

”Optimizing Sales Data Management and Analysis System”

Shraddha Sunil Falle¹, Nikhil Ambati², Jayakrishna Pavuluri³ ¹UB ID: sfalle ²UB ID: nambati ³UB ID: jpavulur

Abstract—Addressing issues like data redundancy, sluggish query performance, and constrained analytical capabilities, this project seeks to address the company’s challenges with managing and analyzing sales data. To meet the organization’s changing data management needs, a transition is necessary from Excel to a database. The project’s significance lies in its potential to enhance sales performance, operational effectiveness, cost savings, and the ability to gain insightful information from customers and also to support the company’s use of data-driven decision-making.

I. PROBLEM STATEMENT

Effective administration and analysis of sales data, particularly inside the organization’s order processing system, is currently a task. There is a lack of normalization in the current data structure which causes problems including data redundancy, sluggish query performance, constrained analytical capabilities, and worries about data integrity. Also, insufficiently representing complex data interactions between consumers, goods, orders, and workers. According to the organization, this creates considerable challenges for data-driven decision-making, hinders productivity, and restricts the potential to gain insightful information from sales data. The project seeks to develop an Entity Relationship Diagram and normalized database schema to address these problems, providing a more complete and effective solution that will give the business simplified operations, accurate reporting, and improved strategic planning capabilities.

Why do you need a database instead of an Excel file?

Companies must switch from Excel to a database for numerous reasons. Excel is a key tool for managing and analyzing basic data, but as the amount and complexity of sales data increases, so do its limits. The advantages of a database, on the other hand, include enforcement of data integrity through constraints, improving query performance, supporting advanced analytics through structured data models, accurately simulating complex data relationships, scaling to accommodate growing data volumes and user concurrency, enhancing data security through user authentication and encryption, and offering effective version control through structured transaction logs. By providing increased data quality, additional analytics capabilities, scalability, security, and efficient version control, this shift to a relational database is in line with the changing requirements of our organization’s data management.

A. Discuss the background of the problem leading to your objectives. Why is it a significant problem?

The changing nature of modern corporate operations serves as the foundation for the problem that led to the project’s aims.

As organizations grow and data complexity rises, a number of major difficulties have emerged. First off, it has become more and more challenging to effectively manage and analyze data within the limitations of conventional technologies due to the growing complexity of sales data, which is driven by different product ranges and expanding client bases. The existing data management method introduces errors, redundancies, and inconsistencies that could ultimately result in poor decision-making. As a result, preserving data consistency and accuracy has become a vital issue. Additionally, as data volume has increased, query and data retrieval performance has substantially slowed down, affecting productivity and decision-making. Advanced analytics are also urgently needed to dive deeper insights from sales data, but the current data structure’s constraints make it difficult to meet these analytical requirements. Additionally, the delicate data relationships between consumers, goods, orders, and staff are crucial to corporate operations; however, the limitations of the current technologies in effectively depicting these complex interactions provide a substantial barrier to data-driven decision-making. The limitations of current tools to scale have also become increasingly obvious as organizations continue to expand, particularly when it comes to supporting multiple users at once and managing expanding data quantities.

B. Explain the potential of your project to contribute to your problem domain. Discuss why this contribution is crucial.

There are many benefits to this project’s successful implementation. Firstly, it enables considerable improvements in sales performance through data-driven decision-making, pricing strategies, and precise customer targeting, all of which increase revenue. Secondly, the project results in remarkable efficiency improvements, which translate into significant cost savings through streamlined operations and quick access to important data, uplifting profitability. The project also gives us the tools to gain insights into customer preferences and behavior, which results in increased customer satisfaction and brand loyalty. The project’s ability to scale up and down is vital for market growth because it opens up new customer segments and increases sales potential. Moreover, the project includes improved data security and compliance measures, significantly reducing the risks of data breaches and enhancing the defense of our reputation. The project also grants the company a competitive edge, fortifying its position in the fast-paced and fiercely competitive business environment. Sales Enhancement: The project supports data-driven decisions, optimizes pricing strategies, and targets customers, resulting

in increased sales and revenue. Cost reductions: Efficiency improvements resulting from streamlined processes and quick data access lead to impressive cost reductions, boosting overall profitability. Customer satisfaction is enhanced, and brand loyalty is fostered by personalized services that are made possible by a deep understanding of customer preferences and behavior. Market Expansion: The project's scalability is the fundamental of our market expansion strategy because it provides the opportunity to reach out to new clients and increase our sales scope. Reduced risk: Better compliance and data security procedures considerably reduce the possibility of data breaches, enhancing our standing and financial security. Competitive Advantage: The efficiency enhancements made possible by the project give our company a competitive edge, securing its place in a fast-paced, competitive market. Better Decision-Making: Because of improved data quality and accurate data representation, the company can make more informed and data-driven decisions. This leads to improved strategic planning and a competitive advantage in the market.

II. TARGET USERS

The Sales and Marketing Teams as well as Customer Support are two of the main user groups who will communicate with the database. Database administrators (DBAs) are in charge of maintaining and improving the databases, which are used by these teams for important operations.

Sales and Marketing Teams: The database's primary users are the Sales and Marketing Teams. They will rely on it to track sales performance and make strategic decisions. For example, sales representatives use the database to access customer order history, identify potential sales opportunities, and track sales goals. The marketing team uses customer data to evaluate the effectiveness of marketing campaigns, customize promotions, and improve their strategies.

Customer Support Representatives: Another essential user group is customer support representatives. They access the database to quickly retrieve customer information, view order history, and effectively respond to inquiries. This means that a customer calling the support line to ask about a recent order can receive prompt assistance because the support agent can access the database and retrieve the customer's order information.

Real-world Scenario: The database plays a crucial role in day-to-day operations in the context of a successful e-commerce business. Salespeople rely on the database to find sales opportunities, as it provides information about customer preferences, order histories, and prospects for upselling and cross-selling. Marketing teams use the database to evaluate the success of their campaigns and optimize product promotions. The customer support team relies on the database to provide accurate and timely assistance to customers. Database administrators are responsible for maintaining the system's health, enforcing security measures, ensuring data accuracy, and improving query performance. Their work on data backups, access control, and software updates keeps the database accessible and secure.

III. ENTITY-RELATIONSHIP DIAGRAM

Understanding the structural foundations of the Sales and Customer Management System requires understanding the Entity-Relationship Diagram (ERD) presented here. It offers a clear visual representation of the data model of the system and provides insightful information about the complex interactions between entities and their attributes.

A. Entities and Attributes

Order Details: The "Order Details" entity comprises attributes, each with significance, representing a crucial part of the system. Each order detail is uniquely identified by the "Order Detail ID," serving as the primary key. The "Order ID," as a foreign key, establishes a connection with the "Order" entity and identifies the order to which each detail relates. Another foreign key, "Product Code," creates a connection with the "Product" entity and identifies the particular product. Additional attributes include "Order Line Number," "Quantity Ordered," "Price Each," "Sales," and "Employee ID," with the latter acting as a foreign key to identify the employee in charge of completing the order.

Product: The "Product" entity has a primary key in the form of "Product Code," which acts as a vital identifier for specific products. It includes the attribute "Product Line" to divide products into distinct product lines. Additionally, "MSRP" (Manufacturer's Suggested Retail Price) clarifies the product's suggested retail pricing.

Customer: The "Customer" entity holds all necessary client-related data and uses the "Customer ID" as its primary key, ensuring that each customer can be uniquely identified. A complete customer profile is provided by the "Customer Name," "Phone," "Address Line1," "Address Line2," "City," "State," "Postal Code," and "Country" attributes.

B. Entities and Attributes

Line2, City, State, and Country: These attributes serve as foreign keys. The "Postal Code" attribute provides information on the postal code associated with the customer's precise location.

Territories: The "Territories" entity uses "Country" as its primary key and associates territories with particular countries.

Order Status: The "Order Status" entity has "Order ID" as the primary key. The connection between "Customer ID," represented as a foreign key, and the "Customer" entity indicates that the customer initiates the order. "Order Status" states each order's status, while the "Order Date" attribute records the specific time of placing an order.

Department: "Department ID" serves as the primary key for the "Department" entity, which is vital for organizational structuring. The "Department Name" attribute clarifies each department's nomenclature.

Employee: The main attributes of the "Employee" entity are "Employee ID" as the primary key, "Employee Name," "Salary," and "Department ID," which, as a foreign key, creates a connection with the "Department" entity to capture each employee's departmental connection.

C. Attributes and Types in Each Table

Order Details: - Order Detail ID (INT) - Order ID (INT) - Product Code (VARCHAR(255)) - Order Line Number (INT) - Ordered Quantity (INT) - Sales (DECIMAL(10,2)) - Employee ID (INT) - Price Each (DECIMAL(10,2))

Product: - MSRP (INT) - Product Line (text) - Product Code (VARCHAR(255))

Customer: - Customer ID (INT) - Customer Name (VARCHAR(255)) - Phone (VARCHAR(255)) - Address Line1 (VARCHAR(255)) - Address Line2 (VARCHAR(255)) - City (VARCHAR(255)) - State (VARCHAR(255)) - Country (VARCHAR(255)) - Postal Code (VARCHAR(255))

Territories: - Territory (VARCHAR(255)) - Country (VARCHAR(255))

Order Status: - Order ID (INT) - Customer ID (INT) - Order Date (DATE) - Order Status (VARCHAR(255))

Department: - Department Name (VARCHAR(255)) - Department ID (INT)

Employee: - Employee Name (VARCHAR(255)) - Salary (DECIMAL(10,2)) - Employee ID (INT) - Department ID (INT)

D. Primary Keys

Order Details: "Order Detail ID" attribute is the primary key for the "Order Details" table.

Product: "Product Code" attribute is the primary key for the "Product" table.

Customer: "Customer ID" attribute is the primary key for the "Customer" table.

Territories: "Country" attribute is the primary key for the "Territories" table.

Order Status: "Order ID" serves as the primary key for the "Order Status" table.

Department: "Department ID" is the primary key for the "Department" table.

Employee: "Employee ID" is the primary key for the "Employee" table.

E. Foreign Keys

Order Details: In the "Order Details" table, the columns "Order ID," "Product Code," and "Employee ID" are designated as foreign keys. "Order ID" connects to "Order Status," "Product Code" connects to "Product," and "Employee ID" connects to "Employee."

Customer: A foreign key is defined between the "Country" column in the "Customer" table and the "Country" primary key in the "Territories" table.

Order Status: The "Customer ID" column in the "Order Status" table serves as a foreign key connecting to the "Customer ID" primary key in the "Customer" table.

Employee: A foreign key connects the "Department ID" primary key in the "Department" table to the "Department ID" foreign key in the "Employee" table.

F. Relationships

Product to Order Details (One-to-Many): There is a one-to-many configuration governing the relational dynamics between the "Product" and "Order Details" entities. Different order details for a single product can be intricately linked together. This connection to the "Product Code" in the "Order Details" entity (Foreign Key) is made possible by the "Product Code" within the "Product" entity, which serves as the Primary Key.

G. Relationships (Continued)

Territories to Customer (One-to-One): A one-to-one correspondence between "Territories" and "Customer" denotes that a specific nation corresponds to a single territory. The "Country" attribute in the "Territories" entity (Primary Key) and the "Country" attribute in the "Customer" entity (Foreign Key) come together to materialize this.

Order Status to Order Details (One-to-Many): A one-to-many relationship is created by the "Order Status" entity's connection to the "Order Details" entity. The "Order ID" in the "Order Status" entity (Primary Key) directly links to the "Order ID" in the "Order Details" entity (Foreign Key), making this connection evident.

Department to Employee (One-to-One): Employees are assigned to one department per employee in a one-to-one relationship with their respective departments. The connection of employees with their organizational units is emphasized by the "Department ID" in the "Department" entity's (Primary Key) smooth coordination with the "Department ID" in the "Employee" entity's (Foreign Key).

Employee to Order Details (One-to-Many): One employee may be in charge of multiple order details because of the one-to-many relationship that exists between the "Employee" entity and the "Order Details" entity. To represent the employees' responsibility during the order fulfillment process, the "Employee ID" in the "Employee" entity (Primary Key) closely connects with the "Employee ID" in the "Order Details" entity (Foreign Key).

H. BOYCE-CODD NORMAL FORM (BCNF) DATABASE DESIGN

Boyce-Codd Normal Form (BCNF) is the framework we used to structure the database design for the Sales and Customer Management System. BCNF is a set of rules for arranging data and assuring data integrity in a database.

Selection of Primary Keys

"Order Detail ID" serves as the primary key in the "Order Details" table.

"Product Code" serves as the primary key for the "Product" table.

"Customer ID" acts as the "Customer" table's primary key.

"Country" serves as the primary key in the "Territories" table.

"Order ID" serves as the "Order Status" table's primary key.

"Department ID" serves as the primary key for the "Department" table.

"Employee ID" has been designated as the primary key in the "Employee" table.

I. Non-Prime Attributes Must Rely on Primary Key

Order Details Table: The primary key, "Order Detail ID," within the "Order Details" table is directly related to attributes like "Order ID," "Product Code," and "Employee ID."

Product Table: The "Product" table ensures that secondary keys like "Product Line" and "MSRP" are closely related to the primary key, "Product Code."

Customer Table: In the "Customer" table, there is a close relationship between "Country" and the primary key, "Customer ID."

Order Status Table: The "Order Status" table shows a definite relationship between "Customer ID" and the primary key, "Order ID."

Employee Table: Within the "Employee" table, values such as "Employee Name" and "Salary" show the relationship with the primary key, "Employee ID."

In order to comply with BCNF specifications, we made sure that every attribute had an obvious connection to the primary key, which decreased complexity and redundancy in the data structure.

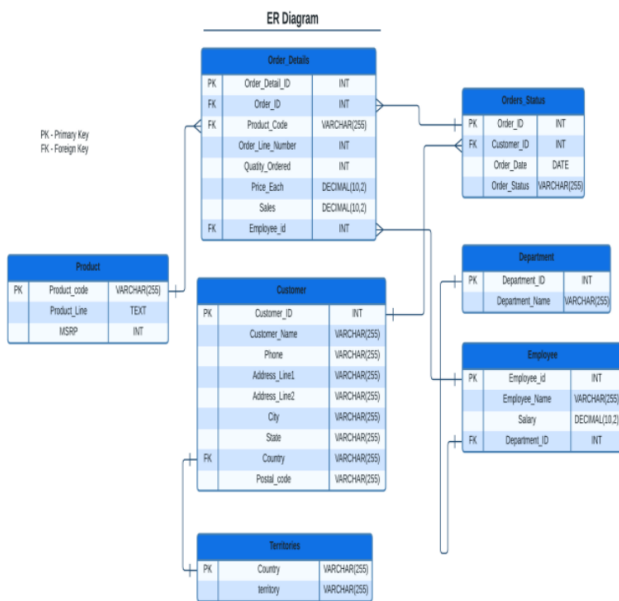


Figure 1

Fig. 1. ER Diagram

effective bulk data loading features. Precisely prepared CSV files allowed for this importation, guaranteeing precision and coherence in the data flow. The efficient procedure created a solid basis for reliable data management and analysis in addition to efficiently populating our tables.

A. Customer Table

customer_id	customer_name	phone	address_line1	address_line2	city	state	postal_code	country
1	Land of Toys Inc.	2125557818	897 Long Airport Avenue		NYC	NY	10022	USA
2	Remo Colodactes	36 47 1005	29 rue de l'Albatre		Paris		91100	France
3	Levi's Jeans	121 146 7242 7223	27 rue de la Couronne		Paris		75008	France
4	TopoDesign.com	6265957265	78934 Hilda Dr		Pasadena	CA	91000	USA
5	Corporate Gift Ideas Co.	6505551986	7734 Strong St		San Francisco	CA	94117	USA
6	Technical Street Inc.	6505555659	9489 Forti Circle		Burlington	CA	94017	USA
7	Swedish Design Imports	20 16 1555	184, chausse de Tourmel		Lille		59000	France
8	Imelia Gilly	447 2347 3215	Damen 121, PR Tui Centrum		Bruges		8 3004	Norway
9	John Street Co.	6505551957	3537 South Parkside Street		San Francisco	CA	94111	USA
10	Auto Care Parts	(71) 47 55 6005	21 rue Lantier		Paris		75016	France
11	Australian Collectors Co.	03 9320 4555	636 19 Klose Road	Level 3	Melbourne	Victoria	3004	Australia
12	Pharmacia Inc.	2125551500	2676 Kingsley Rd	Suite 101	NYC	NY	10022	USA
13	Revo Computers Inc.	2015555560	1476 Nelson Rd		Newark	NJ	08109	USA
14	Gift Report Inc.	2055552070	25593 South Bay Ln		Bridgeport	CT	07626	USA
15	La Rochelle City	40 47 8555	67 rue des Conquêtes Drapeau		Nantes		44000	France
16	Nature Planet Co.	6175555555	20223 Governor Dr		Cambridge	MA	02147	USA
17	Tony of Finland Co.	90 224 8555	Keskustatu 45		Helsinki		00100	Finland
18	Shawn Webb Imports	0716 9555	6719 Chalmers Gate 75		Stavros		41101	Norway
19	Classical Classics Inc.	2125551505	2789 Westmore St		Albany	PA	12207	USA
20	Schubert Collectables	650 9555	Deinweg 14		Saarbrücken		66100	Austria
21	Shawmire and Things Co.	461 2 9445 8555	Martin Money Building, 810 Pacific Hwy	Level 6	Cheswood	NSW	2267	Australia
22	Equilibrium.com	9605550555	1785 First Street		New Bedford	MA	01903	USA
23	UK Collectables Ltd	(171) 555 2282	Bennyway Gardens 12		Liverpool		881 627	UK
24	Euro Shipping Channel	(171) 555 84 44	C/ Alameda 86		Madrid		28014	Spain
25	John Street Business Co.	493 113 9555	Bergsgränd 5		Luleå		9 419 12	Sweden

Fig. 2. Customer table from the sales data database

B. Department Table

department_id	department_name
1	Sales1
2	Sales2
3	Sales3
4	Sales4
5	Sales5

Fig. 3. Department table from the sales data database

IV. DATABASE TABLE CREATION AND DATA IMPORTATION IN POSTGRESQL

We successfully built and organised seven different PostgreSQL tables in our database system to handle different facets of our sales data. In order to guarantee appropriate organisation and data integrity, the procedure involved defining the schema for every table. We created the tables using SQL scripts and loaded the data into the tables using PostgreSQL's

C. Employee Table

Query Query History

```
1 SELECT * FROM sales_data.employee
2 ORDER BY employee_id ASC
```

Data Output Messages Notifications

	employee_id [PK] integer	employee_name character varying (255)	salary numeric (10,2)	department_id integer
1	1	Paul Ramirez	95328.00	4
2	2	Jacqueline Hill	29717.00	5
3	3	Morgan Oneal	20602.00	3
4	4	Peter Bradley	43860.00	2
5	5	Amber Harris	27494.00	5
6	6	Michael Klein	35889.00	5
7	7	Juan Wade	46965.00	5
8	8	Samantha Smith	37158.00	5
9	9	Terri Haney	34895.00	5
10	10	Cathy Morton	13003.00	3
11	11	Misty Roberson	30332.00	3
12	12	Frederick Martinez	22649.00	3
13	13	Trevor Small	12250.00	5
14	14	Michael Farmer	23303.00	3
15	15	Ann Chapman	29766.00	5
16	16	Haley Dillon	28919.00	5
17	17	Wendy Andrews	45301.00	3
18	18	Kristina Stone	49746.00	1
19	19	Cameron Mitchell	22830.00	5
20	20	Alan Johnson	34913.00	3
21	21	Carrie Jones	16676.00	3
22	22	Rachel Walker	45642.00	5
23	23	William Aguirre	49981.00	1
24	24	Wesley Brown	22981.00	3
25	25	Bonnie Bryant	34963.00	5

Total rows: 60 of 60 Query complete 00:00:00.142

Fig. 4. Employee table from the sales data database

D. Order Details Table

Query Query History

```
1 SELECT * FROM sales_data.order_details
2 ORDER BY order_detail_id ASC
```

Data Output Messages Notifications

	order_detail_id [PK] integer	order_id integer	product_code character varying (255)	order_line_number integer	quantity_ordered integer	price_each numeric (10,2)	sales numeric (10,2)	employee_id integer
1	1	10107	S10_1678	2	30	95.70	2871.00	46
2	2	10121	S10_1678	5	34	81.35	2765.90	42
3	3	10134	S10_1678	2	41	94.74	3884.34	55
4	4	10145	S10_1678	6	45	83.26	3746.70	28
5	5	10159	S10_1678	14	49	100.00	4900.00	48
6	6	10168	S10_1678	1	36	96.66	3479.76	26
7	7	10180	S10_1678	9	29	86.13	2497.77	19
8	8	10188	S10_1678	1	48	100.00	4800.00	25
9	9	10201	S10_1678	2	22	98.57	2168.54	34
10	10	10211	S10_1678	14	41	100.00	4100.00	22
11	11	10223	S10_1678	1	37	100.00	3700.00	39
12	12	10237	S10_1678	7	23	100.00	2300.00	53
13	13	10251	S10_1678	2	28	100.00	2800.00	41
14	14	10263	S10_1678	2	34	100.00	3400.00	39
15	15	10275	S10_1678	1	45	92.83	4177.35	36
16	16	10285	S10_1678	6	36	100.00	3600.00	31
17	17	10299	S10_1678	9	23	100.00	2300.00	31
18	18	10309	S10_1678	5	41	100.00	4100.00	13
19	19	10318	S10_1678	1	46	94.74	4358.04	38
20	20	10329	S10_1678	1	42	100.00	4200.00	55
21	21	10341	S10_1678	9	41	100.00	4100.00	42
22	22	10361	S10_1678	13	20	72.55	1451.00	33
23	23	10375	S10_1678	12	21	34.91	733.11	52
24	24	10388	S10_1678	4	42	76.36	3207.12	21
25	25	10403	S10_1678	7	24	100.00	2400.00	29

Total rows: 1000 of 1000 Query complete 00:00:00.141

Fig. 5. Order Details table from the sales data database

E. Product Table

F. Territories Table

G. Order Status Table

Query Query History

```
1 SELECT * FROM sales_data.order_details
2 ORDER BY order_detail_id ASC
```

Data Output Messages Notifications

	order_detail_id [PK] integer	order_id integer	product_code character varying (255)	order_line_number integer	quantity_ordered integer	price_each numeric (10,2)	sales numeric (10,2)	employee_id integer
1	1	10107	S10_1678	2	30	95.70	2871.00	46
2	2	10121	S10_1678	5	34	81.35	2765.90	42
3	3	10134	S10_1678	2	41	94.74	3884.34	55
4	4	10145	S10_1678	6	45	83.26	3746.70	28
5	5	10159	S10_1678	14	49	100.00	4900.00	48
6	6	10168	S10_1678	1	36	96.66	3479.76	26
7	7	10180	S10_1678	9	29	86.13	2497.77	19
8	8	10188	S10_1678	1	48	100.00	4800.00	25
9	9	10201	S10_1678	2	22	98.57	2168.54	34
10	10	10211	S10_1678	14	41	100.00	4100.00	22
11	11	10223	S10_1678	1	37	100.00	3700.00	39
12	12	10237	S10_1678	7	23	100.00	2300.00	53
13	13	10251	S10_1678	2	28	100.00	2800.00	41
14	14	10263	S10_1678	2	34	100.00	3400.00	39
15	15	10275	S10_1678	1	45	92.83	4177.35	36
16	16	10285	S10_1678	6	36	100.00	3600.00	31
17	17	10299	S10_1678	9	23	100.00	2300.00	31
18	18	10309	S10_1678	5	41	100.00	4100.00	13
19	19	10318	S10_1678	1	46	94.74	4358.04	38
20	20	10329	S10_1678	1	42	100.00	4200.00	55
21	21	10341	S10_1678	9	41	100.00	4100.00	42
22	22	10361	S10_1678	13	20	72.55	1451.00	33
23	23	10375	S10_1678	12	21	34.91	733.11	52
24	24	10388	S10_1678	4	42	76.36	3207.12	21
25	25	10403	S10_1678	7	24	100.00	2400.00	29

Total rows: 1000 of 1000 Query complete 00:00:00.141

Fig. 6. Product table from the sales data database

Query Query History

```
1 SELECT * FROM sales_data.territories
2 ORDER BY country ASC
```

Data Output Messages Notifications

	territory character varying (255)	country [PK] character varying (255)
1	APAC	Australia
2	EMEA	Austria
3	EMEA	Belgium
4	NA	Canada
5	EMEA	Denmark
6	EMEA	Finland
7	EMEA	France
8	EMEA	Germany
9	EMEA	Ireland
10	EMEA	Italy
11	Japan	Japan
12	EMEA	Norway
13	Japan	Philippines
14	APAC	Singapore
15	EMEA	Spain
16	EMEA	Sweden
17	EMEA	Switzerland
18	EMEA	UK
19	NA	USA

Total rows: 19 of 19 Query complete 00:00:00.188

Fig. 7. Territories table from the sales data database

Query

Query History

1

SELECT * FROM sales_data.orders_status

2

ORDER BY order_id ASC

Data Output

Messages

Notifications

order_id

[PK] integer

order_date

date

order_status

character varying (255)

customer_id

integer

1

10100

2003-01-06

Shipped

41

2

10101

2003-01-09

Shipped

73

3

10102

2003-01-10

Shipped

12

4

10103

2003-01-29

Shipped

18

5

10104

2003-01-31

Shipped

24

6

10105

2003-02-11

Shipped

49

7

10106

2003-02-17

Shipped

87

8

10107

2003-02-24

Shipped

1

9

10108

2003-03-03

Shipped

67

10

10109

2003-03-10

Shipped

46

11

10110

2003-03-18

Shipped

78

12

10111

2003-03-25

Shipped

9

13

10112

2003-03-24

Shipped

25

14

10113

2003-03-26

Shipped

40

15

10114

2003-04-01

Shipped

63

16

10115

2003-04-04

Shipped

28

17

10116

2003-04-11

Shipped

91

18

10117

2003-04-16

Shipped

27

19

10118

2003-04-21

Shipped

54

20

10119

2003-04-28

Shipped

20

21

10120

2003-04-29

Shipped

11

22

10121

2003-05-07

Shipped

2

23

10122

2003-05-08

Shipped

68

24

10123

2003-05-20

Shipped

47

25

10124

2003-05-21

Shipped

85

Total rows: 308 of 308

Query complete 00:00:00.182

Fig. 8. Order Status table from the sales data database

V. ADDRESSING CHALLENGES AND DRIVING IMPROVEMENTS IN LARGE DATASET MANAGEMENT

As part of our continuous commitment to upholding a reliable and effective database system, we started an effort to assess and improve the performance of our database operations, especially when managing bigger datasets.

A. Performance Challenges with Large Datasets

As our database expanded, we noticed a small but noticeable lag in query response times, primarily due to an increase in the amount of sales data. These delays, only a few seconds, could negatively affect real-time data analysis capabilities and lead to less-than-ideal user experiences.

The issues were particularly evident in the following scenarios:

- **Employee and Customer Name Lookups:** The auto-complete feature and search functionality were not as responsive as required.
- **Order Details Access:** System lags when attempting to retrieve detailed order information affected our order management and reporting procedures.

B. Indexing Strategies Employed

We developed a thorough indexing strategy to address these concerns and improve the efficiency and speed of our database queries.

1) Index Implementation: 1. Employee and Customer Name Indexes:

• SQL Implementation:

```
CREATE INDEX idx_employee_name ON
sales_data.employee (Employee_Name);
CREATE INDEX idx_customer_name ON
sales_data.customer (Customer_Name);
```

- **Outcome:** Significantly reduced execution time for name-based searches.

ANALYSIS OF PERFORMANCE IMPROVEMENT

Our database operations were significantly improved by the addition of these indexes, as the following details demonstrate:

Execution Time of Queries

Large dataset queries executed much more slowly before indexes were added. These queries now run in less than a second thanks to the updated indexing. Faster data retrieval and processing are now possible thanks to this remarkable decrease in query execution time, which has also greatly increased our database operations' throughput and efficiency.

System Scalability

There has been a significant improvement in the database system's capacity to manage growing data volumes. Longer query latency was frequently the outcome of increasing data volume prior to indexing. As of right now, the system shows strong scalability, effectively handling bigger datasets without sacrificing query response time. With this enhancement, our database will be able to function well even as data demands increase.

User Satisfaction

A noticeable improvement in the user experience is a direct benefit of faster query responses. Data analysts and business users, among other internal stakeholders, have reported easier and more productive interactions with the database system. This is particularly significant in situations where prompt data analysis and access are essential to business operations.

VI. OVERVIEW OF DATABASE SQL QUERY TESTING

We give a summary of the comprehensive testing we did using different SQL queries on our database system in this section. The system's performance and capabilities were assessed during this crucial testing phase. The queries were created to cover insertion, deletion, and updating, among other operations, giving a thorough evaluation of the database's functionality. To further show how well the system handled challenging data retrieval scenarios, a range of select queries were run. Several SQL techniques, including joins, order by, group by, and subqueries, were used in these select queries. The results of these searches were carefully documented, and screen grabs were made to visually record the execution outcomes, demonstrating the efficiency and stability of our database system.

A. Select Queries

The select queries we used for testing are the main topic of this subsection. These queries are essential for showcasing how effectively the database can retrieve and handle data. They were thoughtfully created to demonstrate how the database handled a range of intricate querying strategies, such as joins, order by, group by, and subqueries. The outcomes of these select queries demonstrated the database system's ability to handle and analyse massive amounts of data, offering insightful information about its scalability and performance.

Query Query History

```
1  SELECT
2    os.Order_ID,
3    os.Order_Date,
4    os.Order_Status,
5    c.Customer_Name,
6    e.Employee_Name,
7    od.Product_Code,
8    od.Quantity_Ordered,
9    od.Sales
10 FROM sales_data.orders_status os
11 JOIN sales_data.customer c ON os.Customer_ID = c.Customer_ID
12 JOIN sales_data.order_details od ON os.Order_ID = od.Order_ID
13 JOIN sales_data.employee e ON od.Employee_ID = e.Employee_ID
14 LIMIT 10;
15
16
```

Data Output Messages Notifications

	order_id integer	order_date date	order_status character varying (255)	customer_name character varying (255)	employee_name character varying (255)	product_code character varying (255)	quantity_ordered integer	sales numeric (10,2)
1	10107	2003-02-24	Shipped	Land of Toys Inc.	Michelle Gilbert	S10_1678	30	2871.00
2	10121	2003-05-07	Shipped	Reims Collectables	Tammy Bird	S10_1678	34	2765.90
3	10134	2003-07-01	Shipped	Lyon Souvenirs	Paula Hill	S10_1678	41	3884.34
4	10145	2003-08-25	Shipped	Toys4Grownups.com	Brett Ward	S10_1678	45	3746.70
5	10159	2003-10-10	Shipped	Corporate Gift Ideas Co.	Jennifer Mason	S10_1678	49	4900.00
6	10168	2003-10-28	Shipped	Technics Stores Inc.	Olivia Hunter	S10_1678	36	3479.76
7	10180	2003-11-11	Shipped	Daedalus Designs Imports	Cameron Mitchell	S10_1678	29	2497.77
8	10188	2003-11-18	Shipped	Herku Gifts	Bonnie Bryant	S10_1678	48	4800.00
9	10201	2003-12-01	Shipped	Mini Wheels Co.	Elizabeth Miller	S10_1678	22	2168.54
10	10211	2004-01-15	Shipped	Auto Canal Pett	Rachel Walker	S10_1678	41	4100.00

Fig. 9. Select Query 1

Query Query History

```

1 SELECT
2   e.Employee_id,
3   e.Employee_Name,
4   SUM(od.Sales) AS TotalSales
5 FROM sales_data.employee e
6 LEFT JOIN sales_data.order_details od ON e.Employee_id = od.Employee_id
7 GROUP BY e.Employee_id, e.Employee_Name
8 LIMIT 15;
9
10
11
12

```

Data Output Messages Notifications

employee_id	employee_name	totalsales
[PK] integer	character varying (255)	numeric
1	54 Emily Johnson	74790.48
2	29 Mike Norris	46523.96
3	4 Peter Bradley	38543.78
4	34 Elizabeth Miller	50341.50
5	51 Kimberly White	63811.66
6	52 Jennifer Griffin	64090.45
7	10 Cathy Morton	57074.28
8	35 Pamela Bowman	51572.49
9	45 Valerie Davis	54351.35
10	6 Michael Klein	61568.11
11	39 Nathaniel Lopez	49394.95
12	36 Alicia Mccarty	54815.44
13	31 John Conner	55110.10
14	50 Amanda Mccoy	64886.83
15	60 Stephen Lewis	42069.13

Fig. 10. Select Query 2

Query Query History

```

1 SELECT
2   e.Employee_id,
3   e.Employee_Name,
4   AVG(od.Sales) AS AvgSales
5 FROM sales_data.employee e
6 LEFT JOIN sales_data.order_details od ON e.Employee_id = od.Employee_id
7 GROUP BY e.Employee_id, e.Employee_Name
8 LIMIT 15;
9
10
11
12
13
14
15

```

Data Output Messages Notifications

employee_id	employee_name	avgsales
[PK] integer	character varying (255)	numeric
1	54 Emily Johnson	3739.524000000000000000
2	29 Mike Norris	3323.140000000000000000
3	4 Peter Bradley	3211.981666666666666667
4	34 Elizabeth Miller	2796.750000000000000000
5	51 Kimberly White	3545.092222222222222222
6	52 Jennifer Griffin	3373.1815789473684211
7	10 Cathy Morton	3170.793333333333333333
8	35 Pamela Bowman	3438.166000000000000000
9	45 Valerie Davis	3197.1382352941176471
10	6 Michael Klein	2798.5504545454545455
11	39 Nathaniel Lopez	3087.184375000000000000
12	36 Alicia Mccarty	3045.302222222222222222
13	31 John Conner	2624.2904761904761905
14	50 Amanda Mccoy	3604.823888888888888889
15	60 Stephen Lewis	3236.0869230769230769

Fig. 11. Select Query 3

Query Query History

```
1 WITH OrderDetailsWithTotalSales AS (
2   SELECT
3     od.Order_ID,
4     SUM(od.Sales) AS TotalOrderSales
5 FROM sales_data.order_details od
6 GROUP BY od.Order_ID
7 )
8
9 SELECT
10  os.Order_ID,
11  os.Order_Date,
12  os.Order_Status,
13  c.Customer_Name,
14  e.Employee_Name,
15  od.Product_Code,
16  od.Quantity_Ordered,
17  od.Sales,
18  odt.TotalOrderSales
19 FROM sales_data.orders_status os
20 JOIN sales_data.customer c ON os.Customer_ID = c.Customer_ID
21 JOIN sales_data.order_details od ON os.Order_ID = od.Order_ID
22 JOIN sales_data.employee e ON od.Employee_id = e.Employee_id
23 JOIN OrderDetailsWithTotalSales odt ON os.Order_ID = odt.Order_ID
24 LIMIT 5;
25
```

Data Output Messages Notifications

order_id	order_date	order_status	customer_name	employee_name	product_code	quantity_ordered	sales	totalordersales	
integer	date	character varying (255)	character varying (255)	character varying (255)	character varying (255)	integer	numeric (10,2)	numeric	
1	10107	2003-02-24	Shipped	Land of Toys Inc.	Michelle Gilbert	S10_1678	30	2871.00	13622.72
2	10121	2003-05-07	Shipped	Reims Collectables	Tammy Bird	S10_1678	34	2765.90	7765.90
3	10134	2003-07-01	Shipped	Lyon Souvenirs	Paula Hill	S10_1678	41	3884.34	13537.74
4	10145	2003-08-25	Shipped	Toys4Grownups.com	Brett Ward	S10_1678	45	3746.70	19696.40
5	10159	2003-10-10	Shipped	Corporate Gift Ideas Co.	Jennifer Mason	S10_1678	49	4900.00	22624.24

Fig. 12. Select Query 4

Query	Query History
1	WITH OrderDetailsWithTotalSales AS (
2	SELECT
3	od.Order_ID,
4	SUM(od.Sales) AS TotalOrdersSales
5	FROM sales_data.order_details od
6	GROUP BY od.Order_ID
7)
8	
9	SELECT
10	os.Order_ID,
11	os.Order_Date,
12	os.Order_Status,
13	c.Customer_Name,
14	e.Employee_Name,
15	od.Product_Code,
16	od.Quantity_Ordered,
17	od.Sales,
18	od.TotalOrdersSales
19	FROM sales_data.orders_status os
20	JOIN sales_data.customer c ON os.Customer_ID = c.Customer_ID
21	JOIN sales_data.order_details od ON os.Order_ID = od.Order_ID
22	JOIN sales_data.employee e ON od.Employee_ID = e.Employee_ID
23	JOIN OrderDetailsWithTotalSales odt ON os.Order_ID = odt.Order_ID
24	LIMIT 5;
25	

order_id	order_date	order_status	customer_name	employee_name	product_code	quantity_ordered	sales	totalordersales
integer	date	character varying (255)	character varying (255)	character varying (255)	character varying (255)	integer	numeric (10,2)	numeric
1	10107	2003-02-24	Shipped	Land of Toys Inc.	Michelle Gilbert	\$10,1678	30	2871.00
2	10121	2003-05-07	Shipped	Remix Collectables	Tammy Bird	\$10,1678	34	2765.90
3	10134	2003-07-01	Shipped	Lynn Souvenirs	Patsy Hill	\$10,1678	41	3864.34
4	10145	2003-08-25	Shipped	Top4Downfalls.com	Brett Ward	\$10,1678	45	3766.70
5	10159	2003-10-10	Shipped	Corporate Gift Ideas Co.	Jennifer Mason	\$10,1678	49	4900.00

Fig. 13. Select Query 5

Query	Query History
1	WITH EmployeePerformance AS (
2	SELECT
3	e.Employee_id,
4	e.Employee_Name,
5	COUNT(os.Order_ID) AS TotalOrders,
6	SUM(od.Sales) AS TotalSales
7	FROM sales_data.employee e
8	LEFT JOIN sales_data.order_details od ON e.Employee_id = od.Employee_id
9	LEFT JOIN sales_data.orders_status os ON od.Order_ID = os.Order_ID
10	GROUP BY e.Employee_id, e.Employee_Name
11)
12	
13	SELECT
14	ep.Employee_id,
15	ep.Employee_Name,
16	ep.TotalOrders,
17	ep.TotalSales,
18	ROW_NUMBER() OVER (ORDER BY ep.TotalSales DESC) AS SalesRank,
19	RANK() OVER (ORDER BY ep.TotalOrders DESC) AS OrdersRank
20	FROM EmployeePerformance ep
21	LIMIT 5;

employee_id	employee_name	totalorders	totalsales	salesrank	ordersrank
[PK] integer	character varying (255)	bigint	numeric	bigint	bigint
1	54 Emily Johnson	20	74790.48	1	9
2	49 Regina Alvarado	26	73744.23	2	1
3	26 Olivia Hunter	21	71859.43	3	6
4	40 Martin Navarro	22	70227.98	4	3
5	28 Brett Ward	20	68770.73	5	9

Fig. 14. Select Query 6

B. Insert Queries

We will examine the 'Insert Queries' used in our database testing process in this subsection. When assessing the database's ability to smoothly add new data entries, insert queries are essential. These queries were created expressly to evaluate how well the database performed data insertion tasks, from straightforward inserts of a single record to more intricate operations involving several records. When assessing our database system, the effectiveness of these insert operations was crucial, especially in situations where frequent and dynamic data additions were needed. The process of data insertion involves more than just adding new entries; it also involves making sure that these additions don't interfere with concurrent operations or transactional integrity. Thus, our assessment also includes determining how these insertions

Query	Query History
1	SELECT
2	c.Customer_Name,
3	COUNT(*) AS TotalOrders,
4	SUM(od.Quantity_Ordered) AS TotalQuantityOrdered,
5	AVG(od.Price_Each) AS AvgPrice
6	FROM
7	sales_data.customer c
8	JOIN sales_data.orders_status os ON c.Customer_ID = os.Customer_ID
9	JOIN sales_data.order_details od ON os.Order_ID = od.Order_ID
10	WHERE
11	c.Customer_Name LIKE '%Inc%'
12	GROUP BY
13	c.Customer_Name
14	HAVING
15	COUNT(*) > 5
16	ORDER BY
17	TotalOrders DESC, AvgPrice DESC;
18	

customer_name	totalorders	totalquantityordered	avgprice
character varying (255)	bigint	bigint	numeric
1 Land of Toys Inc.	20	613	93.015000000000000000
2 Muscle Machine Inc	17	627	98.320000000000000000
3 Diecast Classics Inc.	15	565	90.164666666666666667
4 Technics Stores Inc.	13	396	87.7276923076923077
5 Collectables For Less Inc.	12	408	90.983333333333333333
6 Gift Depot Inc.	11	416	94.9927272727272727
7 Classic Gift Ideas, Inc	10	303	93.751000000000000000
8 Motor Mint Distributors Inc.	10	325	82.713000000000000000
9 Oulu Toy Supplies, Inc.	9	311	86.7522222222222222
10 Vitachrome Inc.	9	307	77.814444444444444444
11 Super Scale Inc.	8	287	100.000000000000000000
12 Classic Legends Inc.	7	257	88.2871428571428571
13 Tekni Collectables Inc.	6	232	89.535000000000000000

Total rows: 13 of 13	Query complete 00:00:00.113
----------------------	-----------------------------

Fig. 15. Select Query 7

affect the latency of concurrent read and write operations as well as the overall throughput of the system. Particularly when scaling the system to meet our enterprise's ever-growing data demands, these insights are priceless.

C. Update Queries

We'll see the update queries we use in our database testing procedure in this subsection. Update queries are essential for determining how well the database can alter and handle current data. These queries were created to test different aspects of data modification, from straightforward updates of a single record to more intricate modifications based on conditions. When assessing the overall efficacy of our database system, the speed and precision of these update operations were crucial, especially in situations where frequent and dynamic data modifications were necessary.

Query Query History			
1	INSERT INTO sales_data.product (Product_Code, Product_Line, MSRP)		
2	VALUES ('PROD0003', 'Appliances', 1300);		
3			
4	INSERT INTO sales_data.product (Product_Code, Product_Line, MSRP)		
5	VALUES ('PROD0004', 'Laptop', 1300);		
6			
7	SELECT * FROM sales_data.product		
Data Output Messages Notifications			
	product_code [PK] character varying (255)	product_line text	msrp integer
91	S32_4485	Motorcycles	102
92	S50_1341	Vintage Cars	43
93	S50_1392	Trucks and Buses	115
94	S50_1514	Trains	58
95	S50_4713	Motorcycles	81
96	S700_1138	Ships	66
97	S700_1691	Planes	91
98	S700_1938	Ships	86
99	S700_2047	Ships	90
100	S700_2466	Planes	99
101	S700_2610	Ships	72
102	S700_2824	Classic Cars	101
103	S700_2834	Planes	118
104	S700_3167	Planes	80
105	S700_3505	Ships	100
106	S700_3962	Ships	99
107	S700_4002	Planes	74
108	S72_1253	Planes	49
109	S72_3212	Ships	54
110	PROD003	Appliances	1300
111	PROD004	Laptop	1300

Fig. 16. Insert Query 1

Query Query History			
1	SELECT * FROM sales_data.product;		
2			
3			
4	UPDATE sales_data.product		
5	SET MSRP = CASE		
6	WHEN Product_Code = 'PROD0003' THEN 1400		
7	WHEN Product_Code = 'PROD0004' THEN 1500		
8	ELSE MSRP		
9	END;		
10			
11			
12	SELECT * FROM sales_data.product;		
13			
Data Output Messages Notifications			
	product_code [PK] character varying (255)	product_line text	msrp integer
91	S32_4485	Motorcycles	102
92	S50_1341	Vintage Cars	43
93	S50_1392	Trucks and Buses	115
94	S50_1514	Trains	58
95	S50_4713	Motorcycles	81
96	S700_1138	Ships	66
97	S700_1691	Planes	91
98	S700_1938	Ships	86
99	S700_2047	Ships	90
100	S700_2466	Planes	99
101	S700_2610	Ships	72
102	S700_2824	Classic Cars	101
103	S700_2834	Planes	118
104	S700_3167	Planes	80
105	S700_3505	Ships	100
106	S700_3962	Ships	99
107	S700_4002	Planes	74
108	S72_1253	Planes	49
109	S72_3212	Ships	54
110	PROD003	Appliances	1400
111	PROD004	Laptop	1500

Fig. 18. Update Query 1

Query Query History			
1	INSERT INTO sales_data.territories (territory, country)		
2	VALUES ('New Territory', 'XYZ');		
3			
4	SELECT * FROM sales_data.territories		
Data Output Messages Notifications			
	territory character varying (255)	country [PK] character varying (255)	
1	Updated Territory	XYZ	

Fig. 17. Insert Query 2

Query Query History			
1	UPDATE sales_data.territories		
2	SET territory = 'Updated Territory'		
3	WHERE Country = 'XYZ';		
4	SELECT * FROM sales_data.territories		
5	WHERE Country = 'XYZ';		
Data Output Messages Notifications			
	territory character varying (255)	country [PK] character varying (255)	
1	Updated Territory	XYZ	

Fig. 19. Update Query 2

D. Delete Queries

We will examine the 'Delete Queries' used in our database testing process in this subsection. When evaluating the database's capacity to efficiently delete data records, delete queries are crucial. These queries were painstakingly designed to assess how well the database performed when it came to data deletion tasks, which included both the deletion of individual records and more intricate operations like conditional data removal. When assessing the overall performance of our database system, the effectiveness and precision of these delete operations were crucial, especially in situations where the smooth management of data removal was necessary.

Query	Query History
1	DELETE FROM sales_data.product
2	WHERE Product_Code IN ('PROD003', 'PROD004');
3	
4	SELECT * FROM sales_data.product;

Data Output	Messages	Notifications
territory character varying (255)	country [PK] character varying (255)	
1	APAC	Australia
2	EMEA	Austria
3	EMEA	Belgium
4	NA	Canada
5	EMEA	Denmark
6	EMEA	Finland
7	EMEA	France
8	EMEA	Germany
9	EMEA	Ireland
10	EMEA	Italy
11	Japan	Japan
12	EMEA	Norway
13	Japan	Philippines
14	APAC	Singapore
15	EMEA	Spain
16	EMEA	Sweden
17	EMEA	Switzerland
18	EMEA	UK
19	NA	USA

Fig. 20. Delete Query 1

Query	Query History
1	DELETE FROM sales_data.territories
2	WHERE Country = 'XYZ';
3	SELECT * FROM sales_data.territories;

Data Output	Messages	Notifications
territory character varying (255)	country [PK] character varying (255)	
1	APAC	Australia
2	EMEA	Austria
3	EMEA	Belgium
4	NA	Canada
5	EMEA	Denmark
6	EMEA	Finland
7	EMEA	France
8	EMEA	Germany
9	EMEA	Ireland
10	EMEA	Italy
11	Japan	Japan
12	EMEA	Norway
13	Japan	Philippines
14	APAC	Singapore
15	EMEA	Spain
16	EMEA	Sweden
17	EMEA	Switzerland
18	EMEA	UK
19	NA	USA

Fig. 21. Delete Query 2

VII. PROBLEMATIC QUERIES

We were able to identify problematic queries that took a long time to execute based on the queries that were executed. We will go over these questions in-depth and look into possible optimisations in this section.

A. Problematic Query 1

Description and Analysis of the First Problematic Query

Query

Query History

```
1 SELECT
2   c.Customer_Name,
3   COUNT(*) AS TotalOrders,
4   SUM(od.Quantity_Ordered) AS TotalQuantityOrdered,
5   AVG(od.Price_Each) AS AvgPrice
6 FROM
7   sales_data.customer c
8   JOIN sales_data.orders_status os ON c.Customer_ID = os.Customer_ID
9   JOIN sales_data.order_details od ON os.Order_ID = od.Order_ID
10 WHERE
11   c.Customer_Name LIKE '%Inc%'
12 GROUP BY
13   c.Customer_Name
14 HAVING
15   COUNT(*) > 5
16 ORDER BY
17   TotalOrders DESC, AvgPrice DESC;
18
```

Data Output

Messages

Explain

×

Notifications

	customer_name character varying (255)	totalorders bigint	totalquantityordered bigint	avgprice numeric
1	Land of Toys Inc.	20	613	93.0150000000000000
2	Muscle Machine Inc	17	627	98.3200000000000000
3	Diecast Classics Inc.	15	565	90.1646666666666667
4	Technics Stores Inc.	13	396	87.7276923076923077
5	Collectables For Less Inc.	12	408	90.0983333333333333
6	Gift Depot Inc.	11	416	94.9927272727272727
7	Classic Gift Ideas, Inc	10	303	93.7510000000000000
8	Motor Mint Distributors Inc.	10	325	82.7130000000000000
9	Oulu Toy Supplies, Inc.	9	311	86.7522222222222222
10	Vitachrome Inc.	9	307	77.8144444444444444
11	Super Scale Inc.	8	287	100.0000000000000000
12	Classic Legends Inc.	7	257	88.2871428571428571
13	Tekni Collectables Inc.	6	232	89.5350000000000000

Total rows: 13 of 13

Query complete 00:00:00.113

Fig. 22. Problematic Query 1

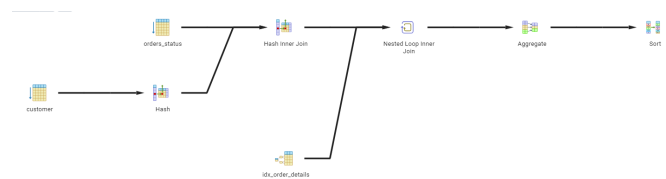


Fig. 23. Explain from PG Admin

Our database testing's first problematic query identified a sizable performance bottleneck. Long execution time for this specific query raised doubts about its effectiveness. We looked into possible optimisations to improve the query's performance as well as the causes of its prolonged execution time.

Optimized Solution

We came up with an optimised solution to deal with Problematic Query 1's performance problems. This approach resulted in a more effective execution process by rewriting the query to include filtering before joining. We were able

to improve overall database performance and drastically cut down on query execution time by putting this optimisation into practice.

Query	Query History
1	WITH filtered_customer_ids AS (
2	SELECT Customer_ID
3	FROM sales_data.customer
4	WHERE Customer_Name LIKE '%Inc%'
5)
6	SELECT
7	c.Customer_Name,
8	COUNT(*) AS TotalOrders,
9	SUM(od.Quantity_Ordered) AS TotalQuantityOrdered,
10	AVG(od.Price_Each) AS AvgPrice
11	FROM
12	filtered_customer_ids fc
13	JOIN sales_data.orders_status os ON fc.Customer_ID = os.Customer_ID
14	JOIN sales_data.customer c ON fc.Customer_ID = c.Customer_ID
15	JOIN sales_data.order_details od ON os.Order_ID = od.Order_ID
16	GROUP BY
17	c.Customer_Name
18	HAVING
19	COUNT(*) > 5
20	ORDER BY
21	TotalOrders DESC, AvgPrice DESC;

customer_name	totalorders	totalquantityordered	avgprice
character varying (255)	bigint	bigint	numeric
1 Land of Toys Inc.	20	613	93.015000000000000000
2 Muscle Machine Inc	17	627	98.320000000000000000
3 Diecast Classics Inc.	15	565	90.164666666666666667
4 Technics Stores Inc.	13	396	87.7276923076923077
5 Collectables For Less Inc.	12	408	90.098333333333333333
6 Gift Depot Inc.	11	416	94.9927272727272727
7 Classic Gift Ideas, Inc	10	303	93.751000000000000000
8 Motor Mint Distributors Inc.	10	325	82.713000000000000000
9 Oulu Toy Supplies, Inc.	9	311	86.7522222222222222
10 Vitachrome Inc.	9	307	77.8144444444444444
11 Super Scale Inc.	8	287	100.0000000000000000
12 Classic Legends Inc.	7	257	88.2871428571428571

Total rows: 13 of 13 Query complete 00:00:00.053

Fig. 24. Optimized Solution for Problematic Query 1

B. Problematic Query 2

Description and Analysis of the Second Problematic Query

1	SELECT
2	od.order_detail_id,
3	od.order_id,
4	os.order_status
5	FROM
6	sales_data.order_details od
7	FULL OUTER JOIN
8	sales_data.orders_status os ON od.order_id = os.order_id;
9	

order_detail_id	order_id	order_status
integer	integer	character varying (255)
1	10107	Shipped
2	10121	Shipped
3	10134	Shipped
4	10145	Shipped
5	10159	Shipped
6	10168	Shipped
7	10180	Shipped
8	10188	Shipped
9	10201	Shipped
10	10211	Shipped
11	10223	Shipped
12	10237	Shipped

Total rows: 1000 of 1042 Query complete 00:00:00.051

Fig. 25. Problematic Query 2

Again long execution time for this specific query raised doubts about its effectiveness. We looked into possible optimisations to improve the query's performance as well as the causes of its prolonged execution time.

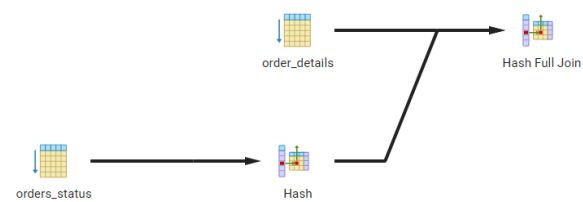


Fig. 26. Explain from PG Admin

Optimized Solution

Our approach to solving the performance problems with Problematic Query 2 was optimised. Applying indexing strategically was required for this. We achieved a significant reduction in query execution time and overall database performance improvement by putting these optimisations into practice.

1	SELECT
2	od.order_detail_id,
3	od.order_id,
4	os.order_status
5	FROM
6	sales_data.order_details od
7	full outer JOIN
8	sales_data.orders_status os ON od.order_id = os.order_id;
9	

order_detail_id	order_id	order_status
integer	integer	character varying (255)
1	10107	Shipped
2	10121	Shipped
3	10134	Shipped
4	10145	Shipped
5	10159	Shipped
6	10168	Shipped
7	10180	Shipped
8	10188	Shipped
9	10201	Shipped
10	10211	Shipped
11	10223	Shipped
12	10237	Shipped

Total rows: 1000 of 1042 Query complete 00:00:00.051

Fig. 27. Optimized Solution for Problematic Query 2

C. Problematic Query 3

Description and Analysis of the Third Problematic Query

1	SELECT
2	os.order_id,
3	os.order_status,
4	os.order_detail_id,
5	c.customer_name,
6	e.employee_name,
7	od.product_code,
8	od.quantity_ordered,
9	od.price_each,
10	FROM
11	sales_data.orders_status os
12	JOIN sales_data.customer c ON os.customer_id = c.customer_id
13	JOIN sales_data.order_details od ON os.order_id = od.order_id
14	JOIN sales_data.employee e ON od.employee_id = e.employee_id
15	WHERE
16	od.quantity_ordered > 5

order_id	order_status	customer_name	employee_name	product_code	quantity_ordered	price_each
integer	character varying (255)	character varying (255)	character varying (255)	integer	integer	numeric
1	10107	Land of Toys Inc	Melissa Knight	516,1678	30	2071.00
2	10121	Land of Toys Inc	Tammy Bird	516,1678	34	2165.90
3	10134	Land of Toys Inc	Paula Hill	516,1678	41	3884.34
4	10145	Land of Toys Inc	Jeffrey Brown	516,1678	45	2516.70
5	10159	Land of Toys Inc	Jeffrey Brown	516,1678	49	4900.00
6	10168	Land of Toys Inc	Jeffrey Brown	516,1678	36	3475.70
7	10180	Land of Toys Inc	Jeffrey Brown	516,1678	29	2487.77
8	10188	Land of Toys Inc	Jeffrey Brown	516,1678	48	4800.00
9	10201	Land of Toys Inc	Jeffrey Brown	516,1678	23	2168.34
10	10211	Land of Toys Inc	Jeffrey Brown	516,1678	41	4100.00

Total rows: 10 of 10 Query complete 00:00:00.333

Fig. 28. Problematic Query 3

Again long execution time almost 333 milli seconds for this specific query raised doubts about its effectiveness. We looked

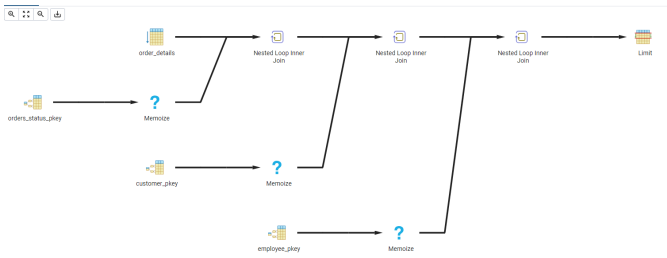


Fig. 29. Explain from PG Admin

into possible optimisations to improve the query's performance as well as the causes of its prolonged execution time.

Optimized Solution

The OrderID, CustomerID, and Employeeid columns in each of the corresponding tables in our database now have indexing implemented. We've significantly improved our query performance by enabling much faster lookup times for join operations by creating these indexes. This optimisation makes sure the database engine can find the required rows quickly rather than having to search through entire tables when our SQL queries involve joining tables on these indexed columns.

```

1 SELECT
2   os.Order_ID,
3   os.Order_Date,
4   os.Order_Status,
5   c.Customer_Name,
6   e.Employee_Name,
7   od.Product_Code,
8   od.Quantity_Ordered,
9   os.Sales
10 FROM
11   (SELECT Order_ID, Order_Date, Order_Status, Customer_ID FROM sales_data.orders_status) os
12 JOIN
13   sales_data.customer c ON os.Customer_ID = c.Customer_ID
14 JOIN
15   sales_data.order_details od ON os.Order_ID = od.Order_ID
16 JOIN
17   sales_data.employee e ON od.Employee_id = e.Employee_id
18 LIMIT 10;

```

Data Output
Messages
Explain
Notifications

	order_id integer	order_date date	order_status character varying (255)	customer_name character varying (255)	employee_name character varying (255)	product_code character varying (255)	quantity_ordered integer	sales numeric (10,2)
1	10107	2003-02-24	Shipped	Land of Toys Inc.	Michelle Gilbert	510,1678	30	2871.00
2	10121	2003-05-07	Shipped	Reims Collectables	Tammy Bird	510,1678	34	2765.90
3	10134	2003-07-01	Shipped	Lyon Souveniers	Paula Hill	510,1678	41	3884.34
4	10145	2003-08-25	Shipped	Toys4Growth.com	Brett Ward	510,1678	45	3746.70
5	10159	2003-10-10	Shipped	Corporate Gift Ideas Co.	Jennifer Mason	510,1678	49	4900.00
6	10168	2003-10-28	Shipped	Technics Stores Inc.	Olivia Hunter	510,1678	36	3479.76
7	10180	2003-11-11	Shipped	Daewidus Designs Imports	Cameron Mitchell	510,1678	29	2497.77
8	10186	2003-11-18	Shipped	Henku Gifts	Bonnie Bryant	510,1678	48	4800.00
9	10201	2003-12-01	Shipped	Mini Wheels Co.	Elizabeth Miller	510,1678	22	2168.54
10	10211	2004-01-15	Shipped	Auto Canal Petit	Rachel Walker	510,1678	41	4100.00

Total rows: 10 of 10 Query complete 00:00:00.042

Fig. 30. Optimized Solution for Problematic Query 3

VIII. DEVELOPMENT OF FLASK WEB APPLICATION

We developed a Flask web application that is tightly integrated with a PostgreSQL database. This application comprises two key routes: the root (/) route and the /search route. The root route is responsible for rendering the HTML template for the index page, while the /search route handles incoming POST requests for executing database searches.

Within the /search route, the application establishes a connection to the PostgreSQL database, utilizing parameters defined in the db_params configuration. To enhance the query generation process, we leveraged the LangChain library.

When a user submits a query through the input form on the index page, the LangChain processes this input and formulates an SQL query. Subsequently, this SQL query is executed

against the connected database. The resulting data, including column names, is dynamically rendered on a dedicated results page.

Our application is designed with robustness in mind. It adeptly handles exceptions that may arise during query execution and ensures that the database connection is closed after each query, maintaining system integrity. Additionally, we have incorporated a conditional block in the script to facilitate running the Flask app in debug mode, simplifying the debugging and development process.

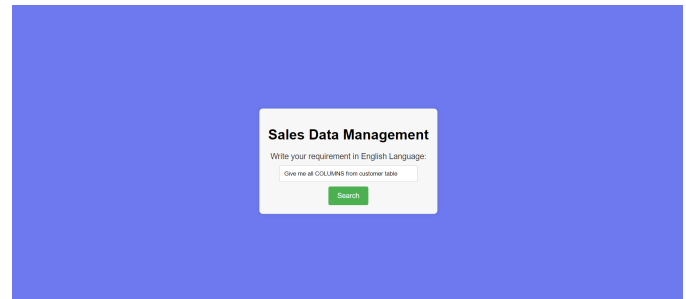


Fig. 31. Website's Home Page

Users can input queries here and accordingly will get the output like:

Query Result

customer_id	customer_name	phone	address_line1	address_line2	city	state	postal_code	country
1	Land of Toys Inc.	2125579118	107 Long Acceptor Avenue	Norae	NYC	NY	10022	USA
2	Reims Collectables	28-67-1555	29 rue de l'Alouette	Norae	Reims	Norae	51100	France
3	Lyon Souveniers	+33 1 46 62 7555	27 rue du Colonel Pierre Avia	Norae	Paris	Norae	75008	France
4	Toys4Growth.com	6285557265	78044 Hillside Dr.	Norae	Pasadena	CA	90063	USA
5	Corporate Gift Ideas Co.	650555286	1744 Strong St.	Norae	San Francisco	CA	Norae	USA
6	Technics Stores Inc.	6505556809	9408 Fourth Circle	Norae	Burlington	CA	94217	USA
7	Daewidus Designs Imports	28-16-1555	184, chemin de l'Entrée	Norae	Lille	Norae	59000	France
8	Henrietta Gifts	+47 2307 2325	Dronningen 11, PB 744 Sentrum	Norae	Bergen	Norae	28 1004	Norway
9	Mini Wheels Co.	6505555787	5557 North Pondera Street	Norae	San Francisco	CA	Norae	USA
10	Auto Canal Petit	11147-555855	25 rue Lantier	Norae	Paris	Norae	75016	France
11	Avantasia Collectables Co.	00 8739 8555	609 W. Kala Road	Norae	Malibu	Norae	90262	USA
12	Vladimir Inc.	2125551000	2678 Kingston Rd.	Norae	NYC	NY	10022	USA
13	Reims Collectables Inc.	2015559150	1476 Main Rd.	Norae	Newark	NJ	07102	USA
14	Karl Duper Inc.	2035553170	2593 South Bay Ln.	Norae	Burlington	VT	05405	USA
15	La Rochelle Gifts	40-67-8555	67, rue des Capucines Champs	Norae	Nantes	Norae	44000	France
16	Henrietta Gifts	6175558555	19523 Spinnaker Dr.	Norae	Cambridge	MA	02147	USA
17	Reims Collectables	906-224-8555	Kerkdalen 45	Norae	Helsinki	Norae	00100	Finland
18	Reims Mini Imports	07-98 9555	Edling Skjolden gate 78	Norae	Stavrum	Norae	4110	Norway
19	Reims Collectables	2125551555	1768 Higgins St.	Norae	Albany	NY	12207	USA
20	Reims Collectables	6505559555	1785 Main Street	Norae	Saltburg	PA	15468	USA
21	Reims Collectables	+61 2 9495 8555	Monter Moore Building, 515 Pacific Hwy	Norae	Cherrywood	NSW	2067	Australia
22	Reims Collectables	5085552555	1785 Main Street	Norae	New Bedford	MA	01905	USA
23	Reims Collectables, Ltd.	0171 555-2382	Bedeley Gardens 12 Brewery	Norae	Liverpool	Norae	06761	UK
24	Reims Collectables, Ltd.	0911 555 84 44	C. Moorland, 86	Norae	London	Norae	20014	Spain
25	Reims Collectables, Ltd.	0911 12 5555	Bergstrasse 8	Norae	Lahr	Norae	6-498 22	Sweden
26	Reims Collectables, Ltd.	0911 555 22 32	C. Avenue 17	Norae	Madrid	Norae	28023	Spain
27	Reims Collectables, Ltd.	+65 221 7555	Bonus Inc., Bonus Art, 3-6 Teviera	Norae	Singapore	Norae	79903	Singapore

Fig. 32. Website's Result Page

ACKNOWLEDGMENT

I sincerely appreciate Professor Dr. Shamsad Parvin for their insightful advice. I also acknowledge all the TAs for their assistance. Both have significantly improved my educational experience.

REFERENCES

- <https://www.kaggle.com/datasets/kyanyoga/sample-sales-data>
- <https://lucid.app/>
- <https://www.lucidchart.com/pages/er-diagrams>
- <https://www.postgresqltutorial.com/>
- <https://www.simplilearn.com/tutorials/sql-tutorial/er-diagram-in-dbms>
- https://python.langchain.com/docs/use_cases/qa_structured/sql
- <https://platform.openai.com/docs/api-reference/making-requests>

<https://flask.palletsprojects.com/en/3.0.x/>
<https://www.codecademy.com/article/sql-indexes>
<https://www.postgresqltutorial.com/postgresql-tutorial/import-csv-file-into-posgresql-table/>