



ECOLE
POLYTECHNIQUE
DE BRUXELLES

2023-2024

INFO-H419 - Data Warehouse

Project 2 TPC-DI Benchmark Using DuckDB

Group:

Jintao Ma(000586557)

Linhan Wang(000585737)

Iyoha Peace Osamuyi(000583313)

Hieu Nguyen Minh(000583782)

Professor

Esteban Zimányi

December 2023



Contents

1	Introduction	1
1.1	Overview	1
1.2	Aim and Objectives	1
1.3	Technology Used	2
1.3.1	DuckDB	2
1.3.2	Simple and portable	2
1.3.3	Python	3
1.3.4	Code Editor: Microsoft Visual Code	4
1.3.5	Tools and materials	4
1.4	Limitations	4
2	ETL and TPC-DI Benchmarking	5
2.1	ETL	5
2.2	Benchmarking	6

2.3	TPC-DI	7
2.4	DIGEN	8
2.5	Phases of operation	8
3	System Setting and Data Generation	11
3.1	Hardware Specification	11
3.2	Setup for Project	11
3.2.1	DuckDB installation	12
3.2.2	TPC-DI kit installation	12
3.3	Database schema creation	12
3.4	Data generation	13
3.4.1	Jre and JVM Setting	13
3.4.2	Data generation	13
4	Benchmark Implementation Process	14
4.1	Implementing TPC-DI on DuckDB	14
4.2	Tasks	16
4.2.1	Task 1	16
4.2.2	Task 2-3	18
4.2.3	Task 4	19
4.2.4	Task 5-6	20
4.2.5	Task 7	21

4.2.6 Task 8-11	21
4.2.7 Task 12-15	22
4.2.8 Task 16	24
4.2.9 Task 17	25
4.2.10 Task 18-20	27
4.2.11 Task 21	29
4.2.12 Task 22	30
4.2.13 Task 23-24	31
4.2.14 Task 25	34
4.2.15 Task 26	35
4.2.16 Task 27	36
4.2.17 Task 28	37
4.2.18 Task 29-30	38

5 Results and Discussion 42

5.1 Total Performance	42
5.2 Comparison of execution times for all tasks	43
5.2.1 The tasks of linear growth	48
5.2.2 The tasks of decreasing tendency	49
5.2.3 The tasks of stable trend	50
5.3 Discussion	50

6 Conclusion	51
6.1 Conclusion	51
6.1.1 Further Works	51
References	52

Listings

2.1	Data generation.	8
3.1	DuckDB installation.	12
3.2	TPC-DS kit installation.	12
3.3	Database setup	12
3.4	Jre and JVM Setting.	13
3.5	Data generation.	13
4.1	Part of Task 1 SQL	16
4.2	Part of Task 2 SQL	18
4.3	Part of Task 4 SQL	19
4.4	Part of Task 5 SQL	20
4.5	Part of Task 8 SQL	22
4.6	Task 14 SQL	22
4.7	Task 15 SQL	23
4.8	Task 16 SQL	24

4.9	btrim_equivalent	25
4.10	Task 17 SQL	25
4.11	Task 19 SQL	28
4.12	Task 20 SQL	29
4.13	Task 20 SQL	30
4.14	Task 23 SQL	31
4.15	Task 24 SQL	33
4.16	Task 25 SQL	34
4.17	Task 26 SQL	35
4.18	Task 27 SQL	36
4.19	Task 28 SQL	37
4.20	Part of Task 29 SQL	38
4.21	Task 30 SQL	41

List of Figures

2.1	ETL Process FLOW	6
4.1	TPC-DI Benchmarking BPMN Diagram	14
4.2	Staging Database Schema	15
4.3	Master Data Warehouse	16
5.1	Line Chart of Total execution time	43
5.2	Bar Chart of Total execution time	43
5.3	execution times for all tasks	44
5.4	Proportion of Execution Time of Various Tasks for SF3	44
5.5	Proportion of Execution Time of Various Tasks for SF5	45
5.6	Proportion of Execution Time of Various Tasks for SF7	45
5.7	Proportion of Execution Time of Various Tasks for SF9	46
5.8	execution times for all tasks after removing two tasks	46
5.9	execution times for all tasks after removing more tasks	47

5.10 execution times for parse_finwire.	47
5.11 execution times for convert_customermgmt_xml_to_csv.	47
5.12 execution times for load_txt_csv_sources_to_staging.	48
5.13 execution times for load_finwire_to_staging.	48
5.14 execution times for transform_load_master_dimaccount.	49
5.15 execution times for transform_load_master_dimtrade.	49
5.16 execution times for create_staging_schema.	49
5.17 execution times for create_master_schema.	50
5.18 execution times for load_master_tradetype.	50

Abstract

In this project, we present an in-depth implementation of the TPC-DI Benchmark using DuckDB, a versatile, free, and open-source database from CWI. Our comprehensive report is organized into six key chapters, each addressing crucial aspects of the benchmarking process:

Chapter 1: Sets the stage by outlining the project's goals, requirements, and introduces the programming tools and databases utilized.

Chapter 2: Delves into the fundamentals of ETL (Extract, Transform, and Load) and the intricacies of TPC-DI benchmarking.

Chapter 3 and 4: Detail our setup and configuration, alongside a comprehensive guide on executing tests via Python scripts. This section includes various assumptions and practical code implementations.

Chapter 5: Focuses on data-driven insights, employing numerous charts for comparison, analysis and interpretation,

Chapter 6: Concludes with a summary of our findings and reflections on the project.

For those interested in a more granular exploration of our experimental and analytical approaches, we have made all related codes and scripts accessible through our GitHub repository at

<https://github.com/facingfrost/tpcdi-duckdb>.

1.1 Overview

The Project entails performing TPC-DI Benchmarking on a selected Database Management System which include DuckDB, PostgreSQL, SparkSQL, MySQL etc. We can execute the benchmarking using a chosen data integration tool such as Python scripts, Apache Airflow, Pentaho Data Integration, Oracle Data Integrator, Talend Data Studio, or SQL scripts. The Benchmarking is done with several scale factors which determine the size of the resulting data warehouse and then implement some metrics to evaluate the performance of the database.

The project is carried out in groups of 4 persons and delivers a report in pdf file explaining the essential aspects of your implementation, and a zip file containing the code of your implementation, with all necessary instructions to be able to replicate your implementation. This report uses Python script for the data integration tool and DuckDB as the preferred choice of Database Management System on which the TPC-DI benchmark is performed.

1.2 Aim and Objectives

The aim of this project is to execute the TPC-DI benchmark, which involves the extraction, transformation, and loading (ETL) of data from an On-Line Transaction Processing (OLTP) system and various other data sources into a chosen data warehouse (DuckDB), with the following objectives include:

- To learn and understand the steps required for a complete ETL (Extract, Transform, Load) implementation.
- Evaluate benchmarking performance and conducting an analysis of the outcomes.

1.3 Technology Used

1.3.1 DuckDB

DuckDB is an in-process SQL OLAP database management system.¹ The first open-source version of it was released by Hannes Mühleisen and Mark Raasveldt, who were researchers in the Database Architectures research group at Centrum Wiskunde & Informatica (CWI) in 2019.[RM19] It was the first purpose built in-process Online Analytical Processing(OLAP)-database management system.² It was designed to be fast, reliable and easy to use. There have been some important features of this tool:

1.3.2 Simple and portable

For DuckDB, there is no DBMS server software to install, update and maintain. It does not run as a separate process, but is completely embedded within a host process. For the analytical use cases that DuckDB targets, this has the additional advantage of high-speed data transfer to and from the database. In some cases, DuckDB can process foreign data without copying. We can use DuckDB without installing, updating, or maintaining separate DBMS server software. It doesn't run as an independent process but is entirely integrated within a host process. This approach provides added benefits for analytical use cases, notably enabling high-speed data transfer to and from the database.

¹<https://duckdb.org/>

²<https://www.cwi.nl/en/news/cwi-spin-off-company-duckdb-labs-provides-solutions-for-fast-database-analytics>

Feature-rich

DuckDB enables users to run complex SQL queries and provides APIs for Java, C, C++, and more. It is also deeply integrated into Python and R, enabling users to conduct efficient interactive data analysis; thus, you can interact with DuckDB from your preferred programming language. There is also access to extra SQL keywords that make SQL queries easier to write, such as EXCLUDE, REPLACE, and ALL.

Fast

DuckDB can support the Online analytical processing (OLAP). It features a query execution engine that employs columnar vectorization. While a few queries may still be interpreted, a substantial batch of values, known as a "vector," is processed in a single operation. This approach significantly minimizes the overhead found in conventional systems like PostgreSQL, MySQL, or SQLite, which handle each row sequentially. The utilization of vectorized query execution significantly enhances performance in OLAP queries. They used an interactive SQL shell to demonstrate the capabilities of an analytical database in the browser[Koh+22].

Free and Extensible

Because the DuckDB is Open Source, all of the source code is freely available on GitHub. You can access it by clicking the link: <https://github.com/duckdb>. They also invite developers from anyone to take issues.

1.3.3 Python

Python is an interpreted, general-purpose, interactive, object-oriented programming language. It incorporates modules, exceptions, dynamic typing, very high level dynamic data types, and classes. It supports multiple programming paradigms beyond object-oriented programming, such as procedural and functional programming. Python combines remarkable power with very clear syntax. It has interfaces to many system calls and libraries, as well as to various window systems, and is extensible in C or C++ [Pyt23]. In this project, we used python to write scripts for creating

the database, loading of tables in our database, querying and performing various evaluation metrics. Also, Jupyter notebook with matplotlib was used to visualize our data. Python version 3.10 was used in this project.

1.3.4 Code Editor: Microsoft Visual Code

Visual Studio Code is a lightweight but powerful source code editor which runs on your desktop and is available for Windows, macOS and Linux. It comes with built-in support for JavaScript, TypeScript and Node.js and has a rich ecosystem of extensions for other languages and runtimes (such as C++, CSharp, Java, Python, PHP, Go, .NET) [Cod23]. It was useful in this project to write our python scripts and running dsdgen and dsqgen on command line.

1.3.5 Tools and materials

The main reading material for conducting our experiments is TPC BenchMark DI Standard Specification Version 1.1.0 [TPC14]. Besides, we follow the tutorial in DuckDB documentation³.

1.4 Limitations

One Limitation was that the SQL Keyword "**ALTER**" hasn't been fully implemented to handle referential integrity constraints. Time was a limiting factor for us to thoroughly expose the objectives stated above.

³<https://duckdb.org/docs/archive/0.9.1/>

ETL and TPC-DI Benchmarking

2.1 ETL

According to Judith Awiti, Alejandro A. Vaisman and Zimányi [JZ20], ETL is defined as the Extraction, Transformation and Load process, in which the data go through three steps: extracting data from internal and external sources of an organization, transforming these data, and loading them into a data warehouse. The Data Integration (DI) should combine data from various sources into a single one to solve possible conflicts of data meaning and data shape.

The ETL process plays a pivotal role in data warehouse maintenance. By making use of ETL on the data, organizations can consolidate information from diverse sources with different formats, and get a unified view of their data. Besides facilitating the data integration, ETL is a foundation for analytics, reporting, and business intelligence, enabling enterprises to extract valuable insights from their huge data. The diagram 2.1 below shows the ETL process flow:

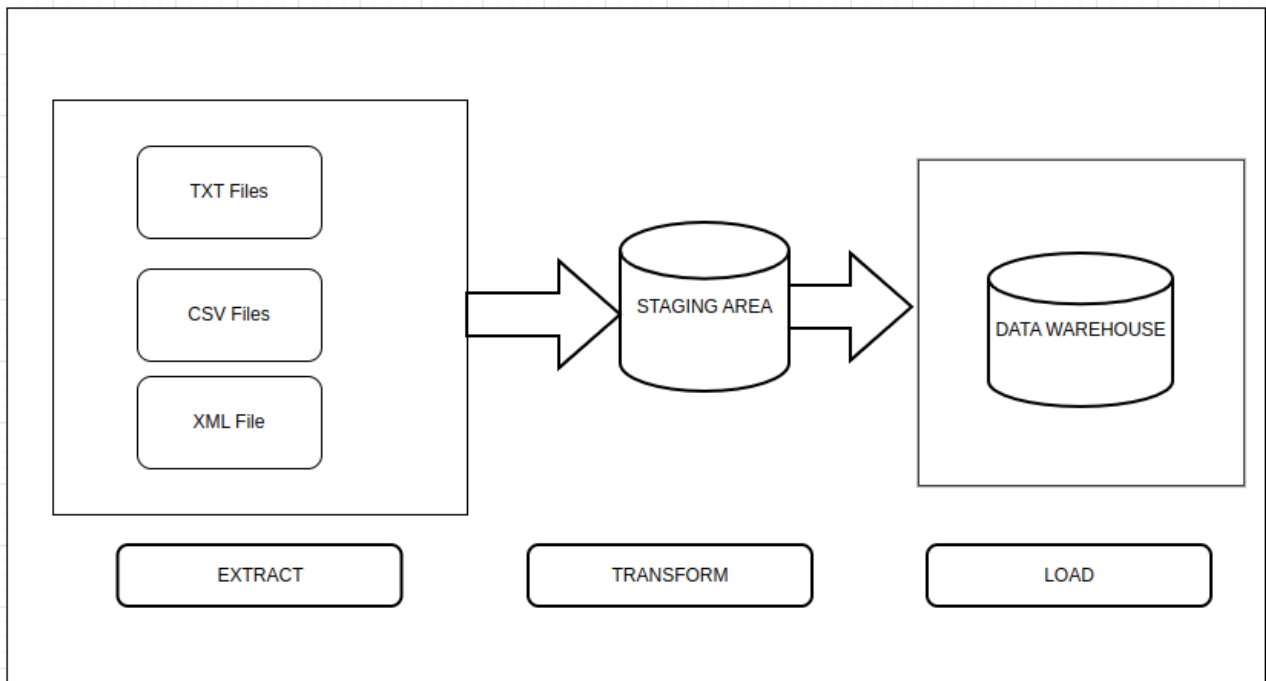


Figure 2.1: ETL Process FLOW

2.2 Benchmarking

Benchmarking is a process of comparing and evaluating an organization's performance against an external standard or best practice model. It involves measuring and analyzing various aspects of an organization's processes, performance, and practices against those of other similar organizations, with the goal of identifying areas for improvement and setting goals for achieving better results [Mah21]. Benchmarking a database is the process of performing well defined tests on that particular database for the purpose of evaluating its performance [Sam21]. Benchmark test results facilitate means for cross platform comparisons of various database management systems by providing valuable information to database professionals on whether to utilise a particular database product. Hence, we are carrying and evaluating the performance of DuckDB against the TPC-DS Benchmarking standard.

Transaction Processing Performance Council (TPC) is a non-profit organization founded to define transaction processing and database benchmarks and to disseminate objective, verifiable TPC performance data to the industry ¹. In terms of data warehouses, TPC provide two benchmarks: TPC-DS and TPC-DI for the evaluation of decision support and data integration, respectively.

¹https://en.wikipedia.org/wiki/Transaction_Processing_Performance_Council

2.3 TPC-DI

The TPC Benchmark™ DI (TPC-DI) serves as a model for data integration benchmarking, covering all aspects of a fundamental data integration setup [TPC14]. This encompasses the complete Extract, Transform, Load (ETL) process and regular data upkeep. TPC-DI offers a detailed framework for data integration, effectively assessing the performance of the System Under Test (SUT). It is a typical data integration platform, aligning with industry standards and distinguished by its thorough method of evaluating data integration proficiency. These components are characterized by:

- Handling and transferring substantial amounts of data.
- Diverse transformations applied to the dataset.
- Incorporation of both historical and incremental data updates.
- Adherence to data consistency standards.
- Integration of various data sources in different formats.
- Management of numerous data tables, each with distinct data types, attributes, and interrelations

Conversely, the operational model for the TPC-DI benchmark is structured in the following manner:

- The **source data** is autonomously generated and stored in flat file formats.
- The **transformation phase** kicks off with the System Under Test (SUT) accessing and reading this source data.
- During transformation, the **source data undergoes validation** and is methodically structured for efficient loading into a Data Warehouse (DWH).
- This **process reaches completion** once all the source data has been transformed and is successfully stored in the DWH.

The **performance metric** of the benchmark is a throughput measure, the **Number of Source data rows processed per second**.

2.4 DIGEN

DIGen, the automated data generator supplied by the TPC-DI benchmark, serves a crucial role in generating Source Data and audit details, and it's the exclusively endorsed method for this purpose. It provides detailed statistics regarding the data generation process.

The data is generated from the TPC-DI tool in the following way:

- **Requirement:** DIGen - Data Generation utility v1.1.0 used for generating source data for the TPC-DI benchmark.
- **Dependencies:**
 - Requires Java SE 8 or above
 - PDGF which is located in the same directory as the DIGen
- **Command line Usage:**

```
1 $ java -jar DIGen.jar -sf <scalefactor: 3,5,7,9> -o <directory:>
```

Listing 2.1: Data generation.

2.5 Phases of operation

The TPC-DI benchmark, reflecting real-world Data Integration (DI) applications, involves two primary processes: a historical load and incremental updates. The historical load typically occurs during the initial creation or restructuring of a data warehouse, while incremental updates represent the ongoing addition of new data. The TPC-DI benchmark specifically models daily updates within a constrained time frame, often overnight, to fit with other maintenance activities. These processes are crucial in handling large volumes of historical data within limited time windows. The benchmark encompasses several phases executed sequentially:

- **Preparation Phase:**
 - **Data Generation:** Utilizing a data generator (as detailed in Clause 6), this phase involves creating data either directly in the Staging Area or elsewhere before transferring it to the Staging Area. This process is not timed for the benchmark.

- **Data Warehouse Creation:** Involves setting up the Data Warehouse database and tables, including disk space allocation. This setup is not considered a DI operation and thus is not timed.
- **Data Integration Preparation:** This step varies among implementations and includes preparing and configuring the data integration software. It is not timed for the benchmark.
- **Historical Load Phase:** This phase, which is timed for the benchmark, involves different transformations than incremental updates. It populates initially empty destination tables with new data. Source files may differ in their ordering properties compared to incremental updates, and the Historical Load typically uses a larger data set. After this phase, a Validation Query is run to collect information for the automated audit phase.
- **Incremental Update Phase:** Distinct from the Historical Load, this phase deals with transformations of CDC extracts from the OLTP database, representing changes since the last extract. It includes two required phases to ensure repeatability, each with a completion time of 30-60 minutes. A Validation Query follows each update for correctness verification in the audit phase.
- **Automated Audit Phase:** This final phase, not timed, involves querying the Data Warehouse to conduct extensive tests on the resultant data and generate a report. All tests must pass for the run to be considered valid.

Each of these phases is integral to the TPC-DI benchmark, emulating the complexities and constraints of real-world data integration tasks. However, for this study, we implement only the Historical load phase which contains the following and is integrated in no particular order:

1. DimDate
2. DimTime
3. StatusType
4. TaxRate
5. TradeType
6. DimBroker

7. DimCompany
8. DimCustomer
9. DimAccount
10. DimSecurity
11. DimTrade
12. FactCashBalances
13. FactMarketHistory
14. FactWatches
15. Industry
16. Financial
17. Prospect

System Setting and Data Generation

3.1 Hardware Specification

Some members of our group used Windows OS, MacOS and Ubuntu OS so we installed DuckDB directly on our computers . The Apple MacBook was the main machine for the benchmark implementation as all the results were obtained from it.

In order to carry out the benchmark experimentation, Apple Macbook and HP EliteBook 840 G5 were used. Table 3.1 shows the hardware specification of this computer:

Name	HP EliteBook 840 G5	Apple Macbook
Operating System	Ubuntu 22.04.3 LTS	MacOS 13.0
System Type	64 Bits	64 Bits
CPU	Intel® Core™ i7-8650U	Apple M1 Pro
CPU Frequency	1.90GHz × 8	3.2 GHz
RAM capacity	24 GB	16 GB
RAM type	DDR4	SDRAM
Hard disk capacity	512 GB	512 GB

Table 3.1: System Specifications.

3.2 Setup for Project

Our implementation was inspired by a partial implementation of TPC-DI in PostgreSQL [SG22].

3.2.1 DuckDB installation

We use `pip`¹ - the package installer for Python - to install the DuckDB.

```
1 pip install duckdb==0.9.1
```

Listing 3.1: DuckDB installation.

3.2.2 TPC-DI kit installation

We install the TPC-DS kit from the `tpcds- toolkit` downloaded from TPC website².

```
1 git clone https://github.com/gregrahn/tpcds-kit.git
2 cd tpcds-kit/tools
3 make OS=MACOS
```

Listing 3.2: TPC-DS kit installation.

The first command creates a folder `tpcds-kit` which is a replicate of the `tpcds-kit` repository. Then we move into the `tools` folder, and run the `make` file to build TPC-DS tools.

3.3 Database schema creation

The first step in the process is to connect to DuckDB and create a database:

```
1 import duckdb
```

Listing 3.3: Database setup

Then we use the scripts in `tpcds-kit/tools/tpcdi.sql` to create the database scheme.

¹<https://duckdb.org/docs/installation/index>

²<https://www.tpc.org/tpcdi>

3.4 Data generation

3.4.1 Jre and JVM Setting

Before generating data, we need to install Jre and JVM³. The version should be higher than 1.8.0.

```
1 C:\Users\Jintao1999>java -version java version "1.8.0_172"  
2 Java(TM) SE Runtime Environment(built 1.8.0_172-b11)  
3 Java HotSpot(TM) 64-Bit Server VM (built 25.172-b11, mixed mode)
```

Listing 3.4: Jre and JVM Setting.

3.4.2 Data generation

We generated data for evaluation with the TPC-DI program dsdgen. At this step, we are in the tpcds-kit/tools folder. We needed to Unzip and change the 'PDGF' to 'pdgf'. We found scale factors must be in range:3-2147. So we set the minimum scale factors as 3.

Then We ran the following command lines to generate data in 3 scale factors and store them in the generated-data folder.

```
1 run java -jar DIGen.jar -o D:\tpcdi\data -sf 3 # SF 3  
2 run java -jar DIGen.jar -o D:\tpcdi\data -sf 5 # SF 5  
3 run java -jar DIGen.jar -o D:\tpcdi\data -sf 7 # SF 7  
4 run java -jar DIGen.jar -o D:\tpcdi\data -sf 9 # SF 9
```

Listing 3.5: Data generation.

³<https://www.java.com/en/download/manual.jsp>

Benchmark Implementation Process

4.1 Implementing TPC-DI on DuckDB

Figure 4.1 shows a business process model depicting a brief rundown of the implementation the TPC-DI benchmark on DuckDB.

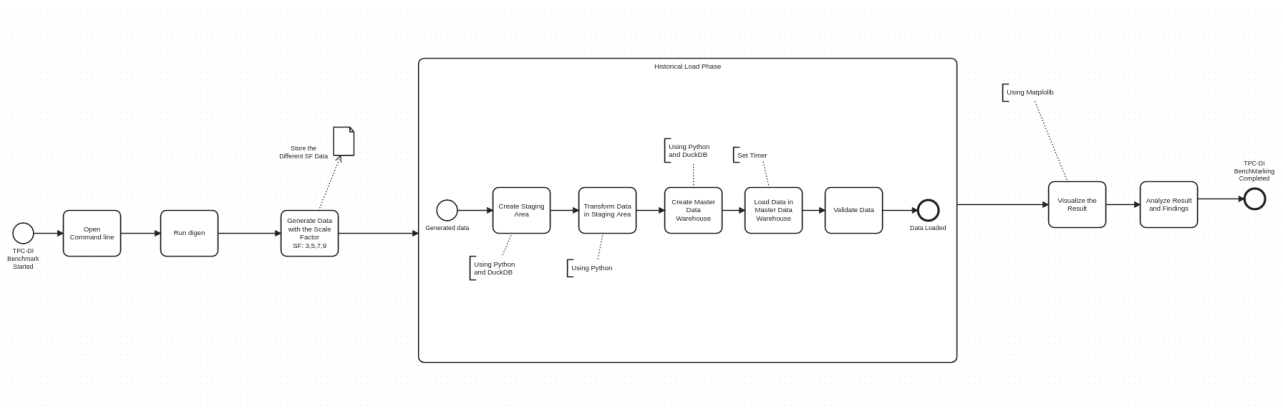


Figure 4.1: TPC-DI Benchmarking BPMN Diagram

From Figure 4.1, some of the process that was done include:

- Using DGen Utility from the TPC-DI tool kit to obtain load data for all our scale factor
- DuckDB is created for each selected SF with 2 schemas:
 - Staging Area: holds the initially transformed data
 - Master data warehouse
- Python Scripts was used for running the ETL tasks

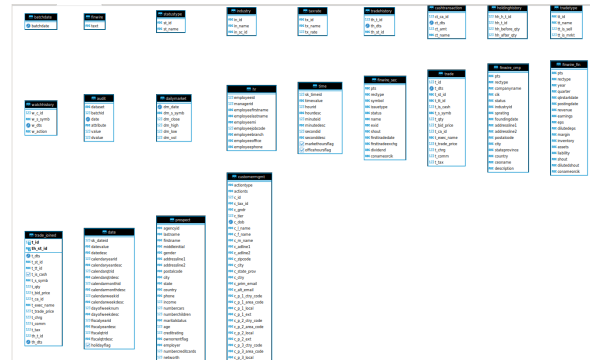
- Visualize all outputs.
- Benchmarking results are obtained and analyzed.

The implementation of the TPC-DI benchmarking were done in different stages, which include:

1. Create staging schema as seen in figure 4.2.
2. Create master schema as seen in figure 4.3.
3. Extract data and populate into staging schema.
4. Transform and load into master schema.

Schema Name:	staging				
Catalog:	tpcdi_3				
Tables	Table Name	Table Type	Catalog	Table Description	Schema
Views	audit	BASE TABLE	tpcdi_3		staging
Indexes	batchdate	BASE TABLE	tpcdi_3		staging
Procedures	cashtransaction	BASE TABLE	tpcdi_3		staging
Data Types	customermgmt	BASE TABLE	tpcdi_3		staging
	dailymarket	BASE TABLE	tpcdi_3		staging
	date	BASE TABLE	tpcdi_3		staging
	finwire	BASE TABLE	tpcdi_3		staging
	finwire_cmp	BASE TABLE	tpcdi_3		staging
	finwire_fin	BASE TABLE	tpcdi_3		staging
	finwire_sec	BASE TABLE	tpcdi_3		staging
	holdinghistory	BASE TABLE	tpcdi_3		staging
	hr	BASE TABLE	tpcdi_3		staging
	industry	BASE TABLE	tpcdi_3		staging
	prospect	BASE TABLE	tpcdi_3		staging
	statustype	BASE TABLE	tpcdi_3		staging
	taxrate	BASE TABLE	tpcdi_3		staging
	time	BASE TABLE	tpcdi_3		staging
	trade	BASE TABLE	tpcdi_3		staging
	trade_joined	BASE TABLE	tpcdi_3		staging
	tradehistory	BASE TABLE	tpcdi_3		staging
	tradetype	BASE TABLE	tpcdi_3		staging
	watchhistory	BASE TABLE	tpcdi_3		staging

(a) Staging Schema



(b) Staging ERD

Figure 4.2: Staging Database Schema

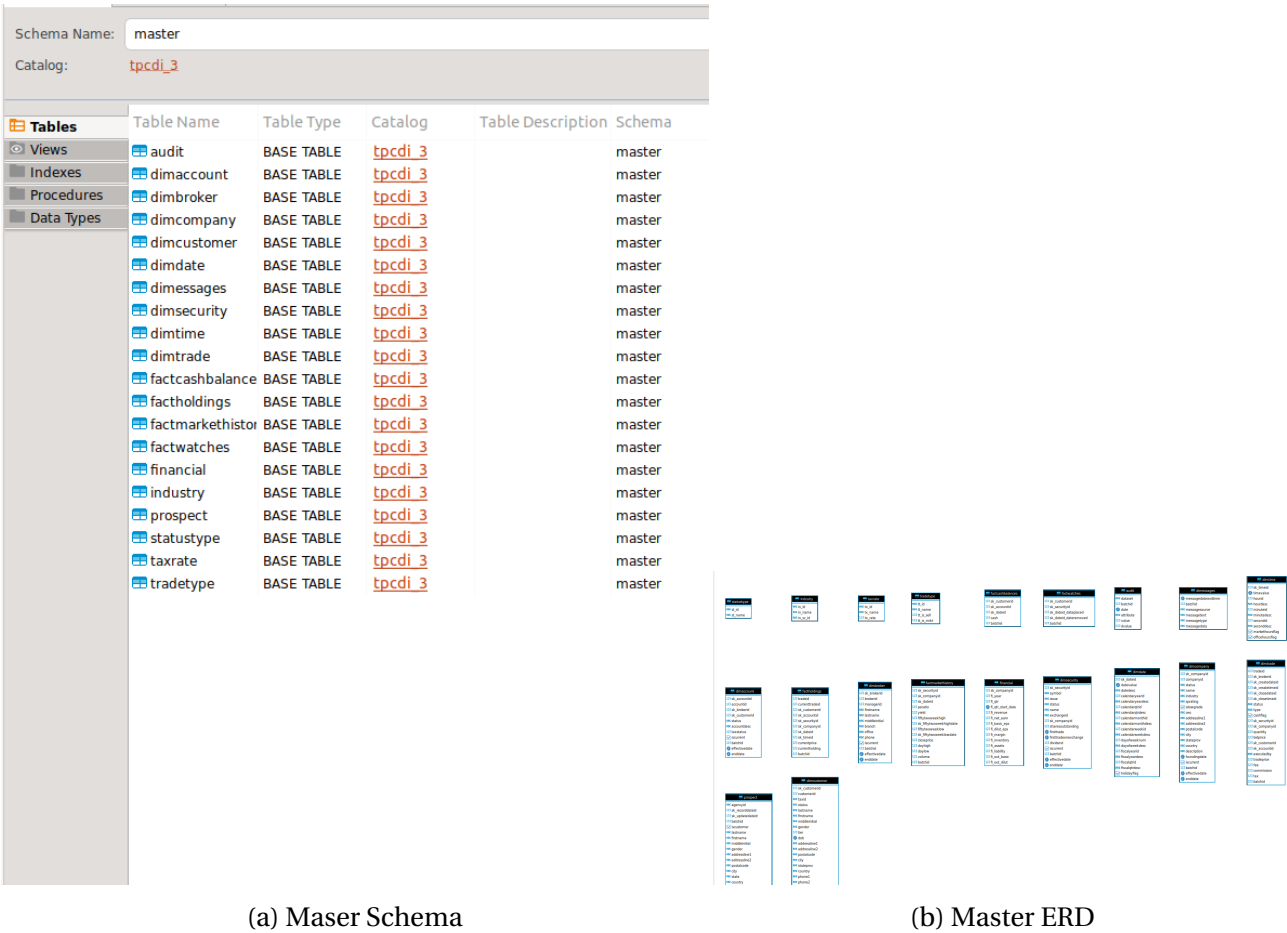


Figure 4.3: Master Data Warehouse

To make it more clear, we follow [tpc-di-benchmark¹](#) to divide the whole ETL process into 30 tasks. They are shown in table 4.1 .

4.2 Tasks

4.2.1 Task 1

Task 1 creates staging schema, storing all data from the source file. The staging schema includes 21 tables: batchdate, cashtransaction, customermgmt, dailymarket, date, finwire, finwire_cmp, finwire_sec, finwire_fin, holdinghistory, hr, industry, prospect, statustype, taxrate, time, tradehistory, trade, tradetype, watchhistory and audit.

Part of Task 1 SQL is shown in listing 4.1.

¹

¹<https://github.com/risg99/tpc-di-benchmark>

Number	Name	Type
1	Create staging schema	Create Schema
2	Load txt and csv sources to staging	Extract & Load
3	Load finwire source to staging	Extract & Load
4	Parse finwire and load to seperate tables	Extract & Load
5	Convert customer management source from xml to csv	Extract
6	Load customer management source to staging	Extract & Load
7	Create master schema	Create schema
8	Direct load master.tradetype	Load
9	Direct load master.statustype	Load
10	Direct load master.taxrate	Load
11	Direct load master.industry	Load
12	Transform & load master.dimdate	Transform & Load
13	Transform & load master.dimtime	Transform & Load
14	Transform & load master.dimcompany	Transform & Load
15	Load master.dimessages with alert from master.dimcompany	Transform & Load
16	Transform & load master.dimbroker	Transform & Load
17	Transform & load master.prospect	Transform & Load
18	Transform & load master.dimcustomer	Transform & Load
19	Load master.dimessages with alert from master.dimcustomer	Transform & Load
20	Update master.prospect	Transform & Load
21	Transform & load master.dimaccount	Transform & Load
22	Transform & load master.dimsecurity	Transform & Load
23	Transform & load master.dimtrade	Transform & Load
24	Load master.dimessages with alert from master.dimtrade	Transform & Load
25	Transform & load master.financial	Transform & Load
26	Transform & load master.factcashbalances	Transform & Load
27	Transform & load master.factholdings	Transform & Load
28	Transform & load master.factwatches	Transform & Load
29	Transform & load master.factmarkethistory	Transform & Load
30	Load master.dimessages with alert from master.factmarkethistory	Transform & Load

Table 4.1: Implementation tasks

```

2 drop table if exists staging.customermgmt;
3 create table staging.customermgmt(
4     --action element
5     actiontype char(9) check(actiontype in ('NEW','ADDACCT','UPDCUST','UPDACCT',
6         '','CLOSEACCT','INACT')),
7     actions varchar check(length(actions) > 0),
8     --action.customer element
9     c_id numeric(11) not null check(c_id >= 0),
10    c_tax_id char(20) check((actiontype = 'NEW' and length(c_tax_id) > 0) or (
11        actiontype != 'NEW')),
12    c_gndr char(1) check(length(c_gndr) > 0),
13    c_tier numeric(1) check(c_tier >= 0)
14    .....
15 );

```

Listing 4.1: Part of Task 1 SQL

4.2.2 Task 2-3

Task 2 and Task 3 simply load data in txt and csv to staging schema, with the delimiter as "|". Part of Task 1 SQL is shown in listing 4.2.

```

1 .....
2 COPY staging.batchdate FROM '/Users/wanglinhan/Desktop/BDMA/ULB/INFO-H419/
3     Project/P2/generated_data/Batch1/BatchDate.txt';
4 COPY staging.cashtransaction FROM '/Users/wanglinhan/Desktop/BDMA/ULB/INFO-
5     H419/Project/P2/generated_data/Batch1/CashTransaction.txt' delimiter '|';
6 COPY staging.dailymarket FROM '/Users/wanglinhan/Desktop/BDMA/ULB/INFO-H419/
7     Project/P2/generated_data/Batch1/DailyMarket.txt' delimiter '|';
8 .....

```

Listing 4.2: Part of Task 2 SQL

4.2.3 Task 4

Task 4 parse finwire and load to separate tables. Because there are three record types in FINWIRE files: 'CMP', 'SEC', and 'FIN', we need to separate them into three tables `finwire_cmp`, `finwire_sec`, `finwire_fin` with selected attributes.

As DuckDB doesn't support Postgresql, we need to define a Python function and register it to the connection. During the parsing process, the function need to connect to database several times and execute sql queries several times, instead of finishing every step in one function.

This feature may show the disadvantage of DuckDB to handle complicated queries and transformation.

Part of the customized function is shown in listing 4.3

```

1 def staging_finwire_split(conn):
2     conn.execute("DELETE FROM staging.finwire_cmp")
3     conn.execute("DELETE FROM staging.finwire_sec")
4     conn.execute("DELETE FROM staging.finwire_fin")
5
6     # Fetch data from staging.finwire
7     result = conn.execute("SELECT * FROM staging.finwire").fetchall()
8     for row in result:
9         rectype = row[0][15:18]
10        if rectype == 'CMP':
11            # Logic for CMP
12            exe_args = [row[0] for i in range(16)]
13            conn.execute("""
14                INSERT INTO staging.finwire_cmp
15                SELECT
16                nullif(trim(both from substring(?,1,15)), '') as pts,
17                nullif(trim(both from substring(?,16,3)), '') as rectype,
18                .....
19                LIMIT 1
20                """, exe_args)
21        elif rectype == 'SEC':
22            # Logic for SEC
23            exe_args = [row[0] for i in range(12)]
24            conn.execute("""
25                INSERT INTO staging.finwire_sec

```

```

26         SELECT
27         nullif(trim(both from substring(?,1,15)), '') as pts,
28         nullif(trim(both from substring(?,16,3)), '') as rectype,
29         .....
30         LIMIT 1
31         """ , exe_args)
32     elif rectype == 'FIN':
33         # Logic for FIN
34         exe_args = [row[0] for i in range(17)]
35         conn.execute("""
36             INSERT INTO staging.finwire_fin
37             SELECT
38             nullif(trim(both from substring(?,1,15)), '') as pts,
39             nullif(trim(both from substring(?,16,3)), '') as rectype,
40             .....
41             LIMIT 1
42             """ , exe_args)

```

Listing 4.3: Part of Task 4 SQL

4.2.4 Task 5-6

Task 5 converts customer management source from xml to csv using a Python function without connecting to database. The output csv has 36 columns.

Part of Task 5 function is shown in listing 4.4

```

1  def customermgmt_convert(filepath):
2      with open(filepath+'/CustomerMgmt.xml') as fd:
3          doc = xmltodict.parse(fd.read())
4          fd.close()
5
6      with open(filepath+"/CustomerData.json", "w") as outfile:
7          outfile.write(json.dumps(doc))
8          outfile.close()
9
10     f = open(filepath+'/CustomerData.json', 'r')
11
12     cust = json.load(f)
13     actions = cust['TPCDI:Actions']

```

```

14     action = actions['TPCDI:Action']
15     cust_df = pd.DataFrame(columns = np.arange(0, 36))
16
17
18     for a in action:
19
20         cust_row = {}
21
22         # action element
23         cust_row.update({0: [f"{a.get('@ActionType')}"]})
24         cust_row.update({1: [f"{a.get('@ActionTS')}"]})
25         .....
26
27         # append to dataframe
28         cust_df = pd.concat([cust_df, pd.DataFrame.from_dict(cust_row)],
29                             axis = 0)
30
31         cust_df.replace(to_replace = np.NaN, value = "", inplace = True)
32         cust_df.replace(to_replace = "None", value = "", inplace = True)
33         cust_df.to_csv(filepath+'CustomerMgmt.csv', index = False)
34     print('Customer Management data converted from XML to CSV')

```

Listing 4.4: Part of Task 5 SQL

Task 6 directly loads the output csv into staging schema.

4.2.5 Task 7

Task 7 creates the master schema, which is the data warehouse we need. The master schema includes 20 tables: tradetype, statustype, taxrate, industry, dimdate, dimtime, dimcompany, dimbroker, dimcustomer, dimaccount, dimsecurity, dimtrade, financial, factcashbalances, factholdings, factmarkethistory, factwatches, prospect, audit and dimessages.

4.2.6 Task 8-11

Task 8-11 direct load master.tradetype, master.statustype, master.taxrate, master.industry.

In these 4 tables, the data is the same in staging schema and master schema.

Part of Task 8 SQL is shown in listing 4.5

```

1 delete from master.tradetype;
2 insert into master.tradetype
3   select * from staging.tradetype;

```

Listing 4.5: Part of Task 8 SQL

4.2.7 Task 12-15

Task 12 loads date related attributes from staging.date to master.dimdate.

Task 13 loads time related attributes from staging.time to master.dimtime.

Task 14 loads company related attributes from staging.finwire_cmp, staging.statustype and staging.industry to master.dimcompany. At the same time, it classifies whether some company is low-grade.

Task 14 SQL is shown in listing 4.6.

```

1 -- dimcompany
2 delete from master.dimcompany;
3 insert into master.dimcompany
4   select
5     row_number() over(order by cik) as sk,
6     cik::numeric(11) as companyid,
7     s.st_name as status,
8     companyname as name,
9     i.in_name as industry,
10    (CASE
11      WHEN sprating not in ('AAA','AA','AA+','AA-','A','A+','A-','BBB','BBB+',
12        'BBB-','BB','BB+','BB-','B','B+','B-','CCC','CCC+','CCC-','CC','C','D')
13      THEN null
14      ELSE f.sprating END) as sprating,
15    (CASE
16      WHEN sprating not in ('AAA','AA','AA+','AA-','A','A+','A-','BBB','BBB+',
17        'BBB-','BB','BB+','BB-','B','B+','B-','CCC','CCC+','CCC-','CC','C','D')
18      THEN null

```



```

17     WHEN f.sprating like 'A%' or f.sprating like 'BBB%'
18     THEN false
19     ELSE
20         true
21 END) as islowgrade,
22 ceoname as ceo,
23 addressline1,
24 addressline2,
25 postalcode,
26 city,
27 stateprovince,
28 country,
29 description,
30 CAST(SUBSTRING(foundingdate, 1, 4) || '-' || SUBSTRING(foundingdate, 5, 2)
    || '-' || SUBSTRING(foundingdate, 7, 2) AS DATE),
31 case when lead( (select batchdate from staging.batchdate) ) over (
    partition by cik order by pts asc ) is null then true else false end as
    iscurrent,
32 1 as batchid,
33 CAST(SUBSTRING(left(f.pts, 8), 1, 4) || '-' || SUBSTRING(left(f.pts, 8),
    5, 2) || '-' || SUBSTRING(left(f.pts, 8), 7, 2) AS DATE) as effectivedate
    ,
34 '9999-12-31'::date as enddate
35 from
36     staging.finwire_cmp f,
37     staging.statustype s,
38     staging.industry i
39 where
40     f.status = s.st_id
41 and f.industryid = i.in_id;

```

Listing 4.6: Task 14 SQL

Task 15 inserts master.dimessages from dimcompany results. A record will be inserted in the DImessages table if a company's SPRating is not one of the valid values for Standard and Poor long-term credit-ratings. The valid values are: AAA, AA[+/-], A[+/-], BBB[+/-], BB[+/-], B[+/-], CCC[+/-], CC, C, D.

Task 15 SQL is shown in listing 4.7.

```

1 -- dimessages alert for dimcompany

```

```

2 delete from master.dimessages;
3 insert into master.dimessages
4     select
5     now(),
6     1 as batchid,
7     'DimCompany' as messagesource,
8     'Invalid SPRating' as messagetext,
9     'Alert' as messagetype,
10    'CO_ID = ' || cik::varchar || ', CO_SP_RATE = ' || sprating::varchar
11    from staging.finwire_cmp
12    where sprating not in ('AAA','AA','AA+','AA-','A','A+','A-','BBB','BBB+','
        BBB-','BB','BB+','BB-','B','B+','B-','CCC','CCC+','CCC-','CC','C','D');

```

Listing 4.7: Task 15 SQL

4.2.8 Task 16

Task 16 transforms and loads master.dimbroker.

Data for DimBroker comes from the HR extract file HR.csv. Those employees from the HR file that are brokers (as indicated by the EmployeeJobCode) will have data copied to the Broker table.

Task 16 SQL is shown in listing 4.8.

```

1 -- dimbroker
2 delete from master.dimbroker;
3 insert into master.dimbroker
4     select
5     row_number() over(order by employeeid) as sk,
6     employeeid as brokerid,
7     managerid,
8     employeefirstname,
9     employeelastname,
10    employeemi,
11    employeebranch,
12    employeeoffice,
13    employeephone,
14    true as iscurrent,
15    1 as batchid,
16    (select min(datevalue) FROM master.dimdate) as effective date,

```

```

17 '9999-12-31'::date as enddate
18 from staging.hr
19 where employeejobcode = 314;

```

Listing 4.8: Task 16 SQL

4.2.9 Task 17

Task 17 transform and load master.prospect.

Prospect table data is obtained from the Prospect file.

In this process, MarketingNameplate is set based on other fields. The tags are defined as 'High-Value', 'Expenses', 'Boomer', 'MoneyAlert', 'Spender' and 'Inherited' according to the documentation².

Because DuckDB doesn't support btrim function in Postgresql, so we need to define custom *btrim_equivalent* function.

```

1 def btrim_equivalent(input_string, trim_characters):
2     start_index = len(input_string) - len(input_string.lstrip(
3         trim_characters))
4     end_index = len(input_string.rstrip(trim_characters)) - 1
5     return input_string[start_index:end_index + 1]
6
7 con = duckdb.connect('tpcdi_'+str(scale_factor)+'.db')
8 con.create_function("btrim_equivalent", btrim_equivalent, ['VARCHAR', '
9     VARCHAR'], 'VARCHAR')

```

Listing 4.9: btrim_equivalent

Part of Task 17 SQL is shown in listing 4.10.

```

1 -- prospect part 1
2 delete from master.prospect;
3 insert into master.prospect
4     with date_record_id as (
5         select
6             dd.sk_dateid

```

²https://www.tpc.org/TPC_Documents_Current_Versions/pdf/TPC-DI_v1.1.0.pdf

```

7      from master.dimdate dd
8      inner join staging.batchdate bd
9          on dd.datevalue = bd.batchdate
10 )
11
12 select
13     p.agencyid
14 , dri.sk_dateid
15 , dri.sk_dateid
16     .....
17 , nullif(btrim_equivalent(btrim_equivalent(btrim_equivalent(
18     btrim_equivalent(btrim_equivalent(
19     case
20     when p.networth > 1000000 or p.income > 200000
21     then 'HighValue'
22     else ''
23     end
24     || '+' ||
25     case
26     when p.numberchildren > 3 or p.numbercreditcards > 5
27     then 'Expenses'
28     else ''
29     end
30     , '+')
31     || '+' ||
32     case
33     when p.age > 45
34     then 'Boomer'
35     else ''
36     end
37     , '+')
38     || '+' ||
39     case
40     when p.income < 50000 or p.creditrating < 600 or p.networth < 100000
41     then 'MoneyAlert'
42     else ''
43     end
44     , '+')
45     || '+' ||
46     case

```

```
46     when p.numbercars > 3 or p.numbercreditcards > 7
47     then 'Spender'
48     else ''
49     end
50     , '+' )
51     || '+' ||
52     case
53     when p.age < 25 and p.networth > 1000000
54     then 'Inherited'
55     else ''
56     end
57     , '+' ), '' )
58 from staging.prospect p
59 cross join date_record_id dri;
```

Listing 4.10: Task 17 SQL

4.2.10 Task 18-20

Task 18 transform and load master.dimcustomer.

DimCustomer data is obtained from the data file CustomerMgmt.xml.

Customer has six ActionType.

- **NEW.** A new customer. A new customer is always created with 1 or more new accounts.
- **ADDACCT.** One or more new accounts for an existing customer. This does not require any change to DimCustomer.
- **UPDACCT.** Updates to the information in one or more existing accounts. No change to DimCustomer is required.
- **UPDCUST.** One or more updates to existing customer's information. Only the identifying data and updated property values are supplied in the Source Data.
- **CLOSEACCT.** Close one or more existing accounts. This does not require any change to DimCustomer.

- **INACT.** Make an existing customer and that customer's currently active accounts inactive. No specific account information is supplied in the Source Data.

Changes are made based on different ActionTypes.

Task 19 load master.dimmessages with alert from master.dimcustomer.

A record will be inserted in the DImessages table if a customer's Tier is not one of the valid values (1,2,3). The MessageSource is "DimCustomer", the MessageType is "Alert" and the MessageText is "Invalid customer tier". The MessageData field is "C_ID = " followed by the natural key value of the record, then ", C_TIER = " and the C_TIER value.

A record will be reported in the DImessages table if a customer's DOB is invalid. A customer's DOB is invalid if $DOB < \text{Batch Date} - 100 \text{ years}$ or $DOB > \text{Batch Date}$ (customer is over 100 years old or born in the future). The MessageSource is "DimCustomer", the MessageType is "Alert" and the MessageText is "DOB out of range". The MessageData field is "C_ID = " followed by the natural key value of the record, then ", C_DOB = " and the C_DOB value.

Task 19 SQL is shown in listing 4.11.

```
1  -- dimessages alert for dimcustomer
2  insert into master.dimmessages
3      select
4          now()
5          , 1
6          , 'DimCustomer'
7          , 'Invalid customer tier'
8          , 'Alert'
9          , 'C_ID = ' || customerid || ', C_TIER = ' || tier
10     from master.dimcustomer
11     where tier not between 1 and 3;
12
13 insert into master.dimmessages
14     select
15         now()
16         , 1
17         , 'DimCustomer'
18         , 'DOB out of range'
19         , 'Alert'
```

```

20 , 'C_ID = ' || customerid || ', C_DOB = ' || dob
21 from master.dimcustomer
22 where dob < (select * from staging.batchdate) - interval '100 years'
23 or dob > (select * from staging.batchdate);

```

Listing 4.11: Task 19 SQL

Task 20 update master.prospect.

If a customer is 'ACTIVE' in master.dimcustomer, the attribute 'iscustomer' should be true.

Task 20 SQL is shown in listing 4.12.

```

1  -- update prospect
2  with current_active_customer as (
3      select p.*
4      from master.prospect p
5      inner join master.dimcustomer c
6      on upper(c.lastname) = upper(p.lastname)
7      and upper(c.firstname) = upper(p.firstname)
8      and upper(c.addressline1) = upper(p.addressline1)
9      and upper(c.addressline2) = upper(p.addressline2)
10     and upper(c.postalcode) = upper(p.postalcode)
11     where c.status = 'ACTIVE'
12     and c.iscurrent = true
13 )
14
15 update master.prospect
16     set iscustomer = true
17     where lastname in (select lastname from current_active_customer)
18     and firstname in (select firstname from current_active_customer)
19     and addressline1 in (select addressline1 from current_active_customer)
20     and addressline2 in (select addressline2 from current_active_customer)
21     and postalcode in (select postalcode from current_active_customer);

```

Listing 4.12: Task 20 SQL

4.2.11 Task 21

Task 21 transform and load master.dimaccount.

DimAccount data is obtained from the data file CustomerMgmt.xml.

Similar to Task 18. Task 21 also do different processing based on different ActionTypes: 'NEW', 'ADDACCT', 'UPDACCT', 'UPDCUST', 'CLOSEACCT' and 'INACT'.

There are no messages written to the DImessages table by this transformation.

4.2.12 Task 22

Task 22 transform and load master.dimsecurity.

DimSecurity data is obtained from the FINWIRE files. All FINWIREyyyyQq files are processed in ascending year and quarter order, and records of type SEC are used.

```

1  -- master.dimsecurity transform and load
2  truncate table master.dimsecurity;
3  insert into master.dimsecurity
4      select
5          row_number() over() as sk_securityid,
6          symbol,
7          issuetype as issue,
8          s.st_name as status,
9          f.name,
10         exid as exchangeid,
11         c.sk_companyid as sk_companyid,
12         shout::numeric(12) as sharesoutstanding,
13         CAST(SUBSTRING(left(firsttradedate, 8), 1, 4) || '-' || SUBSTRING(left(
firsttradedate, 8), 5, 2) || '-' || SUBSTRING(left(firsttradedate, 8), 7,
2) AS DATE),
14         CAST(SUBSTRING(left(firsttradeexchg, 8), 1, 4) || '-' || SUBSTRING(left(
firsttradeexchg, 8), 5, 2) || '-' || SUBSTRING(left(firsttradeexchg, 8),
7, 2) AS DATE),
15         dividend::numeric(10,2),
16         case
17             when lead( (select batchdate from staging.batchdate) ) over (
partition by symbol order by pts asc ) is null
18             then true
19             else false
20         end as iscurrent,

```



```

21 1 as batchid,
22 CAST(SUBSTRING(left(f.pts, 8), 1, 4) || '-' || SUBSTRING(left(f.pts, 8),
    5, 2) || '-' || SUBSTRING(left(f.pts, 8), 7, 2) AS DATE),
23 '9999-12-31'::date as enddate
24 from staging.finwire_sec f,
25 staging.statustype s,
26 master.dimcompany c
27 where f.status = s.st_id
28 and ((ltrim(f.conameorcik, '0') = c.companyid::varchar)
29      or (f.conameorcik = c.name))
30 and CAST(SUBSTRING(left(f.pts, 8), 1, 4) || '-' || SUBSTRING(left(f.pts,
    8), 5, 2) || '-' || SUBSTRING(left(f.pts, 8), 7, 2) AS DATE) >= c.
    effectivedate
31 and CAST(SUBSTRING(left(f.pts, 8), 1, 4) || '-' || SUBSTRING(left(f.pts,
    8), 5, 2) || '-' || SUBSTRING(left(f.pts, 8), 7, 2) AS DATE) < c.enddate;

```

Listing 4.13: Task 20 SQL

4.2.13 Task 23-24

Task 23 transform and load master.dimtrade.

DimTrade data is obtained from the Trade.txt and TradeHistory.txt files. The incoming files may be thought of as logically joined on the T_ID field. When a T_ID encountered does not match a TradeID from the DimTrade table, a new DimTrade record is inserted. When a T_ID is encountered that matches an existing TradeID in the DimTrade table, the DimTrade record is updated.

Because DuckDB is using tmp files to store temporary tables, too much tmp files occupies more space than space left on device, Task 23 is divided into two steps. First, join staging.trade and staging.tradehistory. Second, join and select values.

Task 23 SQL is shown in listing 4.14.

```

1 --step 1
2 insert into staging.trade_joined
3 select * from staging.trade t inner join staging.tradehistory th on t.t_id =
    th.th_t_id;
4
5 --step 2

```

```

6 delete from master.dimtrade;
7 insert into master.dimtrade (sk_createdateid, sk_createtimeid,
    sk_closedateid, sk_closetimeid, tradeid, cashflag,
8     quantity, bidprice, executedby, tradeprice, fee, commission, tax,
    status, type, sk_securityid, sk_companyid,
9     sk_accountid, sk_customerid, sk_brokerid, batchid)
10 SELECT
11 Case When (sth.th_st_id = 'SBMT' and (sth.t_tt_id = 'TMB' or sth.t_tt_id = '
    TMS')) or (sth.th_st_id = 'PNDG')
12     then D.sk_dateid else NULL
13 end,
14 Case When (sth.th_st_id = 'SBMT' and (sth.t_tt_id = 'TMB' or sth.t_tt_id = '
    TMS')) or (sth.th_st_id = 'PNDG')
15     then T.sk_timeid else NULL
16 end,
17 case when (sth.th_st_id = 'CMPT' or sth.th_st_id = 'CNCL') then D.sk_dateid
    else NULL
18 end,
19 case when (sth.th_st_id = 'CMPT' or sth.th_st_id = 'CNCL') then T.sk_timeid
    else NULL
20 end,
21 sth.t_id, sth.t_is_cash, sth.t_qty, sth.t_bid_price, sth.t_exec_name, sth.
    t_trade_price, sth.t_chrg,
22 sth.t_comm, sth.t_tax, st.st_name, tt.tt_name,
23 case when (CAST(sth.th_dts as DATE) >= ds.EffectiveDate) and (CAST(sth.
    th_dts as DATE) <= ds.EndDate)
24     then ds.SK_SecurityID END,
25 case when (CAST(sth.th_dts as DATE) >= ds.EffectiveDate) and (CAST(sth.
    th_dts as DATE) <= ds.EndDate)
26     then ds.SK_CompanyID END,
27 case when (CAST(sth.th_dts as DATE) >= ds.EffectiveDate) and (CAST(sth.
    th_dts as DATE) <= ds.EndDate)
28     then da.SK_AccountID END,
29 case when (CAST(sth.th_dts as DATE) >= ds.EffectiveDate) and (CAST(sth.
    th_dts as DATE) <= ds.EndDate)
30     then da.SK_CustomerID END,
31 case when (CAST(sth.th_dts as DATE) >= ds.EffectiveDate) and (CAST(sth.
    th_dts as DATE) <= ds.EndDate)
32     then da.SK_BrokerID END,
33 1

```

```

34 from staging.trade_joined sth
35 join staging.statustype st on sth.t_st_id = st.st_id
36 join staging.tradetype tt on sth.t_tt_id = tt.tt_id
37 join master.dimsecurity ds on sth.t_s_symb = ds.symbol
38 join master.dimaccount da on sth.t_ca_id = da.accountid
39 JOIN master.dimdate D on CAST(sth.th_dts as DATE) = D.datevalue
40 join master.dimtime T on CAST(sth.th_dts AS TIME) = T.timevalue

```

Listing 4.14: Task 23 SQL

Task24 load master.dimessages with alert from master.dimtrade.

A record will be reported in the DImessages table if a trade's Commission is not null and exceeds TradePrice * Quantity. The MessageSource is "DimTrade", the MessageType is "Alert" and the MessageText is "Invalid trade commission". The MessageData field is T_ID = followed by the key value of the record, then T_COMM = and the T_COMM value.

Task 24 SQL is shown in listing 4.15.

```

1  -- dimessages alert for dimtrade
2  insert into master.dimessages
3      select
4          now()
5          , 1
6          , 'DimTrade'
7          , 'Invalid trade commission'
8          , 'Alert'
9          , 'T_ID = ' || tradeid || ', T_COMM = ' || commission
10 from master.dimtrade
11 where commission is not null
12 and commission > (tradeprice * quantity);
13
14 insert into master.dimessages
15 select
16     now()
17     , 1
18     , 'DimTrade'
19     , 'Invalid trade fee'
20     , 'Alert'
21     , 'T_ID = ' || tradeid || ', T_CHRG = ' || fee
22 from master.dimtrade

```

```

23  where fee is not null
24  and fee > (tradeprice * quantity);

```

Listing 4.15: Task 24 SQL

4.2.14 Task 25

Task25 transform and load master.financial.

Financial data is obtained from the FINWIRE files. All FINWIREyyyyQq files are processed in ascending year and quarter order, and records of type FIN are used. The surrogate key of the associated company must be obtained for the Company dimension reference.

Task 25 SQL is shown in listing 4.16.

```

1  -- financial
2  delete from master.financial;
3  insert into master.financial
4  select
5      c.sk_companyid as sk_companyid,
6      year::numeric as fi_year,
7      quarter::numeric as fi_qtr,
8      CAST(SUBSTRING(qtrstartdate, 1, 4) || '-' || SUBSTRING(qtrstartdate, 5,
9      2) || '-' || SUBSTRING(qtrstartdate, 7, 2) AS DATE) as fi_qtr_start_date,
10     revenue::numeric as fi_revenue,
11     earnings::numeric as fi_net_earn,
12     eps::numeric as fi_basic_eps,
13     dilutedeps::numeric as fi_dilut_eps,
14     margin::numeric as fi_margin,
15     inventory::numeric as fi_inventory,
16     assets::numeric as fi_assets,
17     liability::numeric as fi_liability,
18     shout::numeric as fi_out_basic,
19     dilutedshout::numeric as fi_out_dilut
20 from staging.finwire_fin f,
21     master.dimcompany c
22 where ((f.conameorcik = c.companyid::varchar)
23        or (f.conameorcik = c.name))
24 and CAST(SUBSTRING(left(pts, 8), 1, 4) || '-' || SUBSTRING(left(pts, 8),

```

```

24      5, 2) || '-' || SUBSTRING(left(pts, 8), 7, 2) AS DATE) >= c.effectivedate
      and CAST(SUBSTRING(left(pts, 8), 1, 4) || '-' || SUBSTRING(left(pts, 8),
      5, 2) || '-' || SUBSTRING(left(pts, 8), 7, 2) AS DATE) < c.enddate;

```

Listing 4.16: Task 25 SQL

4.2.15 Task 26

Task26 transform and load master.factcashbalances.

FactCashBalances data is obtained from the data file CashTransaction.txt. The net effect of all cash transactions for a given account on a given day is totaled, and only a single record is generated per account that had changes per day.

Task 26 SQL is shown in listing 4.17.

```

1  -- factcashbalances
2  delete from master.factcashbalances;
3  insert into master.factcashbalances
4      with agg as (
5          select
6              a.sk_customerid as sk_customerid,
7              a.sk_accountid as sk_accountid,
8              d.sk_dateid as sk_dateid,
9              sum(ct_amt) as ct_amt_day
10             from staging.cashtransaction c,
11                  master.dimaccount a,
12                  master.dimdate d
13             where c.ct_ca_id = a.accountid
14                   and ct_dts::date >= a.effectivedate
15                   and ct_dts::date < a.enddate
16                   and ct_dts::date = d.datevalue
17             group by
18                 a.sk_customerid,
19                 a.sk_accountid,
20                 d.sk_dateid
21         )
22     , final_output as (
23         select

```

```

24     sk_customerid,
25     sk_accountid,
26     sk_dateid,
27     sum(ct_amt_day) over(partition by sk_accountid order by sk_dateid rows
    between unbounded preceding and current row) as cash,
28     1 as batchid
29 from agg
30 )
31 select * from final_output;

```

Listing 4.17: Task 26 SQL

4.2.16 Task 27

Task27 transform and load master.factholdings.

Data for FactHoldings comes from the HoldingHistory.txt file and the DimTrade table. The quantity and price values reflect the holdings for a particular security after the most recent trade. The customer can have a positive or negative position (Quantity) as a result of a trade.

Task 27 SQL is shown in listing 4.18.

```

1  -- factholdings
2  delete from master.factholdings;
3  insert into master.factholdings
4      select
5          h.hh_h_t_id as tradeid,
6          t.tradeid as currenttradeid,
7          t.sk_customerid as sk_customerid,
8          t.sk_accountid as sk_accountid,
9          t.sk_securityid as sk_securityid,
10         t.sk_companyid as sk_companyid,
11         t.sk_closedateid as sk_dateid,
12         t.sk_closetimeid as sk_timeid,
13         t.tradeprice as currentprice,
14         h.hh_after_qty as currentholding,
15         1 as batchid
16 from staging.holdinghistory h,
17         master.dimtrade t

```

```
18 where h.hh_t_id = t.tradeid;
```

Listing 4.18: Task 27 SQL

4.2.17 Task 28

Task28 transform and load master.factwatches.

Data for FactWatches comes from the WatchHistory.txt file. Surrogate keys must be obtained for the Customer, Security and Date dimension references. The date keys show the dates the watch was set and removed.

Task 28 SQL is shown in listing 4.19.

```
1 -- factwatches
2 delete from master.factwatches;
3 insert into master.factwatches
4   with watches as (
5     select w1.w_c_id,
6     TRIM(w1.w_s_symb) as w_s_symb,
7     w1.w_dts::date as dateplaced,
8     w2.w_dts::date as dateremoved
9     from staging.watchhistory w1,
10          staging.watchhistory w2
11     where w1.w_c_id = w2.w_c_id
12     and w1.w_s_symb = w2.w_s_symb
13     and w1.w_action = 'ACTV'
14     and w2.w_action = 'CNCL'
15   )
16
17 select
18   c.sk_customerid as sk_customerid,
19   s.sk_securityid as sk_securityid,
20   CAST(EXTRACT(YEAR FROM w.dateplaced) * 10000 + EXTRACT(MONTH FROM w.
dateplaced) * 100 + EXTRACT(DAY FROM w.dateplaced) AS NUMERIC) as
sk_dateid_dateplaced,
21   CAST(EXTRACT(YEAR FROM w.dateremoved) * 10000 + EXTRACT(MONTH FROM w.
dateremoved) * 100 + EXTRACT(DAY FROM w.dateremoved) AS NUMERIC) as
sk_dateid_dateremoved,
```

```

22      1 as batchid
23  from watches w,
24      master.dimcustomer c,
25      master.dimsecurity s,
26      master.dimdate d1,
27      master.dimdate d2
28  where w.w_c_id = c.customerid
29  and w.w_s_symb = s.symbol
30  and w.dateplaced = d1.datevalue
31  and w.dateremoved = d2.datevalue;

```

Listing 4.19: Task 28 SQL

4.2.18 Task 29-30

Task29 transform and load master.factmarkethistory.

FactMarketHistory data is primarily obtained from the file DailyMarket.txt.

Part of Task 29 SQL is shown in listing 4.20.

```

1  -- factmarkethistory
2  delete from master.factmarkethistory;
3  insert into master.factmarkethistory
4  with market_dates_daily as (
5      select
6          dm.dm_s_symb
7          , dm.dm_date
8          , .....
9      from staging.dailymarket dm
10     inner join master.dimdate dd
11         on dm.dm_date = dd.datevalue
12     order by
13         dm.dm_s_symb
14         , dm.dm_date desc
15 )
16
17 , high_low as (
18     select
19         dm_s_symb

```



```

20         .....
21     , max(dm_high) over(partition by dm_s_symb order by dm_date rows between
22       363 preceding and current row) as fiftytwoweekhigh
23     , min(dm_low) over(partition by dm_s_symb order by dm_date rows between
24       363 preceding and current row) as fiftytwoweeklow
25     from market_dates_daily
26 )
27
28 , high_date as (
29     select
30         hl.dm_s_symb
31         .....
32     , max(mdd.dm_date) as sk_fiftytwoweekhighdate
33     from high_low hl
34     inner join market_dates_daily mdd
35         on hl.dm_s_symb = mdd.dm_s_symb
36         and hl.fiftytwoweekhigh = mdd.dm_high
37         and mdd.dm_date <= hl.dm_date
38         and mdd.dm_date >= hl.dm_date - interval '52 weeks'
39     group by
40         hl.dm_s_symb
41         , hl.dm_date
42         .....
43 )
44
45 , low_date as (
46     select
47         hl.dm_s_symb
48         , hl.dm_date
49         .....
50     from high_date hl
51     inner join market_dates_daily mdd
52         on hl.dm_s_symb = mdd.dm_s_symb
53         and hl.fiftytwoweeklow = mdd.dm_low
54         and mdd.dm_date <= hl.dm_date
55         and mdd.dm_date >= hl.dm_date - interval '52 weeks'
56     group by
57         hl.dm_s_symb
58         , hl.dm_date
59         , hl.dm_close

```

```

58         .....
59     )
60
61     , quarters as (
62         select
63             f.sk_companyid
64             .....
65     )
66
67     , final_output as (
68         select
69             s.sk_securityid
70             , s.sk_companyid
71             , CAST(EXTRACT(YEAR FROM ld.dm_date) * 10000 + EXTRACT(MONTH FROM ld.
dm_date) * 100 + EXTRACT(DAY FROM ld.dm_date) AS NUMERIC) as sk_dateid
72             , case
73                 when q.eps_qtr_sum != 0 and q.eps_qtr_sum is not null
74                 then (ld.dm_close / q.eps_qtr_sum)::numeric(10, 2)
75                 else null
76             end as peratio
77             , case
78                 when ld.dm_close != 0
79                 then round((s.dividend / ld.dm_close) * 100, 2)
80                 else null
81             end as yield
82             , ld.fiftytwoweekhigh
83             , CAST(EXTRACT(YEAR FROM ld.sk_fiftytwoweekhighdate) * 10000 + EXTRACT(
MONTH FROM ld.sk_fiftytwoweekhighdate) * 100 + EXTRACT(DAY FROM ld.
sk_fiftytwoweekhighdate) AS NUMERIC) as sk_fiftytwoweekhighdate
84             , ld.fiftytwoweeklow
85             , CAST(EXTRACT(YEAR FROM ld.sk_fiftytwoweeklowdate) * 10000 + EXTRACT(
MONTH FROM ld.sk_fiftytwoweeklowdate) * 100 + EXTRACT(DAY FROM ld.
sk_fiftytwoweeklowdate) AS NUMERIC) as sk_fiftytwoweeklowdate
86             , ld.dm_close as closeprice
87             , ld.dm_high as dayhigh
88             , ld.dm_low as daylow
89             , ld.dm_vol as volume
90             , 1 as batchid
91         from low_date ld
92         inner join master.dimsecurity s

```

```

93         on ld.dm_s_symb = s.symbol
94         and ld.dm_date >= s.effectivedate
95         and ld.dm_date < s.enddate
96     inner join quarters q
97         on s.sk_companyid = q.sk_companyid
98         and q.fi_qtr_start_date <= ld.dm_date
99         and q.next_qtr_start > ld.dm_date
100 )
101 select * from final_output;

```

Listing 4.20: Part of Task 29 SQL

Task30 load master.dimessages with alert from master.factmarkethistory. A record will be reported in the DImessages table if there are no earnings found for a company. The MessageSource is “Fact-MarketHistory”, the MessageType is “Alert” and the MessageText is “No earnings for company”. The MessageData field is DM_S_SYMB = followed by the DM_S_SYMB value of the record.

Task 30 SQL is shown in listing 4.21

```

1  -- dimessages alert for factmarkethistory
2  insert into master.dimessages
3      select
4          now()
5          , 1
6          , 'FactMarketHistory'
7          , 'No earnings for company'
8          , 'Alert'
9          , 'DM_S_SYMB = ' || s.symbol
10 from master.factmarkethistory fmh
11 inner join master.dimsecurity s
12     on fmh.sk_securityid = s.sk_securityid
13 where fmh.peratio is null
14 or fmh.peratio = 0;

```

Listing 4.21: Task 30 SQL

As mentioned in the previous chapter, the benchmark was performed on a local machine, with a total of four different scale factors (3, 5, 7 and 9) and our scale factors increased linearly. In this chapter, we will be focusing on discussing the test result, the evolution of the loading times as the scale factor of the database increases.

5.1 Total Performance

We ran the scripts for 4 scale factors and show the results in Figure 5.1 and Figure 5.2 . Execution times were measured in seconds. The charts provided the execution time for 4 scale factors. The bar chart presents the information about the Load Test Time in seconds for four scale factors. As can be seen from the chart there is a significant difference between Execution Time for each scale factor. The time increases from 268.638778 seconds for SF 3 factor up to 2365.027075 seconds for SF 9. The line chart also shows that the execution time increased almost linearly as the data was increased from sf3 to sf7. But the growth trend suddenly increases from sf7 to sf9.

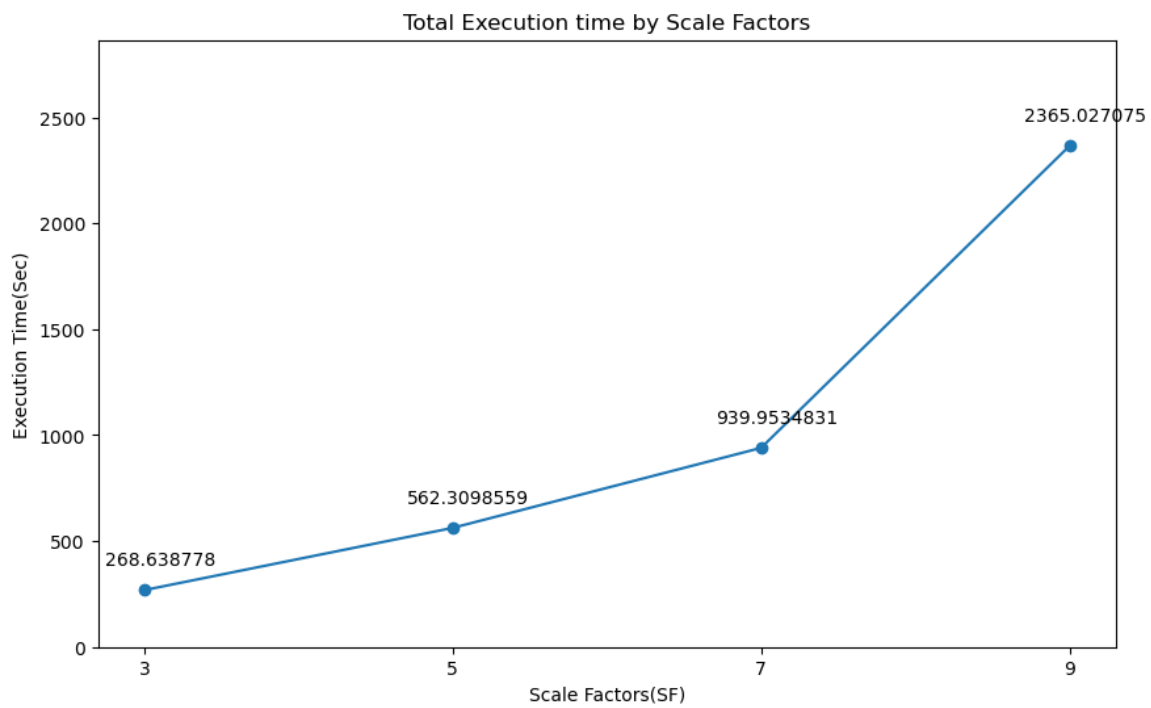


Figure 5.1: Line Chart of Total execution time

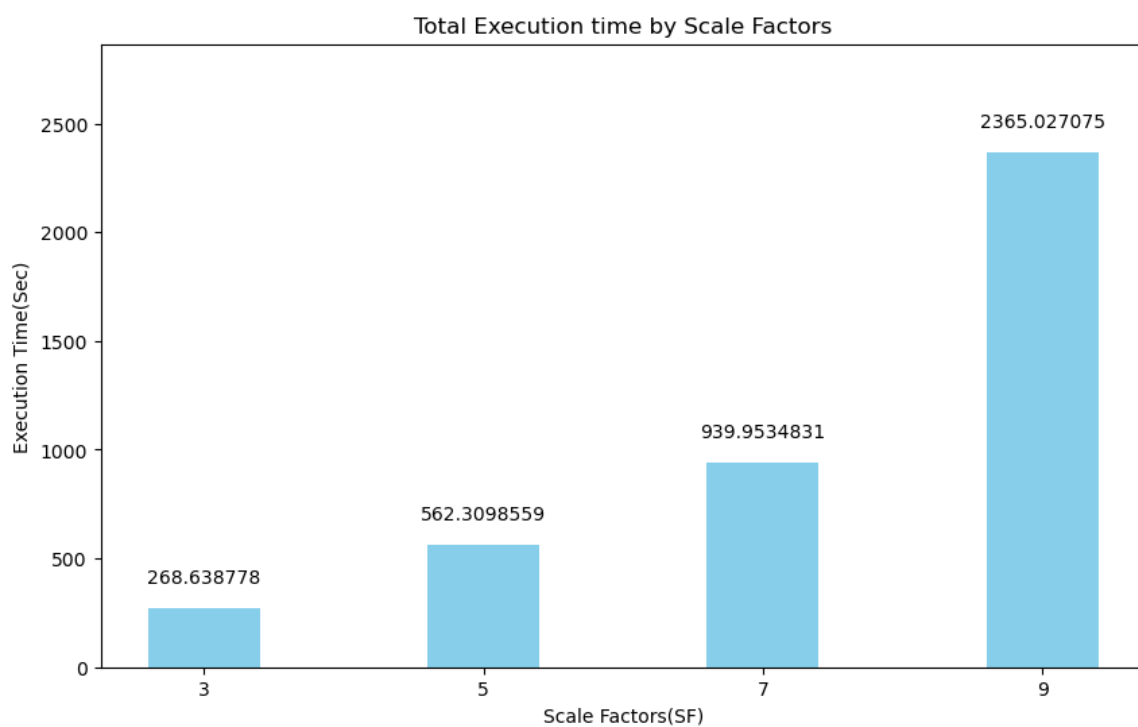


Figure 5.2: Bar Chart of Total execution time

5.2 Comparison of execution times for all tasks

From the Figure 5.3, we can observe the linear tendency of execution times for every tasks.

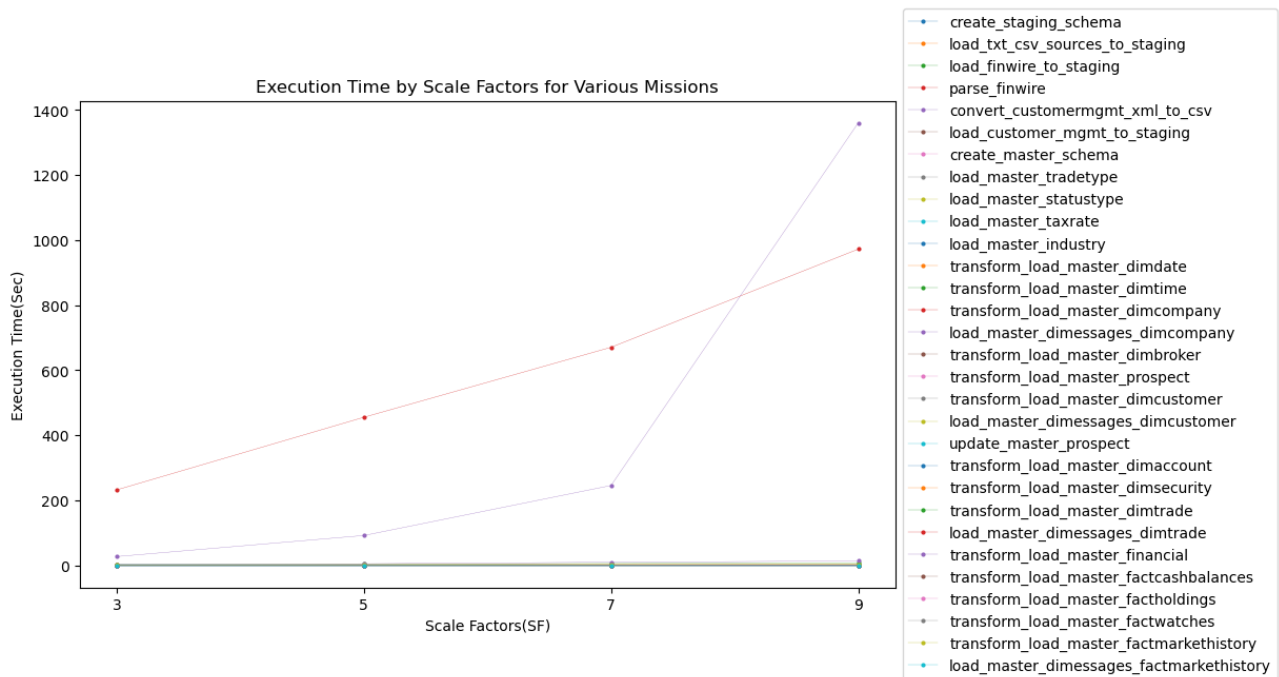


Figure 5.3: execution times for all tasks

From the figure above, we observed a significant disparity in execution times between various tasks. To show this big disparity, We also made the pie charts of Proportion of Execution Time of Various Tasks from SF3 to SF9. Then we found the execution time of parse_finwire and convert_customermgmt_xml_to_csv took up the two largest percentage of time.

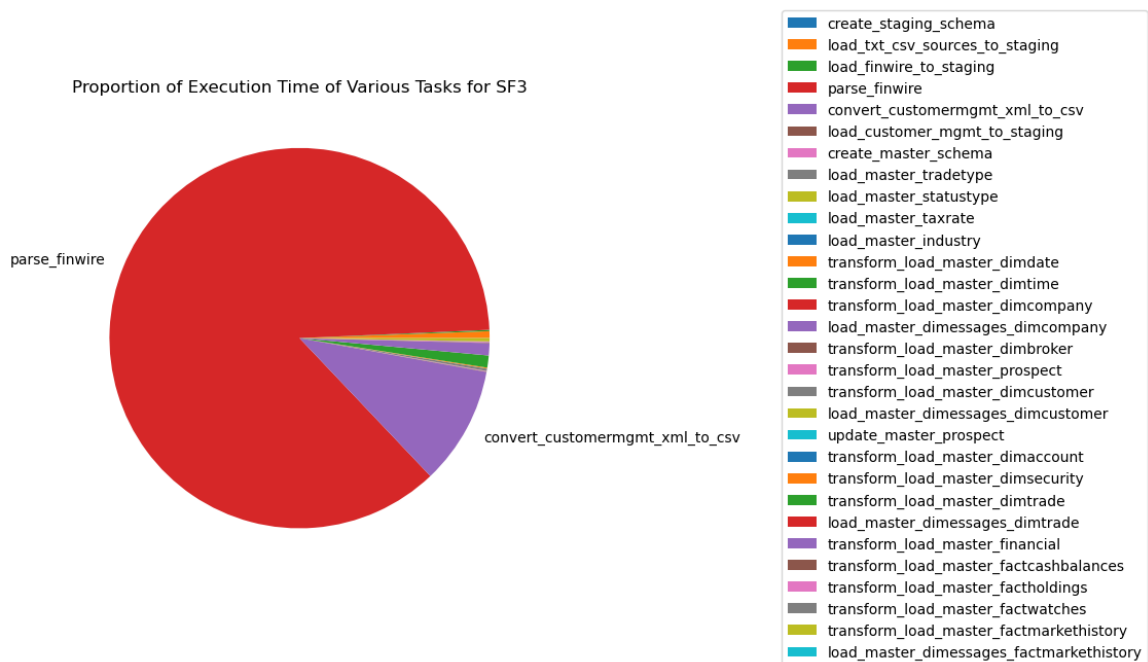


Figure 5.4: Proportion of Execution Time of Various Tasks for SF3

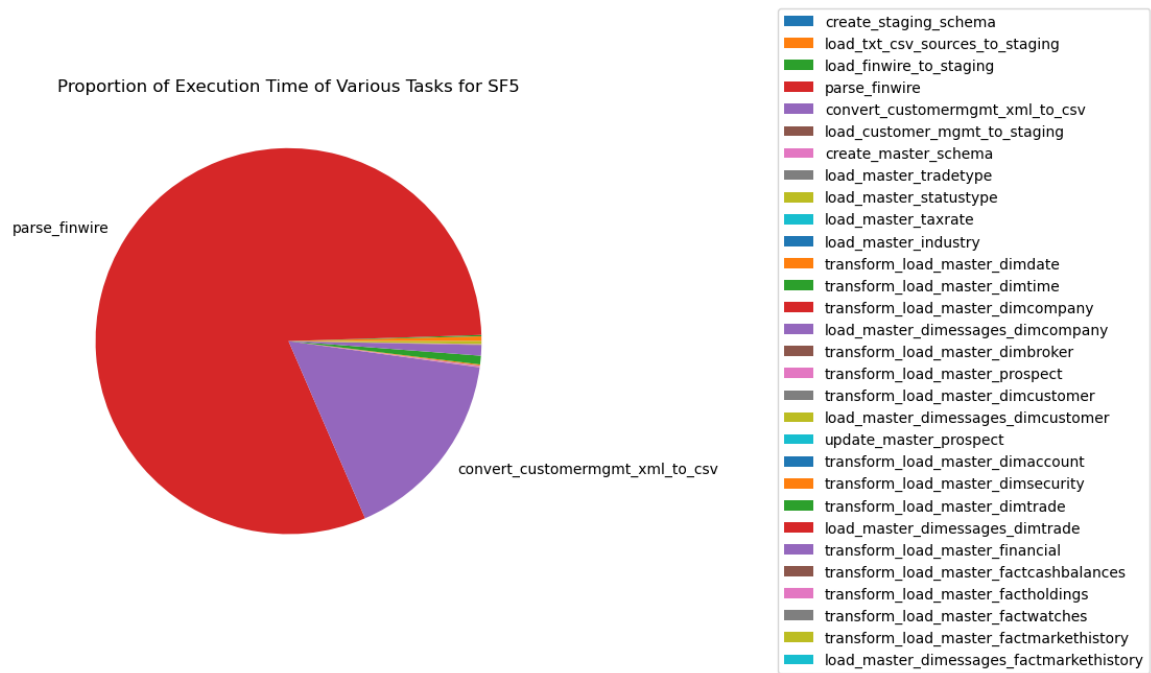


Figure 5.5: Proportion of Execution Time of Various Tasks for SF5

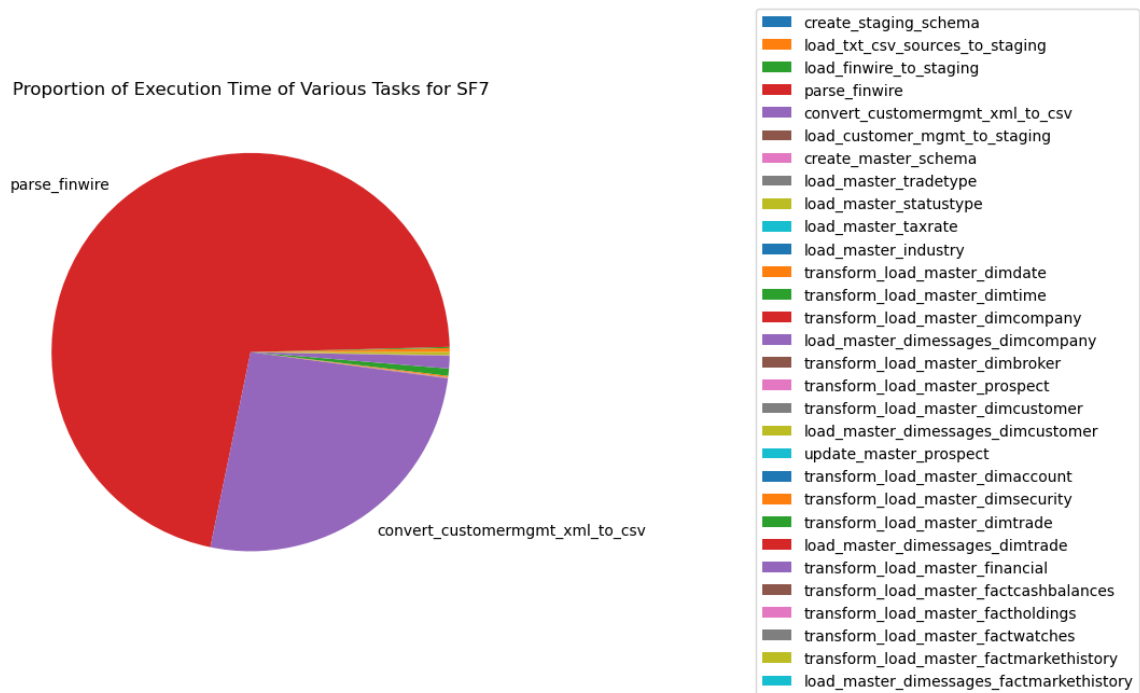


Figure 5.6: Proportion of Execution Time of Various Tasks for SF7

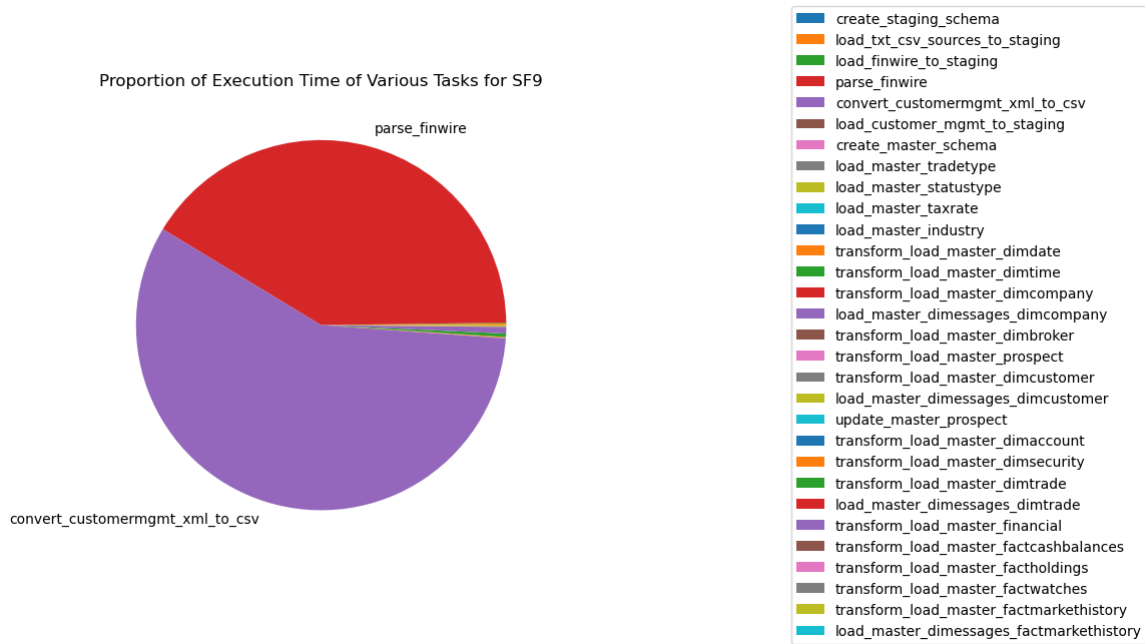


Figure 5.7: Proportion of Execution Time of Various Tasks for SF9

To clearly illustrate the trends of the remaining tasks, we made a strategic decision to exclude the two most time-consuming tasks from our data set. This approach allowed us to focus on the performance of the other tasks more effectively. Figure 5.8 presents the adjusted execution times, offering a clearer view of the underlying trends and patterns

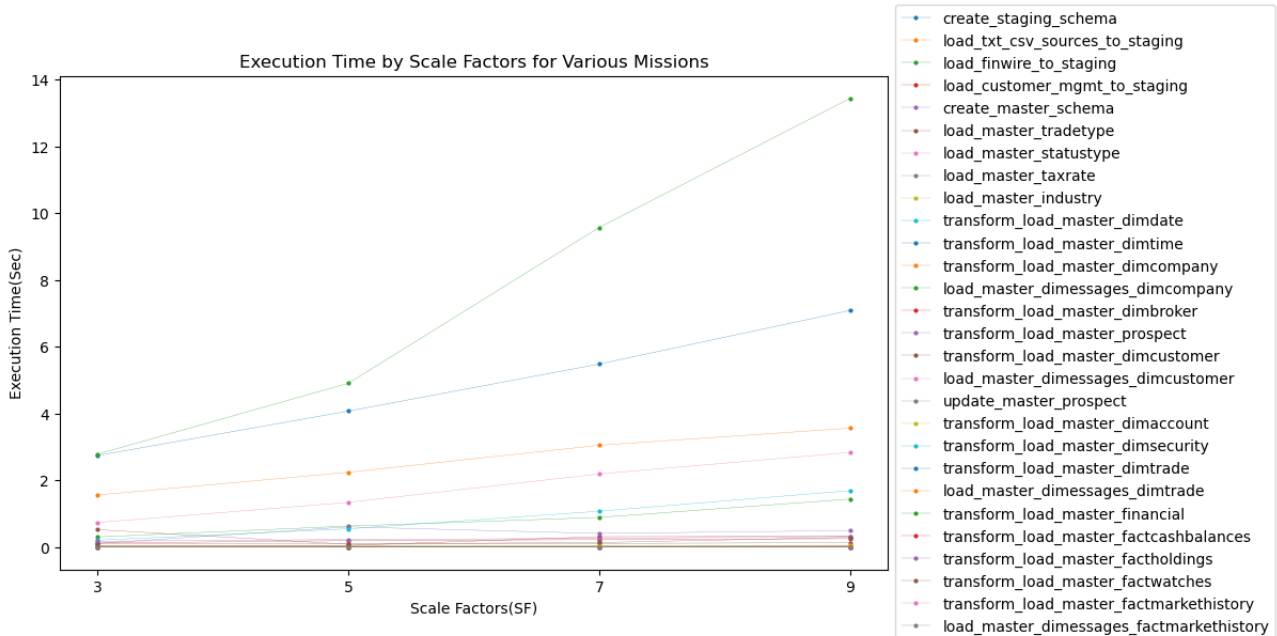


Figure 5.8: execution times for all tasks after removing two tasks

Then we also removed more time-consuming tasks which execution times were above 0.6s as the

Figure 5.9 shown. We can also know that most of the tasks are executed in less than 0.6 seconds.

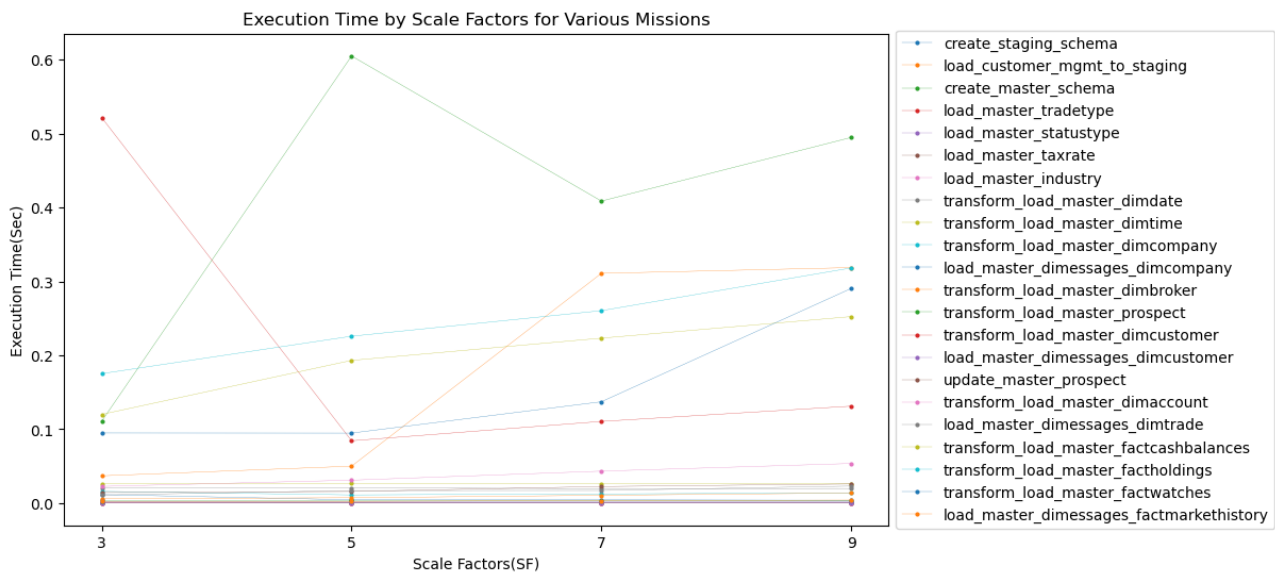


Figure 5.9: execution times for all tasks after removing more tasks

The two most time-consuming tasks were `parse_finwire` and `convert_customermgmt_xml_to_csv`. They are all in that hundreds seconds level compared to other tasks.

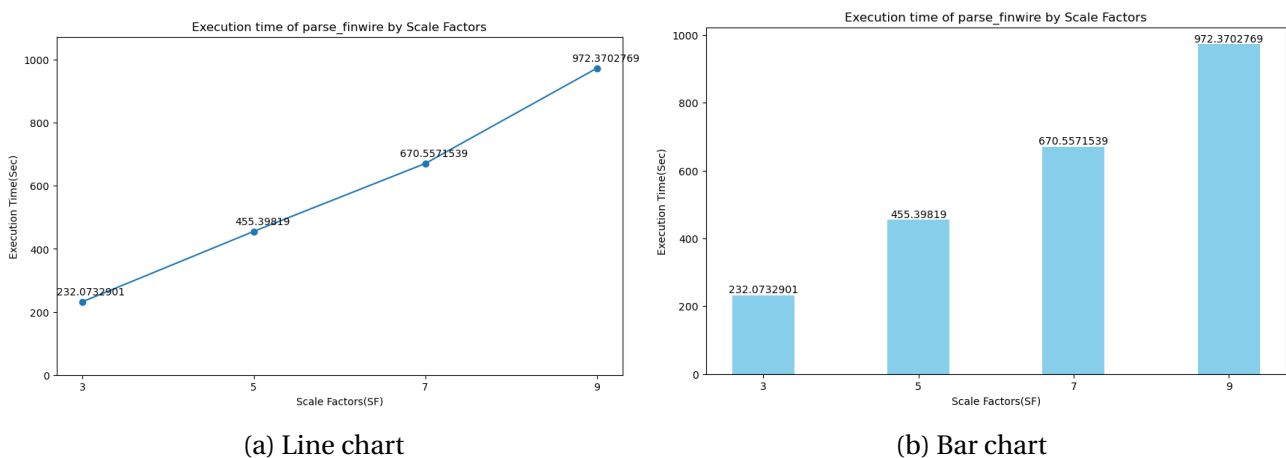


Figure 5.10: execution times for `parse_finwire`.

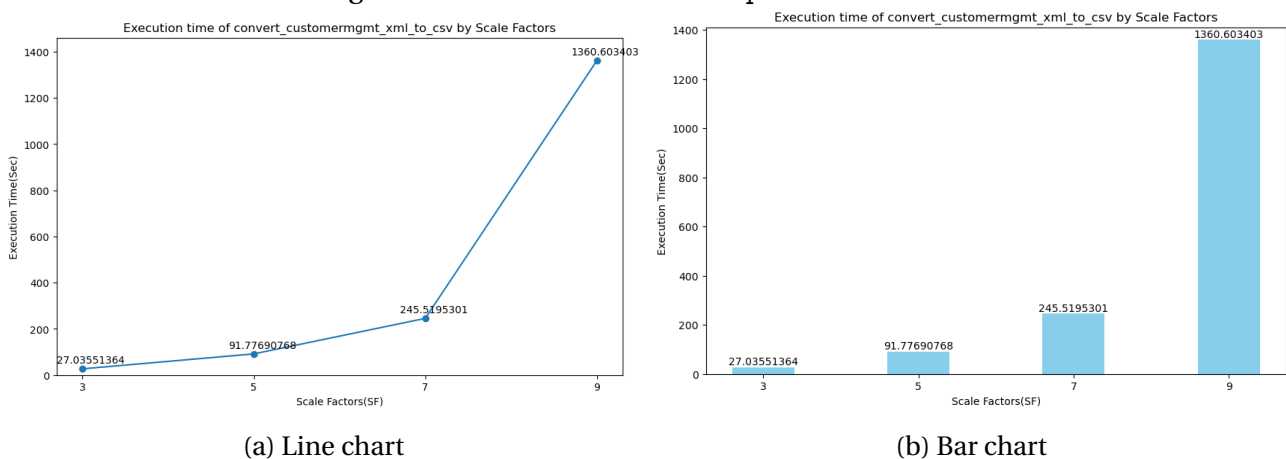


Figure 5.11: execution times for `convert_customermgmt_xml_to_csv`.

5.2.1 The tasks of linear growth

It could be seen that there were many tasks that grew linearly, such as `load_txt_csv_sources_to_staging`, `load_finwire_to_staging`, `transform_load_master_dimaccount` and `transform_load_master_dimtra`. These tasks were mainly load tasks. The charts also showed that the loading time increased almost linearly as the data was increased. We thought the processes used for loading the data may have linear complexity ($O(n)$), meaning that their execution time increases proportionally with the size of the input data. It was same to the load performance we tested in TPC-DS project.

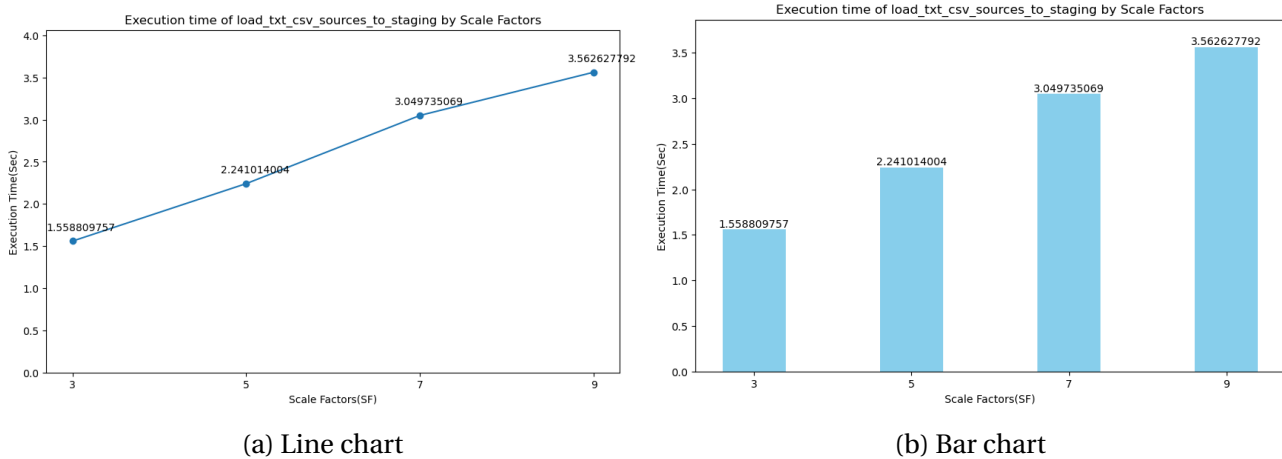


Figure 5.12: execution times for `load_txt_csv_sources_to_staging`.

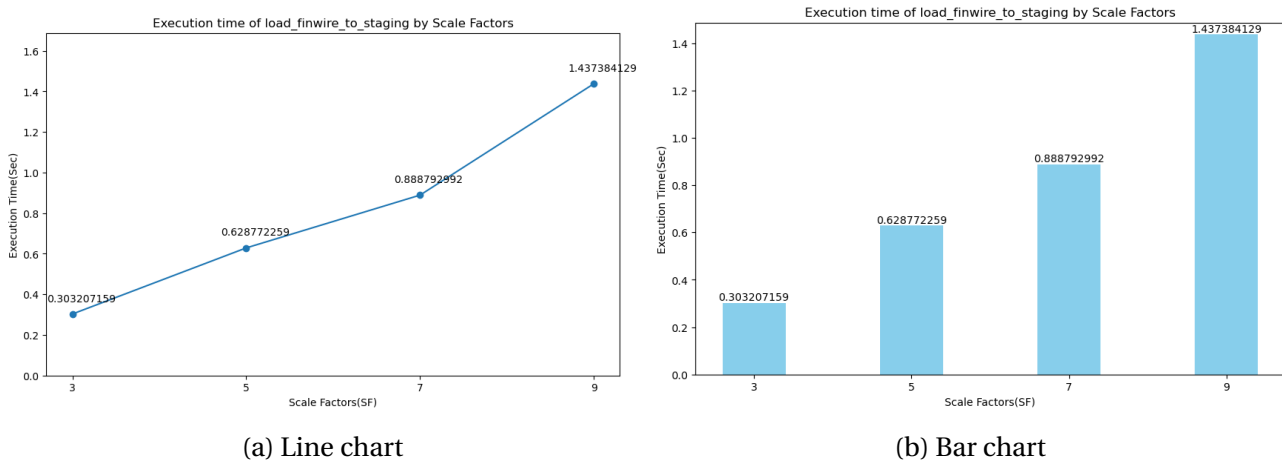
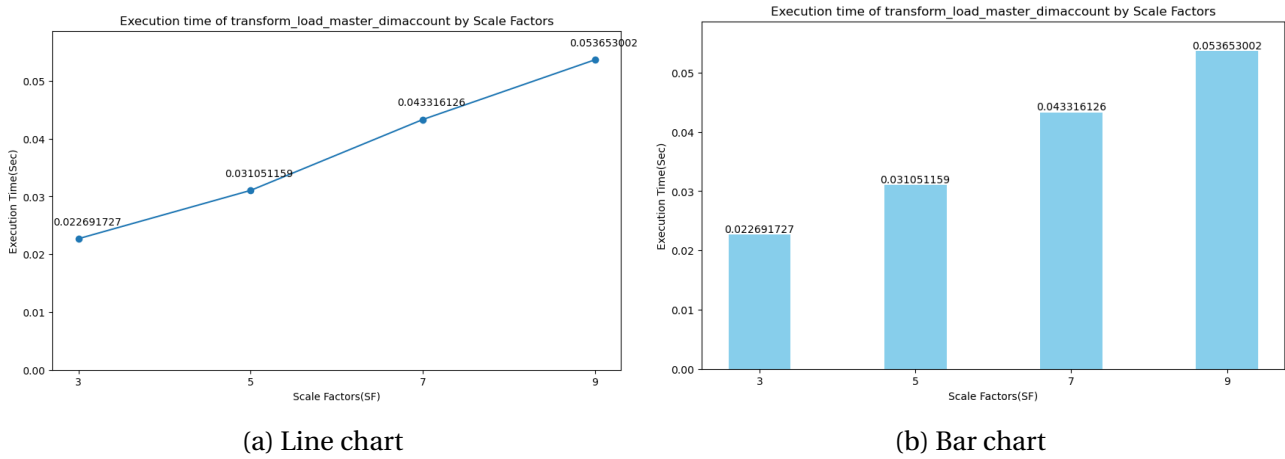
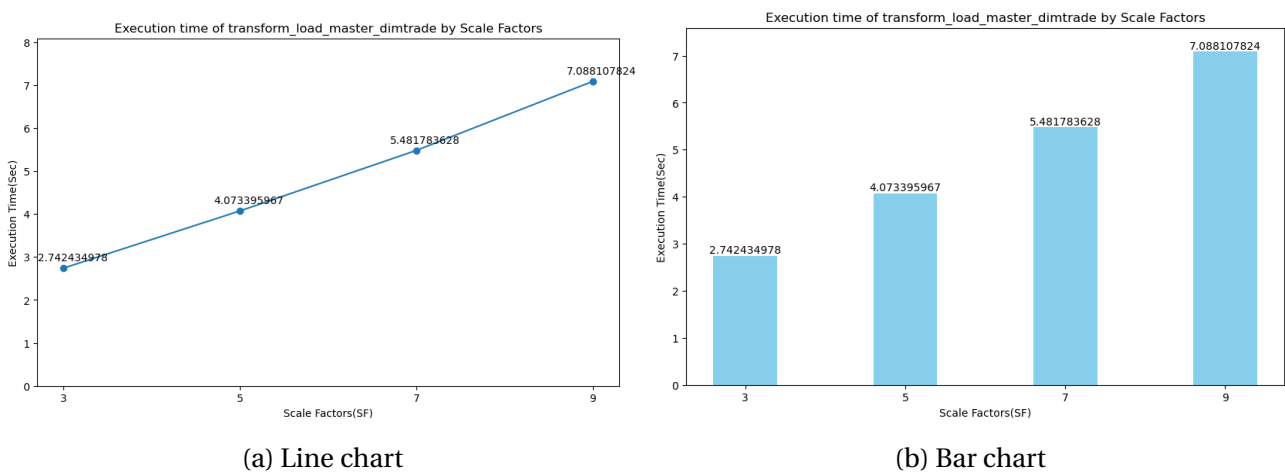
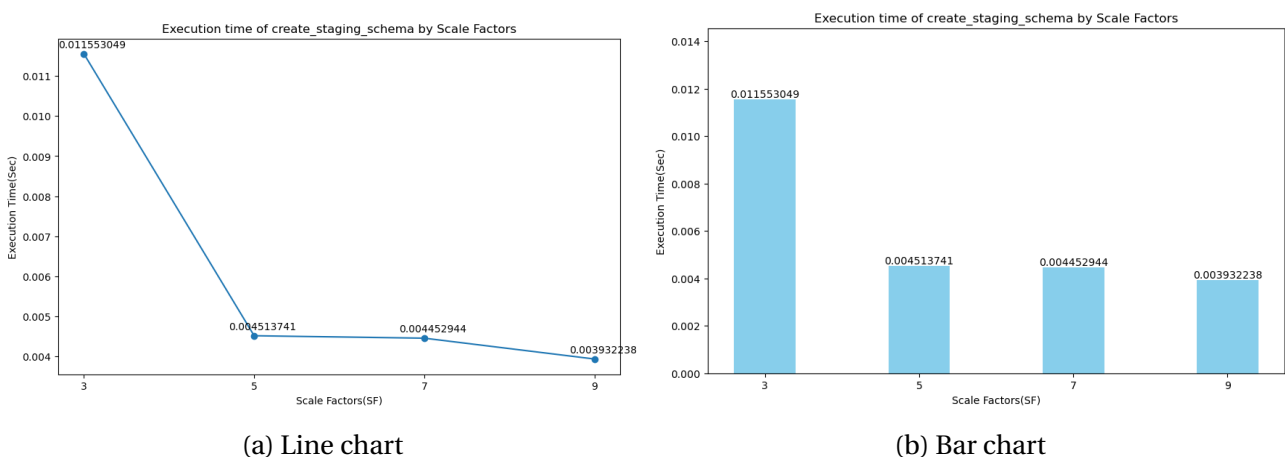


Figure 5.13: execution times for `load_finwire_to_staging`.

Figure 5.14: execution times for `transform_load_master_dimaccount`.Figure 5.15: execution times for `transform_load_master_dimtrade`.

5.2.2 The tasks of decreasing tendency

But there are still some that do not show a linear growth trend. Instead, it was a decrease tendency such as `create_staging_schema`. We thought it was related to Caching: The system might be caching some of the operations required for creating the staging schema, making subsequent operations faster.

Figure 5.16: execution times for `create_staging_schema`.

5.2.3 The tasks of stable trend

In addition, there were some tasks showed a stable up and down trend, such as `create_master_schema` and `load_master_tradetype`. We also thought it was related to caching effects: Repeated executions may benefit from caching, where data or execution plans are stored in memory, speeding up subsequent operations even the scale factors of data was increased.

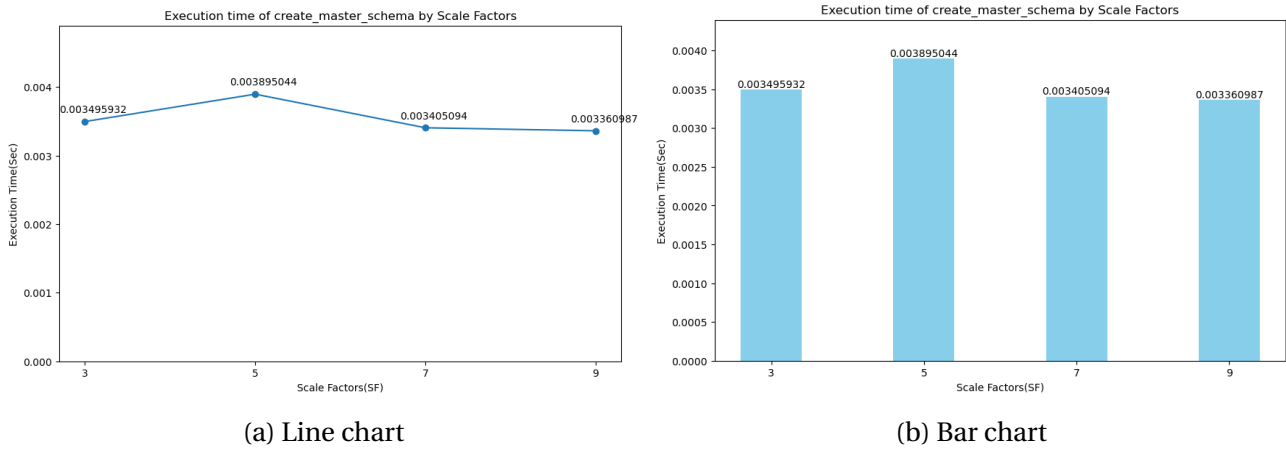


Figure 5.17: execution times for `create_master_schema`.

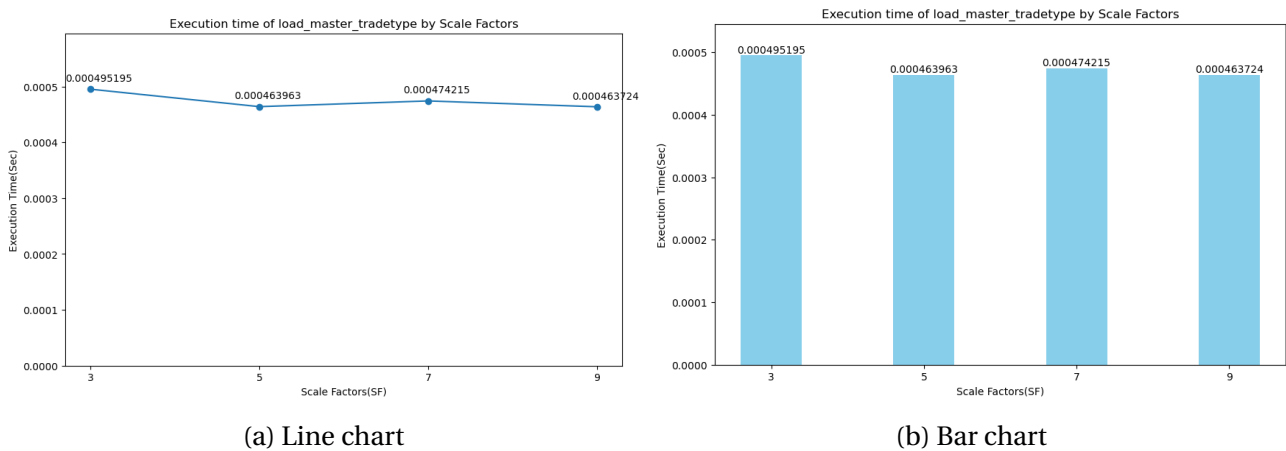


Figure 5.18: execution times for `load_master_tradetype`.

5.3 Discussion

After conducting a benchmarking analysis across a broad spectrum of tasks with several scale factors, we got valuable information on how Duckdb handles escalating data volumes as well as modifications to node steps based on table data size. While the most majority of tasks showed non-linear runtime performance characteristics, there would be some tasks were notable exceptions.

6.1 Conclusion

The TPC-DI framework offers an extensive ETL benchmarking model for DuckDB. This detailed report provides a clear blueprint for replicating our benchmarking results with DuckDB. Through this project, we gained practical insights into Data Integration for systems, learning about the complexities of sourcing data from diverse formats and integrating it into decision support systems through various processing methods.

In this project, We have systematically learned the basic theories and concepts of TPC-DI benchmark and implemented it on DuckDB, which is an open-source and particularly fast database. We finished the data generation. We conducted experiments for four scale factors. It was evident that the only some execution time of load tasks grew linearly with the size of data. Most of the execution time of tasks did not follow the linear tendency. Then we presented and analyzed the results of our experiments in some charts with data visualization.

One significant challenge encountered was the TPC-DI documentation was poorly articulated.

6.1.1 Further Works

To fully complete our TPC-DI implementation, there are key steps yet to be executed:

- Implementing the two phases of Incremental Updates.
- Implementing the Automated Audit Phase

References

- [TPC14] (TPC), Transaction Processing Performance Council (2014). “TPC Benchmark D1 Standard Specification”. In.
- [Cod23] Code, Visual Studio (2023). “Visual Studio Code - Code Editing”. In: URL: <https://code.visualstudio.com/docs>.
- [JZ20] Judith Awiti, Alejandro A. Vaisman and Zimányi (2020). “Design and implementation of ETL processes using BPMN and relational algebra”. In: *Data & Knowledge Engineering*.
- [Koh+22] Kohn, André et al. (2022). “DuckDB-wasm: fast analytical processing for the web”. In: *Proceedings of the VLDB Endowment* 15.12, pp. 3574–3577.
- [Mah21] Mahmud, Md Hedayet (2021). “What is benchmarking and how can it be used to measure our current performance against an ideal or best practice model?” In: URL: <https://www.quora.com/What-is-benchmarking-and-how-can-it-be-used-to-measure-our-current-performance-against-an-ideal-or-best-practice-model>.
- [Pyt23] Python (2023). “General Python FAQ”. In: URL: <https://docs.python.org/3/faq/general.html/>.
- [RM19] Raasveldt, Mark and Hannes Mühleisen (2019). “DuckDB: an Embeddable Analytical Database”. In: *International Conference on Management of Data (SIGMOD '19)*. URL: <https://dl.acm.org/doi/abs/10.1145/3299869.3320212>.
- [SG22] Saiid, Ahmad Abu and Rishika Gupta (2022). “TPC-DI-Benchmark: Data Integration with PostgreSQL & Airflow”. In: URL: <https://github.com/risg99/tpc-di-benchmark>.
- [Sam21] Samy Kabangu, John Ebden (2021). “Benchmarking Databases”. In.