



UNIVERSITÉ LIBRE DE BRUXELLES

ÉCOLE POLYTECHNIQUE DE BRUXELLES

A TCP-DI benchmark implementation using Oracle

INFOH419 - DATA WAREHOUSES

Fall 2022

Authors:

Ivanović, Nikola,
Lorencio Abril, Jose Antonio,
Yusupov, Sayyor,
Živković, Bogdana

Professor: Zimányi, Esteban

Contents

Abstract	iv
1 Data Integration	1
1.1 The ETL process	1
1.2 TPC-DI Benchmark	1
1.2.1 Business and Application environment	2
1.2.2 Summary of operations	3
1.2.3 Transformations	4
1.3 DIGen	4
1.4 Execution Rules & Metrics	4
1.4.1 Execution phases and measurements	4
1.4.2 Calculating throughput	5
1.4.3 Primary metrics	5
1.4.4 Price/Performance metric	5
1.5 Slowly Changing Dimensions	6
2 Implementation Details	7
2.1 Setting up Oracle Database	7
2.2 Generating and Loading the Data	7
2.2.1 Data Generation	7
2.2.2 Data transformation	7
3 Results and discussion	15
4 Conclusion	19
References	20
A Creation script	21
B Control files	30

List of Figures

3.1	Execution Time of Historical Load by Scale Factor.	15
3.2	Cumulative Number of Records by Scale Factor.	16
3.3	Number of Records by Table for each Scale Factor.	17
3.4	Proportion of Records with Alerts by Scale Factor	17

Listings

2.1	DimDate.ctl	8
2.2	DimTime.ctl	8
2.3	FactMarketHistory load query.	12
A.1	oracle-schema.sql	21
B.1	DimDate.ctl	30
B.2	DimTime.ctl	30
B.3	Industry.ctl	31
B.4	StatusType.ctl	31
B.5	TaxRate.ctl	31
B.6	TradeType.ctl	32
B.7	Broker.ctl	32
B.8	Prospect.ctl	32
B.9	Trade.ctl	33
B.10	CashBalances.ctl	33
B.11	Watches.ctl	34
B.12	Holdings.ctl	34

Abstract

In this project we have reproduced the TPC-DI Benchmark specification to execute in an Oracle environment. More precisely, we have used Oracle XE, the free version of Oracle SQL systems.

The present report starts reviewing the ETL process and a summary of the specification of TPC-DI is provided in Chapter 1. In Chapter 2 the implementations details on how we have taken the specification from theory to practice are explained, with emphasis in the different tests performed in the benchmark, whose results are analyzed in Chapter 3. In this chapter we also compare our results with those obtained in [AVZ20].

All the explanations provided in this report refer to the github repository created for the project, which can be accessed with the following link https://github.com/Lorencio/TPC-DI_Oracle.

Chapter 1

Data Integration

In this chapter, we are going to introduce the key concept for which the TPC-DI benchmark is designed: the ETL process. We also present a summary of the TPC-DI specification, trying to represent its main characteristics, phases and objectives.

1.1 The ETL process

Data Integration or, equivalently, the **Extraction, Transformation and Load (ETL) process** refer¹ to the process of extracting data from internal and external sources of an organization, transforming these data, and loading them into a data warehouse.

The ETL process is critical in the data warehouse scheme, and in all the steps of its life cycle: since its design to its maintenance, updating and even redesign to cope with new restrictions and data. The DI should aim to define how all data from different sources are combined into a single one, solving possible conflicts of meaning and shape of the data, as well as minimizing the errors entering the data warehouse.

As can be inferred, the design of the ETL processes must be done before the creation of the data warehouse, because these processes define the scope of the data warehouse and its limits, as they are defining the precise form of the data that will be stored inside it. In addition, one needs to think not only in the population of the data warehouse for the first time, but also how it will be updated when new data is available. This is an important observation, because even if we are able to define the ETL processes in such a way that we handle all possible errors and divergent cases that can arise today, the next month a new type of error that did not exist before can appear, and our ETL processes would need to be consequently modified to be able to tackle it.

1.2 TPC-DI Benchmark

The TPC-DI benchmark² is a performance test of tools that perform the DI process of a system. The benchmark **workload** manipulates a defined volume of data, preparing it for use in a Data Warehouse. The **model** includes data representing an extract from an OLTP system, which needs to be transformed along with data from different sources and loaded into a Data Warehouse.

¹Definition from [VZ22].

²For a completely detailed explanation refer to [TPC14].

This way, the benchmark tests the behavior of system components associated to the DI process. These components are characterized by:

- The manipulation and loading of large volumes of data.
- A mixture of transformation applied to the data.
- Historical loading, as well as incremental updates.
- Consistency requirements.
- Multiple data sources with different formats.
- Multiple data tables with varied data types, attributes and relationships.

On the other hand, the **operations** are modeled as follows:

- **Source data** is automatically generated in flat files.
- Transformation of the data begins with the **System Under Test (SUT)** reading the Source data.
- The **transformations** validate the Source data and properly structure the data for loading into a DWH.
- The process **concludes** when all Source data has been transformed and is available at the DWH.

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

1.2.1 Business and Application environment

The data represents a sample data retrieved from retail brokerage. Five types of data sources are required for the benchmark.

First, there are multiple tables in the OLTP system. They contain information about brokers, trade deals, customers etc. For each table, there are Changed Data Capture extracts that list the changes made to the table. They are used in the Incremental Updates Phase. On the other hand, the full dump of the original tables are made for Historical Load phase.

The HR database contains one table with an information on employees. It is used in the Historical Load and does not contain CDC files for Incremental Updates.

The Prospects file contains information on potential and existing customers with information on their names, demographic data and addresses. The DI process should determine the changes that were made to this file. It is updated daily with a new extract.

In addition, Historical Load Phase uses two other external sources of information. Data on companies and securities comes from Financial Newswire Service (FINWIRE). The file contains three types of records that will be used to fill the tables DimCompany, DimSecurity and Financial during the Historical Load. Also, data with information on new and updated customers and accounts was retrieved from Customer Management System. This data is in the form of .xml file. It is used to fill the tables DimAccount and DimCustomer during Historical Load.

Finally, a collection of reference data is used once during Historical Load.

1.2.2 Summary of operations

It is usually the case that the DI process has two variants. One of them is performing of a historical load, when the DWH is initially created or needs to be refilled with historical data (e.g. a restructuring of the schema). The other variant is adding newly generated data to the data warehouse.

Phases of operation

Preparation phase

Some initial operations are required but their execution times are not used for the benchmark. These include:

- Generating the data and putting it to the Staging Area. The data is generated with an automated generator provided by the Benchmark. The size of the data (scale factor) needs to be specified and as a result, a folder is created, containing three directories *Batch[k]*, where $k = 1, 2, 3$.
- Creating the databases and the tables of the Data Warehouse.

Historical Load phase

Historical Load is done either when the Data Warehouse is initially created or when the Data Warehouse is regenerated using historical records. Transformation rules specific to Historical Load are specified in the benchmark.

This phase uses files from *Batch1* directory of the generated data.

Incremental Update phase

Incremental Updates also require specific transformations. For example, while for Historical Load the data can be loaded in the primary key order, data for Incremental Update phase is loaded in the order the changes to the data were made. Two Incremental Updates are made to confirm repeatability. After each Incremental Update, Validation Query collects information about the execution to use it later for evaluation.

This phase uses CDC extracts from the OLTP system, as mentioned earlier. They contain a column "CDC_FLAG" which contains one of the three characters I, U or D to indicate whether the row was inserted, updated or deleted since the last version. A row can be changed multiple times during a day. The extract also contains column "CDC_DSN" which is used to order the rows by time of the changes. It is specific to each data source.

The tables that are modified during the Incremental Updates, contain not only primary but also surrogate keys. Primary keys are constant but surrogate keys are updated for each recorded change.

This phase uses files from *Batch2* and *Batch3* directories for executing the first and second incremental updates.

Automated Audit phase

This phase is done after all other phases. It runs queries to check for the validity of the resulting data stored in the Data Warehouse. In the end, it generates a report with results of all the tests. Each of

the tests should be passed for the run to be considered valid.

1.2.3 Transformations

In this benchmark, transformation refers to any operation that needs to be done to prepare and load data in the Data Warehouse. In a nutshell, the transformation to be done can be summed up as:

- Convert the string-like data to the correct datatype compatible with the SUT.
- Compute surrogate keys.
- Merge fields or split fields.
- Check data for errors.
- Detect changes in dimension tables and apply the changes to the corresponding records.
- Detect changes in fact data.

Data Dependencies

In the DWH, there exist dependencies between tables, so there is an ordering to the processing of the tables. When a dependent table column refers to a column in a source table, any row in the source table that would change the result of the dependent row must be processed before the latter.

1.3 DIGen

DIGen is the automatic data generator provided by the TPC-DI benchmark. It is used to create Source Data and the audit information and it is the only acceptable way to generate them.

DIGen also produces statistics about the data generation process. This file is used to calculate the metric and for auditing and contains:

- General information about the data generation process.
- Options used.
- Rows counts for each batch.

1.4 Execution Rules & Metrics

1.4.1 Execution phases and measurements

- Preparation phase: it includes the one time tasks needed to prepare the SUT for execution of the benchmark.

Source data is prepared in the staging area for the benchmark using the DIGen tool. A scale factor needs to be chosen.

Also, the DWH needs to be prepared, creating all tables necessary for the execution.

- Benchmark Run: a benchmark run consists of:

1. Initialization phase: provides a SUT prepared to be benchmarked. Not timed.
2. Historical load phase: populates the DWH performing the required transformations, using data from *Batch1*. This phase is timed.
3. Incremental update 1 phase: an incremental insertion of data from *Batch2* directory is performed. It is also timed.
4. Incremental update 2 phase: an incremental insertion of data from *Batch3* directory is performed. It is also timed.
5. Automated audit phase: this phase begins immediately upon completion of the last incremental update phase. It must be executed to confirm the validity of the results.

1.4.2 Calculating throughput

The **historical load throughput** is defined using the elapsed time, in seconds, for the historical load:

$$E_H = CT_1 - CT_0,$$

where CT_i is a timestamp taken after completion of phase i . If the total number of rows of source data is R_H , the throughput is defined as

$$T_H = \frac{R_H}{E_H}.$$

The **incremental update throughput** is defined similarly by

$$T_{I1} = \frac{R_{I1}}{\max(E_{I1}, 1800)}$$

and

$$T_{I2} = \frac{R_{I2}}{\max(E_{I2}, 1800)}.$$

1.4.3 Primary metrics

The **performance metric** is a combined metric using the three throughput values:

$$TPC_DI_RPS = \text{trunc}(\text{geoMean}(T_H, \min(T_{I1}, T_{I2})).$$

1.4.4 Price/Performance metric

This metric is defined as

$$\text{Price} - \text{per} - TPC_DI_RPS = \frac{\$}{TPC_DI_RPS}.$$

Here, \$ is the total 3-year pricing. We have not taken this measure into account.

1.5 Slowly Changing Dimensions

As explained in [AVZ20], **slowly changing dimensions (SCD)** are used in a Data Warehouse to keep the history of changes that occur in data sources. The TPC-DI benchmark focuses on what are called SDCs of type 2 with dependencies. This kind of dimension keeps track of the changes of the values by adding two temporal attributes: *StartDate* and *EndDate*, which indicate the date of insertion of the record in the database and the date when the tuple was deleted from the source table, respectively. Also, '*with dependencies*' means that the SCD table contains at least a surrogate key reference to another table, so when an update in the referenced dimension happens, the referencing table needs to be updated, too.

The dimensions of this type in the TPC-DI benchmark are *DimCustomer* and *DimAccount*, being the latter dependant of the former.

Chapter 2

Implementation Details

In this chapter, a thorough explanation of the implementation of the benchmark is provided, focusing on the underlying ideas, which are complemented by code snippets.

2.1 Setting up Oracle Database

Some members of our group used Windows OS and some Ubuntu OS so we had two types of setups for the project. The Ubuntu users set up a virtual machine using Oracle VM VirtualBox to run Oracle Linux 8.6. Windows users could directly install Oracle DB to their computers. All of the results in Chapter 3 were obtained using a computer running on Windows 64 bits, with CPU 11th Intel Core i7-11800H @ 2.30 GHz and 16 GB of RAM.

All of us used Oracle Database Express Edition (XE) 21c as the tool. It is a free version of Oracle Database. However, it had limitations of up to 12 GB of user data, up to 2 GB of database RAM and up to 2 CPU threads [Dat22].

2.2 Generating and Loading the Data

2.2.1 Data Generation

To generate the data, we have used the data generation tool provided by the TPC-DI developers: *DIGen.jar*.

2.2.2 Data transformation

Our implementation was inspired by a partial implementation of TPC-DI in MySQL [NI18].

To perform the transformations detailed in the specification of the benchmark, we have developed a Python program which makes use of both Python and OracleSQL to perform all transformations needed. In particular, command-line tool *SQL*Loader* was used to load tables with data from files, Oracle Database utility *SQL*Plus* was used to execute the queries and Python Library *python-oracledb* was used to connect to Oracle Database from Python.

For this project, only the transformations for the Historical Load were implemented, due to the time constraints.

The main class in this program is the class `TPCDI_Loader`, which is given the information needed to connect to the database and where to find the flat files of data generated before. Once it has this information, it creates all the tables in the database. Then, the different tables are processed one by one, following the order of their dependencies.

To create all the tables, an OracleSQL script have been defined, it can be seen in Appendix A. It has been made to adhere to the specification of the benchmark, adapting some parts to the OracleSQL syntax and restriction. For example, as OracleSQL does not implement a Boolean type, we have modeled boolean values as *char strings* with only two possible values *'true'* and *'false'*.

Batch Date

To load the Batch Date, we simply read the batch date from *BatchDate.txt* and we enter into the table BatchDate the Batch number and the date.

DimDate

This table is loaded with the *sqlldr* tool, providing the *BatchDate.sql* file and the control file *DimDate.ctl*, which is shown in Listing 2.1.

```

1 LOAD DATA
2 INTO TABLE DimDate
3 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4 (
5     SK_DateID,
6     DateValue DATE "YYYY-MM-DD",
7     DateDesc,
8     CalendarYearID,
9     CalendarYearDesc,
10    CalendarQtrID,
11    CalendarQtrDesc,
12    CalendarMonthID,
13    CalendarMonthDesc,
14    CalendarWeekID,
15    CalendarWeekDesc,
16    DayOfWeeknumeric,
17    DayOfWeekDesc,
18    FiscalYearID,
19    FiscalYearDesc,
20    FiscalQtrID,
21    FiscalQtrDesc,
22    HolidayFlag
23 )

```

Listing 2.1: DimDate.ctl

DimTime

This table is treated very similarly to the previous one: we feed the *sqlldr* with *Time.txt* and *DimTime.ctl*, which can be seen in Listing 2.2.

```

1 LOAD DATA
2 INTO TABLE DimTime
3 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4 (
5     SK_TimeID,
6     TimeValue DATE "HH24:MI:SS",

```

```

7      HourID ,
8      HourDesc ,
9      MinuteID ,
10     MinuteDesc ,
11     SecondID ,
12     SecondDesc ,
13     MarketHoursFlag ,
14     OfficeHoursFlag
15 )

```

Listing 2.2: DimTime.ct1

Industry

Once again, the *sqlldr* tool facilitates our work, just needing to provide the text and control files. From this point on, we will not show more control files in the main text, but the reader is referred to Appendix B. In this case, more precisely to Listing B.3.

StatusType

This table is analogous to the previous ones, we feed *sqlldr* with *StatusType.txt* and *StatusType.ct1*¹.

TaxRate

This table is analogous to the previous ones, we feed *sqlldr* with *TaxRate.txt* and *TaxRate.ct1*².

TradeType

This table is analogous to the previous ones, we feed *sqlldr* with *TradeType.txt* and *TradeType.ct1*³.

DimCompany

This table is more complex than the previous ones. Our approach has been to divide its loading into two steps.

First, we populate a staging table, *S_Company* by means of the function *load_staging_finwire*. This function populates three staging tables: *S_Company*, *S_Financial* and *S_Security*. The idea is simple: we traverse all the *FINWIRE* files in the data directory, and for each of them, we traverse its lines, one by one. The lines can be one out of three types, which is indicated in the characters 15-17 of the line:

- CMP: the record is to be inserted into *DimCompany*.
- SEC: the record is to be inserted into *DimSecurity*.
- FIN: the record is to be inserted into *Financial*.

Now, depending on the type, the different fields are recovered and the insertion query into the staging database is composed. In the case of CMP records, additional checks need to be made, as values can be missing and they need to be inputted as NULL into the database in this case. This does not happen with the other tables.

¹See B.4.

²See B.5.

³See B.6.

An important concern is that *DimCompany* is a SDC, so newly inserted records can make it necessary to modify the previous ones, in terms of time validity: we have to set the column *IsActive* to 'false' and the column *EndDate* to *EffectiveDate* of the new record. This is done with a specific function for filling *DimCompany*.

After *load_staging_finwire* finishes populating the three staging tables, it's time to use these to fill the dimensions. For this task, we developed the function *load_target_dim_company*, which performs the following steps:

1. Load all records of *S_Company* into *DimCompany*, joining with previously loaded *Industry* and *StatusType* tables, and doing some checks on some of the values.
2. Creates a temporal table to perform the SDC changes in the database. This table is a copy of *DimCompany*, with a row number attribute that is useful to order the data by *CompanyID* and effective date.
3. Once we have this order, we update the records in *DimCompany*, comparing each record of each company with the next record for the same company. If there is such a 'next record', then the previous one is updated to deactivate it.
4. The temporal table is dropped.

This function also populates the *DIessages* table with messages generated from *DimCompany* tables, but this matter will be further detailed later.

Financial

Financial also relies on a staging table, *S_Financial*, which is populated as explained before. Once the staging table is filled, a specific function populates *Financial*, function *load_target_financial*.

This function is simpler than the previous one, and its functioning relies on the assumption that *load_target_company* has already been executed. It just joins *S_Financial* with *DimCompany*, checking the correctness of some values and making sure that the records are in the correct dates. Two queries of this type are used to address two different ways to join the values.

DimSecurity

DimSecurity is very similar to the last two tables. First, a staging table *S_Security* is populated with the function *load_staging_finwire*. Then it is populated by two queries treating two distinct cases to join the values of *S_Security*, *StatusType* and *DimCompany*. After this, as *DimSecurity* is a SDC, we do something similar as we did with *DimCompany*, we create a temporary copy of *DimSecurity* which is used to detect changes and perform them in the database.

DimBroker

Broker is also loaded in two parts. First, a staging table is populated using a control file⁴. Then, we execute simple query where only those employees matching the conditions specified in the specification are introduced in the table.

⁴See Listing B.7.

Prospect

Prospect load is similar to *Broker* load: again, a staging table is filled using a control file⁵. After this loading, some messages for the *DimMessages* table are created. Then, two queries are performed to make the appropriate transformations into the data and introduce them into the table *Prospect*. For this purpose, the SQL function *get_marketing_template* is created and used to transform the data from the staging table properly.

Basically, the function classifies the marketing template using different ranges for net worth and income values.

Prospect needs to be updated after *DimCustomer* has been processed. With this objective, the function *update_prospect* was developed. It consists in a simple query that marks 'true' for those prospects for which there is a corresponding customer.

DimCustomer and DimAccount

The table *DimCustomer* is the most complex dimension table in the benchmark, together with *DimAccount*. As they are strongly related and these tables represent the two tables that model the SDC with dependencies tables that we mentioned in the Chapter 1, we are going to explain them together.

The data to fill this tables is first loaded into two staging tables, *S_Customer* and *S_Account*, which are processed at the same time, from the same file, similarly to what we did with *DimCompany*, *Financial* and *DimSecurity*. In this case, we have to parse a XML file, checking for all fields that are to be inserted into the staging tables.

Once the staging tables are filled, the tables start to be treated separately. A series of functions is executed as follows:

1. *load_new_customer*: fills *DimCustomer* table with those records from the staging table that are labeled as 'NEW'. The query joins *Prospect* and *S_Customer*.
2. *load_new_account*: fills *DimAccount* table with those records from the staging table that are labeled as 'NEW' or 'ADDACT'. The query joins *S_Account*, *DimBroker* and *DimCustomer*.
3. *create_trigger_UPD_customer_account*: creates a trigger which automatically updates those records in account that gets affected when the AccountID of a customer is modified.
4. *load_update_customer*: the values from *DimCustomer* table are updated with those records from *S_Customer* labeled as 'UPDCUST'. The basic idea of the query is getting the most recent valid value for each of the fields of each record. To achieve this, an individual query for each field is run. There are two types of queries: those which only need information from *S_Customer* and those which needs a join with the table *Prospect* to obtain the needed information.
5. *load_update_account*: similarly to the last function, here we are updating *DimAccount* table with those records from *S_Account* labeled as 'UPDACCT'.
6. *load_close_account*: in this case, those accounts which are labeled as 'CLOSEACCT' in the staging table are marked as closed in the *DimAccount* table.
7. *load_inact_customer*: this function marks set an EndDate for those customers that are labeled as 'INACT' in the staging table. They are also marked as not valid.

⁵See Listing B.8.

8. *load_inact_account*: this function marks set an *EndDate* for those accounts that are labeled as 'INACT' in the staging table. They are also marked as not valid.

DimTrade

A staging table, *S_Trade* is populated by means of a control file⁶. Then, the function *load_trade* executes a query that makes the pertinent transformations on the data. These transformations basically consist of inserting one value or another into *DimTrade* depending on the values held in *S_Trade*.

FactCashBalances

First, we make use of a control file⁷ to fill a staging table. After that, a query joining *S_Cash_Balances*, *DimAccount* and *DimDate* properly populates the table *FactCashBalances*.

FactWatches

Similarly to *FactCashBalances*, this fact table is filled using a staging table populated by means of a control file⁸ and a simple query which joins appropriately different tables. In this case, distinctions needs to be made using the values of the *W_Action* column in *S_Watches*.

FactHoldings

Again, a control file⁹ and a properly done query makes the work.

FactMarketHistory

This table is the most complex one. First, a staging table with the facts is filled using a control file¹⁰. Then, the *FactMarketHistory* table needs to be populated. The problem at this point is that the query is really complex: it asks to get, for each record, the maximum and minimum value of the security close price for the past year, as well as the dates in which these values occurred.

To do this, we performed an analytic query, based on partitions and windows. As we will see in the results, this query is too complex and not very useful. We have not been able to develop a better solution for this problem. Therefore, we omitted this table from our execution.

The problem of this query is that we need to retrieve the dates when the highest and lowest values occurred. This involves several subqueries and multiple joins. We have not been able to come up with a fast implementation in OracleSQL, because it is not possible to return the values of the row in which we find the maximum value without using subqueries somehow. The fact that Oracle XE only allows for 2GB of RAM used makes it crash at some point. Our query can be seen in Listing 2.3, where one can see the problem: recovering the dates makes it unfeasible.

```

1 INSERT INTO FactMarketHistory (SK_SecurityID, SK_CompanyID, SK_DateID, PERatio, Yield,
2   FiftyTwoWeekHigh,
3   SK_FiftyTwoWeekHighDate, FiftyTwoWeekLow, SK_FiftyTwoWeekLowDate, ClosePrice, DayHigh
4   , DayLow, Volume, BatchID)
5 WITH MAXCLOSE AS (
6   SELECT DM_DATE, DM_S_SYMB, DM_CLOSE, MAX(DM_HIGH) OVER (PARTITION BY DM_S_SYMB ORDER
7   BY DM_DATE ROWS BETWEEN 365 PRECEDING AND CURRENT ROW) AS MAXC

```

⁶See B.9.

⁷See B.10.

⁸See B.11.

⁹See B.12.

¹⁰See ??.

```

5      FROM S_Daily_Market DM),
6  MAXDATE AS (
7      SELECT DM.DM_S_SYMB AS SYMB, MIN(DM.DM_DATE) AS MAXD, MAX(MC.MAXC) AS MAXC, MC.
      DM_DATE AS BASEDATE
8      FROM S_Daily_Market DM
9      JOIN MAXCLOSE MC ON MC.DM_S_SYMB = DM.DM_S_SYMB
10     WHERE DM.DM_DATE >= MC.DM_DATE - 365 AND DM.DM_DATE < MC.DM_DATE AND DM.DM_HIGH = MC.
      MAXC
11     GROUP BY DM.DM_S_SYMB, MC.DM_DATE
12 ),
13 MINCLOSE AS (
14     SELECT DM.DM_DATE, DM.DM_S_SYMB, MIN(DM.DM_LOW) OVER (PARTITION BY DM.DM_S_SYMB ORDER
      BY DM.DM_DATE ROWS BETWEEN 365 PRECEDING AND CURRENT ROW) AS MINC
15     FROM S_Daily_Market DM),
16 MINDATE AS (
17     SELECT DM.DM_S_SYMB AS SYMB, MIN(DM.DM_DATE) AS MIND, MAX(MC.MINC) AS MINC, MC.DM_DATE
      AS BASEDATE
18     FROM S_Daily_Market DM
19     JOIN MINCLOSE MC ON MC.DM_S_SYMB = DM.DM_S_SYMB
20     WHERE DM.DM_DATE >= MC.DM_DATE - 365 AND DM.DM_DATE < MC.DM_DATE AND DM.DM_LOW = MC.
      MINC
21     GROUP BY DM.DM_S_SYMB, MC.DM_DATE
22 ),
23 EPS AS (
24     SELECT DM.DM_DATE, DM.DM_S_SYMB, SUM(F.FI_BASIC_EPS) OVER (PARTITION BY DM.DM_S_SYMB
      ORDER BY DM.DM_DATE ROWS BETWEEN 130 PRECEDING AND CURRENT ROW) AS TOTAL_EPS
25     FROM S_Daily_Market DM
26     JOIN DimSecurity DS ON (DM.DM_S_SYMB = DS.Symbol)
27     JOIN Financial F ON (F.SK_CompanyID = DS.SK_CompanyID))
28 SELECT DS.SK_SecurityID,
29        DS.SK_CompanyID,
30        DD.SK_DateID,
31        CASE
32            WHEN EPS.TOTAL_EPS IS NULL THEN NULL
33            ELSE DM.DM_CLOSE/EPS.TOTAL_EPS
34        END,
35        DS.DIVIDEND/DM.DM_CLOSE*100,
36        MXC.MAXC,
37        DMX.SK_DateID,
38        MNC.MINC,
39        DMN.SK_DateID,
40        DM.DM_CLOSE,
41        DM.DM_HIGH,
42        DM.DM_LOW,
43        DM.DM_VOL,
44        1
45 FROM S_Daily_Market DM INNER JOIN DimDate DD ON (TO_CHAR(DM.DM_DATE, 'YYYY-MM-DD') =
      TO_CHAR(DD.DateValue, 'YYYY-MM-DD'))
46                INNER JOIN DimSecurity DS ON (DM.DM_S_SYMB = DS.Symbol)
47                INNER JOIN MAXCLOSE MXC ON (DM.DM_S_SYMB = MXC.DM_S_SYMB AND
      MXC.DM_DATE = DM.DM_DATE)
48                INNER JOIN MAXDATE MXD ON (MXD.BASEDATE = DM.DM_DATE AND MXD.
      SYMB = DM.DM_S_SYMB)
49                INNER JOIN DimDate DMX ON (TO_CHAR(MXD.MAXD, 'YYYY-MM-DD') =
      TO_CHAR(DMX.DateValue, 'YYYY-MM-DD'))
50                INNER JOIN MINCLOSE MNC ON (DM.DM_S_SYMB = MNC.DM_S_SYMB AND
      MNC.DM_DATE = MNC.DM_DATE)
51                INNER JOIN MINDATE MND ON (MND.BASEDATE = DM.DM_DATE AND MND.
      SYMB = DM.DM_S_SYMB)
52                INNER JOIN DimDate DMN ON (TO_CHAR(MND.MIND, 'YYYY-MM-DD') =
      TO_CHAR(DMN.DateValue, 'YYYY-MM-DD'))

```

```
53         INNER JOIN EPS ON (DM.DM_S_SYMB = EPS.DM_S_SYMB AND DM.DM_DATE
    = EPS.DM_DATE)
54 WHERE DS.EffectiveDate <= DM.DM_DATE AND DM.DM_DATE < DS.EndDate;
```

Listing 2.3: FactMarketHistory load query.

DImessages

The table *DImessages* is used to store alert messages for records, statuses of different execution phases and validation messages.

Alert records are written when columns have values that do not make sense. For example, a record will be inserted *DImessages* if a row in table *DimCustomer* has value of column *Tier* that is not one of the valid values, namely 1, 2 or 3. Such invalid records were deliberately inserted by the system during the Data Generation phase.

DImessages also stores Validation Messages of *Batch Validation Queries* that are run after the Historical Load. They are part of the *Audit Phase* that is described in section 1.2.2. They include the records of the number of records in the resulting tables.

Chapter 3

Results and discussion

Below the results of executing Historical Load of TPC-DI are shown. We decided to run the Load with 4 scale factors, namely 3, 5, 7 and 9. For each scale factor, we saved the records of DIMessages and the execution time of the entire Historical Load.

Figure 3.1 below shows the execution times for each scale factor. The execution time increased at about the same rate, around 3.1 minutes for the next chosen scale factor. The increase from scale factor 5 to 7 was slightly higher, at 4.2 minutes.

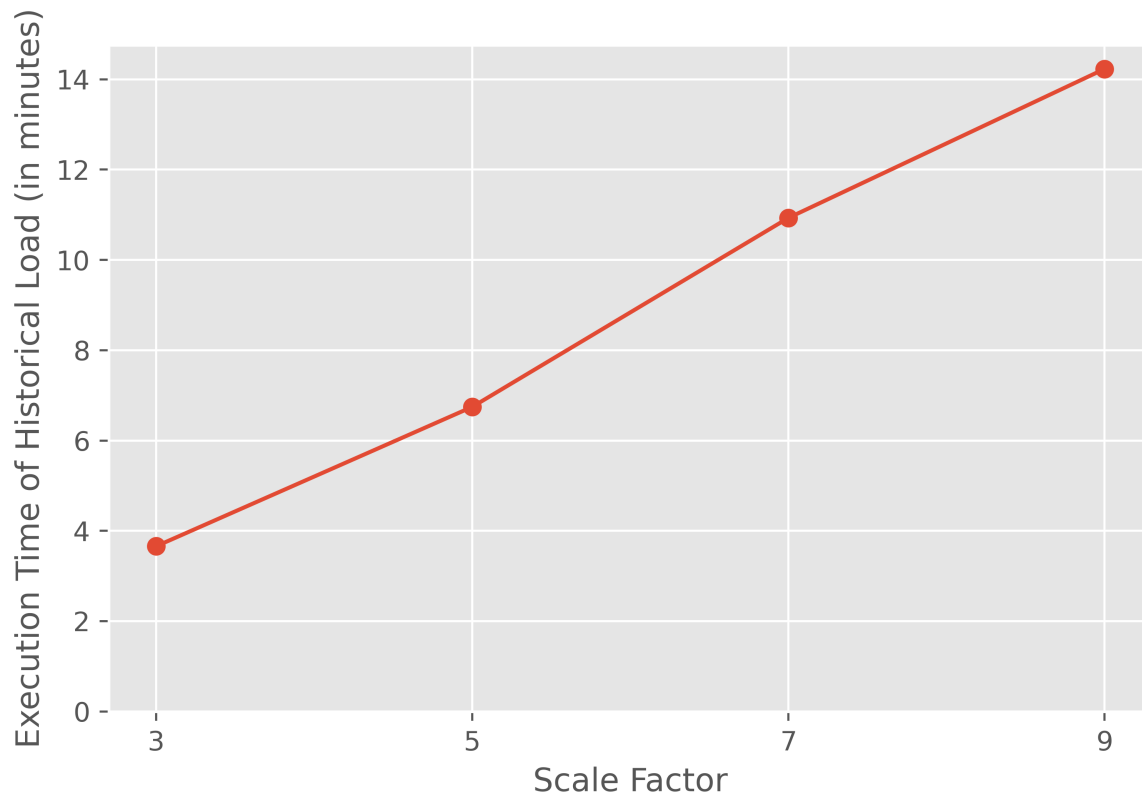


Figure 3.1: Execution Time of Historical Load by Scale Factor.

Figure 3.2 shows the number of records loaded for each scale factor. There were around 850,000 more

records for the next chosen scale factor. However, the rate of this increase decreased slightly. For example, there were around 1 million records more records in scale factor 5 compare to scale factor 3, while there were only 820,000 records more in scale factor 9 compared to scale factor 7.

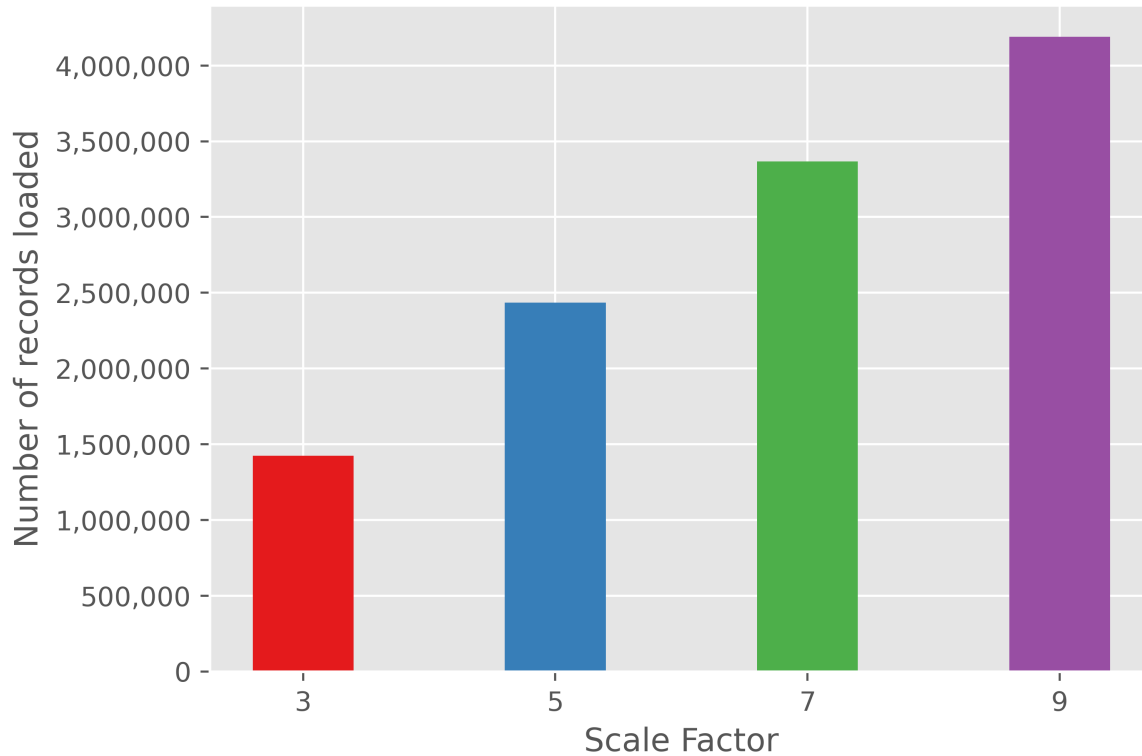


Figure 3.2: Cumulative Number of Records by Scale Factor.

You can see the number of records for each table in Figure 3.3. For all the scale factor, the Dimension Table *DimFactCash* and Fact Tables *Balances*, *FactWatches*, *FactHoldings* were the tables with the greatest number of records in that order, and accounted for around 91% of all the records. All the Reference Tables, including *TradeType*, *StatusType*, *Industry* and *TaxRate*, had a small number of records.

Figure 3.4 shows the number records with alerts which are recorded with *DImessages*, as mentioned previously. The proportion of alerted records are as expected about the same for all the scale factors, with 0.03% of all the records being alerted records of Table *DimCustomer* and around 0.0055% of all the records being alerted records of Table *DimTrade*. From this, we can assume that the records with alerts were not significantly greater than usual and our implementation inserted the records correctly but further assessment is needed.

Now, in Table 3.1 we are comparing the results obtained in [AVZ20] using PostgreSQL and our results using OracleSQL. The only scale factors shared are 3 and 5, and we can see how OracleSQL is performing faster, but we need to take into account that we are not loading *FactMarketHistory*. Nonetheless, in the article the authors show the time taken for each table to load for scale factor 3, so we can subtract the time taken by *FactMarketHistory* and make a more fair comparisons of the times. This is shown in Table 3.2, and we can see how disregarding this table makes OracleSQL a good option for this process. Of course, without that table, the implementation is useless, so there would be work to do in order to be able to load the table in a reasonable time. At the moment we don't know if this is even possible.

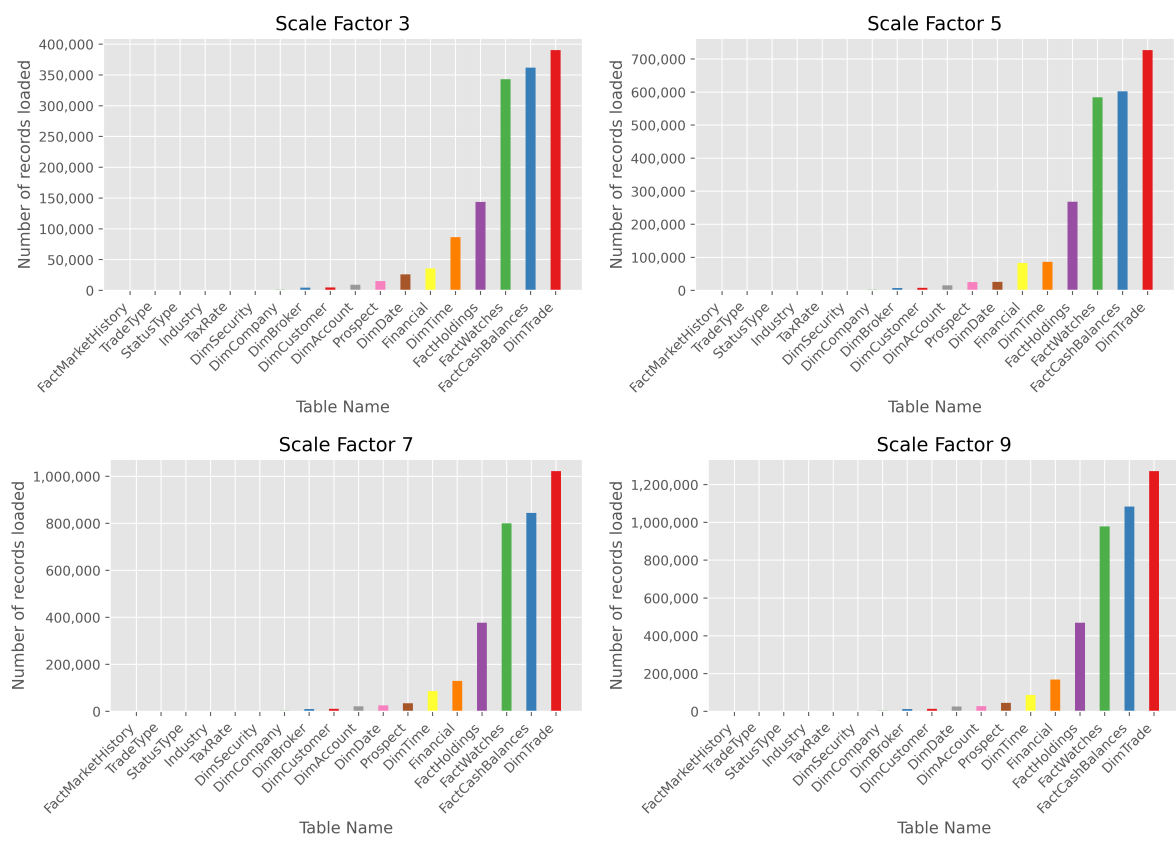


Figure 3.3: Number of Records by Table for each Scale Factor.

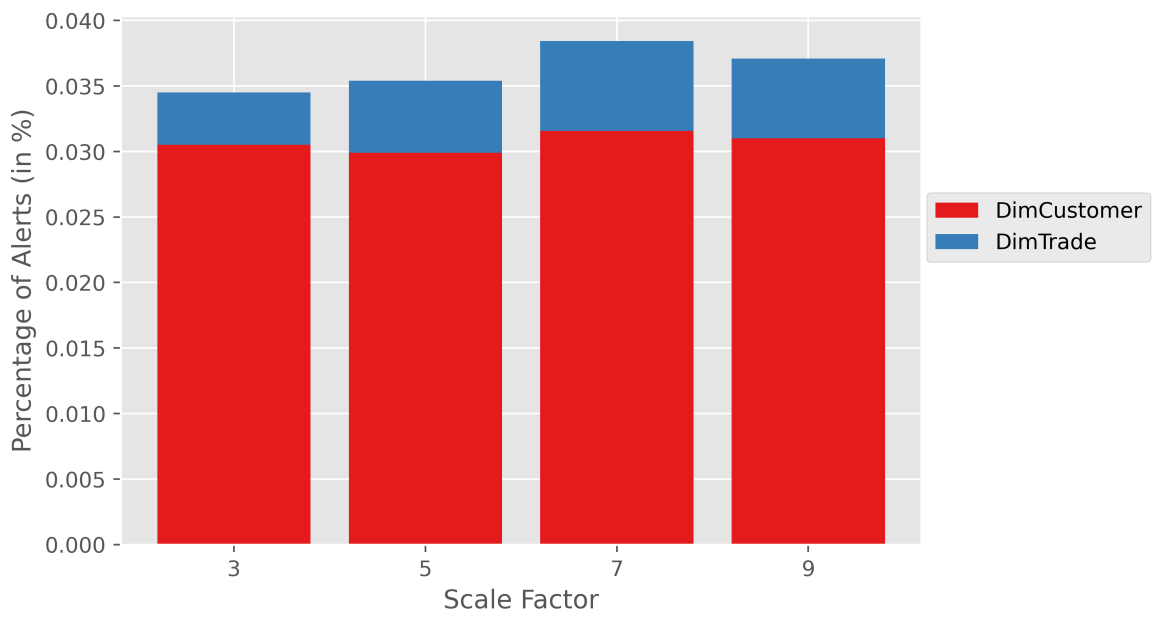


Figure 3.4: Proportion of Records with Alerts by Scale Factor

Technology	SF=3	SF=5
PostgreSQL	00:12:50	00:22:31
OracleSQL	00:03:40*	00:06:45*

Table 3.1: Comparison between PostgreSQL results ([AVZ20]) and OracleSQL.

*: Without the *FactMarketHistory* table.

SF=3	Time
PostgreSQL	00:10:20
OracleSQL	00:03:40

Table 3.2: Comparison between PostgreSQL results ([AVZ20]) and OracleSQL without *FactMarketHistory* table for scale factor 3.

Chapter 4

Conclusion

This project helped us experience first-hand Data Integration for a system. We understood that decision support systems can require input data from various sources with various formats, the amount of detail and constraints. The data needs to be processed in various ways for it to be loaded to a decision support system, often involving combinations of several data sources. It is important to think carefully about the strategy to take in order to create a healthy ETL pipeline and now we are aware of different tools and ideas that are useful for this purpose.

Also, we have experienced the importance of understanding well the operations and transformations used. For instance, we have not been able to come up with an efficient way to recover the values of a row associated with a maximum, which seems like an important operation to have implemented in the system, as it is usual to encounter requirements as those in *FactMarketHistory*.

Further Work

We understand that our implementation of TPC-DI is not complete without the following steps:

- Improving the queries for table *FactMarketHistory*.
- Executing the two phases of *Incremental Updates*.
- Implementing the rest of *Automated Audit Phase* after all the 3 phases are complete.

Bibliography

- [TPC14] Transaction Processing Performance Council (TPC). *TPC BENCHMARK™ DI - Data Integration*. 1.1.0. TPC. Nov. 2014. URL: https://www.tpc.org/tpc_documents_current_versions/current_specifications5.asp.
- [AVZ20] Judith Awiti, Alejandro A. Vaisman, and Esteban Zimányi. “Design and implementation of ETL processes using BPMN and relational algebra”. In: *Data & Knowledge Engineering* 129 (2020), p. 101837. ISSN: 0169-023X. DOI: <https://doi.org/10.1016/j.datak.2020.101837>. URL: <https://www.sciencedirect.com/science/article/pii/S0169023X19306111>.
- [Dat22] Oracle® Database. *Oracle Database Express Edition*. Oracle®. 2022. URL: <https://www.oracle.com/be/database/technologies/appdev/xe.html>.
- [NI18] Dwi Prasetyo Adi Nugroho and Haydar Ali Ismail. *tpcdi-kit*. <https://github.com/haydarai/tpcdi-kit>. 2018.
- [VZ22] Alejandro Vaisman and Esteban Zimányi. *Data Warehouse Systems*. Springer Berlin Heidelberg, 2022. DOI: [10.1007/978-3-662-65167-4](https://doi.org/10.1007/978-3-662-65167-4).

Appendix A

Creation script

In Listing A.1 the script for creating all tables of the database is shown. This script creates the schema explained in TPC-DI documentation adapted to the OracleSQL syntax and restrictions.

```
1 CREATE TABLE BatchDate( BatchNumber NUMBER(3,0),
2     BatchDate DATE NOT NULL
3 );
4
5 CREATE TABLE DimBroker ( SK_BrokerID NUMBER(11,0) GENERATED BY DEFAULT AS IDENTITY
6     PRIMARY KEY,
7     BrokerID NUMBER(11,0) NOT NULL,
8     ManagerID NUMBER(11,0),
9     FirstName CHAR(50) NOT NULL,
10    LastName CHAR(50) NOT NULL,
11    MiddleInitial CHAR(1),
12    Branch CHAR(50),
13    Office CHAR(50),
14    Phone CHAR(14),
15    IsCurrent CHAR(5) NOT NULL CHECK (IsCurrent = 'true' OR IsCurrent = 'false'
16 ),
17    BatchID NUMBER(5,0) NOT NULL,
18    EffectiveDate date NOT NULL,
19    EndDate date NOT NULL
20 );
21
22 CREATE TABLE DimCustomer ( SK_CustomerID NUMBER(11,0) GENERATED BY DEFAULT AS IDENTITY
23     PRIMARY KEY,
24     CustomerID NUMBER(11,0) NOT NULL,
25     TaxID CHAR(20) NOT NULL,
26     Status CHAR(10) NOT NULL,
27     LastName CHAR(30) NOT NULL,
28     FirstName CHAR(30) NOT NULL,
29     MiddleInitial CHAR(1),
30     Gender CHAR(1),
31     Tier NUMBER(1,0),
32     DOB date NOT NULL,
33     AddressLine1 varchar(80) NOT NULL,
34     AddressLine2 varchar(80),
35     PostalCode char(12) NOT NULL,
36     City char(25) NOT NULL,
37     StateProv char(20) NOT NULL,
38     Country char(24),
39     Phone1 char(30),
```

```

37         Phone2    char(30),
38         Phone3    char(30),
39         Email1     char(50),
40         Email2     char(50),
41         NationalTaxRateDesc varchar(50),
42         NationalTaxRate  NUMBER(6,5),
43         LocalTaxRateDesc varchar(50),
44         LocalTaxRate    NUMBER(6,5),
45         AgencyID       char(30),
46         CreditRating   NUMBER(5,0),
47         NetWorth       NUMBER(10),
48         MarketingNameplate varchar(100),
49         IsCurrent     CHAR(5) NOT NULL CHECK (IsCurrent = 'true' OR IsCurrent = 'false'
50     ),
51         BatchID    NUMBER(5,0) NOT NULL,
52         EffectiveDate date NOT NULL,
53         EndDate    date NOT NULL
54 );
55 CREATE TABLE DimAccount ( SK_AccountID  NUMBER(11,0) GENERATED BY DEFAULT AS IDENTITY
56     PRIMARY KEY,
57         AccountID  NUMBER(11,0) NOT NULL,
58         SK_BrokerID  NUMBER(11,0) NOT NULL REFERENCES DimBroker (
59     SK_BrokerID),
60         SK_CustomerID  NUMBER(11,0) NOT NULL REFERENCES DimCustomer (
61     SK_CustomerID),
62         Status          CHAR(10) NOT NULL,
63         AccountDesc      varchar(50),
64         TaxStatus        NUMBER(1,0) NOT NULL CHECK (TaxStatus = 0 OR
65     TaxStatus = 1 OR TaxStatus = 2),
66         IsCurrent       CHAR(5) NOT NULL CHECK (IsCurrent = 'true' OR
67     IsCurrent = 'false'),
68         BatchID         NUMBER(5,0) NOT NULL,
69         EffectiveDate    date NOT NULL,
70         EndDate          date NOT NULL
71 );
72
73 CREATE TABLE DimCompany ( SK_CompanyID  NUMBER(11,0) GENERATED BY DEFAULT AS IDENTITY
74     PRIMARY KEY,
75         CompanyID  NUMBER(11,0) NOT NULL,
76         Status     CHAR(10) Not NULL,
77         Name       CHAR(60) Not NULL,
78         Industry   CHAR(50) Not NULL,
79         SPrating   CHAR(4),
80         isLowGrade CHAR(5) NOT NULL CHECK (isLowGrade = 'true' OR isLowGrade = '
81     false' or isLowGrade = 'NULL'),
82         CEO        CHAR(100) Not NULL,
83         AddressLine1 CHAR(80),
84         AddressLine2 CHAR(80),
85         PostalCode   CHAR(12) Not NULL,
86         City         CHAR(25) Not NULL,
87         StateProv    CHAR(20) Not NULL,
88         Country      CHAR(24),
89         Description  CHAR(150) Not NULL,
90         FoundingDate DATE,
91         IsCurrent    CHAR(5) NOT NULL CHECK (IsCurrent = 'true' OR IsCurrent = 'false'
92     ),
93         BatchID     NUMBER(5,0) Not NULL,
94         EffectiveDate DATE Not NULL,
95         EndDate      DATE Not NULL

```

```

89 );
90
91 CREATE TABLE DimDate ( SK_DateID NUMBER(11,0) PRIMARY KEY,
92     DateValue DATE Not NULL,
93     DateDesc CHAR(20) Not NULL,
94     CalendarYearID NUMBER(4) Not NULL,
95     CalendarYearDesc CHAR(20) Not NULL,
96     CalendarQtrID NUMBER(5) Not NULL,
97     CalendarQtrDesc CHAR(20) Not NULL,
98     CalendarMonthID NUMBER(6) Not NULL,
99     CalendarMonthDesc CHAR(20) Not NULL,
100     CalendarWeekID NUMBER(6) Not NULL,
101     CalendarWeekDesc CHAR(20) Not NULL,
102     DayOfWeeknumeric NUMBER(1) Not NULL,
103     DayOfWeekDesc CHAR(10) Not NULL,
104     FiscalYearID NUMBER(4) Not NULL,
105     FiscalYearDesc CHAR(20) Not NULL,
106     FiscalQtrID NUMBER(5) Not NULL,
107     FiscalQtrDesc CHAR(20) Not NULL,
108     HolidayFlag CHAR(5) NOT NULL CHECK (HolidayFlag = 'true' OR HolidayFlag = '
    false')
109 );
110
111 CREATE TABLE DimSecurity( SK_SecurityID NUMBER(11,0) GENERATED BY DEFAULT AS IDENTITY
    PRIMARY KEY,
112     Symbol CHAR(15) Not NULL,
113     Issue CHAR(6) Not NULL,
114     Status CHAR(10) Not NULL,
115     Name CHAR(70) Not NULL,
116     ExchangeID CHAR(6) Not NULL,
117     SK_CompanyID NUMBER(11,0) Not NULL REFERENCES DimCompany (SK_CompanyID),
118     SharesOutstanding NUMBER(12,0) Not NULL,
119     FirstTrade DATE Not NULL,
120     FirstTradeOnExchange DATE Not NULL,
121     Dividend NUMBER(10,2) Not NULL,
122     IsCurrent CHAR(5) NOT NULL check (IsCurrent = 'false' or IsCurrent = 'true'
    ),
123     BatchID NUMBER(5) Not NULL,
124     EffectiveDate DATE Not NULL,
125     EndDate DATE Not NULL
126 );
127
128 CREATE TABLE DimTime ( SK_TimeID NUMBER(11,0) PRIMARY KEY,
129     TimeValue DATE Not NULL,
130     HourID NUMBER(2) Not NULL,
131     HourDesc CHAR(20) Not NULL,
132     MinuteID NUMBER(2) Not NULL,
133     MinuteDesc CHAR(20) Not NULL,
134     SecondID NUMBER(2) Not NULL,
135     SecondDesc CHAR(20) Not NULL,
136     MarketHoursFlag CHAR(5) NOT NULL check (MarketHoursFlag = 'false' or
    MarketHoursFlag = 'true'),
137     OfficeHoursFlag CHAR(5) NOT NULL check (OfficeHoursFlag = 'false' or
    OfficeHoursFlag = 'true')
138 );
139
140 CREATE TABLE DimTrade ( TradeID NUMBER(11,0) Not NULL,
141     SK_BrokerID NUMBER(11,0) REFERENCES DimBroker (SK_BrokerID),
142     SK_CreateDateID NUMBER(11,0) REFERENCES DimDate (SK_DateID),
143     SK_CreateTimeID NUMBER(11,0) REFERENCES DimTime (SK_TimeID),
144     SK_CloseDateID NUMBER(11,0) REFERENCES DimDate (SK_DateID),

```

```

145         SK_CloseTimeID NUMBER(11,0) REFERENCES DimTime (SK_TimeID),
146         Status CHAR(10) Not NULL,
147         Type CHAR(12) Not NULL,
148         CashFlag CHAR(5) NOT NULL check (CashFlag = 'false' or CashFlag = 'true'),
149         SK_SecurityID NUMBER(11,0) Not NULL REFERENCES DimSecurity (SK_SecurityID),
150         SK_CompanyID NUMBER(11,0) Not NULL REFERENCES DimCompany (SK_CompanyID),
151         Quantity NUMBER(6,0) Not NULL,
152         BidPrice NUMBER(8,2) Not NULL,
153         SK_CustomerID NUMBER(11,0) Not NULL REFERENCES DimCustomer (SK_CustomerID),
154         SK_AccountID NUMBER(11,0) Not NULL REFERENCES DimAccount (SK_AccountID),
155         ExecutedBy CHAR(64) Not NULL,
156         TradePrice NUMBER(8,2),
157         Fee NUMBER(10,2),
158         Commission NUMBER(10,2),
159         Tax NUMBER(10,2),
160         BatchID NUMBER(5) Not Null
161     );
162
163     CREATE TABLE DImessages ( MessageDateAndTime TIMESTAMP Not NULL,
164         BatchID NUMBER(5) Not NULL,
165         MessageSource CHAR(30),
166         MessageText CHAR(50) Not NULL,
167         MessageType CHAR(12) Not NULL,
168         MessageData CHAR(100)
169     );
170
171     CREATE TABLE FactCashBalances ( SK_CustomerID NUMBER(11,0) Not Null REFERENCES
172         DimCustomer (SK_CustomerID),
173         SK_AccountID NUMBER(11,0) Not Null REFERENCES DimAccount (SK_AccountID),
174         SK_DateID NUMBER(11,0) Not Null REFERENCES DimDate (SK_DateID),
175         Cash NUMBER(15,2) Not Null,
176         BatchID NUMBER(5)
177     );
178
179     CREATE TABLE FactHoldings ( TradeID NUMBER(11,0) Not Null,
180         CurrentTradeID NUMBER(11,0) Not Null,
181         SK_CustomerID NUMBER(11,0) Not NULL REFERENCES DimCustomer (SK_CustomerID),
182         SK_AccountID NUMBER(11,0) Not NULL REFERENCES DimAccount (SK_AccountID),
183         SK_SecurityID NUMBER(11,0) Not NULL REFERENCES DimSecurity (SK_SecurityID),
184         SK_CompanyID NUMBER(11,0) Not NULL REFERENCES DimCompany (SK_CompanyID),
185         SK_DateID NUMBER(11,0) Not NULL REFERENCES DimDate (SK_DateID),
186         SK_TimeID NUMBER(11,0) Not NULL REFERENCES DimTime (SK_TimeID),
187         CurrentPrice NUMBER(8,2) CHECK (CurrentPrice > 0) ,
188         CurrentHolding NUMBER(6) Not NULL,
189         BatchID NUMBER(5)
190     );
191
192     CREATE TABLE FactMarketHistory ( SK_SecurityID NUMBER(11,0) Not Null REFERENCES
193         DimSecurity (SK_SecurityID),
194         SK_CompanyID NUMBER(11,0) Not Null REFERENCES DimCompany (SK_CompanyID)
195         ,
196         SK_DateID NUMBER(11,0) Not Null REFERENCES DimDate (SK_DateID),
197         PERatio NUMBER(10,2),
198         Yield NUMBER(5,2) Not Null,
199         FiftyTwoWeekHigh NUMBER(8,2) Not Null,
200         SK_FiftyTwoWeek NUMBER(11,0) Not Null,
201         FiftyTwoWeekLow NUMBER(8,2) Not Null,
202         SK_FiftyTwoWeekL NUMBER(11,0) Not Null,
203         ClosePrice NUMBER(8,2) Not Null,
204         DayHigh NUMBER(8,2) Not Null,
205         DayLow NUMBER(8,2) Not Null,

```

```
203         Volume NUMBER(12) Not Null,
204         BatchID NUMBER(5)
205     );
206
207 CREATE TABLE FactWatches ( SK_CustomerID NUMBER(11,0) Not NULL REFERENCES DimCustomer (
    SK_CustomerID),
208         SK_SecurityID NUMBER(11,0) Not NULL REFERENCES DimSecurity (SK_SecurityID),
209         SK_DateID_DatePlaced NUMBER(11,0) Not NULL REFERENCES DimDate (SK_DateID),
210         SK_DateID_DateRemoved NUMBER(11,0) REFERENCES DimDate (SK_DateID),
211         BatchID NUMBER(5) Not Null
212     );
213
214 CREATE TABLE Industry ( IN_ID CHAR(2) Not NULL,
215         IN_NAME CHAR(50) Not NULL,
216         IN_SC_ID CHAR(4) Not NULL
217     );
218
219 CREATE TABLE Financial ( SK_CompanyID NUMBER(11,0) Not NULL REFERENCES DimCompany (
    SK_CompanyID),
220         FI_YEAR NUMBER(4) Not NULL,
221         FI_QTR NUMBER(1) Not NULL,
222         FI_QTR_START_DATE DATE Not NULL,
223         FI_REVENUE NUMBER(15,2) Not NULL,
224         FI_NET_EARN NUMBER(15,2) Not NULL,
225         FI_BASIC_EPS NUMBER(10,2) Not NULL,
226         FI_DILUT_EPS NUMBER(10,2) Not NULL,
227         FI_MARGIN NUMBER(10,2) Not NULL,
228         FI_INVENTORY NUMBER(15,2) Not NULL,
229         FI_ASSETS NUMBER(15,2) Not NULL,
230         FI_LIABILITY NUMBER(15,2) Not NULL,
231         FI_OUT_BASIC NUMBER(12) Not NULL,
232         FI_OUT_DILUT NUMBER(12) Not NULL
233     );
234
235 CREATE TABLE Prospect ( AgencyID CHAR(30) NOT NULL UNIQUE,
236         SK_RecordDateID NUMBER(11,0) NOT NULL,
237         SK_UpdateDateID NUMBER(11,0) NOT NULL REFERENCES DimDate (SK_DateID),
238         BatchID NUMBER(5) NOT NULL,
239         IsCustomer CHAR(5) NOT NULL check (IsCustomer = 'false' or IsCustomer = 'true
    '),
240         LastName CHAR(30) NOT NULL,
241         FirstName CHAR(30) NOT NULL,
242         MiddleInitial CHAR(1),
243         Gender CHAR(1),
244         AddressLine1 CHAR(80),
245         AddressLine2 CHAR(80),
246         PostalCode CHAR(12),
247         City CHAR(25) NOT NULL,
248         State CHAR(20) NOT NULL,
249         Country CHAR(24),
250         Phone CHAR(30),
251         Income NUMBER(9),
252         numericberCars NUMBER(2),
253         numericberChildren NUMBER(2),
254         MaritalStatus CHAR(1),
255         Age NUMBER(3),
256         CreditRating NUMBER(4),
257         OwnOrRentFlag CHAR(1),
258         Employer CHAR(30),
259         numericberCreditCards NUMBER(2),
260         NetWorth NUMBER(12),
```

```
261         MarketingNameplate CHAR(100)
262     );
263
264     CREATE TABLE StatusType ( ST_ID CHAR(4) Not NULL,
265                                ST_NAME CHAR(10) Not NULL
266     );
267
268     CREATE TABLE TaxRate ( TX_ID CHAR(4) Not NULL,
269                             TX_NAME CHAR(50) Not NULL,
270                             TX_RATE NUMBER(6,5) Not NULL
271     );
272
273     CREATE TABLE TradeType ( TT_ID CHAR(3) Not NULL,
274                               TT_NAME CHAR(12) Not NULL,
275                               TT_IS_SELL NUMBER(1) Not NULL,
276                               TT_IS_MRKT NUMBER(1) Not NULL
277     );
278
279
280     CREATE TABLE Audit_ ( DataSet CHAR(20) Not Null,
281                            BatchID NUMBER(5),
282                            Date_ DATE,
283                            Attribute CHAR(50) not null,
284                            Value NUMBER(15),
285                            DValue NUMBER(15,5)
286     );
287
288     -- staging tables
289
290     CREATE TABLE S_Company (
291         PTS CHAR(15) NOT NULL,
292         REC_TYPE CHAR(3) NOT NULL,
293         COMPANY_NAME CHAR(60),
294         CIK CHAR(10) NOT NULL,
295         STATUS CHAR(4) NOT NULL,
296         INDUSTRY_ID CHAR(2) NOT NULL,
297         SP_RATING CHAR(4),
298         FOUNDING_DATE CHAR(8),
299         ADDR_LINE_1 CHAR(80),
300         ADDR_LINE_2 CHAR(80),
301         POSTAL_CODE CHAR(12),
302         CITY CHAR(25),
303         STATE_PROVINCE CHAR(20),
304         COUNTRY CHAR(24),
305         CEO_NAME CHAR(46),
306         DESCRIPTION CHAR(150)
307     );
308
309     CREATE TABLE S_Security (
310         PTS CHAR(15) NOT NULL,
311         REC_TYPE CHAR(3) NOT NULL,
312         SYMBOL CHAR(15) NOT NULL,
313         ISSUE_TYPE CHAR(6) NOT NULL,
314         STATUS CHAR(4) NOT NULL,
315         NAME CHAR(70) NOT NULL,
316         EX_ID CHAR(6) NOT NULL,
317         SH_OUT CHAR(13) NOT NULL,
318         FIRST_TRADE_DATE CHAR(8) NOT NULL,
319         FIRST_TRADE_EXCHANGE CHAR(8) NOT NULL,
320         DIVIDEN CHAR(12) NOT NULL,
321         COMPANY_NAME_OR_CIK CHAR(60) NOT NULL
```

```
322 );
323
324 CREATE TABLE S_Financial(
325     PTS CHAR(15),
326     REC_TYPE CHAR(3),
327     YEAR CHAR(4),
328     QUARTER CHAR(1),
329     QTR_START_DATE CHAR(8),
330     POSTING_DATE CHAR(8),
331     REVENUE CHAR(17),
332     EARNINGS CHAR(17),
333     EPS CHAR(12),
334     DILUTED_EPS CHAR(12),
335     MARGIN CHAR(12),
336     INVENTORY CHAR(17),
337     ASSETS CHAR(17),
338     LIABILITIES CHAR(17),
339     SH_OUT CHAR(13),
340     DILUTED_SH_OUT CHAR(13),
341     CO_NAME_OR_CIK CHAR(60)
342 );
343
344 CREATE TABLE S_Prospect(
345     AGENCY_ID CHAR(30) NOT NULL,
346     LAST_NAME CHAR(30) NOT NULL,
347     FIRST_NAME CHAR(30) NOT NULL,
348     MIDDLE_INITIAL CHAR(1),
349     GENDER CHAR(1),
350     ADDRESS_LINE_1 CHAR(80),
351     ADDRESS_LINE_2 CHAR(80),
352     POSTAL_CODE CHAR(12),
353     CITY CHAR(25) NOT NULL,
354     STATE CHAR(20) NOT NULL,
355     COUNTRY CHAR(24),
356     PHONE CHAR(30),
357     INCOME NUMBER(9),
358     NUMBER_CARS NUMBER(2),
359     NUMBER_CHILDREN NUMBER(2),
360     MARITAL_STATUS CHAR(1),
361     AGE NUMBER(3),
362     CREDIT_RATING NUMBER(4),
363     OWN_OR_RENT_FLAG CHAR(1),
364     EMPLOYER CHAR(30),
365     NUMBER_CREDIT_CARDS NUMBER(2),
366     NET_WORTH NUMBER(12)
367 );
368
369 CREATE TABLE S_Watches(
370     W_C_ID INTEGER NOT NULL,
371     W_S_SYMB CHAR(15) NOT NULL,
372     W_DTS DATE NOT NULL,
373     W_ACTION CHAR(4) NOT NULL
374 );
375
376 CREATE TABLE S_Cash_Balances(
377     CT_CA_ID INTEGER NOT NULL,
378     CT_DTS DATE NOT NULL,
379     CT_AMT CHAR(20) NOT NULL,
380     CT_NAME CHAR(100) NOT NULL
381 );
382
```



```
383 CREATE TABLE S_Broker(  
384     EmployeeID INTEGER NOT NULL,  
385     ManagerID INTEGER NOT NULL,  
386     EmployeeFirstName CHAR(30) NOT NULL,  
387     EmployeeLastName CHAR(30) NOT NULL,  
388     EmployeeMI CHAR(1),  
389     EmployeeJobCode NUMBER(3),  
390     EmployeeBranch CHAR(30),  
391     EmployeeOffice CHAR(10),  
392     EmployeePhone CHAR(14)  
393 );  
394  
395 CREATE TABLE S_Account ( SK_AccountID NUMBER(11,0),  
396     AccountID NUMBER(11,0),  
397     BrokerID NUMBER(11,0),  
398     CustomerID NUMBER(11,0),  
399     Status CHAR(10) NOT NULL,  
400     AccountDesc VARCHAR(50),  
401     TaxStatus NUMBER(1,0),  
402     BatchID NUMBER(5,0),  
403     EffectiveDate DATE,  
404     EndDate DATE,  
405     ActionType CHAR(10)  
406 );  
407  
408 CREATE TABLE S_Customer ( ActionType CHAR(10),  
409     CustomerID NUMBER(11,0),  
410     TaxID CHAR(20),  
411     Status CHAR(10),  
412     LastName CHAR(30),  
413     FirstName CHAR(30),  
414     MiddleInitial CHAR(1),  
415     Gender CHAR(1),  
416     Tier NUMBER(1,0),  
417     DOB DATE,  
418     AddressLine1 VARCHAR(80),  
419     AddressLine2 VARCHAR(80),  
420     PostalCode CHAR(12),  
421     City CHAR(25),  
422     StateProv CHAR(20),  
423     Country CHAR(24),  
424     Phone1 CHAR(30),  
425     Phone2 CHAR(30),  
426     Phone3 CHAR(30),  
427     Email1 CHAR(50),  
428     Email2 CHAR(50),  
429     NationalTaxRateDesc VARCHAR(50),  
430     NationalTaxRate NUMBER(6,5),  
431     LocalTaxRateDesc VARCHAR(50),  
432     LocalTaxRate NUMBER(6,5),  
433     BatchID NUMBER(5,0),  
434     EffectiveDate DATE,  
435     EndDate DATE  
436 );  
437  
438 CREATE TABLE S_Trade (  
439     cdc_flag CHAR(1),  
440     cdc_dsn NUMBER(12),  
441     t_id NUMBER(15),  
442     t_dts DATE,  
443     t_st_id CHAR(4),
```

```
444     t_tt_id CHAR(3),
445     t_is_cash CHAR(3),
446     t_s_symb CHAR(15) NOT NULL,
447     t_qty NUMBER(6) NOT NULL,
448     t_bid_price NUMBER(8),
449     t_ca_id NUMBER(11),
450     t_exec_name CHAR(49),
451     t_trade_price NUMBER(8),
452     t_chrg NUMBER(10),
453     t_comm NUMBER(10),
454     t_tax NUMBER(10)
455 );
456
457 CREATE TABLE S_Trade_History (
458     th_t_id NUMBER(15),
459     th_dts DATE,
460     th_st_id CHAR(4)
461 );
462
463 exit
```

Listing A.1: oracle-schema.sql

Appendix B

Control files

In this appendix, we show all control files, in the order of appearance in the text.

DimDate

See [2.2.2](#).

```
1 LOAD DATA
2 INTO TABLE DimDate
3 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4 (
5     SK_DateID,
6     DateValue DATE "YYYY-MM-DD",
7     DateDesc,
8     CalendarYearID,
9     CalendarYearDesc,
10    CalendarQtrID,
11    CalendarQtrDesc,
12    CalendarMonthID,
13    CalendarMonthDesc,
14    CalendarWeekID,
15    CalendarWeekDesc,
16    DayOfWeeknumeric,
17    DayOfWeekDesc,
18    FiscalYearID,
19    FiscalYearDesc,
20    FiscalQtrID,
21    FiscalQtrDesc,
22    HolidayFlag
23 )
```

Listing B.1: DimDate.ctl

DimTime

```
1 LOAD DATA
2 INTO TABLE DimTime
3 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4 (
5     SK_TimeID,
6     TimeValue DATE "HH24:MI:SS",
7     HourID,
```

```
8      HourDesc ,
9      MinuteID ,
10     MinuteDesc ,
11     SecondID ,
12     SecondDesc ,
13     MarketHoursFlag ,
14     OfficeHoursFlag
15 )
```

Listing B.2: DimTime.ctf

Industry

See [2.2.2](#).

```
1 LOAD DATA
2 INFILE
3 TRUNCATE
4 INTO TABLE Industry
5 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
6 (
7     IN_ID ,
8     IN_NAME ,
9     IN_SC_ID
10 )
```

Listing B.3: Industry.ctf

StatusType

See [2.2.2](#).

```
1 LOAD DATA
2 INFILE
3 TRUNCATE
4 INTO TABLE StatusType
5 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
6 (
7     ST_ID ,
8     ST_NAME
9 )
```

Listing B.4: StatusType.ctf

TaxRate

See [2.2.2](#).

```
1 LOAD DATA
2 INFILE
3 TRUNCATE
4 INTO TABLE TaxRate
5 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
6 (
7     TX_ID ,
8     TX_NAME ,
9     TX_RATE
10 )
```

Listing B.5: TaxRate.ctf

TradeType

See [2.2.2](#).

```
1 LOAD DATA
2 INFILE
3 TRUNCATE
4 INTO TABLE TradeType
5 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
6 (
7     TT_ID,
8     TT_NAME,
9     TT_IS_SELL,
10    TT_IS_MRKT
11 )
```

Listing B.6: TradeType.ctl

Broker

See [2.2.2](#).

```
1 LOAD DATA
2 INTO TABLE S_Broker
3 FIELDS TERMINATED BY ',' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4 (
5     EmployeeID,
6     ManagerID,
7     EmployeeFirstName,
8     EmployeeLastName,
9     EmployeeMI,
10    EmployeeJobCode,
11    EmployeeBranch,
12    EmployeeOffice,
13    EmployeePhone
14 )
```

Listing B.7: Broker.ctl

Prospect

See [2.2.2](#).

```
1 LOAD DATA
2 INTO TABLE S_Prospect
3 FIELDS TERMINATED BY ',' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4 (
5     AGENCY_ID,
6     LAST_NAME,
7     FIRST_NAME,
8     MIDDLE_INITIAL,
9     GENDER,
10    ADDRESS_LINE_1,
11    ADDRESS_LINE_2,
12    POSTAL_CODE,
13    CITY,
14    STATE,
15    COUNTRY,
16    PHONE,
17    INCOME,
```

```
18  NUMBER_CARS ,
19  NUMBER_CHILDREN ,
20  MARITAL_STATUS ,
21  AGE ,
22  CREDIT_RATING ,
23  OWN_OR_RENT_FLAG ,
24  EMPLOYER ,
25  NUMBER_CREDIT_CARDS ,
26  NET_WORTH
27 )
```

Listing B.8: Prospect.ctl

DimTrade

See [2.2.2](#).

```
1  LOAD DATA
2  INTO TABLE S_Trade
3  FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4  (
5      t_id,
6      t_dts DATE "YYYY-MM-DD HH24:MI:SS",
7      t_st_id,
8      t_tt_id,
9      t_is_cash,
10     t_s_symb,
11     t_qty,
12     t_bid_price,
13     t_ca_id,
14     t_exec_name,
15     t_trade_price,
16     t_chrg,
17     t_comm,
18     t_tax
19 )
```

Listing B.9: Trade.ctl

FactCashBalances

See [2.2.2](#).

```
1  LOAD DATA
2  INTO TABLE S_Cash_Balances
3  FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4  (
5      CT_CA_ID,
6      CT_DTS DATE "YYYY-MM-DD HH24:MI:SS",
7      CT_AMT,
8      CT_NAME
9  )
```

Listing B.10: CashBalances.ctl

FactWatches

See [2.2.2](#).

```
1 LOAD DATA
2 INTO TABLE S_Watches
3 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4 (
5     W_C_ID,
6     W_S_SYMB,
7     W_DTS DATE "YYYY-MM-DD HH24:MI:SS",
8     W_ACTION
9 )
```

Listing B.11: Watches.ctf

FactHoldings

See [2.2.2](#).

```
1 LOAD DATA
2 INTO TABLE S_Holdings
3 FIELDS TERMINATED BY '|' OPTIONALLY ENCLOSED BY '"' TRAILING NULLCOLS
4 (
5     HH_H_T_ID,
6     HH_T_ID,
7     HH_BEFORE_QTY,
8     HH_AFTER_QTY
9 )
```

Listing B.12: Holdings.ctf