

# Midterm team 1 Project

By

Dhrumil Patel

Mittul Sharma

Vijeth Reddy

# Part 1

- Data Profiling

The screenshot shows a data profiling tool interface. At the top, there's a toolbar with options: Sample, Select, Sort, Join, Union, Text To Columns, Summarize, and Comment. Below this is a workflow diagram with several steps: a data source icon, a 'Convert Inspection Date From MM-dd-yy' step, another 'Convert Inspection Date From MM-dd-yy' step, and a 'Container 31' step. The 'Results - Browse (32) - Input' section shows a table with 109 fields and 78,984 records displayed (43 MB). The table has columns: Record, Restaurant Name, Inspection Type, Inspection Score, Street Address, Zip Code, and Viola. The first 9 records are visible.

Record	Restaurant Name	Inspection Type	Inspection Score	Street Address	Zip Code	Viola
1	MICKLE CHICKEN	Routine	100	3203 W CAMP WISDOM RD	75237	N/A
2	TOM THUMB - JUICE BAR	Routine	100	2380 N FIELD ST	75021	N/A
3	BROOKDALE WHITE ROCK	Routine	97	9271 WHITE ROCK TRL	75238	*21 R
4	CHURCH'S CHICKEN #201	Routine	100	10295 FERGUSON RD	75228	N/A
5	PEAK PREPARATORY-PRIMARY SCHOOL	Routine	98	1474 ANNEX AVE	75204	*21 R
6	SAMS CLUB #6482 - DEMO ROOM	Routine	100	5555 S BUCKNER BLVD	75227	N/A
7	AMC THEATRES NORTH PARK 15 (BAR)	Routine	100	8687 N CENTRAL EXPW #3000	75225	N/A
8	SAM'S CLUB #8282 - DEMO ROOM	Routine	100	2900 W WHEATLAND RD	75237	N/A
9	UCHI DALLAS	Routine	98	2817 MAPLE AVE #110	75201	*32 E

- Data analysis and key output:  
Columns- there are total 114 columns in Dallas table out of which restaurant name, inspection type, inspection score, street address, zip code, street, inspection date, inspection month, inspection year, location do not have any null values.

- There is a total of 78984 records in the table.
- Apart from this there is one more data inconsistency such as for violation description, violation details, violation memo, violation points are columns in De-normalize form such as violation description 1, violation description 2, violation description 3 so here is there is violation then it stores in one of column and other columns store null values.
- There are changes of data type of data in inspection date, inspection month, created date in date, zip code, inspection score, violation points in integer and other columns are in string.

Browse (9) - Configuration

78,984 records displayed, 109 fields, 43 MB

Profile

78,984 records displayed, 109 fields, 43 MB

Restaurant Name	
MCDONALDS	533
SUBWAY	312
BURGER KING	293
SONIC DRIVE IN	266
WINGSTOP	264
995 more >	

Inspection Type	
Routine	78,019
Follow-up	935
Complaint	30

Inspection Score	
100	6,564
97	5,249
95	4,964
94	4,724
96	4,675
53 more >	

Street Address	
2201 N STEMMONS FRWY	181
3750 COTTON BOWL PLAZA	171

[Null]	[Null]	[Null]	[Null]	FY2022	2900 W WHEATLAND RD	2021-11-14
[Null]	[Null]	[Null]	[Null]	FY2020	2817 MAPLE AVE #110	2020-03-11

Browse (9) - Configuration

78,984 records displayed, 109 fields, 43 MB

Profile

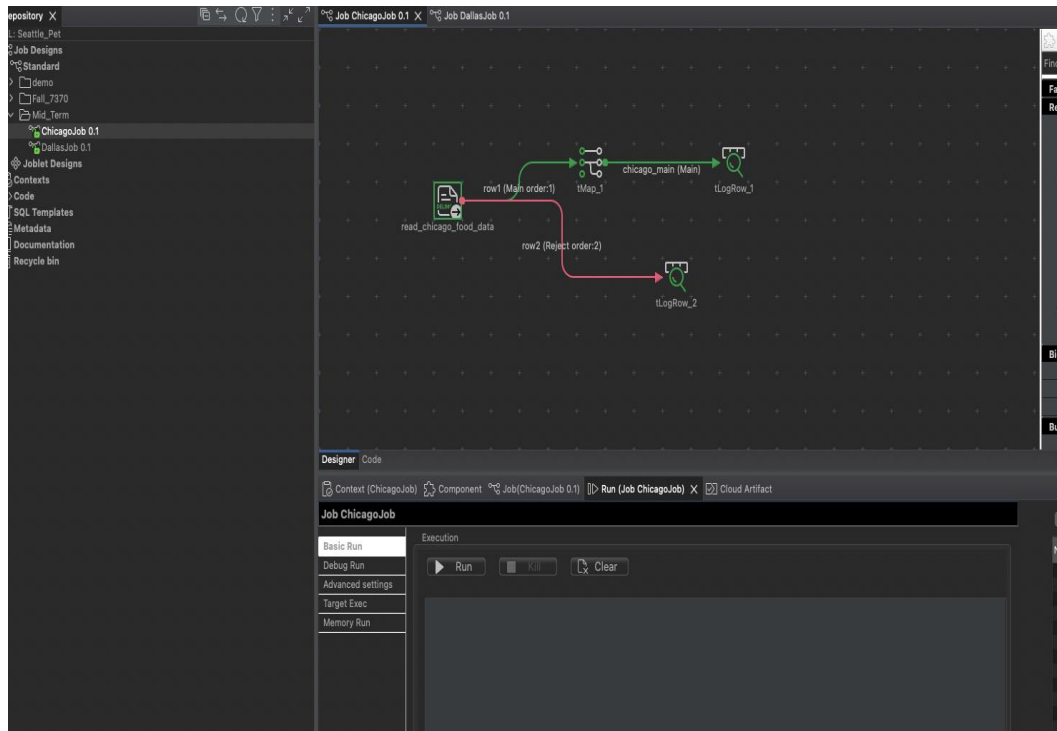
Zip Code	
75201	4,687
75220	3,784
75206	3,373
75243	3,146
75211	3,075
151 more >	
Violation Description - 1	
*10 Clean Sight and Touch	5,113
*21 RFSM - Not On Site	4,421
*02 Cold Hold (41°F/45°F or below)	3,924
*42 Dirty nonfood contact surfaces	2,242
*07 Food safe, good condition, unadulterated, and honestly prese...	2,143
618 more >	
Violation Points - 1	
3	28,550
1	24,947
2	18,907
4	1
Violation Detail - 1	
228.113 Equipment, Utensils, and Linens. Cleaning of equipment a...	5,113
Sec. 17-2.2(c)(1)(D) (c) Registered food service managers. (1) Regis...	4,421
228.75 Food, Time and temperature control. (f) Time/temperature ...	3,924

Browse (9) - Configuration

78,984 records displayed, 109 fields, 43 MB

Profile

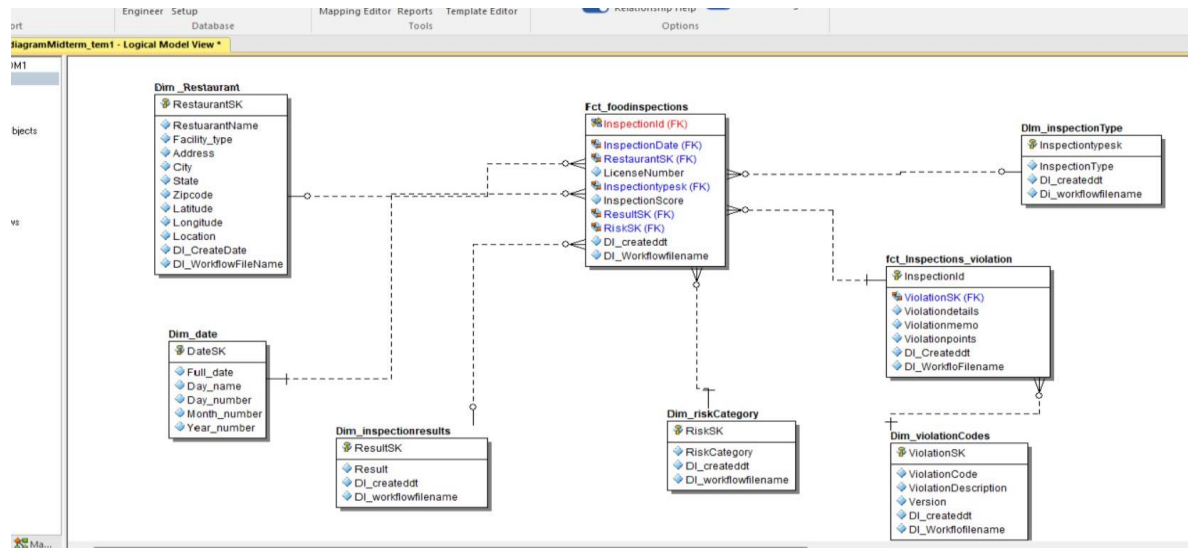
Inspection Year	
FY2018	14,412
FY2017	13,443
FY2019	11,342
FY2023	9,706
FY2021	9,319
3 more >	
Lat Long Location	
650 S GRIFFIN ST (32.772703, -96.799217)	147
2201 N STEMMONS FRWY (32.801595, -96.829026)	126
3750 COTTON BOWL PLAZA (32.778661, -96.760942)	126
555 S LAMAR ST #100 (32.775343, -96.803021)	112
3500 GASTON AVE (32.342322, -86.31818)	97
995 more >	
Inspection_date	
2018-03-21	117
2017-09-21	110
2017-10-12	109
2018-11-28	103
2018-05-31	102
995 more >	
Inspection_month	
2017-12-10	109
2017-06-12	93
2017-12-13	90



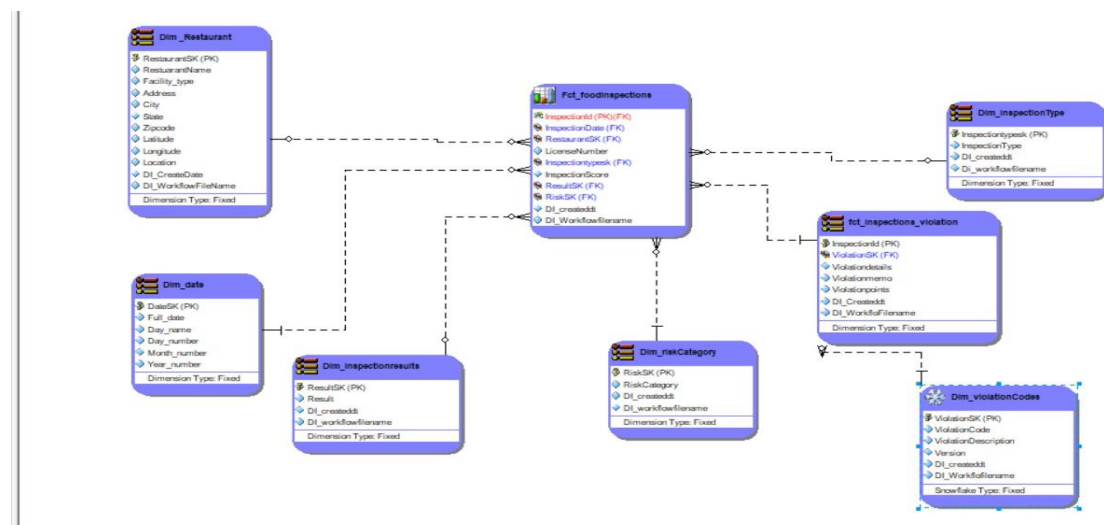
- Compare to Dallas data Chicago has less data inconsistency
- In Chicago data there is not any denormalized data in this data set there is detailed data is created but when it comes to combine both data some of fields are eliminated such as inspection type, license number
- And here we have used similar data type conversion as per Chicago dataset

CHICAGO\_data\_PROFILING.yemd

<



## 2. Physical Model



## 3. DDL Script

```
-- Drop table Dim_InspectionResults;
CREATE TABLE Dim_InspectionResults (
    ResultsSK INT AUTOINCREMENT(1,1) PRIMARY KEY,
    Results VARCHAR(255),
    DI_CreateDate DATETIME,
    DI_WorkflowFileName STRING
);

-- Drop table Dim_InspectionType ;
CREATE TABLE Dim_InspectionType (
    InspectionTypeSK INT AUTOINCREMENT(1,1) PRIMARY KEY,
    InspectionType VARCHAR(255),
    DI_CreateDate DATETIME,
    DI_WorkflowFileName STRING
);

-- Drop table Dim_RiskCategory ;
CREATE TABLE Dim_RiskCategory (
    RiskSK INT AUTOINCREMENT(1,1) PRIMARY KEY,
    Risk VARCHAR(255),
    DI_CreateDate DATETIME,
    DI_WorkflowFileName STRING
);

-- Drop table Dim_Restaurants ;
CREATE TABLE Dim_Restaurants (
    RestaurantSK INT AUTOINCREMENT(1,1) PRIMARY KEY,
    Restaurant_Name VARCHAR(255),
    Facility_Type VARCHAR(255),
    Address VARCHAR(255),
    City VARCHAR(50),
    State CHAR(10),
    Zip INT,
    Latitude CHAR(18),
    Longitude CHAR(18),
    Location CHAR(48),
    DI_CreateDate DATETIME,
    DI_WorkflowFileName STRING
);

-- Drop table Dim_ViolationCodes ;
CREATE TABLE Dim_ViolationCodes (
    ViolationSK INT AUTOINCREMENT(1,1) PRIMARY KEY,
```

ults

Messages

Restaurant_Name	Inspection_Type	Inspection_Date	Inspection_Score	Street_Number	Street_Name	Street_Direction	Street_Type	Street_Unit	Street_Address	Zip_Code	Violation_Description_1
WENDY'S #9780	Routine	2018-09-13	80	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75216	*31 Handwashing lavatory --
WENDY'S #9780	Routine	2021-04-06	96	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75216	*32 Maintain in Good Repair
WENDY'S #9780	Routine	2020-02-24	100	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75216	NULL
WENDY'S #97800	Routine	2021-12-22	92	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75215-4514	*19 Plumbing System Constr
WENDY'S #9780	Routine	2016-12-13	92	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75216	*31 Individual, disposable

```

69 -- Drop table Dim_ViolationCodes ;
70 CREATE TABLE Dim_ViolationCodes (
71     ViolationSK INT AUTOINCREMENT(1,1) PRIMARY KEY,
72     ViolationCode INT,
73     ViolationDescription VARCHAR(1000),
74     DI_CreateDate DATETIME,
75     DI_WorkflowFileName STRING
76 );
77
78 -- Drop table FCT_Inspections_Violations ;
79 CREATE TABLE FCT_Inspections_Violations (
80     InspectionID INT,
81     ViolationSK INT,
82     ViolationDetail VARCHAR(8000) NULL,
83     ViolationMemo VARCHAR(8000) NULL,
84     ViolationPoints Varchar(100) NULL,
85     DI_CreateDate DATETIME,
86     DI_WorkflowFileName STRING,
87     PRIMARY KEY (InspectionID, ViolationSK),
88     FOREIGN KEY (ViolationSK) REFERENCES Dim_ViolationCodes(ViolationSK)
89 );
90
91
92
93
94
95 Select * from FCT_FoodInspections ;
96
97 CREATE TABLE FCT_FoodInspections (
98     FCTInspectionID INT AUTOINCREMENT(1,1) PRIMARY KEY,
99     InspectionID INT,
100     InspectionDate DATE NULL,
101     DateSK INT,
102     RestaurantSK INT NULL,
103     InspectionTypeSK INT NULL,
104     ResultsSK INT NULL,
105     RiskSK INT NULL,
106     Inspection_Score INT NULL,
107     DI_CreateDate DATETIME,
108     DI_WorkflowFileName STRING
```

ults

Messages

Restaurant_Name	Inspection_Type	Inspection_Date	Inspection_Score	Street_Number	Street_Name	Street_Direction	Street_Type	Street_Unit	Street_Address	Zip_Code	Violation_Description_1
WENDY'S #9780	Routine	2018-09-13	80	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75216	*31 Handwashing lavatory --
WENDY'S #9780	Routine	2021-04-06	96	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75216	*32 Maintain in Good Repai
WENDY'S #9780	Routine	2020-02-24	100	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75216	NULL
WENDY'S #97800	Routine	2021-12-22	92	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75215-4514	*19 Plumbing System Constr
WENDY'S #9780	Routine	2016-12-13	92	1507	KIEST	E	BLVD	NULL	1507 E Kiest Blvd	75216	*31 Individual, disposable

## 4. Dallas Stage Table



```
1 CREATE OR REPLACE TABLE stg_DALLAS_FOOD_INSPECTIONS (
2   RECORD_ID NUMBER(38,0) NOT NULL AUTOINCREMENT START 1 INCREMENT 1 NOORDER,
3   INSPECTION_ID NUMBER(38,0),
4   RESTAURANT_NAME VARCHAR(255),
5   INSPECTION_TYPE VARCHAR(100),
6   INSPECTION_DATE DATE,
7   INSPECTION_SCORE NUMBER(5, 2),
8   STREET_NUMBER VARCHAR(50),
9   STREET_NAME VARCHAR(255),
10  STREET_DIRECTION VARCHAR(10),
11  STREET_TYPE VARCHAR(50),
12  STREET_UNIT VARCHAR(50),
13  STREET_ADDRESS VARCHAR(255),
14  ZIP_CODE VARCHAR(10),
15  INSPECTION_MONTH VARCHAR(20),
16  INSPECTION_YEAR VARCHAR(10),
17  LATITUDE VARCHAR(50),
18  LONGITUDE VARCHAR(50),
19  LOCATION VARCHAR(255),
20  -- Violation Columns
21  VIOLATION_DESCRIPTION VARCHAR(1000),
22  VIOLATION_POINTS NUMBER(10, 2),
23  VIOLATION_DETAIL VARCHAR(4000),
24  VIOLATION_MEMO VARCHAR(4000),
25 );
26
```

## 5. Chicago Stage Table

```
1 create or replace TABLE stg_CHICAGO_FOOD_INSPECTIONS (
2   RECORD_ID NUMBER(38,0) NOT NULL autoincrement start 1 increment 1 noorder,
3   INSPECTION_ID NUMBER(38,0),
4   DBA_NAME VARCHAR(255),
5   AKA_NAME VARCHAR(255),
6   LICENSE_NUMBER NUMBER(38,0),
7   FACILITY_TYPE VARCHAR(255),
8   RISK VARCHAR(50),
9   ADDRESS VARCHAR(255),
10  CITY VARCHAR(100),
11  STATE VARCHAR(255),
12  ZIP_CODE VARCHAR(10),
13  INSPECTION_DATE DATE,
14  INSPECTION_TYPE VARCHAR(100),
15  RESULTS VARCHAR(100),
16  LATITUDE VARCHAR(255),
17  LONGITUDE VARCHAR(255),
18  LOCATION VARCHAR(255),
19  VIOLATION_ID NUMBER(38,0),
20  VIOLATION_DESCRIPTION VARCHAR(10000),
21  COMMENTS VARCHAR(10000),
22  DI_CREATEDATE DATE,
23  DI_WORKFLOWFILENAME VARCHAR(255),
24  primary key (RECORD_ID)
25 );
```

## Part 3

# Integration and snowflake

## ● Chicago stage table in snowflake

### 1. Data Integration Workflow to load datasets into Snowflake as staging table

The screenshot displays a Data Integration Workflow in Snowflake. The workflow consists of the following steps:

- FoodInspection**: Input data from ChicagoFoodInspection
- datacleansing**: Creating/updating the columns: Inspection ID, DBA Name, AKA Name, License #, Facility Type, Risk, Address, City, State, Zip
- ViolationArray**: Creating/updating the columns: Inspection ID, DBA Name, AKA Name, License #, Facility Type, Risk, Address, City, State, Zip
- flattenViolationArray**: Unrolling arrays from Violations to with columns: Inspection ID, DBA Name, AKA Name, License #, Facility Type, Risk, Address, City, State, Zip
- SplitViolation**: Creating/updating the columns: Inspection ID, DBA Name, AKA Name, License #, Facility Type, Risk, Address, City, State, Zip
- SnowflakeTable**: Columns: 21 total

Below the workflow, the **Data preview** tab is active, showing a table with 10 columns and 10 rows of data. The columns are: **INSPECTION\_ID**, **DBA\_NAME**, **AKA\_NAME**, **LICENSE\_NUMBER**, **FACILITY\_TYPE**, **RISK**, **ADDRESS**, **CITY**, **STATE**, and **ZIP\_CODE**.

INSPECTION_ID	DBA_NAME	AKA_NAME	LICENSE_NUMBER	FACILITY_TYPE	RISK	ADDRESS	CITY	STATE	ZIP_CODE
2605385	LUCY & M...	LUCY & M...	2896496	Restaurant	Risk 1 (High)	1507 W M...	CHICAGO	IL	60607
2605348	BUTCHER ...	BUTCHER ...	2877649	Restaurant	Risk 1 (High)	11601 W T...	CHICAGO	IL	60666
2605325	3 ASIAN SL...	3 ASIAN SL...	2996879	Restaurant	Risk 1 (High)	818 W FUL...	CHICAGO	IL	60614
2605318	C/O ST AN...	EMPLOYEE...	2334226	Restaurant	Risk 1 (High)	2875 W 19...	CHICAGO	IL	60623
2605306	HANSON P...	HANSON P...	2996860	Grocery St...	Risk 3 (Low)	5424 W FU...	CHICAGO	IL	60639
2605276	KRISPY CH...	KRISPY CH...	2570228	Restaurant	Risk 2 (Me...	1956 W 79...	CHICAGO	IL	60620
2605276	KRISPY CH...	KRISPY CH...	2570228	Restaurant	Risk 2 (Me...	1956 W 79...	CHICAGO	IL	60620
2605276	KRISPY CH...	KRISPY CH...	2570228	Restaurant	Risk 2 (Me...	1956 W 79...	CHICAGO	IL	60620
2605276	KRISPY CH...	KRISPY CH...	2570228	Restaurant	Risk 2 (Me...	1956 W 79...	CHICAGO	IL	60620
2605276	KRISPY CH...	KRISPY CH...	2570228	Restaurant	Risk 2 (Me...	1956 W 79...	CHICAGO	IL	60620

### 2. Load Staging data into Dimensional tables

**Dim Fact Data Load**  
Cluster startup time: 1s 97ms Number of transformations: 10 Data flow status: Success

Refresh Auto refresh On Edit dataflow

Sink	Status	Processing time	Highest processing time	Rows written	Stages	Lineage
DimInspectionType	Succeeded	6s	153ms	98		
DimResults	Succeeded	6s	167ms	7		
DimRisk	Succeeded	6s	144ms	4		
DimViolations	Succeeded	5s	74ms	110		
DimRestaurant	Succeeded	7s	1s 103ms	38827		

### 3. Pipeline to load data into fact table from dimensional table

**Sink Settings**

Output stream name: LoadToFactTable [Learn more](#)

Description: Export data to SnowflakeDimandFact [Reset](#)

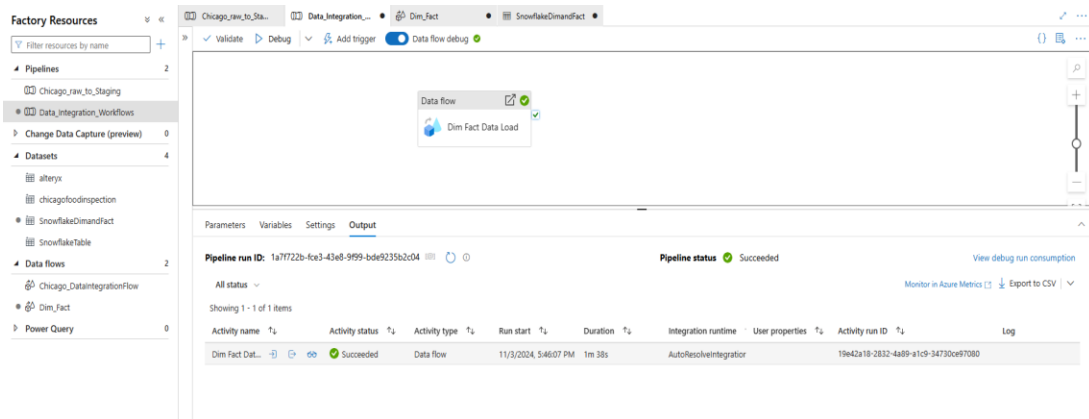
Incoming stream: JoinDate

Sink type: Dataset ☒ Inline ☐ Cache

Dataset: SnowflakeDimandFact [Test connection](#) [Open](#) [New](#)

Options: ☒ Allow schema drift [?](#) ☐ Validate schema [?](#)

### 4. Snapshot of Staging to Dimension tables



STG	12	33320	1354010	FLOUR & STONE	FLOUR & STONE	2220904	Restaurant
Tables	13	33321	1354102	AVALON EAT SHOP	AVALON EAT SHOP	8682	Grocery Store
DIM_DATE	14	33322	1354102	AVALON EAT SHOP	AVALON EAT SHOP	8682	Grocery Store
DIM_INSPECTIONRESULTS	15	33323	1354102	AVALON EAT SHOP	AVALON EAT SHOP	8682	Grocery Store
DIM_INSPECTIONTYPE	16	33324	1369270	HAROLD'S CHICKEN SHACK	HAROLD'S CHICKEN SHACK	20687	Restaurant
DIM_RESTAURANTS	17	33325	1369270	HAROLD'S CHICKEN SHACK	HAROLD'S CHICKEN SHACK	20687	Restaurant
DIM_RISKCATEGORY	18	33326	1369270	HAROLD'S CHICKEN SHACK	HAROLD'S CHICKEN SHACK	20687	Restaurant
DIM_VIOLATIONCODES	19	33327	1369270	HAROLD'S CHICKEN SHACK	HAROLD'S CHICKEN SHACK	20687	Restaurant
FCT_FOODINSPECTIONS	20	33328	1235912	HOPLEAF BAR	HOPLEAF BAR	2152669	Restaurant
FCT_INSPECTIONS_VIOLATIO...	21	33329	1235912	HOPLEAF BAR	HOPLEAF BAR	2152669	Restaurant
SNOWFLAKE	22	33330	1235912	HOPLEAF BAR	HOPLEAF BAR	2152669	Restaurant
SNOWFLAKE_SAMPLE_DATA	23	33331	1235912	HOPLEAF BAR	HOPLEAF BAR	2152669	Restaurant

To merge the stg\_DALLAS\_FOOD\_INSPECTIONS and stg\_CHICAGO\_FOOD\_INSPECTIONS datasets, here's the plan:

- Align Column Names:** Map similar columns, like RESTAURANT\_NAME (Dallas) to DBA\_NAME (Chicago) and consolidate address components into a single ADDRESS field.
- Include Unique Columns:** Add Chicago-specific fields (e.g., LICENSE\_NUMBER, FACILITY\_TYPE, RISK, RESULTS) to the merged schema, with null values allowed for Dallas entries.
- Standardize Violation Details:** Combine violation columns across datasets, such as VIOLATION\_DESCRIPTION and COMMENTS.
- Unified Schema:** Create a schema that captures fields from both datasets, allowing flexibility in columns exclusive to either Dallas or Chicago.
- Merge similar columns:** Map VIOLATION\_DESCRIPTION from both datasets into a unified VIOLATION\_DESCRIPTION field.

6. Include VIOLATION\_POINTS (Dallas-only) and VIOLATION\_ID (Chicago-only) in the merged schema, allowing nulls where they don't apply.
7. Consolidate VIOLATION\_DETAIL (Dallas) and COMMENTS (Chicago) into a single VIOLATION\_DETAIL field, storing specific notes or remarks related to the violation.
8. Include VIOLATION\_MEMO (Dallas-only) for additional detail where available.
9. Design a unified violation section in the merged schema to capture all aspects of violation data, ensuring any city-specific fields have null values for records from the other city.