# CE/CZ4123/SC4023 BIG DATA MANAGEMENT

# SEMESTER GROUP PROJECT

| Name | Matriculation Number |
|---|---|
| HENDY | U2122559J |
| GWEE JIA XIANG | U2122287H |
| GOH PENG AIK | U2022363E |

# Table Of Contents

# Data Storage: Column-Oriented

### 1. Column Store Implementation

Column store approach was implemented through the ColumnStore class that takes in pandas dataframe obtained from input data source from CSV file located in the directory "Data/ResalePricesSingapore.csv", the program would then iterate over the dataframe columns, extracting each attribute (column) and storing it individually into separate .store files on the disk. Refer to Figure 1 in the appendix.

### 2. Data Column Design for Efficient Processing

This approach enables efficient storage and retrieval of data that allows direct and selective access to specific attributes by reading the necessary .store file instead of performing full table scans during the query process. The program also allows parallel access to different data attributes, increasing the speed during query execution.

### 3. Reading and Writing Input/Output Files

The load_data_from_disk(file_path) method allows the reading of input files from the disk and It employs caching using functools.lru_cache(maxsize=None) decorator, which stores previously read data in memory to minimize disk access, this method helps avoid redundant reads, enhancing efficiency in data retrieval. By utilizing caching, the program reduces the overhead associated with repetitive file reads, leading to faster data retrieval.

The save_to_files(directory) method allows the writing of output files to the disk. It organises and stores each column of the dataset in separate files within the specified directory. If the existing attribute files exist, it will overwrite it with the new file.

### 4. Data anomalies and exception handling

In addressing exceptions like empty qualified entries, the program demonstrated a robust approach to handling data anomalies by making use of a try-catch algorithm during the storing and loading of the data to make it easier to troubleshoot. The program also replaces found empty entries in the table and replace it with "@#NULL@#" and highlights to the user the existence of empty entries. Figure 2 in the appendix shows examples of empty entries in 3 columns "floor_area_sqm", "flat_model", and "lease_commence_date", those entries are replaced and the user will be notified.

The program also checks for the existence of incorrect data format among the data entries of attributes and highlights it to the user. Figure 2 in the appendix shows different data formats in the last column of "resale_price" which have the value of the string, this value will still be stored to ensure the same number of items is the same across all the attributes. however, the user will be notified.

The solution implemented to address empty qualified entries and incorrect data format could result in errors in query computation if those specific entries are selected for computation for the result. However, this solution ensures the correctness of column storage and enables users to easily track the issues in the data entry.

# Data Processing: Scanning Process

Firstly, the last 3 numerical digits of the student matriculation number are mapped to the town, month, and year respectively. After mapping, the month and year are appended such that they fit the data format in the month column (e.g. 2014-10). The program will proceed to scan the month column for data with the desired date value, saving the indices into a variable called matching_indices. For each value in matching_indices, the program scans the town column to filter for indices that match the desired town before saving the indices into matching_indices again.

Next, for each index in matching_indices, the program will scan the resale_price column and read each value associated with the index. For each value, the program will compare them to save the minimum price under the variable min_price as well as store the sum of all the values under total_price. Each individual value is also saved under the matching_price_data variable. This step is repeated for the floor_area_sqm column, returning min_area, total_area and matching_area_data. To obtain avg_price and avg_area, total_price and total_area are divided using the length of matching_indices which is obtained via *len(matching_indices)*. Standard Deviation is obtained using the stdev function from the statistics module in Python on matching_price_data or matching_area_data. All the results are rounded to 2 decimal points.

## Improving the Efficiency in Column Scans

As the data was already sorted based on the date in the month column and the initial scanning method was done using indices, we did not implement any enhancement with index and sorting. Instead, we implemented Zone Mapping and Shared Scans for enhancements.

### 1. Zone Mapping

After the program has stored all the values in the database, it calls for the zone mapping function which splits the month column into zones based on the year. (e.g. Zone 1: 2014; min1: 1; max1:16096). Instead of scanning the entire month column for values that fit the query, the program uses the year retrieved from the mapping to get the starting and ending indices of the Zone that the year is associated with. This reduces the number of indices that need to be scanned by the program (Refer to Figure 3 in the appendix).

### 2. Shared Scanning

The program originally scanned through either the floor_area_sqm and resale_price column three times, once for each query (Min, Avg, S.D). Using Shared Scanning, the program would instead scan through the columns and store all the required variables in a single scan, reducing the time taken (Refer to Figure 4 in the appendix).

### 3. Zone Mapping + Shared Scanning

As Zone Mapping and Shared Scanning enhance different parts of the program, the final enhancement combines both methods together to reduce the time taken for column scans by reducing I/O operations, enhancing data locality, and enabling parallel processing (Refer to Figure 5 in the appendix).

# Experiment Result

In our solution, query function and result construction are stored in "Output.py" and program flow in "main.py". The program starts by running the Python file "main.py" and it will perform column storage first. It would then perform queries using: Default scan, Shared scan, Zone Map, Zone Map + Shared scan.

Here is a screenshot of the program execution on one of the IDs using our Python file.



These are screenshots of the output files and comparisons with Excel pivot table results.

U2022363E: Town= "CLEMENTI"  Month=6  Year=2023.

| U2022363E | | | | | | |
|---|---|---|---|---|---|---|
| month | 2023-06 | | | | | |
| town | CLEMENTI | | | | | |
| | | | | | | |
| Min of resale_price | Min of floor_area_sqm | Average of resale_price | Average of floor_area_sqm | StdDev of resale_price | StdDev of floor_area_sqm | |
| 330000 | 67 | 625089.5484 | 83.41935484 | 227788.8549 | 16.54846255 | |

U2122287H: Town= "BUKIT BATOK"  Month=8  Year=2017.

| U2122287H | | | | | | |
|---|---|---|---|---|---|---|
| month | 2017-08 | | | | | |
| town | BUKIT BATOK | | | | | |
| | | | | | | |
| Min of resale_price | Min of floor_area_sqm | Average of resale_price | Average of floor_area_sqm | StdDev of resale_price | StdDev of floor_area_sqm | |
| 250000 | 60 | 384553.6271 | 91.86440678 | 133012.5176 | 25.15500747 | |

```
ScanResult_U2122287H.csv  ×

Result >  ScanResult_U2122287H.csv
    1    Year,Month,Town,Category,Value
    2    2017,08,BUKIT BATOK,Minimum Price,250000.0
    3    2017,08,BUKIT BATOK,Minimum Area,60.0
    4    2017,08,BUKIT BATOK,Average Price,384553.63
    5    2017,08,BUKIT BATOK,Average Area,91.86
    6    2017,08,BUKIT BATOK,Standard Deviation of Price,133012.52
    7    2017,08,BUKIT BATOK,Standard Deviation of Area,25.16
```

U2122559J: Town= "HOUGANG"  Month=5  Year=2019.

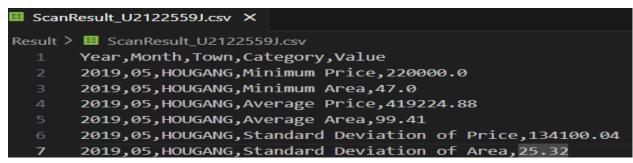| U2122559J | | | | | | |
|---|---|---|---|---|---|---|
| month | 2019-05 | | | | | |
| town | HOUGANG | | | | | |
| | | | | | | |
| Min of resale_price | Min of floor_area_sqm | Average of resale_price | Average of floor_area_sqm | StdDev of resale_price | StdDev of floor_area_sqm | |
| 220000 | 47 | 419224.8846 | 99.41346154 | 134100.0399 | 25.31763673 | |

```
ScanResult_U2122559J.csv  ×

Result >  ScanResult_U2122559J.csv
    1    Year,Month,Town,Category,Value
    2    2019,05,HOUGANG,Minimum Price,220000.0
    3    2019,05,HOUGANG,Minimum Area,47.0
    4    2019,05,HOUGANG,Average Price,419224.88
    5    2019,05,HOUGANG,Average Area,99.41
    6    2019,05,HOUGANG,Standard Deviation of Price,134100.04
    7    2019,05,HOUGANG,Standard Deviation of Area,25.32
```

We also tested the program with the possible samples, for example, U212214H covers both "2014-01" and "2024-01"

The program handles the query by combining different years that have the same last digit (refers to our matriculation ID).

U2122214H: Town= "BUKIT BATOK"  Month=1  Year=2014 and 2024.

| U2122214H | | | | | | |
|---|---|---|---|---|---|---|
| month | (Multiple Items) | | | | | |
| town | BUKIT BATOK | | | | | |
| | | | | | | |
| Min of resale_price | Min of floor_area_sqm | Average of resale_price | Average of floor_area_sqm | StdDev of resale_price | StdDev of floor_area_sqm | |
| 295000 | 59 | 466680.7458 | 92.47457627 | 132471.3485 | 23.84241073 | |

```
ScanResult_U2122214H.csv  ×

Result >  ScanResult_U2122214H.csv
    1    Year,Month,Town,Category,Value
    2    "2014,2024",01,BUKIT BATOK,Minimum Price,295000.0
    3    "2014,2024",01,BUKIT BATOK,Minimum Area,59.0
    4    "2014,2024",01,BUKIT BATOK,Average Price,466680.75
    5    "2014,2024",01,BUKIT BATOK,Average Area,92.47
    6    "2014,2024",01,BUKIT BATOK,Standard Deviation of Price,132471.35
    7    "2014,2024",01,BUKIT BATOK,Standard Deviation of Area,23.84
```

## Discussion

We verified using the Excel pivot table tool that the implementation of the column stores and results of the queries are accurate as shown in the figures above.

The results demonstrate considerable improvements in query performance through Zone Mapping and Shared Scans, with the combined approach yielding the best results. The default method took 0.1280 seconds for queries, whereas the optimized methods reduced this significantly. Shared Scans dropped it to 0.1057 seconds by consolidating multiple scans into a single operation, and Zone Mapping cut it down to 0.0994 seconds by limiting scans to relevant data segments pre-mapped by year. The integration of both techniques further reduced the time to 0.0988 seconds, showcasing a reduction of about 30% from the default scan.

These improvements can be attributed to more efficient use of computational resources and reduced I/O operations, where Shared Scans decrease the number of redundant data retrieval operations and Zone Mapping minimizes the data scanned by strategically segmenting the database. This optimized scanning approach not only enhances performance by reducing the amount of data to be processed but also speeds up access and retrieval.

## Conclusion

In conclusion, our project demonstrated data storage through column-based storage that can handle data anomalies and store each column into individual separate files. Our project also suggested multiple different efficient data scanning and access methods through the use of shared scans and zone maps to retrieve data from the information that we had stored. The implementation of these techniques significantly enhanced query performance, reducing both the time and computational resources required for data retrieval. By leveraging shared scans, we minimized the number of read operations needed, while zone mapping allowed us to focus these operations on relevant segments of data, thus avoiding unnecessary scans of the entire database. These optimizations not only improved the speed and efficiency of our data retrieval processes but also demonstrated the potential for scalable solutions in environments with larger datasets.
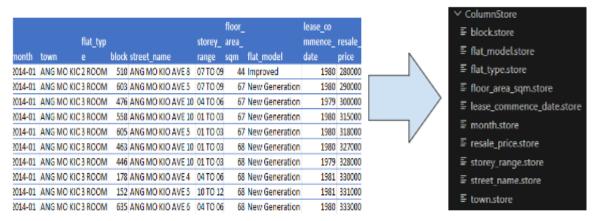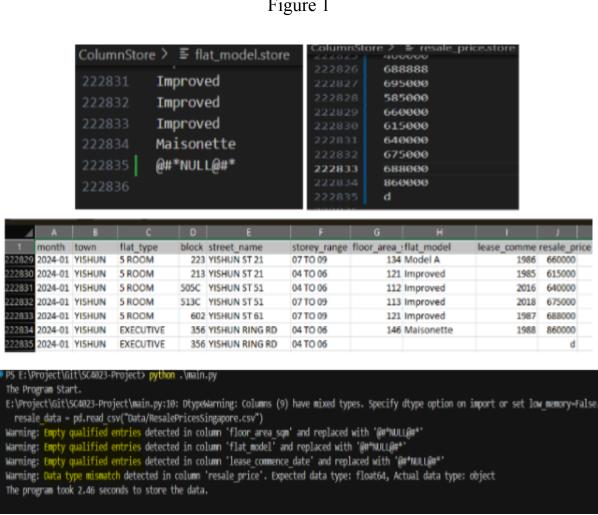
# Appendix



Figure 1



Figure 2

**Scanning Process with Zone Mapping**

Start Query → Map Last 3 Numerical Digits → Scan Month Column within Associated Zone → Scan Town Column → Return Matching Indices

Scan Price Column — For each query (min, avg, s.d)
Scan Area Column — For each query (min, avg, s.d)
Return Query Results

Figure 3

**Scanning Process with Shared Scan**

Start Query → Map Last 3 Numerical Digits → Scan Month Column → Scan Town Column → Return Matching Indices

Scan Price and Area Columns together → Return Query Results

Figure 4

**Scanning Process with Zone Mapping and Shared Scan**

Start Query → Map Last 3 Numerical Digits → Scan Month Column within Associated Zone → Scan Town Column → Return Matching Indices

Scan Price and Area Columns together → Return Query Results

Figure 5

# Contribution Form

| Name | Matriculation Number | Detailed Individual Contribution | Percentage (100% in total) |
|---|---|---|---|
| Hendy | U2122559J | <ul><li>Data Storage system implementation</li><li>Convert querying results into statistics tables and output it as a file</li><li>Report writing</li></ul> | 40% |
| Jia Xiang | U2122287H | <ul><li>Column Scan Implementation</li><li>Scan Enhancement</li><li>Report writing</li></ul> | 40% |
| Peng Aik | U2022363E | <ul><li>Scan enhancement</li><li>Results verification</li><li>Report writing</li></ul> | 20% |

**Name and Signature from all group members:**

Hendy                                                                Jia Xiang

Peng Aik