mazda 3 2.0 SKYACTIV-G High Plus Hatchback NO PROCESSING FEE LOW HARGA",
  "Car Brand": "Mazda",
  "Car Model": "3",
  "Manufacture Year": 2019,
  "Body Type": "Hatchback",
  "Fuel Type": "Petrol - Unleaded (ULP)",
  "Mileage (K KM)": "50 - 55",
  "Transmission": "Automatic",
  "Color": "Grey",
  "Price (RM)": 77800
}

```

Figure 6: JSON Output of Polars-Processed Dataset in MongoDB

## Modin (listings\_modin\_cleaned):

Modin enabled parallel processing with minimal code changes from Pandas, benefiting from scalability. Figures 7 and 8 show its processed outputs.

|    | Car Name   | String | Car Brand       | String | Car Model       | String | Manufacture Year | Int32 | Body Type   | String | Fuel Type                 |
|----|--|--------|-----------------|--------|-----------------|--------|------------------|-------|-------------|--------|---------------------------|
| 1  | "2023 Lexus RX350 2.4 F Sport SUV"   |        | "Lexus"         |        | "RX350"         |        | 2023             |       | "SUV"       |        | "Petrol - Unleaded (ULP)" |
| 2  | "2010 Toyota Estima 2.4 Aeras MPV Hot Mpv Car In Market / Tip Top Condition" |        | "Toyota"        |        | "Estima"        |        | 2010             |       | "MPV"       |        | "Petrol - Unleaded (ULP)" |
| 3  | "2020 Porsche Cayenne Coupe"   |        | "Porsche"       |        | "Cayenne"       |        | 2020             |       | "Coupe"     |        | "Petrol - Unleaded (ULP)" |
| 4  | "2021 Honda City 1.5 V i-Matic"  |        | "Honda"         |        | "City"          |        | 2021             |       | "Hatchback" |        | "Petrol - Unleaded (ULP)" |
| 5  | "2022 Toyota Corolla Cross"  |        | "Toyota"        |        | "Corolla Cross" |        | 2022             |       | "SUV"       |        | "Petrol - Unleaded (ULP)" |
| 6  | "2018 Mazda CX-5 2.0 SKYACTIV-G 4WD"   |        | "Mazda"         |        | "CX-5"          |        | 2018             |       | "SUV"       |        | "Petrol - Unleaded (ULP)" |
| 7  | "2020 Porsche Cayenne Coupe"   |        | "Porsche"       |        | "Cayenne"       |        | 2020             |       | "Coupe"     |        | "Petrol - Unleaded (ULP)" |
| 8  | "2020 Porsche Cayenne Coupe"   |        | "Porsche"       |        | "Cayenne"       |        | 2020             |       | "Coupe"     |        | "Petrol - Unleaded (ULP)" |
| 9  | "2012 Proton Preve 1.6 G-Power"  |        | "Proton"        |        | "Preve"         |        | 2012             |       | "Sedan"     |        | "Petrol - Unleaded (ULP)" |
| 10 | "2015 Nissan X-Trail 2.0 AWD"  |        | "Nissan"        |        | "X-Trail"       |        | 2015             |       | "SUV"       |        | "Petrol - Unleaded (ULP)" |
| 11 | "2015 Lamborghini Huracan"   |        | "Lamborghini"   |        | "Huracan"       |        | 2015             |       | "Coupe"     |        | "Petrol - Unleaded (ULP)" |
| 12 | "2020 Porsche Cayenne 4.0 V8 GTS"  |        | "Porsche"       |        | "Cayenne"       |        | 2020             |       | "SUV"       |        | "Petrol - Unleaded (ULP)" |
| 13 | "2024 Toyota Alphard 2.4 AWD"  |        | "Toyota"        |        | "Alphard"       |        | 2024             |       | "MPV"       |        | "Petrol - Unleaded (ULP)" |
| 14 | "2020 Bentley Flying Spur V8"  |        | "Bentley"       |        | "Flying Spur"   |        | 2020             |       | "Sedan"     |        | "Petrol - Unleaded (ULP)" |
| 15 | "2014 Honda Jazz 1.3 Hybrid"   |        | "Honda"         |        | "Jazz"          |        | 2014             |       | "Hatchback" |        | "Hybrid - Unleaded (ULP)" |
| 16 | "2022 Porsche 911 3.0 Carrera S"   |        | "Porsche"       |        | "911"           |        | 2022             |       | "Coupe"     |        | "Petrol - Unleaded (ULP)" |
| 17 | "2023 Mercedes-Benz C300 2.0 AWD"  |        | "Mercedes-Benz" |        | "C300"          |        | 2023             |       | "Sedan"     |        | "Petrol - Unleaded (ULP)" |
| 18 | "2022 Honda HR-V 1.5 V SUV"  |        | "Honda"         |        | "HR-V"          |        | 2022             |       | "SUV"       |        | "Petrol - Unleaded (ULP)" |

Figure 7: Table-Format View of Modin-Processed Dataset in MongoDB

```

[{"_id": "5f3e0d000000000000000001", "Car Name": "2023 Lexus RX350 2.4 F Sport SUV", "Car Brand": "Lexus", "Car Model": "RX350", "Manufacture Year": 2023, "Body Type": "SUV", "Fuel Type": "Petrol - Unleaded (ULP)", "Mileage (K KM)": "n/a - 10", "Transmission": "Automatic", "Color": "Black", "Price (RM)": 37880, "Instalment (RM)": 4862, "Condition": "Refurbished", "Location": "Selangor, Klang ", "Sales Channel": "Sales Agent", "Seat Capacity": 5}, {"_id": "5f3e0d000000000000000002", "Car Name": "2010 Toyota Estima 2.4 Aeras MPV Hot Mpv Car In Market / Tip Top Condition", "Car Brand": "Toyota", "Car Model": "Estima", "Manufacture Year": 2010, "Body Type": "MPV", "Fuel Type": "Petrol - Unleaded (ULP)", "Mileage (K KM)": "115 - 120", "Transmission": "Automatic", "Color": "White", "Price (RM)": 55999, "Instalment (RM)": 726, "Condition": "Used", "Location": "Selangor, Kajang ", "Sales Channel": "Sales Agent"}]

```

Figure 8: JSON Output of Modin-Processed Dataset in MongoDB

## Dask (listings\_dask\_cleaned):

Dask efficiently handled large datasets by processing them in partitions. Figures 9 and 10 show its outputs stored in MongoDB.

|    | Car Name  | Car Brand       | Car Model       | Manufacture Year | Body Type   | Fuel Ty |
|----|---|-----------------|-----------------|------------------|-------------|---------|
| 1  | "2021 Honda City 1.5 V i-VTEC Hatchback"  | "Honda"         | "City"          | 2021             | "Hatchback" | "Petrol |
| 2  | "2022 Toyota Corolla Cross 1.8 V SUV Full Service Record Under Warranty By Toyota Malaysia" | "Toyota"        | "Corolla Cross" | 2022             | "SUV"       | "Petrol |
| 3  | "2014 Inokom Elantra 1.6 M CVT Sedan"   | "Inokom"        | "Elantra"       | 2014             | "Sedan"     | "Petrol |
| 4  | "2023 Toyota Veloz 1.5 MPV"   | "Toyota"        | "Veloz"         | 2023             | "MPV"       | "Petrol |
| 5  | "2019 Perodua Myvi 1.3 AVT Sedan"   | "Perodua"       | "Myvi"          | 2019             | "Hatchback" | "Petrol |
| 6  | "2020 Porsche Cayenne Coupe 4.0 V8 GTS SUV"   | "Porsche"       | "Cayenne"       | 2020             | "Coupe"     | "Petrol |
| 7  | "2019 Perodua Bezza 1.0 G MT Sedan"   | "Perodua"       | "Bezza"         | 2019             | "Sedan"     | "Petrol |
| 8  | "2018 Honda BR-V 1.5 V i-VTEC SUV"  | "Honda"         | "BR-V"          | 2018             | "SUV"       | "Petrol |
| 9  | "2011 Toyota Vios 1.5 G Sedan"  | "Toyota"        | "Vios"          | 2011             | "Sedan"     | "Petrol |
| 10 | "2021 BMW X3 2.0 xDrive30i SUV"   | "BMW"           | "X3"            | 2021             | "SUV"       | "Petrol |
| 11 | "2011 Mazda CX-9 3.7 Gateau Sedan"  | "Mazda"         | "CX-9"          | 2011             | "SUV"       | "Petrol |
| 12 | "2006 Mercedes-Benz S350 L Sedan"   | "Mercedes-Benz" | "S350"          | 2006             | "Sedan"     | "Petrol |
| 13 | "2023 Proton Saga 1.3 Sedan"  | "Proton"        | "Saga"          | 2023             | "Sedan"     | "Petrol |
| 14 | "2016 Proton Perdana 2.0 Sedan"   | "Proton"        | "Perdana"       | 2016             | "Sedan"     | "Petrol |
| 15 | "2024 Lexus RX350 2.4 Luxury SUV"   | "Lexus"         | "RX350"         | 2024             | "SUV"       | "Petrol |
| 16 | "2012 Volkswagen Passat 1.8 TSI Sedan"  | "Volkswagen"    | "Passat"        | 2012             | "Sedan"     | "Petrol |
| 17 | "2016 Toyota Vios 1.5 TRD Sedan"  | "Toyota"        | "Vios"          | 2016             | "Sedan"     | "Petrol |
| 18 | "2015 Proton Preve 1.6 Executive Sedan"   | "Proton"        | "Preve"         | 2015             | "Sedan"     | "Petrol |

Figure 9: Table-Format View of Dask-Processed Dataset in MongoDB

```

[{"id": "2021 Honda City 1.5 V i-VTEC Hatchback", "Car Name": "2021 Honda City 1.5 V i-VTEC Hatchback", "Car Brand": "Honda", "Car Model": "City", "Manufacture Year": 2021, "Body Type": "Hatchback", "Fuel Type": "Petrol - Unleaded (ULP)", "Mileage (K KM)": "99 - 95", "Transmission": "Automatic", "Color": "Silver", "Price (RM)": 67600, "Instalment (RM)": 669, "Condition": "Used", "Location": "Dutamas, Ulu Tiram ", "Sales Channel": "Sales Agent", "Seat Capacity": 5}, {"id": "2022 Toyota Corolla Cross 1.8 V SUV Full Service Record Under Warranty By Toyota Malaysia", "Car Name": "2022 Toyota Corolla Cross 1.8 V SUV Full Service Record Under Warranty By Toyota Malaysia", "Car Brand": "Toyota", "Car Model": "Corolla Cross", "Manufacture Year": 2022, "Body Type": "SUV", "Fuel Type": "Petrol - Unleaded (ULP)", "Mileage (K KM)": "80 - 85", "Transmission": "Automatic", "Color": "White", "Price (RM)": 98999, "Instalment (RM)": 1283, "Condition": "Used", "Location": "Selangor, Kajang ", "Sales Channel": "Sales Agent"}]

```

Figure 10: JSON Output of Dask-Processed Dataset in MongoDB

## 4.3 Transformation and Formatting

### 4.3.1 Combining Redundant Columns

**Target Columns:** 'Seat Capacity' and 'Seating Capacity'

**Objective:** To eliminate semantic redundancy by merging similar columns.

**Process:** In instances where the 'Seat Capacity' column was null but 'Seating Capacity' contained a value, the value from 'Seating Capacity' was used to fill in the missing data. If 'Seat Capacity' already had a value, it was retained. After the merging process, the 'Seating Capacity' column was dropped to avoid redundancy. Figure 11 illustrates the Pandas code implementation used to perform this column combination and removal.

```
# === Data Cleaning Step: Combine 'Seat Capacity' and 'Seating Capacity' ===
total_rows_combinecol_pandas = df_pandas.shape[0]

# Combine 'Seat Capacity' and 'Seating Capacity' into one column
df_pandas['Seat Capacity'] = df_pandas['Seat Capacity'].combine_first(df_pandas['Seating Capacity'])

# Drop 'Seating Capacity' column
df_pandas.drop(columns=['Seating Capacity'], inplace=True)
```

Figure 11: Pandas Code for Combining Redundant Columns: 'Seat Capacity' and 'Seating Capacity'

### 4.3.2 Removing Noisy and Irrelevant Columns

**Target Column:** 'URL'

**Objective:** To eliminate noisy and irrelevant data.

**Method:** The 'URL' column was identified as analytically irrelevant and subsequently removed to reduce dataset dimensionality and eliminate noise. This operation is illustrated in Figure 12.

```

# === Data Cleaning Step: Remove 'URL' column ===
total_rows_removecol_pandas = df_pandas.shape[0]

# =====

# Remove the 'URL' column
df_pandas.drop(columns=['URL'], inplace=True)

# =====

```

Figure 12: Pandas Code for Removing Noisy Column: 'URL'

### 4.3.3 String Cleaning and Formatting

This section details the preprocessing steps applied to key textual columns to ensure consistency and enable numerical analysis. It covers cleaning monetary strings in the Installment column, extracting categorical labels from schema URLs in the Condition column, standardising sales source identifiers in the Sales Channel column, and binning mileage values into defined intervals. Code examples demonstrating these transformations are provided in Figures 13 through 16.

#### 4.3.3.1 Installment Column

**Objective:** To clean monetary values for numeric operations.

**Method:** Monetary string values were cleaned by removing the suffix '/month', the prefix 'RM', and any commas. The cleaned strings were then converted to numeric format to enable mathematical operations. The code for this conversion is presented in Figure 13.

```

# === Data Cleaning: Clean 'Installment' Column ===
total_rows_cleaninstallment_pandas = df_pandas.shape[0]

# Cleaning operations
df_pandas['Installment'] = df_pandas['Installment'].str.replace('RM', "", regex=False)
df_pandas['Installment'] = df_pandas['Installment'].str.replace(',', "", regex=False)
df_pandas['Installment'] = df_pandas['Installment'].str.replace('/month', "", regex=False)
df_pandas['Installment'] = pd.to_numeric(df_pandas['Installment'], errors='coerce')

```

Figure 13: Pandas Code for Cleaning the 'Installment' Column

#### 4.3.3.2 Condition Column

**Objective:** To extract simple categorical labels from schema URLs.

**Method:** Schema URLs in the Condition column were parsed to extract simplified categorical labels such as 'Used' or 'New' (e.g., from <https://schema.org/UsedCondition>). This transformation of schema-based condition values into human-readable categories is reflected in Figure 14.

```
# === Data Cleaning: Extract Clean Values from 'Condition' Column ===
total_rows_condition_pandas = df_pandas.shape[0]

# Extract value from "...schema.org/XCondition"
df_pandas['Condition'] = df_pandas['Condition'].apply(
    lambda x: re.search(r'\.org/([A-Za-z]+)Condition', x).group(1) if isinstance(x, str) and 'schema.org/' in x else x
)

# Strip 'Condition' if it remains (for post-processed cases)
df_pandas['Condition'] = df_pandas['Condition'].str.replace('Condition$', '', regex=True)
```

Figure 14: Pandas Code for Extracting Categorical Values from the 'Condition' Column

#### 4.3.3.3 Sales Channel Column

**Objective:** To simplify and standardise sales source identifiers.

**Method:** Textual content in the *Sales Channel* column was parsed and categorised using regular expressions to standardise valid entries such as 'Sales Agent' or 'Dealer'. This transformation process is demonstrated in Figure 15.

```
# === Data Cleaning: Clean 'Sales Channel' Column ===
total_rows_saleschannel_pandas = df_pandas.shape[0]

# Extract Sales Channel (either 'Sales Agent' or 'Dealer')
df_pandas['Sales Channel'] = df_pandas['Sales Channel'].str.extract(r'^Sales Agent|Dealer')
```

Figure 15: Pandas Code for Standardising the 'Sales Channel' Column

#### 4.3.3.4 Mileage Column

**Objective:** To standardise mileage ranges into bins.

**Method:** Suffixes such as 'K KM' were removed, and mileage ranges (e.g., '5 - 10') were parsed. Single values were binned into predefined intervals (e.g., 12 was binned into '10 - 15'). The binning logic is presented in Figure 16.

```

# === Data Cleaning: Clean 'Mileage' Column ===
total_rows_mileage_pandas = df_pandas.shape[0]

# Remove "K KM" and clean mileage into 5K step ranges
def process_mileage_pandas(mileage):
    if not isinstance(mileage, str) or mileage.strip() == '':
        return None
    try:
        if '-' in mileage:
            start, end = mileage.split('-')
            return f"{int(start.strip())} - {int(end.strip())}"
        else:
            val = int(mileage.strip())
            lower = (val // 5) * 5
            upper = lower + 5
            return f"{lower} - {upper}"
    except:
        return None

df_pandas['Mileage'] = df_pandas['Mileage'].str.replace('K KM', '', regex=False)
df_pandas['Mileage'] = df_pandas['Mileage'].apply(process_mileage_pandas)

```

Figure 16: Pandas Code for Binning Values in the 'Mileage' Column

#### 4.3.4 Handling Missing Data

**Objective:** To retain only rows with complete critical data.

**Critical Columns Checked:** 'Body Type', 'Fuel Type', 'Mileage', 'Sales Channel', and 'Seat Capacity'

**Method:** Rows containing null values in any of the critical columns were removed using the dropna() function. This approach was taken to preserve data integrity without introducing imputation bias. The implementation is detailed in Figure 17.

```

# =====
# === Data Cleaning Step: Drop rows with missing values ===
total_rows_before_cleaning_null_pandas = df_pandas.shape[0]

# Drop rows with missing values in specific columns
df_pandas.dropna(subset=['Body Type', 'Fuel Type', 'Mileage', 'Sales Channel', 'Seat Capacity'], inplace=True)

# =====

```

Figure 17: Pandas Code for Handling Missing Data in Critical Columns

#### 4.3.5 Renaming Columns for Clarity

**Objective:** To add units to numeric columns for readability.

**Renaming Map:**

'Mileage' → 'Mileage (K KM)'

'Price' → 'Price (RM)'

'Installment' → 'Installment (RM)'

**Method:** Figure 18 shows rename columns of the dataset by appending units to avoid ambiguity: For instance, Price to Price (RM), and Mileage to Mileage (K KM). This kind of detailing unquestionably refines the understanding of the data by specifying the units of measurement.

```
# =====
df_pandas.rename(columns={
    'Mileage': 'Mileage (K KM)',
    'Price': 'Price (RM)',
    'Installment': 'Installment (RM)'
}, inplace=True)
# =====
```

Figure 18: Pandas Code for Renaming Columns for Clarity

#### 4.3.6 Removing Duplicate Rows

**Objective:** To prevent duplicate entries from distorting analysis.

**Method:** Duplicate entries were identified using `.duplicated()` and removed using `.drop_duplicates()`, ensuring that exact duplicates were eliminated while retaining valid lookalikes. The applied logic is captured in Figure 19.

```
# == Drop Duplicate Rows ==
df_pandas = df_pandas.drop_duplicates()
# =====
```

Figure 19: Pandas Code for Removing Duplicate Rows

## **5.0 Optimisation Techniques**

In data preprocessing, optimisation techniques are considered crucial for efficiently handling large-scale datasets. The performance of a data pipeline can be improved by selecting libraries that offer optimisation capabilities such as parallel computing and multi-threading. In this project, the performance of four Python-based data processing libraries—Pandas, Polars, Modin, and Dask—is investigated by applying equivalent cleaning and transformation operations to a combined car listings dataset. Each library is evaluated based on execution time (wall time), CPU usage and CPU time, memory usage (current and peak), and processing throughput (rows per second). Through this evaluation, the objective is to determine how different libraries optimise resource usage and processing time when identical data cleaning operations—such as column combination, string manipulation, and handling of missing values—are performed.

### **5.1 Methods Used**

To ensure a fair comparison, identical data cleaning steps were applied across all four libraries. The performance optimisation techniques specific to each library are detailed in Appendix E (Table E1).

### **5.2 Code Overview**

This section presents the code implementations used to apply optimisation techniques in each data processing library. While the data cleaning and transformation steps were standardised across all libraries, the code was written separately to accommodate the unique syntax and execution model of each library. The content is organised into three parts: the task performed, the library used, and the corresponding code snippet. Screenshots are included for reference, and a summary of all tasks and associated libraries is provided in Appendix E (Table E2).

## 6.0 Performance Evaluation

Performance evaluation is conducted to assess the effectiveness of the optimisation techniques applied in this project. This evaluation focuses specifically on the data processing component by comparing how different libraries perform identical data cleaning tasks under various optimisation methods. The objectives of the performance evaluation are to achieve:

- Reduced processing time
- Lower memory consumption
- Improved CPU utilisation
- Higher throughput (i.e., more records processed per minute)

### 6.1 Before vs After Optimisation

This section presents a comparison of the performance before and after applying optimisation techniques using different libraries: Pandas, Polars, Modin, and Dask.

- Pandas is used as the baseline, representing a non-optimised, single-threaded processing approach.
- In contrast, Polars, Modin, and Dask represent optimised solutions.
  - Polars utilises multi-threaded execution and is implemented in Rust to achieve high performance.
  - Modin introduces parallel execution by distributing tasks across multiple CPU cores using backend engines such as Dask.
  - Dask employs task-based parallelism with lazy evaluation, enabling it to efficiently scale data processing workflows.

By comparing the results obtained from performing the same data cleaning tasks, the impact of these optimisation techniques on overall system performance can be assessed.

### 6.2 Time, Memory, CPU Usage, Throughput

To evaluate performance, six metrics were collected during data processing tasks using each library. These metrics were consistently measured using the Python libraries time, psutil, and tracemalloc. The metrics include:

- **Elapsed Time (Wall Time) (s):** The real-world time taken to complete the task from start to finish.
- **CPU Time (ms):** The actual CPU processing time consumed, measured in milliseconds.
- **CPU Usage (%):** The intensity of CPU usage during code execution.
- **Throughput (rows/sec):** The number of records processed per second, calculated as total rows divided by elapsed time.
- **Current Memory Usage (MB):** The amount of memory in use at the end of execution.
- **Peak Memory Usage (MB):** The highest amount of memory used at any point during execution.

The compiled results for these metrics across all libraries are summarised in Appendix F (Table F1).

### 6.3 Charts and Graphs

Figure 20 presents a comparison of four Python data-processing libraries—Pandas, Polars, Modin, and Dask—across a typical ETL pipeline. Polars is observed to perform the fastest overall, particularly in the parsing, bucketing, and MongoDB export stages. Pandas demonstrates the quickest performance in data import and simple column operations but shows moderate efficiency on more complex string-cleaning tasks. Higher overheads are incurred by Modin and Dask, especially during import, null-filling, and memory optimisation steps, rendering them more beneficial primarily at larger scales or within distributed environments. Consequently, Polars is identified as the preferred choice for accelerating most in-memory workflows, whereas Pandas remains advantageous for straightforward, single-machine operations.

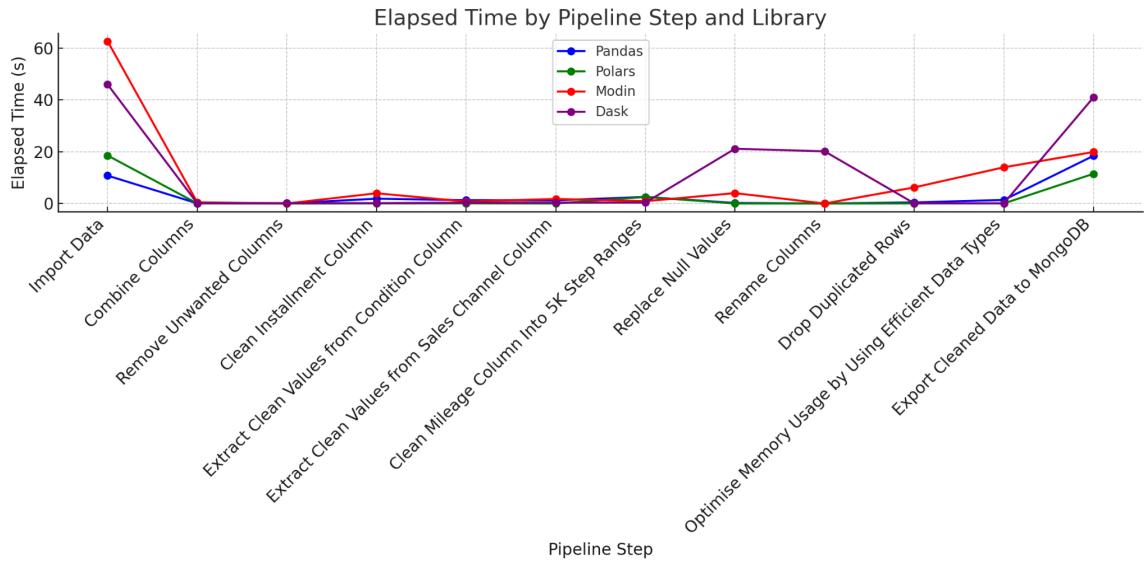


Figure 20: Line Graph of Elapsed Time by Pipeline Step and Library

Figure 21 shows that Polars is by far the most CPU-efficient, using the fewest milliseconds for virtually every ETL step, while Pandas is positioned in the middle, with low overhead on simple operations but increasing on heavier parsing and export. Huge CPU cycles are consumed by Modin on import (~180,000 ms) and memory tuning, after which the usage levels off to a similar but still higher amount than Pandas on string-cleaning and MongoDB writes. Significant CPU overhead is also observed in Dask on broadcast-style tasks (null-replacement, memory optimisation) and export, which reflects its chunked and distributed scheduling costs. In summary, the leanest CPU footprint is delivered by Polars, Pandas is considered reasonable for small-scale in-memory work, and Modin and Dask only prove advantageous when distributed throughput is required.

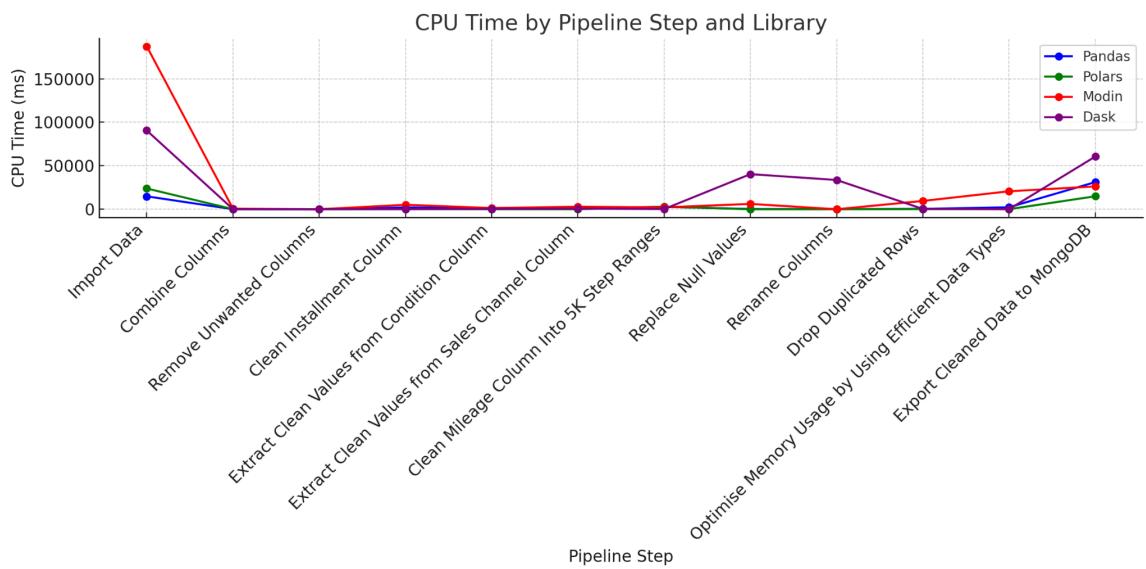


Figure 21: Line Graph of CPU Time by Pipeline Step and Library

Figure 22 demonstrates that Polars consistently maintains the smallest and steadiest CPU footprint, with utilisation pegged between approximately 6% and 16% throughout the pipeline, which is attributed to its efficient, multi-threaded Rust core. Pandas is positioned in the middle, consuming about 7%–20% of CPU with no dramatic spikes, reflecting its straightforward, single-threaded execution. By contrast, utilisation is driven above 25% by Modin’s parallel scheduler on simple combines and approaches 45% when dropping duplicates (with 12%–20% elsewhere), while Dask likewise increases to around 20% for broadcast-style tasks such as null replacement and memory down-casting (and otherwise remains near 8%–15%). In practice, Polars is identified as the most CPU-efficient choice, Pandas offers stable, moderate usage, and Modin and Dask justify their additional overhead only when parallel or distributed throughput is truly required.

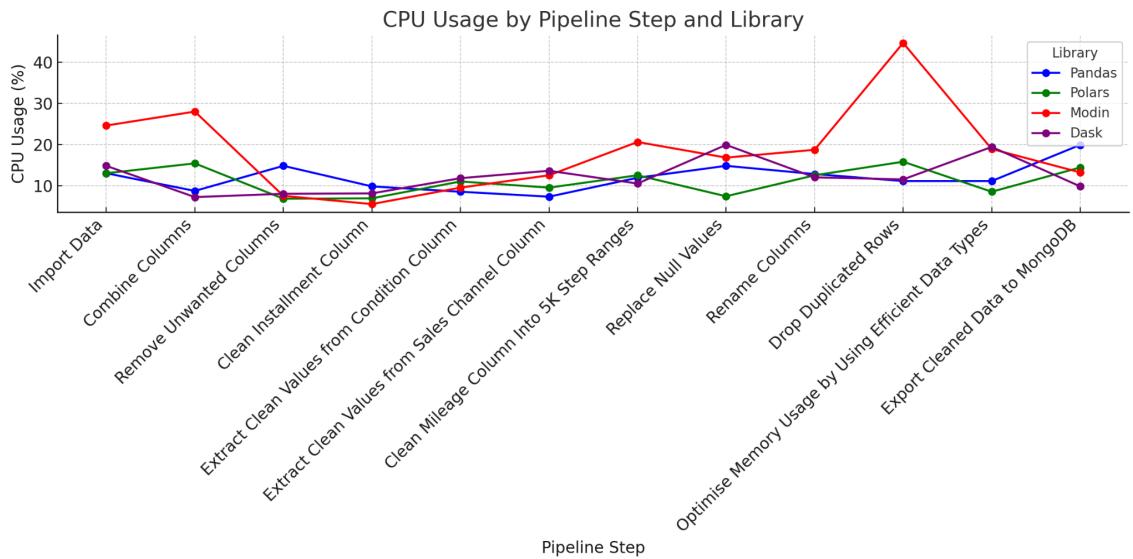


Figure 22: Line Graph of CPU Usage by Pipeline Step and Library

Figure 23 illustrates the number of rows processed per second by each library at various ETL steps. Polars significantly outperforms the others, achieving approximately 150 million rows per second during string-cleaning and renaming tasks and maintaining 80–90 million rows per second during bucketing and channel-parsing. Pandas ranks second, peaking at around 120 million rows per second during the efficient dtype-conversion step and processing 30–50 million rows per second for most cleaning operations. Moderate performance improvements are observed with Modin’s parallel engine, which processes approximately 30–40 million rows per second on intensive parsing and renaming tasks, though it does not reach the performance levels of Polars. Dask, when executed on a single node, exhibits significantly lower throughput (generally under 1 million rows per second) due to the overhead imposed by its chunked scheduling approach on these in-memory workloads. In summary, Polars demonstrates the highest throughput, Pandas performs respectably, Modin provides some additional speed, and Dask proves advantageous primarily when scaled across multiple machines.

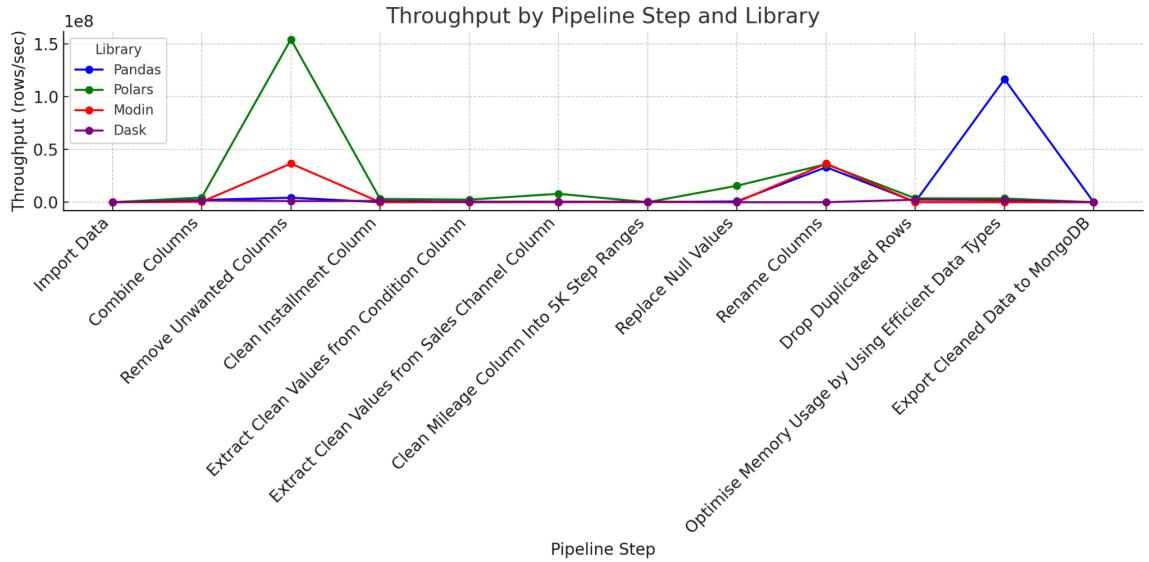


Figure 23: Line Graph of Throughput by Pipeline Step and Library

Figure 24 presents the current memory usage by pipeline step for each library. After loading the CSV file, all four libraries occupy approximately 420–450 MB of RAM; however, their steady-state memory footprints diverge markedly during subsequent transformations. Polars promptly releases its buffers and maintains an additional resident memory usage close to 0 MB, consistent with its streaming architecture. Pandas sustains approximately 20–30 MB of memory usage throughout the cleaning steps, whereas Modin generally holds a slightly lower footprint of around 10–15 MB during most string parsing and renaming operations, increasing to roughly 15 MB during export. Dask’s memory usage fluctuates, alternating between near-zero usage on simple column operations and a per-chunk overhead, ranging from approximately 5 MB during null replacement to up to 25 MB during the final MongoDB write. Overall, Polars demonstrates the most minimal working set, Pandas and Modin maintain moderate memory buffers, and Dask’s footprint varies in accordance with its chunked execution model.

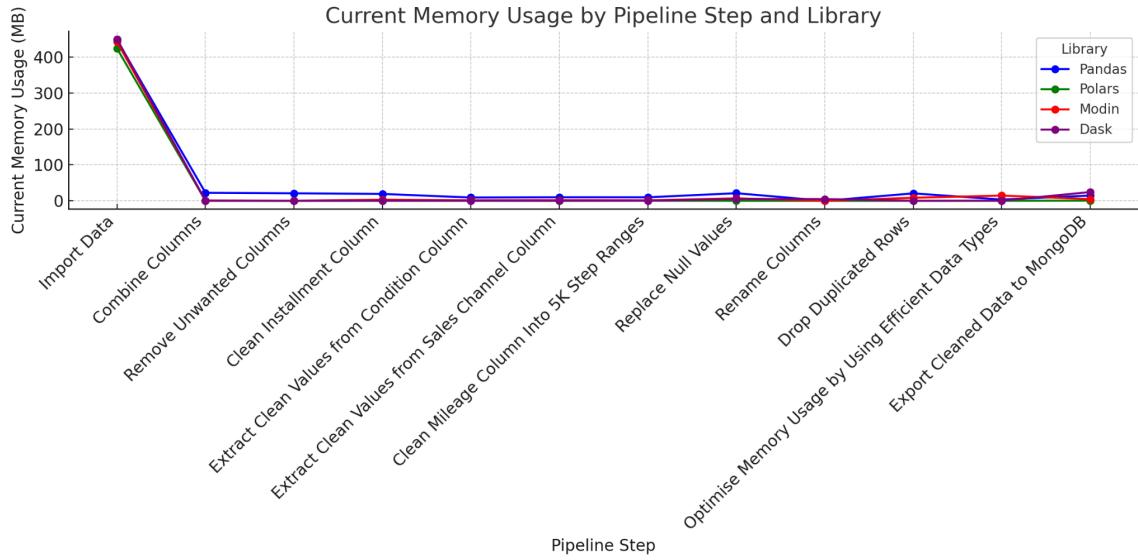


Figure 24: Line Graph of Current Memory Usage by Pipeline Step and Library

Figure 25 illustrates the peak memory usage by pipeline step for each library. It is evident that reading the CSV file constitutes the largest single memory allocation, with Dask exhibiting a peak of approximately 800 MB, Modin 550 MB, Pandas 520 MB, and Polars 430 MB. Following this initial load, Polars demonstrates the smallest additional memory peaks, ranging from single-digit to low-teens MB during string parsing and bucketing operations. In contrast, Pandas and Modin reach peak usages between 15 and 35 MB throughout their various cleaning steps. Dask, however, maintains a substantial peak of 340 MB during broadcast-style operations such as null replacement and memory down-casting, before decreasing to near zero during simpler column operations. During the MongoDB export phase, notable transient peaks are again observed, with Dask and Polars reaching approximately 340 MB and 270 MB, respectively, while Modin and Pandas peak at 200 MB and 180 MB. All in all, Polars minimises additional memory consumption following the initial load, Pandas and Modin incur moderate buffer peaks, and Dask exhibits the highest transient overhead within this single-node ETL workflow.

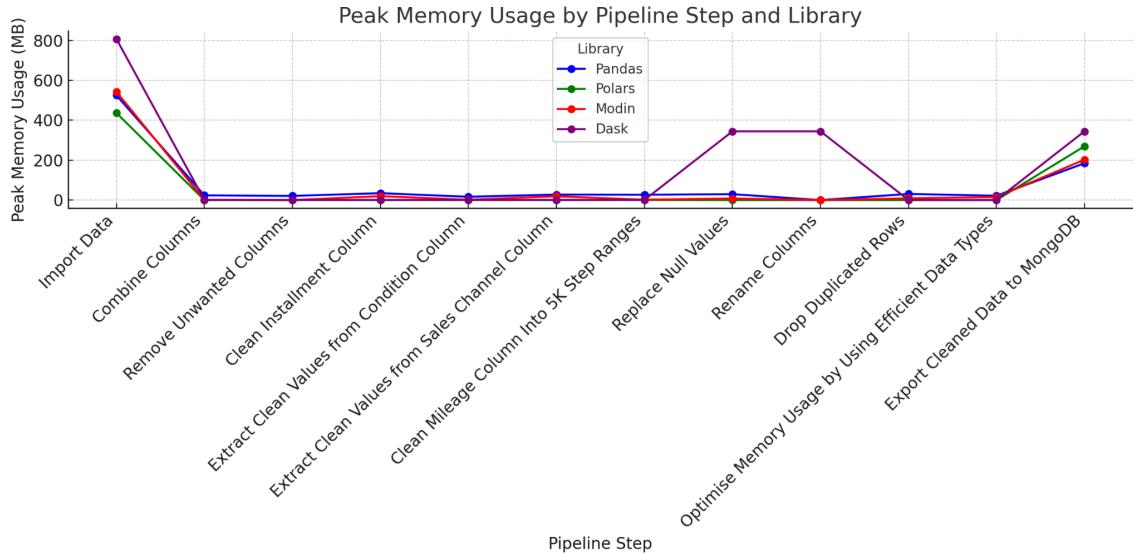


Figure 25: Line Graph of Peak Memory Usage by Pipeline Step and Library

## **7.0 Challenges and Limitations**

Several difficulties were encountered when scraping automobile listings using Visual Studio Code (VSCode). One major issue was that the script would unexpectedly stop, causing VSCode to freeze—especially during attempts to scrape large volumes of data from the website. This inability to complete the task efficiently significantly hampered both progress and overall productivity.

In addition, web scraping was limited by restrictions in robots.txt and the availability of accessible records. The robots.txt file may block certain paths, rendering it illegal or unethical to scrape specific types of content. Furthermore, some websites implement poor pagination or display only a limited number of items per page, making it difficult to collect large datasets effectively.

Another challenge stemmed from the incompatibility of some Python libraries with the Google Colab environment. Libraries such as Playwright and Selenium encountered installation or runtime issues due to Colab's limited support for headless browsers and restricted system-level operations. Additionally, Colab's memory and session time limits imposed further constraints when processing large volumes of scraped data.

## **7.1 Reasons of Challenges and Limitations**

The freezing issue in VSCode was likely caused by the high resource demands of running the scraping script directly within the IDE. The integrated terminal in VSCode struggled to handle the large amount of data, especially when real-time output was printed extensively. As a result, VSCode became unresponsive during execution.

The restriction imposed by robots.txt exists to manage web traffic, protect user data, and prevent server overload. Website owners use this file to define which bots may access their content. Limited accessible records may also result from dynamic content loading via JavaScript, API constraints, or intentional design choices aimed at discouraging scraping.

Google Colab, although effective for quick data analysis and learning machine learning models, does not perform well when high computing power or large-scale data processing is required. Its lack of support for running long background jobs and modifying system-wide settings limits the capabilities of tools that require persistent browser control. In addition, Colab's RAM is typically restricted to 12–16 GB, and execution time is limited, which complicates long-running scraping and data-gathering tasks.

## 7.2 Solutions

To address the freezing issue encountered in Visual Studio Code (VSCode), several adjustments were implemented. The scraping process was divided into two smaller tasks. Each request was delayed by approximately two to five seconds to reduce the workload on both the target server and the local system. Data was incrementally saved between each batch to manage memory usage and prevent overload. The resulting CSV files from each segment were subsequently merged to create a single comprehensive dataset. These modifications effectively mitigated system failures and improved overall performance.

Ethical scraping practices were observed by adhering to the website's usage policies. In instances where direct access to resources was restricted, publicly available alternatives and open APIs were utilised. Automation tools such as Playwright and Selenium were employed to emulate user interactions and load dynamically generated content, even when data availability was limited. To avoid detection or triggering of anti-scraping mechanisms, data collection was evenly distributed across multiple categories and pages, with deliberate time delays between requests.

To overcome the limitations of Google Colab, the `httpx` library was combined with lightweight parsers such as Selectolax or BeautifulSoup. This approach reduced dependency overhead and offered more stable performance within Colab's constrained environment. For large-scale data processing, local servers or cloud services such as Google Cloud Platform, Amazon Web Services (AWS), or Microsoft Azure were recommended. Intermediate results were stored using platforms like Google Drive or BigQuery, and the scraping tasks were divided into smaller batches to optimise resource allocation.

## **8.0 Conclusion and Future Work**

This section presents key findings derived from the data collection and processing activities conducted on various websites. It evaluates the performance of each library and outlines the major challenges encountered during their use. Besides, it offers recommendations and practical suggestions to improve the efficiency, reliability, and scalability of future implementations.

### **8.1 Summary of Findings**

From an ETL pipeline perspective, Polars demonstrates superior performance in terms of speed, CPU efficiency, and memory usage compared to Pandas, Modin, and Dask. This makes it a highly suitable option for in-memory processing of large datasets. In terms of overall performance, Polars consistently outperforms the other libraries, particularly in tasks involving data parsing, transformation, and migration to MongoDB. It recorded the lowest CPU usage (6%–16%) and the smallest memory footprint throughout the processing stages. These advantages are largely attributed to its Rust-based backend, which is optimised for multi-threading and capable of handling up to 150 million rows per second.

Although slower in handling intensive operations, Pandas remains effective for simpler tasks. It maintains a low to moderate CPU load (7%–20%) and exhibits stable memory consumption across most ETL stages. For small-scale workflows executed on a single machine, Pandas continues to serve as a reliable and widely adopted solution.

Modin and Dask, which are designed for distributed computing across multiple machines, incur higher initialisation overhead, particularly during data import and missing value imputation. Modin tends to consume higher CPU resources during large-scale operations, yielding performance gains primarily when true parallelism is required. Dask, due to its distributed task scheduler, demonstrates suboptimal performance on single-machine setups, with relatively low throughput and notable memory fluctuations during specific operations.

Upon initial CSV loading, all libraries register similar memory usage in the range of 420–450 MB. However, their memory behaviour diverges during processing. Polars, benefiting from a streaming architecture, exhibits the most efficient memory handling. In contrast, Pandas and Modin manage data using smaller buffers, while Dask shows the highest variability—peaking at 800 MB during loading and settling at approximately 340 MB during transformations.

In conclusion, for a standalone ETL workflow, Polars offers the best overall performance in terms of processing speed, resource efficiency, and scalability when compared to Pandas, Modin, and Dask.

## 8.2 Improvement for Future Works

Several improvements should be implemented to ensure that future web scraping and data processing can be conducted reliably and at scale. Firstly, establishing a stable and flexible computing environment is essential. It is recommended to use local machines equipped with sufficient resources or cloud services such as AWS EC2, Azure, or DigitalOcean, as they support a broader range of tools and offer greater flexibility. Tools such as Playwright and Selenium are better suited for these environments, since platforms like Google Colab do not fully support headless browsers and impose runtime limitations (Mitchell, 2021; Amazon Web Services, 2023).

Robust error-handling mechanisms are also critical to maintaining continuous operation during intensive scraping tasks. Continuity and ease of debugging can be improved by incorporating structured try-except blocks, persistent logging practices, and checkpointing strategies (Ryan, 2020). Additionally, breaking down scraping tasks into smaller, manageable units can prevent memory overflows and increase system fault tolerance.

Adherence to ethical and legal standards in web scraping must also be ensured. Integrating robots.txt compliance into the scraping workflow reduces the risk of violating

site policies. Whenever possible, accessing data through official APIs or through data-sharing agreements is strongly preferred (Singh & Sharma, 2022).

The quality of data can be enhanced by addressing limitations associated with datasets derived from a small number of pages. Future work should explore sources that provide structured or fragmented listings or aggregate data from multiple websites to expand the dataset without overburdening any single server (Kumar & Bansal, 2021).

When processing large volumes of data, effective memory management is essential. Streaming techniques, temporary disc storage, and efficient data formats such as Parquet are recommended. Libraries such as Polars are particularly suitable, as they are optimised for performance and exhibit low memory consumption, enabling multiple ETL tasks to be executed in memory simultaneously (van der Maas, 2023).

## 9.0 References

- Amazon Web Services. (2023). Amazon EC2 features.  
<https://aws.amazon.com/ec2/features/>
- Faiz, A. (2022, December 12). Selectolax — A better Python web scraping library. Medium.  
<https://medium.com/@afiffaiz/selectolax-a-better-python-web-scraping-library-2a8c1ee5551b>
- Fischer, B. (2024, December 23). Python advanced: 10 things you can do with Polars (and didn't know about it). Medium.  
<https://captain-solaris.medium.com/python-advanced-10-things-you-can-do-with-polars-and-didnt-know-about-it-cb8c071227ba>
- Kolli, A. (2024, February 8). Accelerating pandas with Modin: A comprehensive guide. Medium.  
<https://aravindkolli.medium.com/accelerating-pandas-with-modin-a-comprehensive-guide-90697ca4173d>
- Kumar, R., & Bansal, S. (2021). Efficient web scraping approaches and tools for data extraction: A comprehensive review. International Journal of Web & Semantic Technology, 12(2), 45–59. <https://doi.org/10.5121/ijwest.2021.12204>
- Kumawat, K. (2024, November 23). What is Parsel in Python? Medium.  
<https://medium.com/@kuldeepkumawat195/what-is-parsel-in-python-b76237cb7e6d>
- Mitchell, R. (2021). Web scraping with Python: Collecting data from the modern web (2nd ed.). O'Reilly Media.
- Omisola, I. (2025, February 27). 7 best Python web scraping libraries in 2025. ZenRows.  
<https://www.zenrows.com/blog/python-web-scraping-library>
- Ryan, D. (2020). Practical Python web scraping: Best practices and techniques. Apress.
- Seng, S. A. C. (2025, May 15). HTTPX [Python] | Fast-performing and robust Python library. Apidog: An integrated platform for API design, debugging, development, mock, and testing. <https://apidog.com/blog/what-is-htpx/>

- Shahizan, D. (n.d.). drshahizan [GitHub repository]. GitHub.  
<https://github.com/drshahizan/drshahizan>
- Singh, A., & Sharma, M. (2022). Ethical issues in web data extraction: Legal challenges and best practices. *Journal of Information Ethics*, 31(1), 50–65.  
<https://doi.org/10.3172/JIE.31.1.050>
- van der Maas, R. (2023). Polars user guide: Fast and memory-efficient dataframes in Python. <https://pola-rs.github.io/polars-book/>
- Vishnoi, P. (2023, August 20). What is Dask in Python. Medium.  
<https://premvishnoi.medium.com/what-is-dask-in-python-bc86730ec69e>
- Whelan, L. (2025, April 8). Web crawling vs. web scraping – What's the difference? SOAX.  
<https://soax.com/blog/web-crawling-vs-web-scraping#:~:text=Web%20crawlers%20visit%20the%20target,the%20information%20to%20be%20scraped>

## 10.0 Appendices

### 10.1 Appendix A

Table A1 below lists the data attributes extracted from Carlist.my, which form the basis of our data collection and analysis.

Table A1: Overview of Car Data Attributes from Carlist.my

| Attribute        | Description  |
|------------------|--|
| Car Name         | The full name of the car, including model year, brand, and trim/version details. |
| Car Brand        | The manufacturer or brand of the vehicle.  |
| Car Model        | The specific model of the car.   |
| Manufacture Year | The year in which the car was made.  |
| Body Type        | The style or design of the car.  |
| Fuel Type        | The type of fuel the car uses.   |
| Mileage          | The total distance the car has been driven, measured in kilometres (KM).         |
| Transmission     | The type of transmission the car uses.   |
| Color            | The colour of the car's exterior.  |
| Price            | The cost of the car, usually in the local currency.                              |
| Installment      | The monthly payment if the car is purchased through an instalment plan.          |
| Condition        | The state of the car, either new, used or reconditioned.                         |

|               |  |
|---------------|--|
| Seat Capacity | The number of people the car can accommodate, including the driver and passengers. |
| Location      | The geographic location where the car is being sold.                               |
| Sales Channel | The type of seller or platform through which the car is being sold.                |

## 10.2 Appendix B

Table B1 provides detailed descriptions of the components involved in the system architecture, which includes the web scraping process and big data tools used for data processing, cleaning, transformation, and performance evaluation.

Table B1: Web Scraping and Big Data Processing Architecture Components

| Component                                       | Description   |
|---|---|
| Crawler   | The web scraping process was initiated by navigating the targeted website, following links, and collecting raw HTML content (Whelan, 2025).   |
| Parser  | The raw HTML content was processed to extract key elements such as car name, car model, and car price for further analysis (Whelan, 2025).  |
| Data Collection                                 | The parsed data from individually assigned page ranges was stored in a structured format (CSV file) and later saved in MongoDB for further export and processing.   |
| Data Cleaning, Transformation, and Optimisation | The scraped data was cleaned by handling missing values, duplicates, and inconsistencies. It was then transformed into a suitable format for analysis. During this process, optimisation techniques such as parallel processing and multi-threading were implemented using different libraries (Pandas, Polars, Modin, Dask). |
| Performance Evaluation                          | The performance of the entire system was measured and evaluated based on speed, CPU usage, memory usage, and throughput.  |
| End-To-End Flow                                 | A complete pipeline was established, where data flowed from scraping to storage, cleaning, transformation, optimisation, and performance evaluation for efficient data analysis.  |

Table B2 summarises the libraries and frameworks utilised by each team member during the web scraping and big data processing stages of the project. The tools are categorised based on their respective functionalities to provide an overview of the software stack employed for efficient data acquisition and analysis.

Table B2: Tools and Frameworks Used

| <b>Web Crawling &amp; Web Scraping</b>  |                              |                              |   |
|---|------------------------------|------------------------------|---|
| <b>Name</b>                             | <b>Library<br/>(Crawler)</b> | <b>Library<br/>(Scraper)</b> | <b>Description</b>  |
| Marcus Joey Sayner                      | httpx                        | parsel                       | httpx is a modern, async-capable HTTP client used to fetch web content (Seng, 2025), whereas parsel is used to extract data from HTML/XML using CSS or XPath selectors (Kumawat, 2024). |
| Muhammad Luqman Hakim bin Mohd Rizaudin | urllib                       | Selectolax                   | urllib is a built-in Python module to handle HTTP request (Omisola, 2025), and selectolax is a fast, lightweight HTML parser optimised for speed and large documents.                   |
| Camily Tang Jia Lei                     | urllib                       | BeautifulSoup                | Similarly, urllib handles the retrieval of web pages, and BeautifulSoup parses and extracts structured data from HTML content (Omisola, 2025).  |
| Goh Jing Yang                           | requests                     | BeautifulSoup                | requests is a simple and widely used HTTP library for making web requests, and BeautifulSoup is used to parse and navigate the HTML content for data                                    |

|   |         |   | extraction (Omisola, 2025). |
|---|---------|---|-----------------------------|
| <b>Big Data Processing</b>              |         |   |                             |
| Name                                    | Library | Description   |                             |
| Marcus Joey Sayner                      | Polars  | A fast DataFrame library for data manipulation, optimised for performance and scalability using multi-threaded execution (Fischer, 2024).                                     |                             |
| Muhammad Luqman Hakim bin Mohd Rizaudin | Modin   | A library designed to speed up pandas workflows using built-in parallelisation (Kolli, 2024).   |                             |
| Camily Tang Jia Lei                     | Pandas  | A popular Python library for data manipulation and analysis, offering efficient data structures like DataFrame and Series to handle and process large datasets (Kolli, 2024). |                             |
| Goh Jing Yang                           | Dask    | A flexible parallel computing library that extends pandas-like APIs for analytics that can handle larger-than-memory computations (Vishnoi, 2023).                            |                             |

## 10.3 Appendix C

Figures C1 to C4 present the source code implementations for the crawling and scraping methods used by each team member. Each figure corresponds to a specific combination of libraries used for web data extraction from the Carlist.my website.



```
 1 import httpx
 2 import csv
 3 import time
 4 import json
 5 from parsel import Selector
 6
 7 base_url = "https://www.carlist.my/cars-for-sale/malaysia?page_number={}&page_size=25"
 8 start_page = 1
 9 end_page = 2
10 car_listings = []
11 output_file = 'carlist_data.csv'
12
13 # Extract from the <script type="application/ld+json"> block
14 def extract_json_ld(selector):
15     scripts = selector.css('script[type="application/ld+json"]::text').getall()
16     for script in scripts:
17         try:
18             data = json.loads(script)
19             if isinstance(data, list):
20                 for d in data:
21                     if 'itemListElement' in d:
22                         return d['itemListElement']
23             except json.JSONDecodeError:
24                 continue
25     return None
26
27 # Get mileage from <i class="icon--meter">
28 def get_mileage(selector):
29     meter_icon = selector.css('i.icon--meter')
30     if meter_icon:
31         return meter_icon.xpath('following-sibling::text()[1]').get(default='').strip()
32     return ''
33
34 # Scraping function to get the page content
35 def fetch_page(url):
36     response = httpx.get(url)
37     return response.text
38
39 # Correct full location (e.g., Petaling Jaya, Selangor)
40 def get_location(selector, car):
41     try:
42         # Extract the locality and region from the car JSON data
43         address_locality = car.get('seller', {}).get('homeLocation', {}).get('address', {}).get('addressLocality', '')
44         address_region = car.get('seller', {}).get('homeLocation', {}).get('address', {}).get('addressRegion', '')
45
46         # Combine locality and region to form the location
47         if address_locality and address_region:
48             location = f"(address_locality), {address_region}"
49         else:
50             location = ', '.join([address_locality, address_region]).strip(', ')
51
52     return location
53     except Exception as e:
54         return ''
55
56 # Correct sales agent (e.g., Sales Agent)
57 def get_sales_channel(selector):
58     try:
59         # Use CSS selector to find the sales agent text
60         sales_info = selector.css('span.c-chip--icon::text').get(default='').strip()
61         # Replace "Sales" with "Sales Agent"
62         sales_channel = sales_info.replace("Sales", "Sales Agent") if "Sales" in sales_info else sales_info
63         return sales_channel
64     except Exception as e:
65         return ''
66
67 # Start timing
68 start_time = time.time()
```

```

70 # Start scraping
71 for page in range(start_page, end_page + 1):
72     url = base_url.format(page)
73     print(f"\n🌐 Scraping page {page}: {url}")
74
75     response = httpx.get(url, headers={"User-Agent": "Mozilla/5.0"})
76     selector = Selector(response.text)
77
78     articles = selector.css('article.listing')
79     json_ld = extract_json_ld(selector)
80
81     if not json_ld:
82         print("⚠️ No JSON-LD found.")
83         continue
84
85     for i, (article_sel, json_car) in enumerate(zip(articles, json_ld)):
86         car = json_car['item']
87
88         name = article_sel.attrib.get('data-title', '')
89         brand = article_sel.attrib.get('data-make', '')
90         model = article_sel.attrib.get('data-model', '')
91         body = article_sel.attrib.get('data-body-type', '')
92         transmission = article_sel.attrib.get('data-transmission', '')
93         installment = article_sel.attrib.get('data-installment', '')
94
95         year = car.get('vehicleModelDate', '')
96         fuel = car.get('fuelType', '')
97         color = car.get('color', '')
98         price = car.get('offers', {}).get('price', '')
99         condition = car.get('itemCondition', '')
100        seats = car.get('seatingCapacity', '')
101
102        mileage = get_mileage(article_sel)
103        location = get_location(article_sel, car)
104        sales_channel = get_sales_channel(article_sel)
105
106        car_listings.append({
107            'Car Name': name,
108            'Car Brand': brand,
109            'Car Model': model,
110            'Manufacture Year': year,
111            'Body Type': body,
112            'Fuel Type': fuel,
113            'Mileage': mileage,
114            'Transmission': transmission,
115            'Color': color,
116            'Price': price,
117            'Installment': installment,
118            'Condition': condition,
119            'Seating Capacity': seats,
120            'Location': location,
121            'Sales Channel': sales_channel
122        })
123
124    print(f"✅ Extracted {len(articles)} listings from page {page}")
125    time.sleep(2) # Be kind to the server
126
127    # Save to CSV
128    if car_listings:
129        with open(output_file, 'w', newline='', encoding='utf-8') as f:
130            writer = csv.DictWriter(f, fieldnames=car_listings[0].keys())
131            writer.writeheader()
132            writer.writerows(car_listings)
133        print(f"\n🌐 Saved {len(car_listings)} listings to '{output_file}'")
134    else:
135        print("⚠️ No data extracted.")
136
137    # End timing
138    end_time = time.time()
139    print(f"\n⌚ Total scraping time: {end_time - start_time:.2f} seconds")

```

Figure C1: Crawling and scraping code using httpx + parsel

```

● ● ●

1 import urllib.request
2 from selectolax.parser import HTMLParser
3 import csv
4 import time
5 import json
6
7 # Configuration
8 base_url = "https://www.carlist.my/cars-for-sale/malaysia?page_number={}&page_size=25"
9 start_page = 1745      # Starting page number to scrape
10 end_page = 3488       # Ending page number to scrape
11 car_listings = []     # List to store car listings data
12
13 output_file = 'car_listings.csv'
14
15 def get_attr(element, attr_name):
16     return element.attributes.get(attr_name, '') if element else ''
17
18 def get_text(node):
19     return node.text().strip() if node else ''
20
21 def get_sales_channel(article):
22     dealer_div = article.css_first('div.listing__spec--dealer')
23     if dealer_div:
24         channel = dealer_div.text().strip()
25         # Convert "Sales" to "Sales Agent"
26         return channel.replace("Sales", "Sales Agent") if "Sales" in channel else channel
27     return ''
28
29 def get_mileage(article):
30     icon = article.css_first('i.icon--meter')
31     if icon and icon.next:
32         return icon.next.text().strip()
33     return ''
34
35 def get_location(article):
36     icon = article.css_first('i.icon--location')
37     if not icon:
38         return ''
39     text = ''
40     current = icon.next
41     while current:
42         if current.tag == '-text':
43             text += current.text().strip()
44         elif current.tag == 'span':
45             text += current.text().strip()
46         else:
47             break
48     current = current.next
49     return text.strip()
50
51 def extract_json_ld(parser):
52     for script in parser.css('script[type="application/ld+json"]'):
53         try:
54             data = json.loads(script.text())
55             if isinstance(data, list):
56                 for d in data:
57                     if 'itemListElement' in d:
58                         return d['itemListElement']
59             except:
60                 continue
61     return None
62
63 start_time = time.time()
64

```

```

65  for page_num in range(start_page, end_page + 1):
66      url = base_url.format(page_num)
67      print(f"\n🌐 Extracting page {page_num} - {url}")
68
69      req = urllib.request.Request(url, headers={"User-Agent": "Mozilla/5.0"})
70      html = urllib.request.urlopen(req).read()
71      parser = HTMLParser(html)
72
73      articles = parser.css('article.listing')
74      ld_json = extract_json_ld(parser)
75      if not ld_json:
76          print("⚠️ No JSON-LD found on this page.")
77          continue
78
79      page_count = 0
80      for article, item in zip(articles, ld_json):
81          car = item['item']
82
83          name = get_attr(article, 'data-title')
84          brand = get_attr(article, 'data-make')
85          model = get_attr(article, 'data-model')
86          body = get_attr(article, 'data-body-type')
87          transmission = get_attr(article, 'data-transmission')
88          installment = get_attr(article, 'data-installment')
89
90          mileage = get_mileage(article)
91          sales_channel = get_sales_channel(article) # This now returns "Sales Agent" if "Sales" was found
92          location = get_location(article)
93
94          year = car.get('vehicleModelDate', '')
95          fuel = car.get('fuelType', '')
96          color = car.get('color', '')
97          price = car.get('offers', {}).get('price', '')
98          condition = car.get('itemCondition', '')
99          seats = car.get('seatingCapacity', '')
100
101         car_listings.append({
102             'Car Name': name,
103             'Car Brand': brand,
104             'Car Model': model,
105             'Manufacture Year': year,
106             'Body Type': body,
107             'Fuel Type': fuel,
108             'Mileage': mileage,
109             'Transmission': transmission,
110             'Color': color,
111             'Price': price,
112             'Installment': installment,
113             'Condition': condition,
114             'Seating Capacity': seats,
115             'Location': location,
116             'Sales Channel': sales_channel # Will show "Sales Agent" instead of "Sales"
117         })
118         page_count += 1
119
120     time.sleep(2)
121
122 if car_listings:
123     with open(output_file, 'w', newline='', encoding='utf-8') as f:
124         writer = csv.DictWriter(f, fieldnames=car_listings[0].keys())
125         writer.writeheader()
126         writer.writerows(car_listings)
127         print(f"\n✅ Saved {len(car_listings)} cars to '{output_file}'")
128 else:
129     print("\n⚠️ No car listings found.")
130
131 end_time = time.time()
132 execution_time = end_time - start_time
133 print(f"\n⌚ Total execution time: {execution_time:.2f} seconds")

```

Figure C2: Crawling and scraping code using urllib + Selectolax

```
1 import urllib.request
2 from bs4 import BeautifulSoup, NavigableString
3 import csv
4 import time
5 import json
6
7 # Configuration
8 base_url = "https://www.carlist.my/cars-for-sale/malaysia?page_number={}&page_size=25"
9 start_page = 5233    # Starting page number to scrape
10 end_page = 6977     # Ending page number to scrape
11 car_listings = []   # List to store car listings data
12
13 output_file = 'D:/camily_carlist_listings_1.csv'
14
15 # Get HTML attribute or empty string if attribute doesn't exist
16 def get_attr(element, attr_name):
17     return element.get(attr_name, '') if element else ''
18
19 # Get the first text node inside a <div> (ignores non-text elements)
20 def get_first_text_node(div):
21     for node in div.contents if div else []:
22         if isinstance(node, NavigableString) and node.strip():
23             return node
24     return '' # Return empty string if no text node is found
25
26 # Extract raw mileage information from the article (after the meter icon)
27 def get_mileage(article):
28     icon = article.find('i', class_='icon--meter') # Find meter icon
29     return str(icon.next_sibling) if icon and icon.next_sibling else '' # Return mileage text
30
31 # Extract raw location information from the article (after the location icon)
32 def get_location(article):
33     icon = article.find('i', class_='icon--location') # Find location icon
34     if not icon:
35         return ''
36     text = ''
37     for sib in icon.next_siblings: # Loop through siblings after icon
38         if isinstance(sib, NavigableString):
39             text += str(sib)
40         elif getattr(sib, 'name', None) == 'span':
41             text += sib.get_text()
42         else:
43             break
44     return text # Return extracted location text
45
46 # Pull JSON-LD block containing the car listings from the page
47 def extract_json_ld(soup):
48     for script in soup.find_all('script', type='application/ld+json'):
49         try:
50             data = json.loads(script.string) # Parse the JSON-LD data
51             if isinstance(data, list):
52                 for d in data:
53                     if 'itemListElement' in d:
54                         return d['itemListElement'] # Return list of car listings
55         except:
56             continue
57     return None # Return None if no valid JSON-ID is found
58
59 # Start time
60 start_time = time.time()
61
62 # Main scraping loop
63 for page_num in range(start_page, end_page + 1):
64     url = base_url.format(page_num) # Generate URL for the current page
65     print(f"\nExtracting page {page_num} - {url}")
```

```

67     # Load the page and parse it with BeautifulSoup
68     req = urllib.request.Request(url, headers={"User-Agent": "Mozilla/5.0"})
69     html = urllib.request.urlopen(req).read()
70     soup = BeautifulSoup(html, 'html.parser')
71
72     # Find all articles (car listings) and extract JSON-LD data
73     articles = soup.find_all('article', class_='listing')
74     ld_json = extract_json_ld(soup)
75     if not ld_json:
76         print("⚠️ No JSON-LD found on this page.")
77         continue
78
79     page_count = 0
80     # Loop through articles and corresponding JSON-LD data
81     for article, item in zip(articles, ld_json):
82         car = item['item'] # Extract JSON details for the car
83
84         # Extract fields from the <article> element's attributes (HTML-based fields)
85         name = get_attr(article, 'data-title')
86         brand = get_attr(article, 'data-make')
87         model = get_attr(article, 'data-model')
88         body = get_attr(article, 'data-body type')
89         transmission = get_attr(article, 'data-transmission')
90         installment = get_attr(article, 'data-installment')
91
92         # Extract exactly what is displayed on the page (visible text-based fields)
93         mileage = get_mileage(article)
94         sales_channel = get_first_text_node(article.find('div', class_-'listing_spec_dealer'))
95         location = get_location(article)
96
97         # Extract additional fields from the JSON-LD data
98         year = car.get('vehicleModelDate', '')
99         fuel = car.get('fuelType', '')
100        color = car.get('color', '')
101        price = car.get('offers', {}).get('price', '')
102        condition = car.get('itemCondition', '')
103        seats = car.get('seatingCapacity', '')
104
105        # Append the extracted data to the car_listings list
106        car_listings.append({
107            'Car Name': name,
108            'Car Brand': brand,
109            'Car Model': model,
110            'Manufacture Year': year,
111            'Body Type': body,
112            'Fuel Type': fuel,
113            'Mileage': mileage,
114            'Transmission': transmission,
115            'Color': color,
116            'Price': price,
117            'Installment': installment,
118            'Condition': condition,
119            'Seating Capacity': seats,
120            'Location': location,
121            'Sales Channel': sales_channel
122        })
123        page_count += 1
124
125    print(f"\n✅ Found {page_count} cars on page {page_num}")
126    print(f"\n⚠️ Total scraped: {len(car_listings)}")
127    time.sleep(2) # Sleep to avoid making too many requests in a short time
128
129    # Save the extracted data to a CSV file
130    if car_listings:
131        with open(output_file, 'w', newline='', encoding='utf-8') as f:
132            writer = csv.DictWriter(f, fieldnames=car_listings[0].keys())
133            writer.writeheader() # Write the column headers
134            writer.writerows(car_listings) # Write the car listings data
135            print(f"\n✅ Saved {len(car_listings)} cars to '{output_file}'")
136    else:
137        print("\n⚠️ No car listings found.")
138
139    # End time and elapsed time calculation
140    end_time = time.time()
141    execution_time = end_time - start_time
142    print(f"\n⌚ Total execution time: {execution_time:.2f} seconds"

```

Figure C3: Crawling and scraping code using urllib + BeautifulSoup

```
● ● ●
1 import requests
2 from bs4 import BeautifulSoup
3 import csv
4 import json
5 import time
6 import random
7
8 # Configuration
9 base_url = "https://www.carlist.my/cars-for-sale/malaysia?page_number={}&page_size=25"
10 start_page = 3471
11 end_page = 5205
12 car_listings = []
13 output_file = 'car_listings.csv'
14 max_retries = 3
15
16 # Start time
17 start_time = time.time()
18
19 # Main scraping loop
20 try:
21     for page_num in range(start_page, end_page + 1):
22         url = base_url.format(page_num)
23         print(f"\n● Extracting page {page_num} - {url}")
24
25         # Retry mechanism for network errors
26         retries = 0
27         while retries < max_retries:
28             try:
29                 headers = {"User-Agent": "Mozilla/5.0"}
30                 response = requests.get(url, headers=headers, timeout=10)
31                 response.raise_for_status()
32                 soup = BeautifulSoup(response.content, 'html.parser')
33                 break
34             except requests.exceptions.RequestException as e:
35                 retries += 1
36                 print("▲ Network error: {} . Retrying ({}/{})...")
37                 time.sleep(2)
38         else:
39             print(f"✖ Failed to fetch page {page_num} after {} retries. Skipping...")
40             continue
41
42         # Find all articles (car listings) and extract JSON-LD data
43         articles = soup.find_all('article', class_='listing')
44         ld_json = None
45         for script in soup.find_all('script', type='application/ld+json'):
46             try:
47                 data = json.loads(script.string)
48                 if isinstance(data, list):
49                     for d in data:
50                         if 'itemListElement' in d:
51                             ld_json = d['itemListElement']
52                             break
53             except Exception as e:
54                 print("▲ Error parsing JSON-LD: {}")
55                 continue
56
57             if not ld_json:
58                 print("▲ No JSON-LD found on page {}. Skipping...".format(page_num))
59                 continue
60
61             page_count = 0
62             for article, item in zip(articles, ld_json):
63                 car = item['item']
```

```

65         # Extract fields directly without utility functions
66         name = article.get('data-title', '')
67         brand = article.get('data-make', '')
68         model = article.get('data-model', '')
69         body = article.get('data-body-type', '')
70         transmission = article.get('data-transmission', '')
71         installment = article.get('data-installment', '')
72         mileage = ""
73         location = ""
74         for icon in article.find_all('i'):
75             if 'icon--meter' in icon.get('class', []):
76                 mileage = str(icon.next_sibling.strip())
77             elif 'icon--location' in icon.get('class', []):
78                 location_text = []
79                 for sib in icon.next_siblings:
80                     text = sib.get_text(strip=True) if hasattr(sib, 'get_text') else str(sib).strip()
81                     if text:
82                         location_text.append(text)
83                 location = ' '.join(location_text)
84
85         # Extract additional fields from the JSON-LD data
86         year = car.get('vehicleModelDate', '')
87         fuel = car.get('fuelType', '')
88         color = car.get('color', '')
89         price = car.get('offers', {}).get('price', '')
90         condition = car.get('itemCondition', '').lower()
91         condition = "New" if "new" in condition else "Used"
92         seats = car.get('seatingCapacity', '')
93         sales_channel = ""
94         dealer_div = article.find('div', class_=['listing_spec--dealer'])
95         if dealer_div:
96             sales_channel = dealer_div.get_text(strip=True)
97
98         # Append data to the car_listings list
99         car_listings.append({
100             'Car Name': name,
101             'Car Brand': brand,
102             'Car Model': model,
103             'Manufacture Year': year,
104             'Body Type': body,
105             'Fuel Type': fuel,
106             'Mileage': mileage,
107             'Transmission': transmission,
108             'Color': color,
109             'Price': price,
110             'Installment': installment,
111             'Condition': condition,
112             'Seat Capacity': seats,
113             'Location': location,
114             'Sales Channel': sales_channel,
115             'URL': car.get('url', '')
116         })
117         page_count += 1
118
119     print(f"\n✅ Found {page_count} cars on page {page_num}")
120     print(f"\n💡 Total scraped: {len(car_listings)}")
121
122     # Randomized delay to prevent IP blocking
123     delay = random.uniform(2, 5)
124     print(f"\n🕒 Waiting {delay:.2f} seconds before next page...")
125     time.sleep(delay)
126
127 except KeyboardInterrupt:
128     print("\n⚠ Script interrupted by the user. Saving progress...")
129
130 # Save the extracted data to a CSV file
131 if car_listings:
132     with open(output_file, 'w', newline='', encoding='utf-8') as f:
133         writer = csv.DictWriter(f, fieldnames=car_listings[0].keys())
134         writer.writeheader()
135         writer.writerows(car_listings)
136     print(f"\n✅ Saved {len(car_listings)} cars to '{output_file}'")
137 else:
138     print("\n⚠ No car listings found.")
139
140 # End time and elapsed time calculation
141 end_time = time.time()
142 execution_time = end_time - start_time
143 print(f"\n⌚ Total execution time: {execution_time:.2f} seconds")

```

Figure C4: Crawling and scraping code using requests + BeautifulSoup

## 10.4 Appendix D

Table D1 provides a comprehensive overview of the data cleaning procedures applied to the car listings dataset. Each entry in the table outlines the specific objective and detailed steps followed to ensure the dataset's accuracy, consistency, and analytical readiness.

Table D1: Summary of Data Cleaning Methods

| <b>Combine Columns: ‘Seat Capacity’ and ‘Seating Capacity’</b> |   |
|--|---|
| Objective  | To merge redundant or synonymous columns for uniformity of the data.  |
| Steps  | <ol style="list-style-type: none"> <li>1. The dataset was observed to contain two columns, namely ‘Seat Capacity’ and ‘Seating Capacity’, both representing the same attribute.</li> <li>2. To eliminate redundancy and ensure consistency, these columns were consolidated into a single column labelled ‘Seat Capacity’, which was designated as the primary reference for seating information.</li> <li>3. The following procedure was applied:           <ul style="list-style-type: none"> <li>• For each row, if the value in ‘Seat Capacity’ was null but ‘Seating Capacity’ contained a valid entry, the null value was filled with the corresponding value from ‘Seating Capacity’.</li> <li>• If a value already existed in ‘Seat Capacity’, it was retained and given precedence.</li> </ul> </li> <li>4. After completing the merge, the ‘Seating Capacity’ column was removed from the dataset.</li> </ol> |
| <b>Remove Unwanted Column(s): ‘URL’</b>                        |   |
| Objective  | To eliminate unnecessary or noisy data that does not contribute to the  |

|  |   |
|--|---|
|  | analysis.   |
| Steps  | <p>1. The column ‘URL’ was identified as lacking analytical relevance.</p> <p>2. It was observed that the column contained lengthy strings referring to external web resources, which were:</p> <ul style="list-style-type: none"> <li>• Not utilised in analysis or visualizations,</li> <li>• Not consistently present across all entries,</li> <li>• Potentially misleading or distracting.</li> </ul> <p>3. Consequently, the ‘URL’ column was removed from the dataset to reduce dimensionality and streamline the analytical process.</p> |
| <b>String Cleaning and Formatting (Clean ‘Installment’ Column)</b> |   |
| Objective  | To remove currency symbols and units from the values to convert them into numeric form.   |
| Steps  | <p>1. The “RM” currency symbol was removed.</p> <p>2. Commas (,) used as thousand separators were eliminated.</p> <p>3. The “/month” suffix was stripped from the values.</p> <p>4. The cleaned strings were converted into numeric data.</p> <p>5. As a result, the ‘Installment’ column now contains standardised numeric values suitable for aggregation and statistical analysis</p>  |
| <b>String Cleaning and Formatting (Clean ‘Condition’ Column)</b>   |   |
| Objective  | To extract and simplify metadata-style URLs into clear categorical labels.  |
| Steps  | <p>1. Strings were parsed to extract the main condition label from URLs such as "https://schema.org/UsedCondition", retaining only values like "Used" or "New".</p>   |

|  |   |
|--|---|
|  | <p>2. Any trailing "Condition" suffixes were removed.</p> <p>3. The column was finalised with clean categorical values suitable for filtering and classification.</p>   |
| <b>String Cleaning and Formatting (Clean 'Sales Channel' Column)</b> |   |
| Objective  | To standardise and extract meaningful labels from long-form sales channel strings.  |
| Steps  | <p>1. Regular expressions were applied to extract relevant categories such as "Sales Agent" or "Dealer".</p> <p>2. Entries not matching the expected pattern were either removed or replaced with null values.</p> <p>3. The column now contains standardised sales channel identifiers for reliability.</p>                                |
| <b>String Cleaning and Formatting (Clean 'Mileage' Column)</b>       |   |
| Objective  | To normalise mileage values into consistent categorical bins for analysis.  |
| Steps  | <p>1. The "K KM" unit suffix was removed from all entries.</p> <p>2. Range values (e.g., "5 - 10") were parsed and normalised.</p> <p>3. Single numeric values (e.g., "12") were binned into 5K-step ranges (e.g., "10 - 15").</p> <p>4. The column now presents uniform mileage categories for improved aggregation and visualisation.</p> |
| <b>Replace / Drop Null Values</b>                                    |   |
| Objective  | To ensure data completeness by removing rows with missing values in key columns.  |

|   |   |
|---|---|
| Steps   | 1. The following primary columns were selected: 'Body Type', 'Fuel Type', 'Mileage', 'Sales Channel', and 'Seat Capacity'.  |
|   | 2. These columns were checked for missing (NaN) values  |
|   | 3. Rows containing nulls in any of the selected columns were removed using .dropna().   |
|   | 4. Only complete rows were retained, eliminating potential bias from imputation or scattered missing data.  |
| <b>Rename Columns (Add Units for Clarity)</b> |   |
| Objective                                     | To improve readability by appending measurement units to column names.  |
| Steps   | 1. Key numeric columns were identified and renamed as follows: <ul style="list-style-type: none"> <li>• 'Mileage' → 'Mileage (K KM)'</li> <li>• 'Price' → 'Price (RM)'</li> <li>• 'Installment' → 'Installment (RM)'</li> </ul> |
|   | 2. The .rename() method with inplace=True was used to apply changes directly to the original DataFrame.   |
| <b>Drop Duplicated Rows</b>                   |   |
| Objective                                     | To eliminate redundant records and avoid distortion in the analysis.  |
| Steps   | 1. Duplicate rows were identified using .duplicated() across all columns.   |
|   | 2. Exact duplicates were removed using .drop_duplicates().  |
|   | 3. Care was taken to retain semantically valid duplicates, ensuring meaningful data was preserved.  |

## 10.5 Appendix E

Table E1 summarises the key performance optimisation strategies used by each data processing library.

Table E1: Overview of Data Processing Libraries

| Library | Explanation   |
|---------|---|
| Pandas  | <ul style="list-style-type: none"><li>• Operates in a single-threaded manner by default.</li><li>• No built-in parallelism is provided.</li><li>• Act as the control group to evaluate the performance improvements of other libraries.</li></ul>                                       |
| Polars  | <ul style="list-style-type: none"><li>• Operates using <b>multi-threading</b> to parallelise data operations at the thread level.</li><li>• Built in Rust, providing faster performance and lower memory overhead compared to Pandas.</li></ul>   |
| Modin   | <ul style="list-style-type: none"><li>• Enables automatic <b>parallelism</b> using the Dask execution engines.</li><li>• Spreads DataFrame operations across all CPU cores to speed up processing on multi-core systems.</li><li>• Requires minimal code changes from Pandas.</li></ul> |
| Dask    | <ul style="list-style-type: none"><li>• Applies task-based <b>parallelism</b> and lazy computation graphs to break down large operations into smaller, manageable tasks.</li><li>• Optimises memory usage by chunking data and only loading necessary partitions into memory.</li></ul> |

Table E2 provides a summary of optimisation tasks and the corresponding code snippets implemented across the data processing libraries.

Table E2: Code Overview of Each Library

| Library                             | Code   |
|-------------------------------------|--|
| <b>Combine Dataset</b>              | <p>Figure E1 shows the code used to combine multiple scraped datasets from group members into a single DataFrame using pandas. Four CSV files, each scraped separately, are read and merged into one dataset. The combined result is then saved as a new CSV file. Subsequently, a connection is established to a local MongoDB database, and the combined data is inserted into a collection named “listings”.</p> <pre> import pandas as pd from pymongo import MongoClient  # List of group member scraping files files = [     "camily_carlist_listings.csv",     "luqman_carlist_listings.csv",     "jingyang_carlist_listings.csv",     "marcus_carlist_listings.csv" ]  # Read and combine datasets dfs = [pd.read_csv(file) for file in files] combined_df = pd.concat(dfs, ignore_index=True)  # Save combined dataset combined_df.to_csv("carlist_combined.csv", index=False)  # Output shape and preview print("Combined dataset shape:", combined_df.shape) print(combined_df.head(5))  # Connect to MongoDB client = MongoClient("mongodb://localhost:27017") db = client["carlist_db"] collection = db["listings"]  # Insert combined data into MongoDB collection.insert_many(combined_df.to_dict(orient="records")) print("✓ Data inserted into MongoDB successfully.") </pre> |
| <b>Install and Import Libraries</b> |  |

|        |  |
|--------|--|
| Pandas | <p>Figure E2 shows the essential libraries for the Pandas implementation where:</p> <ul style="list-style-type: none"> <li>• <b>pandas:</b> For reading, merging, and manipulating structured tabular data.</li> <li>• <b>numpy:</b> To support numerical operations.</li> <li>• <b>re:</b> For performing regular expression-based string cleaning and extraction.</li> <li>• <b>time:</b> To record execution time for performance evaluation.</li> <li>• <b>psutil:</b> A system-level information such as memory and CPU usage.</li> <li>• <b>os:</b> For interacting with the operating system environment.</li> <li>• <b>tracemalloc:</b> To trace memory allocation during code execution.</li> <li>• <b>Pymongo:</b> Communicate with a MongoDB database to store the processed data.</li> </ul> <pre style="background-color: black; color: cyan; padding: 10px;">import pandas as pd import numpy as np import re import time import psutil import os import tracemalloc from pymongo import MongoClient</pre> |
| Polars | <p>Figure E3 shows the essential libraries for the Polars implementation where:</p> <ul style="list-style-type: none"> <li>• <b>Polars:</b> For fast and memory-efficient DataFrame operations, optimised for performance on large datasets.</li> <li>• <b>numpy, re, time, psutil, os, and tracemalloc:</b> Same purposes as in the Pandas version, enabling numerical operations, system monitoring, and memory tracking.</li> <li>• <b>pymongo:</b> For inserting the final dataset into MongoDB.</li> </ul>  |

```

import polars as pl
import numpy as np
import re
import time
import psutil
import os
import tracemalloc
from pymongo import MongoClient

```

Figure E3: Polars Code for Install and Import Libraries

|       |  |
|-------|--|
|       | <pre> import polars as pl import numpy as np import re import time import psutil import os import tracemalloc from pymongo import MongoClient </pre>   |
| Modin | <p>Figure E4 shows the essential libraries for the Modin implementation where:</p> <ul style="list-style-type: none"> <li>• The environment variable MODIN_ENGINE is set to "dask" to specify Dask as the backend engine for distributed execution.</li> <li>• <b>modin.pandas</b>: Acts as a drop-in replacement for Pandas but enables parallel processing.</li> <li>• <b>numpy, psutil, re, tracemalloc, and time</b>: For data manipulation, system monitoring, and performance tracking.</li> <li>• <b>Os</b>: To set the environment configuration.</li> </ul> <pre> import os os.environ["MODIN_ENGINE"] = "dask"  import modin.pandas as md import numpy as np import psutil import re import tracemalloc import time </pre> |
|       | <p>Figure E4: Modin Code for Install and Import Libraries</p>  |
| Dask  | <p>Figure E5 shows the essential libraries for the Dask implementation where:</p> <ul style="list-style-type: none"> <li>• <b>dask.dataframe</b>: To perform scalable DataFrame operations on large datasets.</li> </ul>   |

- **dask.array:** Allow distributed processing of numerical and unstructured data respectively.
- **dask:** Provides additional task scheduling and parallelism utilities.
- **pandas:** Used for fallback operations or conversion between formats.
- **numpy, time, psutil, re, os, and tracemalloc:** Imported for standard data handling, timing, system monitoring, and memory tracking.

```

import dask.dataframe as dd
import numpy as np
import dask.array as da
import dask.bag as db
import dask
import pandas as pd
import time
import psutil
import re
import os
import tracemalloc

```

Figure E5: Dask Code for Install and Import Libraries

### General Code

Figure E6 and E7 illustrates the general performance tracking code used consistently throughout the project for all libraries during the import, cleaning and export phases. At the beginning of execution, the process ID is retrieved using psutil.Process(os.getpid()), and memory tracking is initiated with tracemalloc.start(). The start time, CPU time, and CPU usage percentage are recorded to establish a baseline. After the operation (e.g., import or cleaning), the end time and CPU time are captured, and both current and peak memory usage are measured using tracemalloc.get\_traced\_memory() before stopping the tracker.

Performance metrics are then calculated:

- **Elapsed time:** The total real-world (wall clock) time taken to complete the operation, measured from start to finish.
- **CPU time:** The actual amount of CPU processing time consumed during execution, calculated as the difference in user CPU time before and after the task and converted into milliseconds.
- **CPU usage:** Average percentage of CPU resources utilised during the process, measured over a one-second sampling interval to reflect system load.
- **Throughput:** The processing efficiency by calculating the number of rows handled per second, obtained by dividing the total rows processed by the elapsed time.
- **Current memory usage:** The amount of memory occupied by the Python process at the time of measurement.
- **Peak memory usage:** The maximum memory consumption during execution, highlighting the highest resource demand encountered.

```
# Setup
process_import_pandas = psutil.Process(os.getpid())
tracemalloc.start()

# Record start time and CPU usage
start_time_import_pandas = time.perf_counter()
start_cpu_time_import_pandas = psutil.cpu_times().user # Start CPU time
start_cpu_percent_import_pandas = psutil.cpu_percent(interval=None)
```

Figure E6: General Code Used Across all Libraries (Initial Setup)

```

# Record end time, memory, and CPU usage
end_time_import_pandas = time.perf_counter()
end_cpu_time_import_pandas = psutil.cpu_times().user # End CPU time
current_mem_import_pandas, peak_mem_import_pandas = tracemalloc.get_traced_memory()
tracemalloc.stop()

# CPU percent used during execution
cpu_usage_import_pandas = psutil.cpu_percent(interval=1) # sampled over 1 second

# Calculate metrics
elapsed_time_import_pandas = end_time_import_pandas - start_time_import_pandas
cpu_time_import_pandas = (end_cpu_time_import_pandas - start_cpu_time_import_pandas) * 1000 # CPU time in ms
total_rows_import_pandas = len(df_pandas)
throughput_import_pandas = total_rows_import_pandas / elapsed_time_import_pandas if elapsed_time_import_pandas > 0 else 0

# Output preview
print("✅ Data loaded from MongoDB")
print(df_pandas.head(3))

# Performance summary
print("\n===== Performance Summary =====")
print(f"Total rows loaded : {total_rows_import_pandas}")
print(f"Wall Time (Elapsed) : {elapsed_time_import_pandas:.4f} seconds")
print(f"CPU Time : {cpu_time_import_pandas:.0f} ms")
print(f"CPU Usage : {cpu_usage_import_pandas:.2f}%")
print(f"Throughput : {throughput_import_pandas:.2f} rows/sec")
print(f"Current Memory (Python): {current_mem_import_pandas / 1e6:.2f} MB")
print(f"Peak Memory (Python) : {peak_mem_import_pandas / 1e6:.2f} MB")
print("====")

```

Figure E7: General Code Used Across all Libraries (Calculation)

## Import Data

Pandas

Figure E8 shows how data is imported from MongoDB into a Pandas DataFrame. The data is first retrieved as a list of dictionaries from the MongoDB collection. This list is then converted directly into a Pandas DataFrame. Afterward, the MongoDB-specific `_id` column is removed since it is not needed for further processing.

```

# === Read from MongoDB into pandas DataFrame ===
data_import_pandas = list(collection.find())
df_pandas = pd.DataFrame(data_import_pandas)
df_pandas.drop(columns=[ '_id' ], inplace=True)
# =====

```

Figure E8: Pandas Code for Import Data

Polars

Figure E9 illustrates the import of data from MongoDB into a Polars DataFrame. Similar to Pandas, the data is fetched as a list of dictionaries. The Polars DataFrame is then created by inferring the schema from a large

|  |  |
|--|--|
|  | <p>sample of rows to ensure proper data typing. The <code>_id</code> column generated by MongoDB is dropped to clean the dataset.</p> <pre># === Read from MongoDB into Polars DataFrame === data_import_polars = list(collection.find()) df_polars = pl.DataFrame(data_import_polars, infer_schema_length=100000) df_polars = df_polars.drop("_id") # =====</pre>   |
|  | <p>Figure E9: Polars Code for Import Data</p>  |
| Modin  | <p>The procedure for loading data into a Modin DataFrame is shown in Figure E10. In order to speed up Pandas operations through parallel processing, the MongoDB documents are first gathered as a list and then transformed into a Modin DataFrame. After that, the <code>_id</code> column is eliminated to get the data ready for analysis.</p> <pre># === Read from MongoDB into Modin DataFrame === data_import_modin = list(collection.find()) df_modin = md.DataFrame(data_import_modin) df_modin.drop(columns=['_id'], inplace=True) # =====</pre>   |
|  | <p>Figure E10: Modin Code for Import Data</p>  |
| Dask   | <p>The process of reading data into a Dask DataFrame is shown in Figure E11. Initially, MongoDB data is retrieved and transformed into a Pandas DataFrame. Dask is then used to divide this Pandas DataFrame into several smaller pieces so that parallel computation is possible. To maintain the data's relevance for upcoming tasks, the <code>_id</code> column is removed.</p> <pre># === Read from MongoDB into Dask DataFrame === data = list(collection.find()) # Fetch all documents as a list of dicts df_dask = dd.from_pandas(pd.DataFrame(data), npartitions=4) df_dask = df_dask.drop(columns=['_id']) # Drop the '_id' column # =====</pre> |
| <p>Figure E11: Dask Code for Import Data</p> |  |
| <p><b>Combine Columns</b></p>                |  |

|        |   |
|--------|---|
| Pandas | <p>Figure E12 shows the Pandas approach to combining the Seat Capacity and Seating Capacity columns. It uses the <code>combine_first()</code> method, which efficiently fills missing values in the first column using values from the second column.</p> <pre># Combine 'Seat Capacity' and 'Seating Capacity' into one column df_pandas['Seat Capacity'] = df_pandas['Seat Capacity'].combine_first(df_pandas['Seating Capacity'])  # Drop 'Seating Capacity' column df_pandas.drop(columns=['Seating Capacity'], inplace=True)</pre>   |
| Polars | <p>Figure E13 demonstrates the process in Polars, where the <code>coalesce()</code> function is used after converting "NaN" strings into proper nulls. This method is optimised with built-in multithreading, allowing it to perform faster on large datasets without any extra configuration.</p> <pre># Replace string "NaN" with nulls df_polars = df_polars.with_columns([     pl.when(pl.col("Seat Capacity") == "NaN").then(None).otherwise(pl.col("Seat Capacity")).alias("Seat Capacity"),     pl.when(pl.col("Seating Capacity") == "NaN").then(None).otherwise(pl.col("Seating Capacity")).alias("Seating Capacity") ])  # Combine the two columns using coalesce df_polars = df_polars.with_columns([     pl.coalesce(["Seat Capacity", "Seating Capacity"]).alias("Seat Capacity") ])  # Drop column df_polars = df_polars.drop("Seating Capacity")</pre> |
| Modin  | <p>Figure E14 displays the Modin approach, using the <code>fillna()</code> function to fill missing values in Seat Capacity with values from Seating Capacity. This is similar to Pandas but benefits from Modin's parallel execution for faster performance on large data.</p> <pre># Combine 'Seat Capacity' and 'Seating Capacity' into one column df_modin['Seat Capacity'] = df_modin['Seat Capacity'].fillna(df_modin['Seating Capacity'])  # Drop 'Seating Capacity' column df_modin.drop(columns=['Seating Capacity'], inplace=True)</pre>  |

|                                  |   |
|----------------------------------|---|
| Dask                             | <p>Figure E15 shows the operation using Dask, also using <code>combine_first()</code> like Pandas. In this context, it runs in parallel on local cores, giving moderate speed improvements.</p> <pre># Combine 'Seat Capacity' and 'Seating Capacity' into one column df_dask['Seat Capacity'] = df_dask['Seat Capacity'].combine_first(df_dask['Seating Capacity'])  # Drop 'Seating Capacity' column df_dask = df_dask.drop(columns=['Seating Capacity'])</pre> |
| <b>Remove Unwanted Column(s)</b> |   |
| Pandas                           | <p>Figure E16 shows the removal of the URL column using Pandas. The <code>drop()</code> function is used with <code>inplace=True</code>, and the operation runs on a single thread.</p> <pre># Remove the 'URL' column df_pandas.drop(columns=['URL'], inplace=True)</pre>  |
|                                  | Figure E16: Pandas Code for Remove Unwanted Column(s)   |
| Polars                           | <p>Figure E17 demonstrates the same column removal in Polars. The <code>drop()</code> method is applied, and the operation benefits from Polars' built-in multithreading, making it faster on larger datasets.</p> <pre># Remove the 'URL' column df_polars = df_polars.drop("URL")</pre>   |
|                                  | Figure E17: Polars Code for Remove Unwanted Column(s)   |
| Modin                            | <p>Figure E18 illustrates how the URL column is removed in Modin using the same syntax as Pandas. However, unlike Pandas, Modin executes the operation in parallel across CPU cores, providing performance gains.</p> <pre># Remove the 'URL' column df_modin.drop(columns=['URL'], inplace=True)</pre>   |
|                                  | Figure E18: Modin Code for Remove Unwanted Column(s)  |

|  |  |
|--|--|
| Dask   | <p>Figure E19 shows the URL column being dropped in Dask. The operation runs in parallel on partitions.</p> <pre># Remove the 'URL' column df_dask = df_dask.drop(columns=['URL'])</pre>   |
| <b>String Cleaning and Formatting (Clean Installment Column)</b> |  |
| Pandas   | <p>Figure E20 shows the cleaning process of the Installment column using Pandas. The string replacement operations and numeric conversion are performed sequentially.</p> <pre># Cleaning operations df_pandas['Installment'] = df_pandas['Installment'].str.replace('RM', '', regex=False) df_pandas['Installment'] = df_pandas['Installment'].str.replace(',', '', regex=False) df_pandas['Installment'] = df_pandas['Installment'].str.replace('/month', '', regex=False) df_pandas['Installment'] = pd.to_numeric(df_pandas['Installment'], errors='coerce')</pre> |
|  | <p>Figure E20: Pandas Code for Clean Installment Column</p>  |
| Polars   | <p>Figure E21 shows how Polars cleans the Installment column with optimised execution. The string replacements and type casting are combined within a single expression using the with_columns() method.</p> <pre># Clean the 'Installment' column: remove 'RM', ',', '/month', strip spaces, and convert to int df_polars = df_polars.with_columns(     pl.col("Installment")         .str.replace_all("RM", "")         .str.replace_all(",", "")         .str.replace_all("/month", "")         .str.strip_chars()         .cast(pl.Int64) )</pre>                  |
|  | <p>Figure E21: Polars Code for Clean Installment Column</p>  |
| Modin  | <p>Figure E22 shows the Modin implementation for cleaning the Installment column. Here, the string manipulations and numeric conversion benefit from Modin's parallel execution framework, which splits the data and processes it across multiple CPU cores.</p>   |

|  |  |
|--|--|
|  | <pre> # Cleaning operations df_modin['Installment'] = df_modin['Installment'].str.replace('RM', '', regex=False) df_modin['Installment'] = df_modin['Installment'].str.replace(',', '', regex=False) df_modin['Installment'] = df_modin['Installment'].str.replace('/month', '', regex=False) df_modin['Installment'] = pd.to_numeric(df_modin['Installment'], errors='coerce') </pre>   |
|  | Figure E22: Modin Code for Clean Installment Column  |
| <b>String Cleaning and Formatting (Extract Clean Values from Condition Column)</b> |  |
| Pandas   | <p>Figure E23 shows the cleaning of the Installment column using Dask. The DataFrame is partitioned, and string replacements along with the numeric conversion are applied in parallel over these partitions.</p> <pre> # Convert to string before applying str methods df_dask['Installment'] = df_dask['Installment'].astype(str) df_dask['Installment'] = df_dask['Installment'].str.replace('RM', '', regex=False) df_dask['Installment'] = df_dask['Installment'].str.replace(',', '', regex=False) df_dask['Installment'] = df_dask['Installment'].str.replace('/month', '', regex=False) df_dask['Installment'] = dd.to_numeric(df_dask['Installment'], errors='coerce') </pre> |
|  | Figure E23: Dask Code for Clean Installment Column   |
| Polars   | <p>Figure E24 shows the Pandas method where the regex extraction and string replacement are applied using .apply() and .str.replace() respectively. Both operations run sequentially without parallelisation, so no specific optimisation is applied in this code.</p> <pre> # Extract value from '...schema.org/Condition' df_pandas['Condition'] = df_pandas['Condition'].apply(     lambda x: re.search(r'\.\org/([A-Za-z]+)Condition', x).group(1) if isinstance(x, str) and 'schema.org/' in x else x )  # Strip 'Condition' if it remains (for post-processed cases) df_pandas['Condition'] = df_pandas['Condition'].str.replace(r'Condition\$', '', regex=True) </pre>          |
|  | Figure E24: Pandas Code for Clean Condition Column   |
| Polars   | <p>Figure E25 demonstrates Polars' approach where string extraction and replacement are done via fully vectorised expressions inside .with_columns().</p>  |

```

# Vectorized: Extract text before 'Condition' after the last '/' (e.g., 'NewCondition' → 'New')
df_polars = df_polars.with_columns(
    pl.col("Condition")
        .str.extract(r"/([A-Za-z]+)Condition", 1)
        .fill_null(
            pl.col("Condition").str.replace(r"Condition$", "", literal=False)
        )
    .alias("Condition")
)

```

Figure E25: Polars Code for Clean Condition Column

|       |   |
|-------|---|
|       | <p>Figure E26 uses Modin to perform regex extraction with <code>.apply()</code> and string replacement with <code>.str.replace()</code>. The optimisation occurs because Modin distributes these Pandas-like operations across multiple CPU cores, enabling parallel execution transparently.</p> <pre> # First: Extract value from '...schema.org/XCondition' df_modin['Condition'] = df_modin['Condition'].apply(     lambda x: re.search(r'\.org/([A-Za-z]+)Condition', x).group(1) if isinstance(x, str) and 'schema.org/' in x else x )  # Second: Strip 'Condition' if it remains (for post-processed cases) df_modin['Condition'] = df_modin['Condition'].str.replace(r'Condition\$', '', regex=True) </pre> |
| Modin | <p>Figure E26 uses Modin to perform regex extraction with <code>.apply()</code> and string replacement with <code>.str.replace()</code>. The optimisation occurs because Modin distributes these Pandas-like operations across multiple CPU cores, enabling parallel execution transparently.</p> <pre> # First: Extract value from '...schema.org/XCondition' df_modin['Condition'] = df_modin['Condition'].apply(     lambda x: re.search(r'\.org/([A-Za-z]+)Condition', x).group(1) if isinstance(x, str) and 'schema.org/' in x else x )  # Second: Strip 'Condition' if it remains (for post-processed cases) df_modin['Condition'] = df_modin['Condition'].str.replace(r'Condition\$', '', regex=True) </pre> |

Figure E26: Modin Code for Clean Condition Column

|      |  |
|------|--|
|      | <p>Figure E27 applies regex extraction via <code>.apply()</code> with the <code>meta</code> parameter to define output type, allowing Dask to parallelise this operation on each DataFrame partition. The subsequent string replacement runs in parallel across these partitions, providing moderate speedup over sequential processing.</p> <pre> # First: Extract value from '...schema.org/XCondition' df_dask['Condition'] = df_dask['Condition'].apply(     lambda x: re.search(r'\.org/([A-Za-z]+)Condition', x).group(1) if isinstance(x, str) and 'schema.org/' in x else x,     meta=('Condition', 'object') )  # Second: Strip 'Condition' if it remains (for post-processed cases) df_dask['Condition'] = df_dask['Condition'].str.replace(r'Condition\$', '', regex=True) </pre> |
| Dask | <p>Figure E27 applies regex extraction via <code>.apply()</code> with the <code>meta</code> parameter to define output type, allowing Dask to parallelise this operation on each DataFrame partition. The subsequent string replacement runs in parallel across these partitions, providing moderate speedup over sequential processing.</p> <pre> # First: Extract value from '...schema.org/XCondition' df_dask['Condition'] = df_dask['Condition'].apply(     lambda x: re.search(r'\.org/([A-Za-z]+)Condition', x).group(1) if isinstance(x, str) and 'schema.org/' in x else x,     meta=('Condition', 'object') )  # Second: Strip 'Condition' if it remains (for post-processed cases) df_dask['Condition'] = df_dask['Condition'].str.replace(r'Condition\$', '', regex=True) </pre> |

Figure E27: Dask Code for Clean Condition Column

| <b>String Cleaning and Formatting (Clean Sales Channel)</b> |  |
|---|--|
| Pandas  | <p>Figure E28 uses Pandas' <code>.str.extract()</code> method to apply regex extraction on the "Sales Channel" column.</p> <pre> # Extract Sales Channel (either 'Sales Agent' or 'Dealer') df_pandas['Sales Channel'] = df_pandas['Sales Channel'].str.extract(r'^Sales Agent Dealer') </pre> |

|   |  |
|---|--|
|   | Figure E28: Pandas Code for Clean Sales Channel Column   |
| Polars  | Figure E29 shows Polars performing regex extraction using the vectorised .str.extract() expression within .with_columns().   |
|   | <pre># Extract either "Sales Agent" or "Dealer" using regex df_polars = df_polars.with_columns(     pl.col("Sales Channel")         .str.extract(r'^Sales Agent Dealer", 1)         .alias("Sales Channel") )</pre>  |
|   | Figure E29: Polars Code for Clean Sales Channel Column   |
| Modin   | Figure E30 applies .str.extract() in Modin, which internally parallelises the operation across multiple CPU cores by distributing the Pandas-like computation, resulting in faster execution compared to vanilla Pandas.   |
|   | <pre># Extract Sales Channel (either 'Sales Agent' or 'Dealer') df_modin['Sales Channel'] = df_modin['Sales Channel'].str.extract(r'^Sales Agent Dealer')</pre>  |
|   | Figure E30: Modin Code for Clean Sales Channel Column  |
| Dask  | Figure E31 demonstrates Dask's use of .str.extract() on its partitioned DataFrame, where the extraction runs in parallel on each partition, leveraging Dask's partitioned task scheduling for improved speed on large data.  |
|   | <pre># Extract Sales Channel (either 'Sales Agent' or 'Dealer') df_dask['Sales Channel'] = df_dask['Sales Channel'].str.extract(r'^Sales Agent Dealer', expand=False)</pre>  |
|   | Figure E31: DaskCode for Clean Sales Channel Column  |
| <b>String Cleaning and Formatting (Clean Mileage)</b> |  |
| Pandas  | Figure E32 shows Pandas cleaning mileage by first removing the substring "K KM" with a vectorised .str.replace() and then applying a custom Python function via .apply(). The .apply() runs sequentially on the entire column, so the performance is limited by single-thread execution. |

```

# Remove "K KM" and clean mileage into 5K step ranges
def process_mileage_pandas(mileage):
    if not isinstance(mileage, str) or mileage.strip() == '':
        return None
    try:
        if '-' in mileage:
            start, end = mileage.split('-')
            return f"{int(start.strip())} - {int(end.strip())}"
        else:
            val = int(mileage.strip())
            lower = (val // 5) * 5
            upper = lower + 5
            return f"{lower} - {upper}"
    except:
        return None

df_pandas['Mileage'] = df_pandas['Mileage'].str.replace('K KM', '', regex=False)
df_pandas['Mileage'] = df_pandas['Mileage'].apply(process_mileage_pandas)

```

Figure E32: Pandas Code for Clean Mileage Column

|        |   |
|--------|---|
| Polars | <p>Figure E33 uses Polars' <code>.map_elements()</code> to apply a custom cleaning function on each element of the "Mileage" column. Polars optimises this by executing the function in a parallel and compiled manner internally, giving it better speed than pure Python <code>.apply()</code> in Pandas.</p> <pre> # Define a function to clean mileage and convert into 5K step range def clean_mileage(val):     if val is None or str(val).strip() == '':         return None     val = str(val).replace('K KM', '').strip()     try:         if '-' in val:             start, end = val.split('-')             return f"{int(start.strip())} - {int(end.strip())}"         else:             num = int(val)             lower = (num // 5) * 5             upper = lower + 5             return f"{lower} - {upper}"     except:         return None  # Apply the cleaning using map_elements df_polars = df_polars.with_columns(     pl.col("Mileage")     .map_elements(clean_mileage, return_dtype=pl.Utf8)     .alias("Mileage") ) </pre> |
|--------|---|

|       |  |
|-------|--|
|       | Figure E33: Polars Code for Clean Mileage Column   |
| Modin | <p>Figure E34 applies the same approach as Pandas but leverages Modin's parallelised .apply(). Modin partitions the DataFrame and distributes the .apply() function across CPU cores, enabling concurrent execution and improved performance over standard Pandas.</p> <pre># Remove "K KM" and clean mileage into 5K step ranges def process_mileage_modin(mileage):     if not isinstance(mileage, str) or mileage.strip() == '':         return None     try:         if '-' in mileage:             start, end = mileage.split('-')             return f"{int(start.strip())} - {int(end.strip())}"         else:             val = int(mileage.strip())             lower = (val // 5) * 5             upper = lower + 5             return f"{lower} - {upper}"     except:         return None  df_modin['Mileage'] = df_modin['Mileage'].str.replace('K KM', '', regex=False) df_modin['Mileage'] = df_modin['Mileage'].apply(process_mileage_modin)</pre> |
| Dask  | <p>Figure E35 performs a two-step process: a vectorised .str.replace() to remove "K KM", followed by a .apply() with a custom function executed on each Dask partition. Dask's task scheduler runs these .apply() functions in parallel across partitions, accelerating the processing on large datasets.</p>  |

```

# Remove "K KM"
df_dask['Mileage'] = df_dask['Mileage'].str.replace('K KM', '', regex=False)

# Clean and standardize mileage into 5K step ranges
def process_mileage_dask(mileage):
    if not isinstance(mileage, str) or mileage.strip() == '':
        return None
    try:
        if '-' in mileage:
            start, end = mileage.split('-')
            return f'{int(start.strip())} - {int(end.strip())}'
        else:
            val = int(mileage.strip())
            lower = (val // 5) * 5
            upper = lower + 5
            return f'{lower} - {upper}'
    except:
        return None

df_dask['Mileage'] = df_dask['Mileage'].apply(process_mileage_dask, meta=('x', 'object'))

```

Figure E35: Dask Code for Clean Mileage Column

## Replace Null Values

Pandas

Figure E36 shows Pandas replacing empty strings and "-" with np.nan using vectorisation. .replace() calls are efficient built-in string operations. The subsequent .dropna() call removes rows with missing values in specified columns, operating in memory and sequentially.

```

# Convert to np.nan
df_pandas.replace(r'^\s*$', np.nan, regex=True, inplace=True) # Empty and whitespace-only strings
df_pandas.replace('-', np.nan, inplace=True) # "-"

# Check null
df_pandas.isnull().sum()

# Drop rows with missing values in specific columns
df_pandas.dropna(subset=['Body Type', 'Fuel Type', 'Mileage', 'Sales Channel', 'Seat Capacity'], inplace=True)

```

Figure E36: Pandas Code for Replace Null Values

Polars

Figure E37 uses Polars' lazy evaluation and expression API to replace empty, "-", or "NaN" strings with null values in all string columns via vectorisation .when().then().otherwise() expressions. The .drop\_nulls() method then drops rows with nulls efficiently using Polars' optimised execution engine, which is highly parallelised and memory efficient.

```

string_cols = [col for col, dtype in zip(df_polars.columns, df_polars.dtypes) if dtype == pl.Utf8]

df_polars = df_polars.with_columns([
    pl.when(
        [
            pl.col(col).str.strip_chars().is_in(["", "-", "NaN"])
        ]
    ).then(None)
    .otherwise(pl.col(col))
    .alias(col)
    for col in string_cols
])

# Show null counts in vertical format
for col in df_polars.columns:
    null_counts_before = df_polars.select(pl.col(col).null_count()).item()
    print(f"{col:<25} : {null_counts_before}")

# Drop rows with missing values in specific columns
columns_to_check = ['Body Type', 'Fuel Type', 'Mileage', 'Sales Channel', 'Seat Capacity']
df_polars = df_polars.drop_nulls(subset=columns_to_check)

```

Figure E37: Polars Code for Replace Null Values

Modin

Figure E38 applies Pandas-like `.replace()` with regex support and `.dropna()` on Modin DataFrame. Modin optimises these operations by splitting data into partitions and executing replacements and null drops concurrently across multiple cores, speeding up these typically expensive row filtering steps.

```

# Convert to np.nan
df_modin.replace(r'^\s*$', np.nan, regex=True, inplace=True) # Empty and whitespace-only strings
df_modin.replace('-', np.nan, inplace=True) # "-"

# Check null
df_modin.isnull().sum()

# Drop rows with missing values in specific columns
df_modin.dropna(subset=['Body Type', 'Fuel Type', 'Mileage', 'Sales Channel', 'Seat Capacity'], inplace=True)

```

Figure E38: Modin Code for Replace Null Values

Dask

Figure E39 performs `.replace()` operations to convert blanks and `"-"` to `np.nan`, then `.dropna()` to remove incomplete rows on Dask DataFrame. Dask parallelises these steps by partitioning data and executing the replacements and null drops across partitions concurrently, then triggers computation with `.compute()` for aggregated results.

```

# Convert to np.nan
df_dask = df_dask.replace(r'^\s*$', np.nan, regex=True) # Empty and whitespace
df_dask = df_dask.replace('-', np.nan) # -

# Check for null values in each column
df_dask.isnull().sum().compute()

# Drop rows with missing values in specific columns
df_dask = df_dask.dropna(subset=['Body Type', 'Fuel Type', 'Mileage', 'Sales Channel', 'Seat Capacity'])

```

Figure E39: Dask Code for Replace Null Values

## Rename Columns

Pandas

Figure E40 illustrates the column renaming in Pandas. The .rename() method with inplace=True is applied, which modifies the DataFrame directly in memory. This method is efficient for smaller to medium datasets, as it avoids creating a copy of the DataFrame and quickly updates column labels.

```

df_pandas.rename(columns={
    'Mileage': 'Mileage (K KM)',
    'Price': 'Price (RM)',
    'Installment': 'Installment (RM)'
}, inplace=True)

```

Figure E40: Pandas Code for Rename Columns

Polars

Figure 5.2.41 shows the renaming process in Polars. Since Polars operates on immutable DataFrames, .rename() returns a new DataFrame with the updated column names. This fits Polars' optimised execution model, allowing the operation to be part of a fast, chainable query pipeline suitable for large datasets.

```

df_polars = df_polars.rename({
    'Mileage': 'Mileage (K KM)',
    'Price': 'Price (RM)',
    'Installment': 'Installment (RM)'
})

```

Figure E41: Polars Code for Rename Columns

|                             |  |
|-----------------------------|--|
| Modin                       | <p>Figure E42 demonstrates the renaming method used by Modin. It applies <code>.rename()</code> with <code>inplace=True</code>, similar to Pandas, but internally Modin distributes the task across multiple cores or nodes, enabling parallelised column renaming. This parallelism improves performance when handling large-scale data.</p> <pre data-bbox="595 496 1281 707"><code>df_modin.rename(columns={     'Mileage': 'Mileage (K KM)',     'Price': 'Price (RM)',     'Installment': 'Installment (RM)' }, inplace=True)</code></pre>          |
| Dask                        | <p>Figure E43 presents the renaming step in Dask. The <code>.rename()</code> method returns a lazy DataFrame with new column names, deferring actual computation. This lazy evaluation allows Dask to optimise the execution plan and efficiently handle renaming as part of a larger distributed workflow before triggering <code>.compute()</code>.</p> <pre data-bbox="612 1108 1264 1320"><code>df_dask = df_dask.rename(columns={     'Mileage': 'Mileage (K KM)',     'Price': 'Price (RM)',     'Installment': 'Installment (RM)' })</code></pre> |
| <b>Drop Duplicated Rows</b> |  |
| Pandas                      | <p>Figure E44 illustrates how Pandas handles duplicate detection and removal using in-memory operations optimised for single-threaded performance. Pandas applies vectorised methods but can be limited by memory when working with very large datasets.</p> <pre data-bbox="465 1742 1428 1869"><code># Check for duplicate rows duplicates_before_pandas = df_pandas.duplicated().sum() print(f"Number of duplicate rows: {duplicates_before_pandas}")</code></pre>  |

|        |  |
|--------|--|
|        | <pre># === Drop Duplicate Rows === df_pandas = df_pandas.drop_duplicates() # =====</pre>   |
|        | <p>Figure E44: Pandas Code for Drop Duplicated Rows</p>  |
| Polars | <p>Figure E45 shows Polars using its built-in unique() method, which is implemented in Rust and optimised for parallel execution. This allows Polars to efficiently drop duplicates with lower memory usage and faster speed compared to pandas.</p> <pre># === Drop Duplicate Rows in Polars === df_polars = df_polars.unique() # Drops duplicate rows # =====</pre>  |
|        | <p>Figure E45: Polars Code for Drop Duplicated Rows</p>  |
| Modin  | <p>Figure E46 depicts Modin's approach, which distributes the duplicate check and removal tasks across multiple cores or nodes using Ray or Dask backends. This parallelisation provides scalability and improved runtime on large datasets relative to pandas.</p> <pre># Check for duplicate rows duplicates_before_modin = df_modin.duplicated().sum() print(f"Number of duplicate rows: {duplicates_before_modin}")  # === Drop Duplicate Rows === df_modin = df_modin.drop_duplicates() # =====</pre> |
|        | <p>Figure E46: Modin Code for Drop Duplicated Rows</p>   |
| Dask   | <p>Figure E47 demonstrates Dask's method, which computes the lazy DataFrame into pandas to leverage pandas' duplicate detection capabilities.</p>  |

```

print("Computing Dask DataFrame to pandas...")
df_pandas_duplicates_before = df_dask.compute()

# Use pandas' duplicated method
duplicates_before_dask = df_pandas_duplicates_before.duplicated().sum()
print(f"Number of duplicate rows: {duplicates_before_dask}")

# === Drop Duplicate Rows ===
df_dask = df_dask.drop_duplicates()

```

Figure E47: Dask Code for Drop Duplicated Rows

### Optimise Memory Usage by Using Efficient Data Types

Pandas

Figure E48 illustrates the numeric downcasting and categorical conversion process applied to the pandas DataFrame. The numeric columns ‘Manufacture Year’, ‘Price (RM)’, ‘Installment (RM)’, and ‘Seat Capacity’ were downcast to smaller integer types to reduce memory usage. Additionally, categorical columns such as ‘Car Brand’, ‘Body Type’, and ‘Fuel Type’ were converted to the category dtype to optimise storage and improve performance in subsequent operations.

```

# === Numeric Downcasting and Categorical Conversion for pandas ===
df_pandas['Manufacture Year'] = df_pandas['Manufacture Year'].astype('int16')
df_pandas[['Price (RM)', 'Installment (RM)']] = df_pandas[['Price (RM)', 'Installment (RM)']].astype('int32')
df_pandas['Seat Capacity'] = pd.to_numeric(df_pandas['Seat Capacity'], errors='coerce').astype('int8')

# Categorical conversion
categorical_columns_numeric_downcasting_pandas = ['Car Brand', 'Car Model', 'Body Type', 'Fuel Type',
| 'Transmission', 'Color', 'Condition', 'Sales Channel', 'Mileage (K KM)']
df_pandas[categorical_columns_numeric_downcasting_pandas] = df_pandas[categorical_columns_numeric_downcasting_pandas].apply(lambda x: x.astype('category'))
# *****

```

Figure E48: Pandas Code for Optimise Memory Usage by Using Efficient Data Types

Polars

Figure E49 depicts the implementation of numeric downcasting and categorical conversion in the Polars DataFrame. The relevant numeric columns were cast to smaller integer types using Polars’ casting methods. For categorical data, each specified column was cast to the categorical type, facilitating more efficient memory consumption and faster query performance.

|       |  |
|-------|--|
|       | <pre> # === Numeric Downcasting and Categorical Conversion for Polars === df_polars = df_polars.with_columns([     pl.col("Manufacture Year").cast(pl.Int16),     pl.col("Price (RM)").cast(pl.Int32),     pl.col("Installment (RM)").cast(pl.Int32),     pl.col("Seat Capacity").cast(pl.Int8), ])  # Categorical conversion categorical_columns_numeric_downcasting_polars = ['Car Brand', 'Car Model', 'Body Type', 'Fuel Type',     'Transmission', 'Color', 'Condition', 'Sales Channel', 'Mileage (K KM)'] for col in categorical_columns_numeric_downcasting_polars:     df_polars = df_polars.with_columns([         pl.col(col).cast(pl.Categorical)     ]) </pre>  |
|       | <p>Figure E49: Polars Code for Optimise Memory Usage by Using Efficient Data Types</p>   |
| Modin | <p>Figure E50 shows the numeric downcasting and categorical conversion carried out on the Modin DataFrame. Similar to pandas, key numeric columns were downcasted to reduced integer types, and selected categorical columns were converted to the category datatype.</p> <pre> # === Numeric Downcasting and Categorical Conversion for Modin === df_modin['Manufacture Year'] = df_modin['Manufacture Year'].astype('int16') df_modin[['Price (RM)', 'Installment (RM)']] = df_modin[['Price (RM)', 'Installment (RM)']].astype('int32') df_modin['Seat Capacity'] = pd.to_numeric(df_modin['Seat Capacity'], errors='coerce').astype('int8')  # Categorical conversion categorical_columns_numeric_downcasting_modin = ['Car Brand', 'Car Model', 'Body Type', 'Fuel Type',     'Transmission', 'Color', 'Condition', 'Sales Channel', 'Mileage (K KM)'] df_modin[categorical_columns_numeric_downcasting_modin] = df_modin[categorical_columns_numeric_downcasting_modin].apply(lambda x: x.astype('category')) # ===== </pre> |
|       | <p>Figure E50: Modin Code for Optimise Memory Usage by Using Efficient Data Types</p>  |
| Dask  | <p>Figure E51 presents the numeric downcasting and categorical conversion performed on the Dask DataFrame. The numeric columns were downcasted to smaller integer types, while categorical columns were individually converted to the category dtype.</p> <pre> # Numeric downcasting for Dask df_dask['Manufacture Year'] = df_dask['Manufacture Year'].astype('int16') df_dask[['Price (RM)', 'Installment (RM)']] = df_dask[['Price (RM)', 'Installment (RM)']].astype('int32') df_dask['Seat Capacity'] = dd.to_numeric(df_dask['Seat Capacity'], errors='coerce').astype('int8')  # Categorical conversion categorical_columns_numeric_downcasting_dask = ['Car Brand', 'Car Model', 'Body Type', 'Fuel Type',     'Transmission', 'Color', 'Condition', 'Sales Channel', 'Mileage (K KM)'] # Convert each column to category individually for col in categorical_columns_numeric_downcasting_dask:     df_dask[col] = df_dask[col].astype('category') </pre>   |
|       | <p>Figure E51: Dask Code for Optimise Memory Usage by Using Efficient Data Types</p>   |

## Export Cleaned Data to MongoDB

Pandas

Figure E52 shows the cleaned pandas DataFrame being converted to dictionaries and inserted into MongoDB. The process used direct conversion for efficient data export.

```
# === Insert cleaned DataFrame into MongoDB ===
client_export_pandas = MongoClient("mongodb://localhost:27017")
db_export_pandas = client_export_pandas["carlist_db"]

df_pandas_carlist = df_pandas.copy() # Rename cleaned DataFrame

# Create a new collection for the cleaned data
cleaned_collection_export_pandas = db_export_pandas["listings_pandas_cleaned"]

# Insert cleaned data
cleaned_collection_export_pandas.insert_many(df_pandas_carlist.to_dict(orient="records"))
# =====
```

Figure E52: Pandas Code for Export Cleaned Data to MongoDB

Polars

Figure E53 illustrates the cleaned Polars DataFrame cloned and converted to dictionaries before insertion into MongoDB, leveraging Polars' fast data handling.

```
# === Insert cleaned DataFrame into MongoDB ===
client_export_polars = MongoClient("mongodb://localhost:27017")
db_export_polars = client_export_polars["carlist_db"]

df_polars_carlist = df_polars.clone() # Rename cleaned DataFrame

# Create a new collection for the cleaned data
cleaned_collection_export_polars = db_export_polars["listings_polars_cleaned"]

# Insert cleaned data
cleaned_collection_export_polars.insert_many(df_polars_carlist.to_dicts())
# =====
```

Figure E53: Polars Code for Export Cleaned Data to MongoDB

Modin

Figure E54 depicts the cleaned Modin DataFrame copied and transformed to dictionary format for batch insertion into MongoDB, benefiting from Modin's parallel processing.

```

# === Insert cleaned DataFrame into MongoDB ===
client_export_modin = MongoClient("mongodb://localhost:27017")
db_export_modin = client_export_modin["carlist_db"]

df_modin_carlist = df_modin.copy() # Rename cleaned DataFrame

# Create a new collection for the cleaned data
cleaned_collection_export_modin = db_export_modin["listings_modin_cleaned"]

# Insert cleaned data
cleaned_collection_export_modin.insert_many(df_modin_carlist.to_dict(orient="records"))
# =====

```

Figure E54: Modin Code for Export Cleaned Data to MongoDB

|      |   |
|------|---|
| Dask | <p>Figure E55 presents the cleaned Dask DataFrame being computed to pandas, then converted and inserted into MongoDB, accounting for Dask's lazy evaluation.</p> <pre> # === Insert cleaned Dask DataFrame into MongoDB ===  client_export_dask = MongoClient("mongodb://localhost:27017") db_export_dask = client_export_dask["carlist_db"]  # Compute Dask DataFrame to Pandas DataFrame for MongoDB insertion df_dask_carlist = df_dask.compute() # Convert Dask DataFrame to Pandas DataFrame  # Create a new collection for the cleaned data cleaned_collection_export_dask = db_export_dask["listings_dask_cleaned"]  # Insert cleaned data into MongoDB cleaned_collection_export_dask.insert_many(df_dask_carlist.to_dict(orient="records")) # ===== </pre> |
|------|---|

Figure E55: Dask Code for Export Cleaned Data to MongoDB

## 10.6 Appendix F

Table F1 summarises the performance metrics collected during data processing tasks for each library, including time, CPU usage, throughput, and memory consumption.

Table 6.2.1: Performance Metrics Comparison Across Data Processing Libraries

|                           | Metrics       | Elapsed Time (Wall Time) (s) | CPU Time (ms) | CPU Usage (%) | Throughput (rows/sec) | Current Memory Usage (MB) | Peak Memory Usage (MB) |
|---------------------------|---------------|------------------------------|---------------|---------------|-----------------------|---------------------------|------------------------|
| Steps                     | Library       |                              |               |               |                       |                           |                        |
| Import Data               | <b>Pandas</b> | 10.7775                      | 14828         | 13.00         | 16158.60              | 448.21                    | 524.84                 |
|                           | <b>Polars</b> | 18.5478                      | 23859         | 13.00         | 9389.27               | 424.51                    | 436.65                 |
|                           | <b>Modin</b>  | 62.6945                      | 187359        | 24.60         | 2777.75               | 442.99                    | 542.54                 |
|                           | <b>Dask</b>   | 45.9962                      | 90578         | 14.80         | 3786.18               | 450.60                    | 807.31                 |
| Combine Columns           | <b>Pandas</b> | 0.0824                       | 62            | 8.70          | 2112196.89            | 22.34                     | 23.71                  |
|                           | <b>Polars</b> | 0.0393                       | 94            | 15.40         | 4429832.37            | 0.04                      | 0.06                   |
|                           | <b>Modin</b>  | 0.4683                       | 641           | 28.00         | 371837.94             | 0.83                      | 1.24                   |
|                           | <b>Dask</b>   | 0.1055                       | 78            | 7.20          | 1650515.34            | 0.20                      | 0.65                   |
| Remove Unwanted Column(s) | <b>Pandas</b> | 0.0405                       | 31            | 14.80         | 4298248.86            | 20.92                     | 20.94                  |
|                           | <b>Polars</b> | 0.0011                       | 0             | 6.80          | 154333569.63          | 0.00                      | 0.02                   |
|                           | <b>Modin</b>  | 0.0048                       | 0             | 7.50          | 36599974.78           | 0.01                      | 0.03                   |
|                           | <b>Dask</b>   | 0.1541                       | 78            | 8.00          | 1130332.57            | 0.13                      | 0.65                   |
| Clean Installme           | <b>Pandas</b> | 1.8746                       | 2047          | 9.80          | 92898.42              | 19.13                     | 34.51                  |

|  |               |        |      |       |             |       |       |
|--|---------------|--------|------|-------|-------------|-------|-------|
| nt<br>Column   | <b>Polars</b> | 0.0556 | 94   | 6.90  | 3129715.50  | 0.02  | 0.04  |
|  | <b>Modin</b>  | 3.9595 | 5094 | 5.50  | 43983.06    | 3.11  | 19.68 |
|  | <b>Dask</b>   | 0.1376 | 109  | 8.10  | 1265861.43  | 0.26  | 0.66  |
| Extract<br>Clean<br>Values<br>from<br>Condition<br>Column        | <b>Pandas</b> | 1.2767 | 1375 | 8.50  | 13604.88    | 9.13  | 16.64 |
|  | <b>Polars</b> | 0.0710 | 62   | 11.00 | 2453134.74  | 0.02  | 0.29  |
|  | <b>Modin</b>  | 0.7658 | 1391 | 9.50  | 227422.13   | 1.31  | 1.75  |
|  | <b>Dask</b>   | 0.3295 | 375  | 11.80 | 528492.96   | 0.28  | 0.78  |
| Extract<br>Clean<br>Values<br>from<br>Sales<br>Channel<br>Column | <b>Pandas</b> | 1.2360 | 1578 | 7.30  | 140902.18   | 9.84  | 27.68 |
|  | <b>Polars</b> | 0.0217 | 16   | 9.50  | 8010594.42  | 0.02  | 0.04  |
|  | <b>Modin</b>  | 1.7670 | 2922 | 12.50 | 98556.84    | 1.67  | 19.02 |
|  | <b>Dask</b>   | 0.3289 | 703  | 13.60 | 529517.81   | 0.20  | 0.72  |
| Clean<br>Mileage<br>Column<br>Into 5K<br>Step<br>Ranges          | <b>Pandas</b> | 2.5527 | 2719 | 11.90 | 68222.64    | 9.75  | 26.75 |
|  | <b>Polars</b> | 2.5342 | 2922 | 12.50 | 68719.17    | 0.28  | 0.54  |
|  | <b>Modin</b>  | 0.8994 | 2109 | 20.60 | 193626.21   | 1.28  | 1.89  |
|  | <b>Dask</b>   | 0.3569 | 203  | 10.50 | 487891.87   | 0.26  | 0.75  |
| Replace<br>Null<br>Values  | <b>Pandas</b> | 0.2154 | 281  | 14.80 | 808379.11   | 21.13 | 29.60 |
|  | <b>Polars</b> | 0.0112 | 0    | 7.40  | 15585287.27 | 0.01  | 0.03  |
|  | <b>Modin</b>  | 4.0161 | 6125 | 16.80 | 43362.71    | 6.92  | 8.30  |

|  |               |         |       |       |             |       |        |
|--|---------------|---------|-------|-------|-------------|-------|--------|
|  | <b>Dask</b>   | 21.1509 | 40359 | 19.90 | 8233.69     | 5.05  | 344.34 |
| <b>Rename Columns</b>                                      | <b>Pandas</b> | 0.0050  | 156   | 12.80 | 33195271.89 | 0.00  | 0.02   |
|  | <b>Polars</b> | 0.0046  | 0     | 12.50 | 36049846.76 | 0.01  | 0.03   |
|  | <b>Modin</b>  | 0.0045  | 0     | 18.70 | 36660828.65 | 0.02  | 0.04   |
|  | <b>Dask</b>   | 20.1589 | 33641 | 12.00 | 8160.95     | 4.83  | 344.34 |
| <b>Drop Duplicate Rows</b>                                 | <b>Pandas</b> | 0.4250  | 406   | 11.10 | 378294.51   | 20.65 | 30.84  |
|  | <b>Polars</b> | 0.0435  | 141   | 15.80 | 3696134.72  | 0.01  | 0.03   |
|  | <b>Modin</b>  | 6.2243  | 9531  | 44.70 | 25833.47    | 8.46  | 8.97   |
|  | <b>Dask</b>   | 0.0668  | 344   | 11.50 | 2407074.39  | 0.09  | 0.33   |
| <b>Optimise Memory Usage by Using Efficient Data Types</b> | <b>Pandas</b> | 1.3797  | 2297  | 11.10 | 116544046   | 3.53  | 21.75  |
|  | <b>Polars</b> | 0.0452  | 94    | 8.50  | 3553856.58  | 0.02  | 0.04   |
|  | <b>Modin</b>  | 13.9932 | 20641 | 18.90 | 11490.86    | 14.53 | 15.20  |
|  | <b>Dask</b>   | 0.0840  | 78    | 19.40 | 1913473.14  | 0.34  | 0.36   |
| <b>Export Cleaned Data to MongoDB</b>                      | <b>Pandas</b> | 18.4362 | 31219 | 19.90 | 8721.67     | 14.62 | 184.51 |
|  | <b>Polars</b> | 11.4875 | 14891 | 14.40 | 13997.33    | 0.50  | 269.53 |
|  | <b>Modin</b>  | 19.8898 | 26250 | 13.20 | 8084.24     | 4.73  | 202.58 |
|  | <b>Dask</b>   | 41.0053 | 60609 | 9.80  | 3921.30     | 24.40 | 344.57 |

## **10.7 Appendix G**

The full project code and dataset are available at the following GitHub repository:  
<https://github.com/Jingyong14/HPDP02/tree/main/2425/project/p1/Group%201>