

NANYANG TECHNOLOGICAL UNIVERSITY

SCHOOL OF COMPUTER SCIENCE

AND ENGINEERING



CZ4031 Group 31 Project 2 Report

CHIA SONGCHENG U2021738F

CHIN KIN MUN U2020779H

LIU LIWEN N2202357F

LUO MINJUAN N2202356J

Table of Contents

1. Introduction	3
2. project.py	3
2.1. Database:	3
2.2. main function:	3
3. interface.py	4
4. explain.py	5
4.1. Retrieval of query plan.....	5
4.2. Check depth of query plans	6
4.3. Differences between two query plans.....	6
4.4. Checking for parallel execution and limit	6
4.5. Checking for similar operations	6
4.6. Check for any other unique operations	6
5. Features and limitations	7
5.1. Features	7
5.2. Limitation	7
Appendix A: Installation Guide	8
Appendix B: Examples.....	10

1. Introduction

The report is going to explain details of the software including formal descriptions of the key algorithms with examples base on the submission.

Hence the content should be:

1. project.py
2. interface.py
3. explain.py

2. project.py

The software provided is a Python-based application that allows users to execute and visualize query execution plans for a database. It is composed of three main Python files. The program uses the PostgreSQL database and connects to a specified database by providing the necessary login credentials. The main program file, project.py, is responsible for invoking all the necessary procedures from the other two files.

2.1. Database:

The Database contains a class that initializes the connection to the database and has methods for executing SQL queries and returning the results.

The class constructor takes the name of the database as an argument and initializes the necessary login credentials such as the hostname, username, password, and port id. The psycopg2 library is used to establish a connection to the database.

The method get_query_results(query) take an SQL query as an argument, executes the query using the established connection and returns the results.

The method print_explanation_results(queryresult) print out the explanation of the query execution plan while print_query_results(queryresult) print out the query results. This is used for debugging.

2.2. main function:

The main function in the project.py file will invoke all the necessary procedures from the other three files. It initializes a database object and a user interface object and connects to the database.

When the user inputs an SQL query and presses the execute button, the program executes the query using the database object and displays the results in the table widget of the user interface.

3. interface.py

The file interface.py contains the code for the GUI.

The interface.py file provides a graphical user interface for executing SQL queries. The file uses the PyQt5 library for building the GUI. The file consists of a class called MyWindow that is a subclass of QMainWindow. The following are the key components of the GUI:

The class constructor initializes the user interface window and creates the necessary widgets including text boxes, labels, and tables to display the query results. It also defines the necessary functions to execute the SQL query, retrieve the query results, and display them in the table widget.

1. Textboxes for inputting SQL queries: The GUI has two textboxes for inputting SQL queries. These textboxes are named queryTextbox1 and queryTextbox2.
2. Textbox for inputting the database name: The GUI has a textbox for inputting the name of the database. This textbox is named dbNameTextbox.
3. Labels for displaying output: The GUI has three labels for displaying output. These labels are named queryOutput1, queryOutput2, and queryExplain.
 - a. The queryOutput1 label displays the output of the first SQL query
 - b. The queryOutput2 label displays the output of the second SQL query
 - c. The queryExplain label displays the user-friendly explanation of the changes to the two inputted query execution plans that take place during data exploration.
4. Button for executing SQL queries: The GUI has a button named submitButton for executing the SQL queries.
5. Error message box: The GUI has an error message box for displaying error messages.

The program uses the following algorithms:

1. onClick() method: This method is executed when the submitButton is clicked. The method first checks if the database name entered by the user is the same as the current database name. If not, it establishes a connection with the new database. Next, it calls the process () method from the project.py file to execute the SQL queries and obtain the output. It then calls the explain () method from the explain.py file to obtain the execution plan of the first SQL query. Finally, it displays the output and execution plan in the appropriate labels.
2. ScrollableLabel class: This is a custom class that extends the QScrollArea class. The ScrollableLabel class is used for displaying the output and execution plan in the GUI. The class provides a setText() method that sets the text to be displayed in the label.

Conclusion:

The interface.py file provides a GUI for executing SQL queries. The file uses the PyQt5 library for building the GUI. The program uses the onClick() method and the ScrollableLabel class for executing the SQL queries and displaying the output and execution plan in the GUI.

4.explain.py

The explain.py contains the code for generating the explanation.

The program uses the PostgreSQL EXPLAIN command to get the execution plans for the two queries, and then compares the node types and depth of the plans to identify any changes. The program outputs a report describing the differences between the two plans. This output is shown in Appendix B. This result is written in the text box named "Explain for Output Query:".

```
def explain(self, database, query1, query2):  
    print("working on explain")
```

Figure 1: explain function

The program defines a class called 'explain' and a method called 'explain' which takes three arguments: 'database', 'query1', and 'query2'. The database argument is an object representing the database to execute the queries on, while 'query1' and 'query2' are strings representing the two SQL queries to compare.

4.1. Retrieval of query plan

The "explain" method first retrieves the query execution plans for the two SQL queries using the "get_query_results" method of the database object. It then uses a helper function called "get_node_types" to extract the node types and execution times from the raw explanation of the query execution plan. The "get_node_types" function recursively traverses the query execution plan and extracts the node types and execution times of each operation.

```
def get_node_types(raw_explanation):  
    node_types_time = []  
    node_types = []  
    for Plans in raw_explanation[0][0]:  
        node_types_time, node_types=get_node_helper(Plans['Plan'])  
    return node_types_time, node_types
```

Figure 2: Function for get_node_types

This function takes raw explanation generated by PostgreSQL explain command as input and gets information about which operations are performed and the execution time for each operation. It uses a get_node_helper function to recursively get information about all operations and return two lists: the operations with their execution time, and the name of operations that are used to compare the difference between two queries.

```
def get_node_helper(raw_explanation):
```

Figure 3: Function for get_node_helper

The get_node_helper function is a recursive function that recursively goes through each depth of the query execution plan tree to get information about operations and adds them into lists. The get_node_helper function is called by the get_node_types function.

4.2. Check depth of query plans

```
def check_depth(query):
```

Figure 4: Function for check_depth

The "explain" method then compares the node types of the two-query execution plans to identify any differences. It uses another helper function called "check_depth" to calculate the depth of each plan and checks if they are the same. If the depth is different, the method adds an explanation stating the difference.

4.3. Differences between two query plans

Next, the "explain" method compares the node types of the two-query execution plans to identify any differences in the operations performed. It uses a nested loop to compare the node types of each plan and identify any nodes that are present in one plan but not the other.

4.4. Checking for parallel execution and limit

The program then checks if a "Gather" node was added or removed and adds an explanation to indicate whether parallel execution was enabled or disabled. It also checks for limit by checking for a "Limit" node in the query plan.

4.5. Checking for similar operations

The "explain" method identifies any nodes that were present in both query execution plans but had a different depth in each plan. It generates an explanation to highlight these nodes, indicating that the operations were the same but were executed at a different depth.

It stores these nodes in a list called 'same_operation_diff_depth', removes them from the 'nodes_in_query1' and 'nodes_in_query2' lists, and then proceeds to compare the remaining nodes in the same way as before.

4.6. Check for any other unique operations

```
def compare_query(query1, query2):
```

Figure 5: Function for compare_query

The compare_query function checks for any other unique operations in both queries so as to map the differences in the query plan. It removes similar parts from the queries and gives an array of three arrays. Each of the three arrays contains another two arrays of strings to show the differences in the select clause, where clause and the from clause. It then adds the explanation for the unique operations done in one query compared to the other based on the differences in the queries.

Overall, the explain method aims to provide a detailed and informative comparison of the query execution plans of query1 and query2. The method achieves this by recursively traversing the execution plan trees of the two queries, extracting relevant information, and identifying differences and similarities between the two.

5. Features and limitations

5.1. Features

1. Able to identify whether parallel execution is used for the query.
2. Able to identify if there is a limit to the number of tuples in the query plan.
3. Able to show the complexity of the query plan:
 - a. Able to show the depth of the query plan.
 - b. Able to show the different number of operations in the query plan.
4. Gives reasons why the query plans are different.
5. Looks up at the query and check for differences to explain the differences to the query plan.

5.2. Limitation

1. Only up to two query plans are generated and get compared.
2. The query execution system will work slow if the query output contains many results.
3. Unable to map each difference in the query plan to a difference in a query.
4. If the query run into an error, the program must be terminated and start again to work.

Appendix A: Installation Guide

Steps to install the python environment:

1. Install python from <https://www.python.org/downloads/>

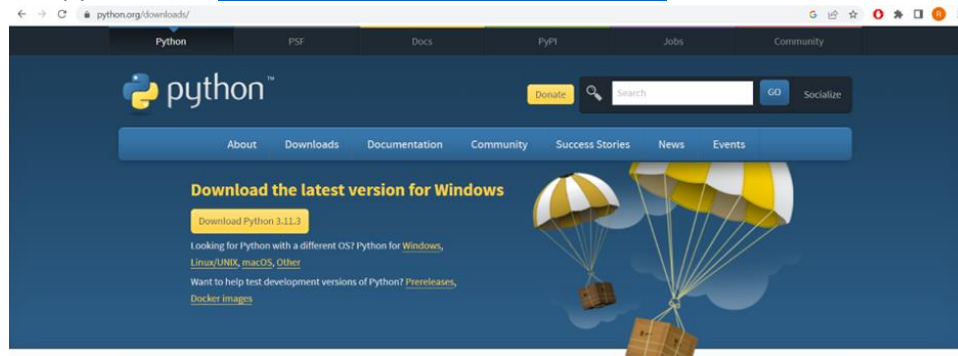


Figure 6: Figure of website to install python

2. Set up virtual environment:
 - a. Open Python terminal
 - b. Enter command:

```
Pip install virtualenv
python -m venv venv
```

Which would create a virtual environment with venv as its name.
 - c. To activate the virtual environment enter the following in the terminal: `.\venv\Scripts\activate`
 - d. Libraries installation

After the virtual environment is activated, install the following libraries by entering the following:

Pip install `psycpg2==2.9.1`

Pip install `PyQt5==5.15.9`

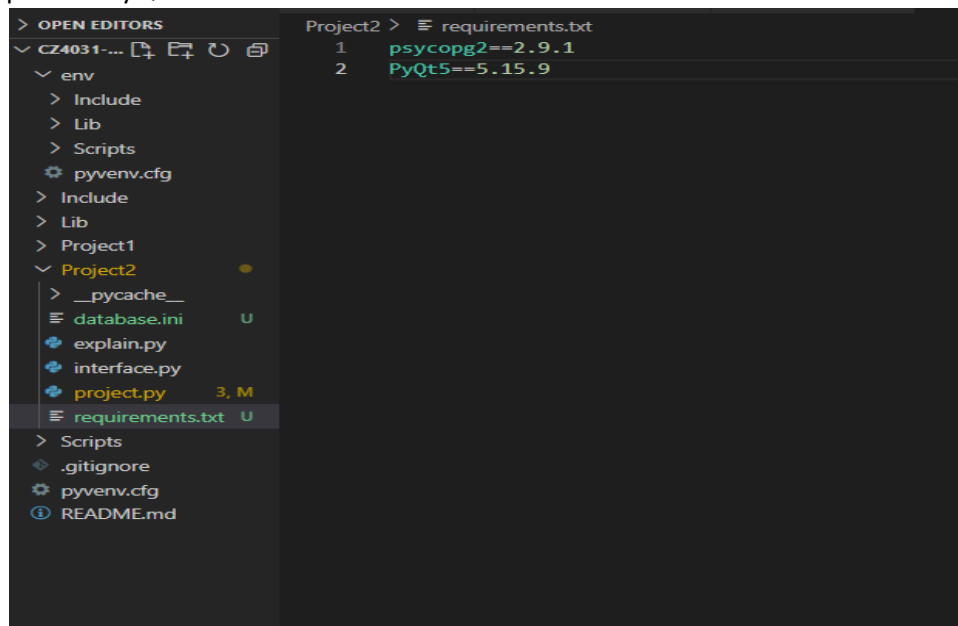


Figure 7: Figure of requirements.txt

Virtual environment is now set up.

3. PostgreSQL extension is enabled/installed in VisualStudioCode(VSC)

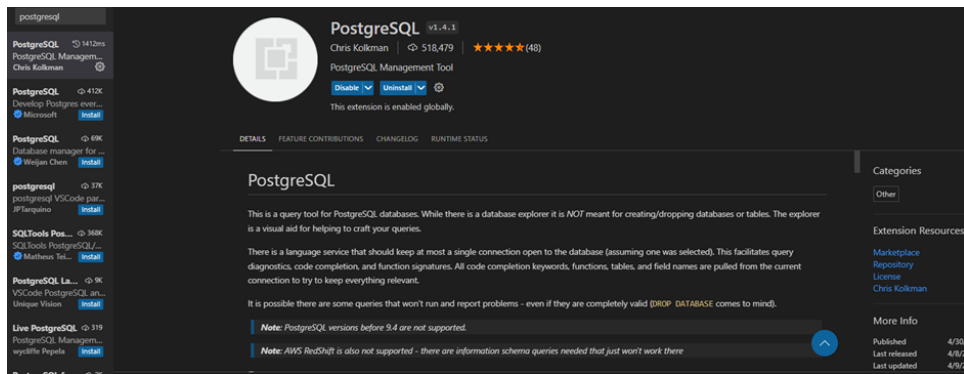


Figure 8: Figure of PostgreSQL extension in VSC

4. Modify database.ini with login details for database

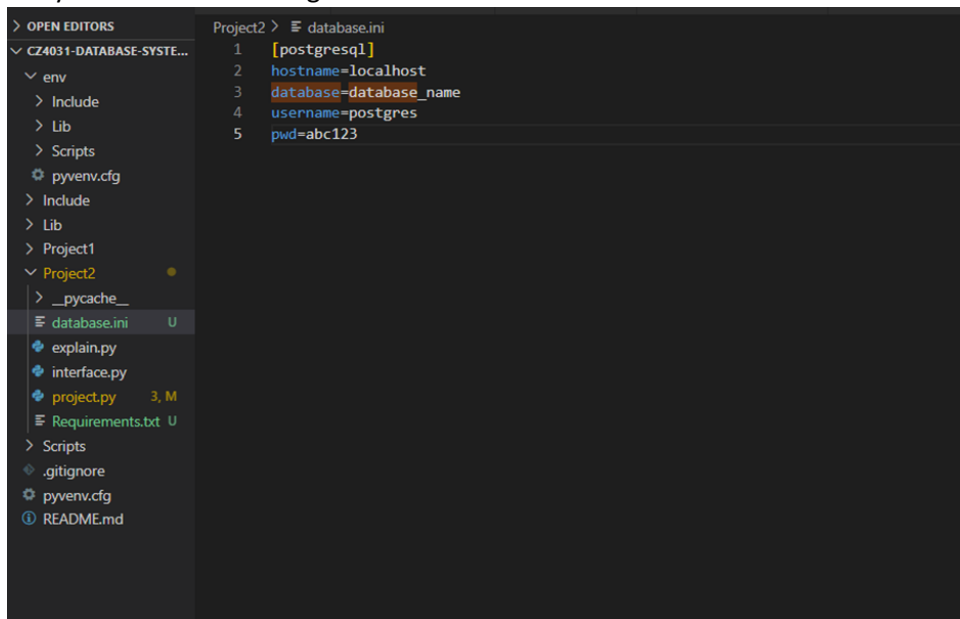


Figure 9: Figure of database.ini in the project folder

Appendix B: Examples

CZ4031-Assignment2-Group 31

Insert SQL Query Q1 Here :

Insert SQL Query Q2 Here :

Insert Database Name Here :

Output Query Q1 :

Output Query Q2 :

Current Database Name:

Explain for Output Query :

Submit Query

Refresh

Figure 10: UI for showing the explanation and output of the query.

CZ4031-Assignment2-Group 31

select * from lineitem as l, orders as o,part as p where l.l_partkey=p.p_partkey and l.l_orderkey=o.o_orderkey and o.o_orderpriority='S-LOW' and p.p_mfgr='Manufacturer#5'

select * from lineitem as l, orders as o,part as p where l.l_partkey=p.p_partkey and l.l_orderkey=o.o_orderkey and o.o_orderpriority='S-LOW' and p.p_mfgr='Manufacturer#5' order by p.p_partkey

CZ4031

2440688406865771141 0043134 4610 0810 08101995-12-021995-10-221995-12-16NONE
2501511142309233546 0052634 5810 0310 08101997-03-111996-12-241997-04-10TAKE I
250177525887599145 0069326 1010 0810 08101997-11-011997-11-081997-11-22COLLE
250177399277437335 0065342 2010 0810 01101997-11-221997-12-161997-12-08TAKE I
250305336996403243 0078614 2710 0810 051993-07-051993-06-151993-08-01NONE
250305388853886513 0023320 4410 0810 041993-10-121993-07-171993-10-17DELIVE
2503051246347147626 0043124 3810 0410 071993-06-221993-07-291993-09-08COLL
2510111286873712126 0044607 6810 0710 051998-06-211998-06-111998-06-28TAKE
2510111520909639921 0023983 6910 0310 051998-07-241998-09-221998-08-09COLL
2511118572789116 0028124 3210 0010 031992-07-211992-08-221992-07-23DELIVE
2511111536586174547 0080447 5510 0310 011992-08-281992-07-291992-09-03NONE
2513351129922993123 0046114 7710 0710 001997-06-051997-05-141997-06-11TAKE
2515551308423356710 0018728 4010 0410 081995-06-271995-07-301995-07-16TAKE
25164823254826345 005229 4010 1010 051998-06-251998-05-311998-07-08COLLE
252517121042105111 001126 1010 0910 021994-08-311994-08-311994-09-08NONE
2525171932028241247 0060874 4010 0410 071994-10-171994-09-231994-10-25DELA
2525171885061025418 0028701 0010 1010 001994-09-061994-09-201994-09-27NONE
252550643716878119 0012018 3310 0210 001995-04-191995-02-271995-05-07COLLE
333488311015011533 0030030 3310 0810 011992-06-261992-07-071992-07-22COLLEC
2629221110251123 002730 0310 0110 031997-05-281997-04-291997-06-08NONE
353526910111144 0040040 4410 0810 071998-09-231998-07-081998-10-21DELIVER IN
56951111011134 003640 0410 0210 001993-06-241993-05-241993-07-16DELIVER IN FI
13144641111627 0024570 2710 0810 011995-09-221996-07-271996-10-08DELIVER IN
47647751102511532 0029120 3210 0710 051992-10-081992-08-221992-10-23DELIVER
2883074132514337 0033781 3710 0810 021995-05-211995-04-251995-06-11TAKE BA
545921137514224 0021912 2410 0310 031995-06-231995-08-311995-07-02NONE
17725141325146947 0042911 4710 0510 031996-10-021996-11-211996-10-10COLLEC
28192021350144446 0043824 4810 0110 081992-09-111992-11-081992-09-12DELIVER
447108137514650 0045650 5010 0010 001997-12-271997-11-161998-01-21NONE
2842339247525547 0043428 9410 0010 081997-10-221997-08-251997-11-15NONE
7762952425246 0042504 9210 0810 051994-05-291994-06-191994-06-04COLLECT C
19964812450253523 0021252 4810 0410 021995-05-211995-04-051995-06-07COLLEC
33315842450294033 0030492 6610 0810 071997-10-031997-08-231997-10-31TAKE BA
3954851247525129 0026796 5810 1010 001997-03-281997-05-251997-04-26NONE
24741442425240 0036960 8010 0010 081994-06-261994-07-051994-07-06NONE
5949380242514 003696 0810 0110 021998-07-301998-06-211998-07-31TAKE BACK F

These operations are performed on query 1:
Gather (2201.844miliseconds) on
->Hash Join (2049.492miliseconds) on
->Hash Join (1175.593miliseconds) on
->Seq Scan on l (477.549miliseconds)
->Hash (21.525miliseconds) on
->Seq Scan on p (16.934miliseconds)
->Hash (175.73miliseconds) on
->Seq Scan on o (135.871miliseconds)

These operations are performed on query 2:
Gather Merge (2326.97miliseconds) on
->Sort (2129.718miliseconds) on
->Hash Join (1974.496miliseconds) on
->Hash Join (1143.055miliseconds) on
->Seq Scan on l (462.722miliseconds)
->Hash (20.136miliseconds) on
->Seq Scan on p (15.784miliseconds)
->Hash (172.239miliseconds) on

Submit Query

Refresh

Figure 11: Execution of Query Pair A part 1

Figure 12 shows the execution of Query Pair A part 2. The interface displays two SQL queries side-by-side, both selecting from a table named 'lineitem' based on 'l_orderkey' and 'o_orderkey' conditions. The left query includes a filter for 'p_partkey' and 'p_mfgr'. The right query includes a filter for 'p_partkey' and 'p_mfgr'. Below the queries, the execution plan for the left query is shown, detailing the sequence of operations: a nested loop join, followed by a hash join, and then a seq scan on the 'lineitem' table. The execution time for the left query is 19.268 milliseconds. The right query's execution plan is also shown, detailing the sequence of operations: a nested loop join, followed by a hash join, and then a seq scan on the 'lineitem' table. The execution time for the right query is 19.268 milliseconds. The interface includes a 'Submit Query' button and a 'Refresh' button.

Figure 12: Execution of Query Pair A part 2

Figure 13 shows the execution of Query Pair B part 1. The interface displays two SQL queries side-by-side, both selecting from a table named 'customer' based on 'c_nation' and 'n_nationkey' conditions. The left query includes a filter for 'c_nation' and 'n_nationkey'. The right query includes a filter for 'c_nation' and 'n_nationkey'. Below the queries, the execution plan for the left query is shown, detailing the sequence of operations: a nested loop join, followed by a hash join, and then a seq scan on the 'customer' table. The execution time for the left query is 19.268 milliseconds. The right query's execution plan is also shown, detailing the sequence of operations: a nested loop join, followed by a hash join, and then a seq scan on the 'customer' table. The execution time for the right query is 19.268 milliseconds. The interface includes a 'Submit Query' button and a 'Refresh' button.

Figure 13: Execution of Query Pair B part 1

