# TPC—DI Project

| | |
|---|---|
| 📅 Dates | @December 24, 2024 |
| ⊙ Course | INFO-H419 |
| ☑ Status | ☐ |
| ≔ Task | |

## DATA GENERATION

我不知道為什麼但是只有退到用java8才能run

sf 3 =

AuditTotalRecordsSummaryWriter - TotalRecords for Batch1: 4539962
AuditTotalRecordsSummaryWriter - TotalRecords for Batch2: 19900
AuditTotalRecordsSummaryWriter - TotalRecords for Batch3: 19855
AuditTotalRecordsSummaryWriter - TotalRecords all Batches: 4579717 695794.14 records/second

Overall time    0h:00m:06s:581ms
Generated       267.1 MiB
Speed           40.6 MiB/s

sf 6 =

AuditTotalRecordsSummaryWriter - TotalRecords for Batch1: 9443149
AuditTotalRecordsSummaryWriter - TotalRecords for Batch2: 40317
AuditTotalRecordsSummaryWriter - TotalRecords for Batch3: 40283
AuditTotalRecordsSummaryWriter - TotalRecords all Batches: 9523749 1074308.97 records/second

Overall time    0h:00m:08s:867ms
Generated       564.2 MiB
Speed           63.6 MiB/s

sf 12 =

AuditTotalRecordsSummaryWriter - TotalRecords for Batch1: 19242651
AuditTotalRecordsSummaryWriter - TotalRecords for Batch2: 81117
AuditTotalRecordsSummaryWriter - TotalRecords for Batch3: 81095
AuditTotalRecordsSummaryWriter - TotalRecords all Batches: 19404863
1466621.04 records/second

Overall time    0h:00m:13s:232ms
Generated       1.1 GiB
Speed           87.6 MiB/s

sf 24 =

AuditTotalRecordsSummaryWriter - TotalRecords for Batch1: 38842886
AuditTotalRecordsSummaryWriter - TotalRecords for Batch2: 162844
AuditTotalRecordsSummaryWriter - TotalRecords for Batch3: 162628
AuditTotalRecordsSummaryWriter - TotalRecords all Batches: 39168358
1780218.07 records/second

Overall time    0h:00m:22s:001ms
Generated       2.3 GiB
Speed           107.0 MiB/s

以下都是手動ETL只針對sf3的的code，airflow automate的時候注意看path

## CREATE DATABASE AND SCHEMAS

```
MariaDB [(none)]> CREATE DATABASE tpc_di;
```

```
MariaDB [(none)]> USE tpc_di;
```

```
MariaDB [tpc_di]> CREATE SCHEMA staging;
```

```
MariaDB [tpc_di]> CREATE SCHEMA master;
```

# DESIGN TABLES

```
MariaDB [master]> SHOW TABLES;
+------------------+
| Tables_in_master |
+------------------+
| audit            |
| dimaccount       |
| dimbroker        |
| dimcompany       |
| dimcustomer      |
| dimdate          |
| dimessages       |
| dimsecurity      |
| dimtime          |
| dimtrade         |
| factcashbalances |
| factholdings     |
| factmarkethistory |
| factwatches      |
| financial        |
| industry         |
| prospect         |
| statustype       |
| taxrate          |
| tradetype        |
+------------------+
20 rows in set (0.001 sec)
```

```
MariaDB [staging]> show tables;
+-------------------+
| Tables_in_staging |
+-------------------+
| audit             |
| batchdate         |
| cashtransaction   |
| customermgmt      |
| dailymarket       |
| date              |
| finwire           |
| finwire_cmp       |
| finwire_fin       |
| finwire_sec       |
| holdinghistory    |
| hr                |
| industry          |
| prospect          |
| statustype        |
| taxrate           |
| time              |
| trade             |
| tradehistory      |
| tradetype         |
| watchhistory      |
+-------------------+
21 rows in set (0.000 sec)
```

# LOAD DATA INTO STAGING

Error populating staging tables: (1292, "Incorrect datetime value: '2013-03-31 02:18:26' for column `staging`.`cashtransaction`.`ct_dts` at row 39188")

The datetime value `2013-03-31 02:18:26` might fall into a range that is invalid due to Daylight Saving Time (DST) changes or input formatting. For example, in many regions, `2013-03-31 02:18:26` would be invalid because clocks jumped from `01:59:59` to `03:00:00`.

MariaDB may reject certain datetime values if they fall into invalid ranges due to DST. Address this by:

```
SET GLOBAL sql_mode = '';
SET SESSION sql_mode = '';
```

run load_staging.py + populate_staging_customermgmt.py

# TRANSFORM AND LOAD TO MASTER

1. Historical load phase: Transforms and loads initial data from Batch1 (timed).

2. Incremental update 1: Loads Batch2 data incrementally (timed).

3. Incremental update 2: Loads Batch3 data incrementally (timed).

4. Automated audit: Validates results after final update.

我在github看了一圈，沒有任何組曾經做過234，這後面的工作量超級無敵不合理的恐怖，我就默認我們不需要做了。https://github.com/QasimKhan5x/TPC-DI/blob/main/scripts/incremental.py 這有個兩千多行的code可以改著用,可以直接在report説我們refer to這個，如果我們打算做的話。。。(如果我沒有update說明我放棄了（99%我會放棄）)

## static tables = These tables are loaded during the Historical Load phase and are not modified afterward:

run 1.load_master_static.sql

- **StatusType**: From `StatusType.txt`

- **TaxRate**: From `TaxRate.txt`

- **TradeType**: From `TradeType.txt`

- **DimDate**: From `Date.txt`

- **DimTime**: From `Time.txt`

- **Industry**: From `Industry.txt`

## non-static tables = These tables are updated dynamically through history tracking or incremental updates:

- **DimCustomer**: Updated based on incoming data and history tracking mechanisms.

- **DimAccount**: Implements history tracking updates.

- **DimBroker**: History tracked but rarely updated in the benchmark.

- **DimSecurity**: Obtains updates from FINWIRE files through history tracking.

- **DimCompany**: Updated as needed from FINWIRE files.

- **Fact Tables**: Such as `FactHoldings` , `FactMarketHistory` , `FactCashBalances` , and `FactWatches` are regularly updated based on transactions and other inputs

**other Intermediate tables:**

- **Prospect**
    - Temporary staging table for potential customer data.
    - Transforms data into `DimCustomer` through periodic updates.
- **DimMessages**
    - Logging table for system events.
    - Stores ETL process messages and alerts.
- **financial**
    - serves as a bridge between staging data ( `staging.finwire_fin` ) and fact tables ( `factmarkethistory` )

run all 2-18 sql files in Historical Transformation file.

**Ordered Execution Steps:**

1. **Load static tables**:
2. **Load** `master.dimcompany` :
    - Required for `dimsecurity` , `financial` , `factmarkethistory` , and `dimessages` `dimcompany` .
    - **Records: 4500**
    - **manually: 0.258 sec**
3. **Load** `master.dimessages dimcompany` :
    - Uses `dimcompany` .
    - **Records: 450**
    - **manually: 0.010 sec**
4. **Transform and load** `master.dimbroker` :
    - Depends on `dimdate` .

- **Records: 4293**
- **manually: 0.044 sec**

5. **Transform and load** `master.prospect` :

   - Depends on `dimdate` .
   - **Records: 14981**
   - **manually: 1 min 22.277 sec**

6. **Transform and load** `master.dimcustomer` :

   - Depends on `taxrate` and `prospect` .
   - **Records: 6440**
   - **manually: 1.223 sec**

7. **Load** `master.dimessages dimcustomer` :

   - Depends on `dimcustomer` .
   - **Records: 1720**
   - **manually:0.016 sec**

8. **Update** `master.prospect` :

   - Depends on `dimcustomer` .
   - **Rows matched: 3377  Changed: 3377**
   - **manually: 16.741 sec**

9. **Transform and load** `master.dimaccount` :

   - Depends on `dimbroker` and `dimcustomer` .
   - **Records: 11261**
   - **manually: 2 min 9.531 sec**

10. **Transform and load** `master.dimsecurity` :

    - Depends on `dimcompany` .
    - **Records: 8838**
    - **manually: 6.715 sec**

11. **Transform and load** `master.dimtrade` :

- Depends on `dimaccount` , `dimsecurity` , `statustype` , and `tradetype` .

- **Records: 3721368**

- **manually: 3 min 40.046 sec**

- query optimisation from 1st(takes more than 30mins) to 2nd(index, reduced row processed, avoid full table scan on trade)

```
+----+-------------+-----------+------+---------------+------+---------+----------------+--------+-----------------------------------------------------+
| id | select_type | table     | type | possible_keys | key  | key_len | ref            | rows   | Extra                                               |
+----+-------------+-----------+------+---------------+------+---------+----------------+--------+-----------------------------------------------------+
|  1 | PRIMARY     | tt        | ALL  | NULL          | NULL | NULL    | NULL           | 5      |                                                     |
|  1 | PRIMARY     | st        | ALL  | NULL          | NULL | NULL    | NULL           | 6      | Using join buffer (flat, BNL join)                  |
|  1 | PRIMARY     | s         | ALL  | NULL          | NULL | NULL    | NULL           | 8963   | Using join buffer (incremental, BNL join)           |
|  1 | PRIMARY     | a         | ALL  | NULL          | NULL | NULL    | NULL           | 11394  | Using join buffer (incremental, BNL join)           |
|  1 | PRIMARY     | t         | ALL  | NULL          | NULL | NULL    | NULL           | 391100 | Using where; Using join buffer (incremental, BNL join) |
|  1 | PRIMARY     | <derived3>| ref  | key0          | key0 | 8       | staging.t.t_id | 10     | Using where                                         |
|  1 | PRIMARY     | <derived2>| ref  | key0          | key0 | 8       | staging.t.t_id | 10     |                                                     |
|  3 | DERIVED     | th        | ALL  | NULL          | NULL | NULL    | NULL           | 982497 | Using temporary                                     |
|  2 | DERIVED     | th        | ALL  | NULL          | NULL | NULL    | NULL           | 982497 | Using temporary; Using filesort                     |
+----+-------------+-----------+------+---------------+------+---------+----------------+--------+-----------------------------------------------------+
```

```
+----+-------------+-----------+------+-----------------------------+---------------------+---------+------------------------+--------+------------------------------------------+
| id | select_type | table     | type | possible_keys               | key                 | key_len | ref                    | rows   | Extra                                    |
+----+-------------+-----------+------+-----------------------------+---------------------+---------+------------------------+--------+------------------------------------------+
|  1 | PRIMARY     | tt        | ALL  | NULL                        | NULL                | NULL    | NULL                   | 5      |                                          |
|  1 | PRIMARY     | st        | ALL  | NULL                        | NULL                | NULL    | NULL                   | 6      | Using join buffer (flat, BNL join)       |
|  1 | PRIMARY     | s         | ALL  | idx_dimsecurity_symbol      | NULL                | NULL    | NULL                   | 8963   | Using join buffer (incremental, BNL join)|
|  1 | PRIMARY     | t         | ref  | idx_trade_t_ca_id,idx_trade_t_s_symb | idx_trade_t_s_symb | 60    | master.s.symbol        | 247    | Using where                              |
|  1 | PRIMARY     | a         | ref  | idx_dimaccount_accountid    | idx_dimaccount_accountid | 5  | staging.t.t_ca_id      | 4      |                                          |
|  1 | PRIMARY     | <derived3>| ref  | key0                        | key0                | 8       | staging.t.t_id         | 10     | Using where                              |
|  1 | PRIMARY     | <derived2>| ref  | key0                        | key0                | 8       | staging.t.t_id         | 10     |                                          |
|  3 | DERIVED     | th        | ALL  | idx_tradehistory_th_t_id,idx_t_id | NULL          | NULL    | NULL                   | 982497 | Using where; Using temporary             |
|  2 | DERIVED     | th        | ALL  | idx_tradehistory_th_t_id,idx_t_id | NULL          | NULL    | NULL                   | 982497 | Using where; Using temporary; Using filesort |
+----+-------------+-----------+------+-----------------------------+---------------------+---------+------------------------+--------+------------------------------------------+
```

12. **Load** `master.dimessages dimtrade` :

- Depends on `dimtrade` .

- **Records: 312**

- **manually: 2.311 sec**

13. **Transform and load** `master.financial` :

- Depends on `dimcompany` .

- **Records: 496617**

- **manually: 2 min 25.753 sec**

14. **Transform and load** `master.factcashbalances` :

- Depends on `dimaccount` and `dimdate` .

- **Records: 564461**

- **manually: 18 min 55.160 sec**

15. **Transform and load** `master.factholdings` :

   - Depends on `dimtrade` .

   - **Records: 3485814**

   - **manually: 10.837 sec**

   - query optimisation from 1st(takes more than 1h) to 2nd(index, avoid full table scan on holdinghistory)

```
+------+-------------+-------+------+---------------+------+---------+------+---------+-----------------------------------------------+
| id   | select_type | table | type | possible_keys | key  | key_len | ref  | rows    | Extra                                         |
+------+-------------+-------+------+---------------+------+---------+------+---------+-----------------------------------------------+
|    1 | SIMPLE      | h     | ALL  | NULL          | NULL | NULL    | NULL | 362529  |                                               |
|    1 | SIMPLE      | t     | ALL  | NULL          | NULL | NULL    | NULL | 3721227 | Using where; Using join buffer (flat, BNL join) |
+------+-------------+-------+------+---------------+------+---------+------+---------+-----------------------------------------------+
```

```
ON h.hh_t_id = t.tradeid;
+------+-------------+-------+------+---------------+-----------+---------+-----------------+---------+-----------------------+
| id   | select_type | table | type | possible_keys | key       | key_len | ref             | rows    | Extra                 |
+------+-------------+-------+------+---------------+-----------+---------+-----------------+---------+-----------------------+
|    1 | SIMPLE      | t     | ALL  | idx_tradeid   | NULL      | NULL    | NULL            | 3721311 |                       |
|    1 | SIMPLE      | h     | ref  | idx_hh_t_id   | idx_hh_t_id | 7     | master.t.tradeid | 1       | Using index condition |
+------+-------------+-------+------+---------------+-----------+---------+-----------------+---------+-----------------------+
```

16. **Transform and load** `master.factwatches` :

   - Depends on `dimcustomer` , `dimsecurity` , and `dimdate` .

   - **Records: 1532010**

   - **manually: 4.846 sec (** optimised with indexes on `w_c_id` , `w_s_symb` , `w_action` in `staging.watchhistory` for faster filtering and joins + Indexes on `customerid` , `symbol` , and `datevalue` in related tables)

17. **Transform and load** `master.factmarkethistory` :

   - Depends on `dimdate` , `financial` , `dimcompany` , and `dimsecurity` .

   - **Records: 6240582**

   - **manually:  5 min 42.005 sec**  (Added indexes to `dailymarket` , `dimsecurity` , `dimdate` , and `financial` for join optimization . tables are filter to only include recent years data because it still takes damn long for mariadb to process)

```
+----+-------------------+------------+------+--------------------------+--------------------------+---------+-----------------------+---------+--------------------------------+
| id | select_type       | table      | type | possible_keys            | key                      | key_len | ref                   | rows    | Extra                          |
+----+-------------------+------------+------+--------------------------+--------------------------+---------+-----------------------+---------+--------------------------------+
|  1 | PRIMARY           | s          | ALL  | idx_dimsecurity_symbol   | NULL                     | NULL    | NULL                  | 8963    |                                |
|  1 | PRIMARY           | <derived5> | ref  | key0                     | key0                     | 6       | master.s.sk_companyid | 55      |                                |
|  1 | PRIMARY           | <derived4> | ref  | key0                     | key0                     | 63      | master.s.symbol       | 1431    | Using where                    |
|  4 | DERIVED           | dm         | ALL  | idx_dailymarket_symbol_date | NULL                  | NULL    | NULL                  | 1282927 | Using temporary; Using filesort|
|  4 | DERIVED           | dd         | ref  | idx_dimdate_datevalue    | idx_dimdate_datevalue    | 3       | staging.dm.dm_date    | 1       |                                |
|  4 | DERIVED           | <derived3> | ref  | key0                     | key0                     | 63      | staging.dm.dm_s_symb  | 10      | Using where                    |
|  3 | DERIVED           | dm         | ALL  | NULL                     | NULL                     | NULL    | NULL                  | 1282927 | Using temporary                |
|  3 | DERIVED           | dd         | ref  | idx_dimdate_datevalue    | idx_dimdate_datevalue    | 3       | staging.dm.dm_date    | 1       | Using index                    |
|  5 | DERIVED           | f          | ALL  | idx_financial_company_qtr| NULL                     | NULL    | NULL                  | 496713  | Using temporary                |
|  6 | DEPENDENT SUBQUERY| f2         | ref  | idx_financial_company_qtr| idx_financial_company_qtr| 5       | master.f.sk_companyid | 81      | Using where; Using index       |
+----+-------------------+------------+------+--------------------------+--------------------------+---------+-----------------------+---------+--------------------------------+
```

18. **Load** `master.dimessages factmarkethistory` :

- Depends on `factmarkethistory` .

- **Records: 6642**

- **manually: 4.231 sec**

"The performance metric reported for TPC-DI is a throughput measure, the number of **Source**

**Data**

rows processed per second. Conceptually, it is calculated by dividing the total rows

processed by the elapsed time of the run.

The completion timestamp (CT) for each phase is found in DImessages.MessageDateAndTime where MessageType='PCR' and BatchID matches the phase. Example query for **Historical Load**:

select MessageDateAndTime from DImessages where BatchID = 1 and MessageType = 'PCR'

For **Historical Load** throughput calculation:

1. Elapsed time (EH) = CT1 - CT0

2. Total rows (RH) = Batch1 row count from **DIGen**

3. Throughput (TH) = RH / EH"

TPC-DI specification是這麼寫的，但是我看到的組也都只看airflow measure execution time？

過程referred to: https://github.com/risg99/tpc-di-benchmark/tree/main airflow的 code改改應該就能用了