

Prepared by group 21

Traffic incidents in Chicago

Matilde Polezzi, Gina Santoro

13 February, 2025

TABLE OF CONTENTS

01 Missing Values Management

02 DW Schema

03 Python Functions

04 DW Data Loading

05 Visual Studio

06 Multidimensional Cube

07 MDX Queries

08 PowerBI Dashboards

INTRODUCTION

Our project aims to simulate a Decision Support System for an Insurance Company. To achieve this, we developed a Data Warehouse that enabled various analyses. Based on this database, we then created a data cube to perform multidimensional analysis.



MISSING VALUES MANAGEMENT

File Crashes

- REPORT TYPE -> "Not on Scene".
- STREET DIRECTION: Imputed cross-checked values: "S" (South) and "N" (North).
- STREET NAME -> The missing value was imputed as "76TH ST" on the basis of relationships with other columns.
- BEAT OF OCCURRENCE -> Replacement based on STREET NAME using consistent matches.
- MOST SEVERE INJURY -> "No Indication of Injuries".
- LATITUDE, LONGITUDE, LOCATION -> Retrieved missing values using Geopy.

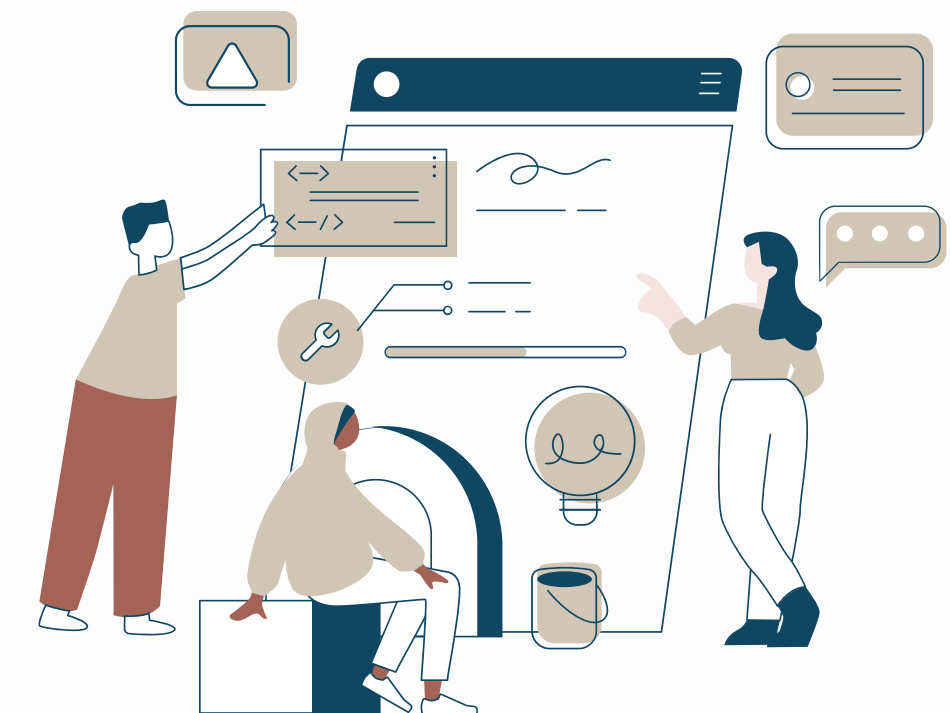
File Vehicles

- UNIT TYPE: Deleted the only record without matches with RD NO in files People or Crashes.
- VEHICLE ID -> empty or 0.0.
- MAKE, MODEL, VEHICLE DEFECT, VEHICLE TYPE, VEHICLE USE, MANEUVER -> left blank. or "Unknown".
- LIC PLATE STATE -> left blank or "XX".
- VEHICLE YEAR -> left blank or 0.0.
- TRAVEL DIRECTION -> "Unknown".
- OCCUPANT CNT -> left blank.
- FIRST CONTACT POINT -> left blank or the column was manually integrated for some specific RD NO.

MISSING VALUES MANAGEMENT

File People

- VEHICLE ID -> 0.0
- CITY, STATE, SEX -> "Unknown", "XX", and "U" respectively.
- AGE: Inconsistencies were corrected or assigning values of 0.0 to represent null values.
- SAFETY EQUIPMENT -> "Usage Unknown".
- AIRBAG DEPLOYED:
 - For Pedestrian, Bicycle, and Non-Motor Vehicle: "NOT APPLICABLE".
 - For Non-Contact Vehicle: "DID NOT DEPLOY".
 - For Drivers: "Deployment Unknown".
- EJECTION:
 - For Pedestrian, Bicycle, and Non-Motor Vehicle: left blank.
 - For others: replaced with "Unknown".
- INJURY CLASSIFICATION -> "No Indication of Injuries".
- DRIVER VISION: Missing values estimated using weather and lighting conditions.
- DAMAGE -> 500
- DRIVER ACTION, PHYSICAL CONDITION -> "Unknown".
- BAC RESULT -> Left blank.



DW SCHEMA

Creation on SQL Server Management Studio.

Dimension

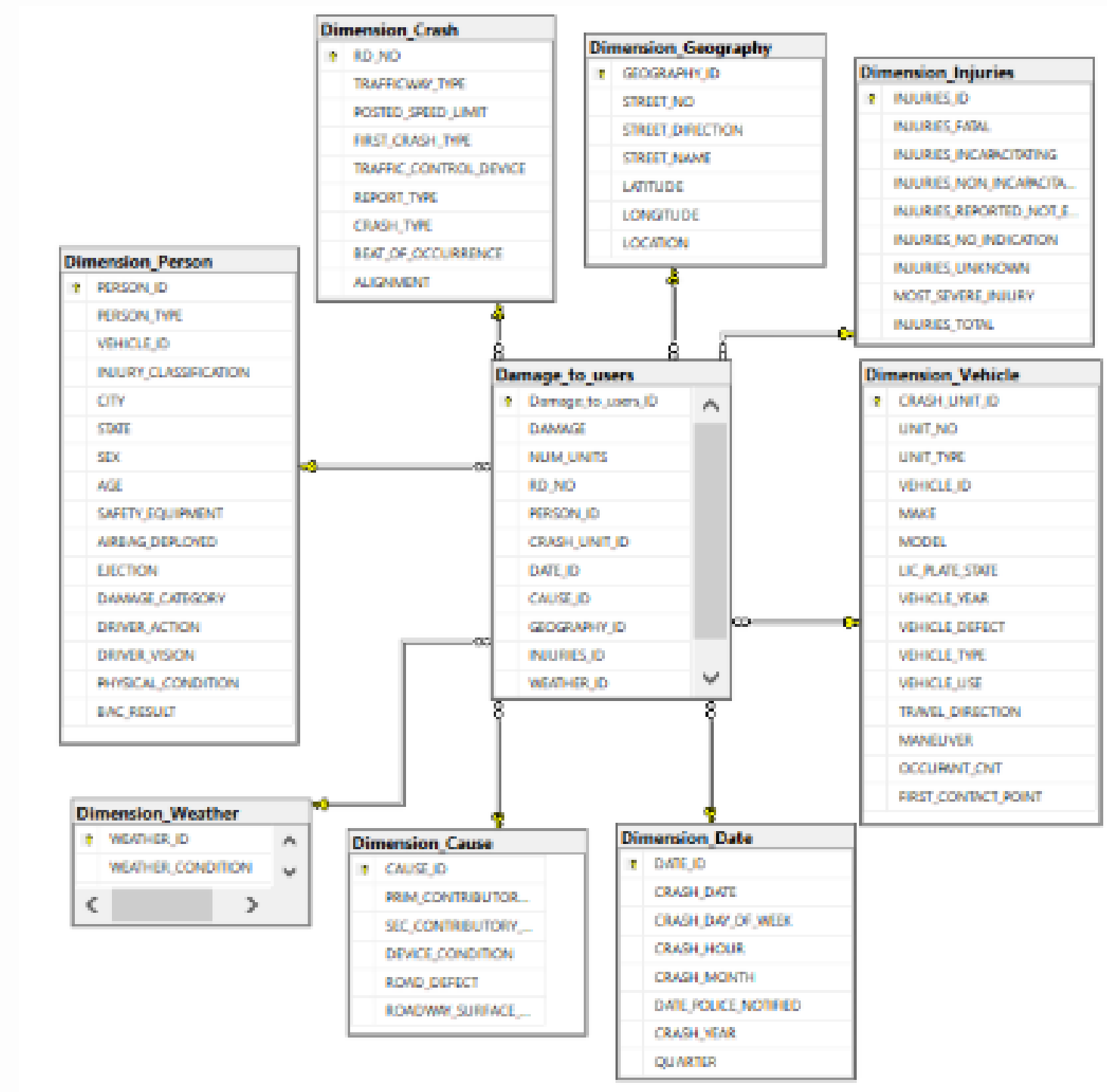
- Crash
- Person
- Vehicle
- Geography
- Injuries
- Cause
- Date
- Weather

FACT TABLE :

- Damage_to_users: Each row concerns the damage cost for each individual associated with a crash.

Measures:

- NUM_UNITS
- DAMAGE



DATA PREPARATION

Functions implemented before populating the data warehouse

FILE PYTHON	PURPOSE
Dimensioni	<ul style="list-style-type: none">• Each dimension is saved in a separate CSV file.• <i>create_dimension</i> function merges data from different sources.• Ensures unique value combinations using a dictionary as a cache.• Assigns a unique incremental ID (starting from 1).
Dimensioni_csv_originali	<ul style="list-style-type: none">• Generates tables by removing duplicates.• Uses <i>csv_no_duplicates</i> function.• Outputs structured in separate CSV files.
Damage_to_users	<ul style="list-style-type: none">• Creates the Damage to users dimension by linking 'RD NO' to primary keys of other dimensions.• Maps source table data to dimensions for clear relationships.• Integrates pre-generated dimensional tables, matching keys, and adding corresponding IDs.• Repeats the process for all relevant dimensions and measurements, updating the fact table sequentially.

POPULATING THE DATA WAREHOUSE

To populate, the `bulk_insert` function has been used:

- Row insertion
- Batches of size 10.000

Memory usage and transaction overhead balance to ensure efficient data insertion.

Commit after each successful batch:

- Minimize data loss and reduce

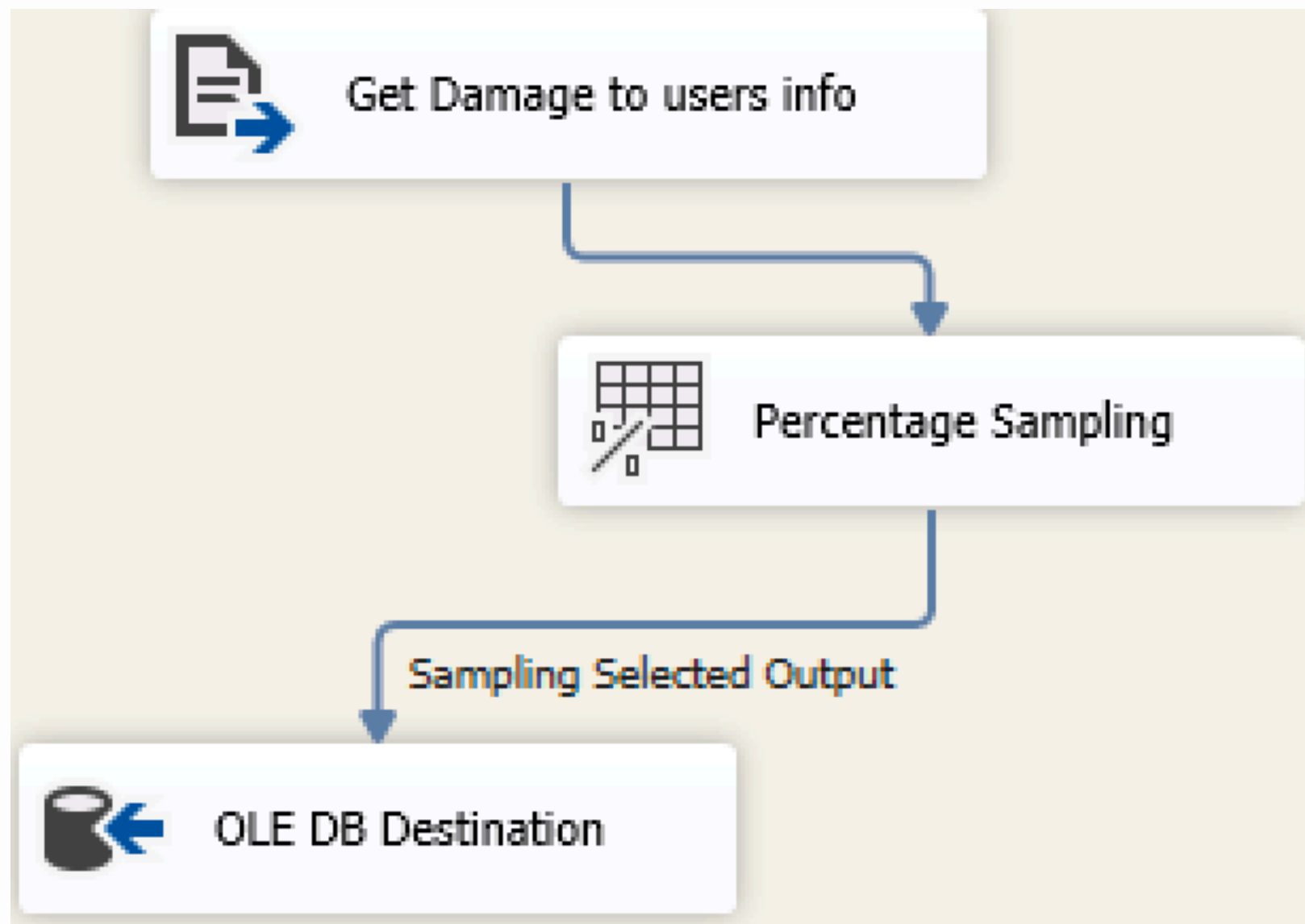
In case of errors:

- Failed rows are saved into a separate file
- One retry for inserting the failed rows

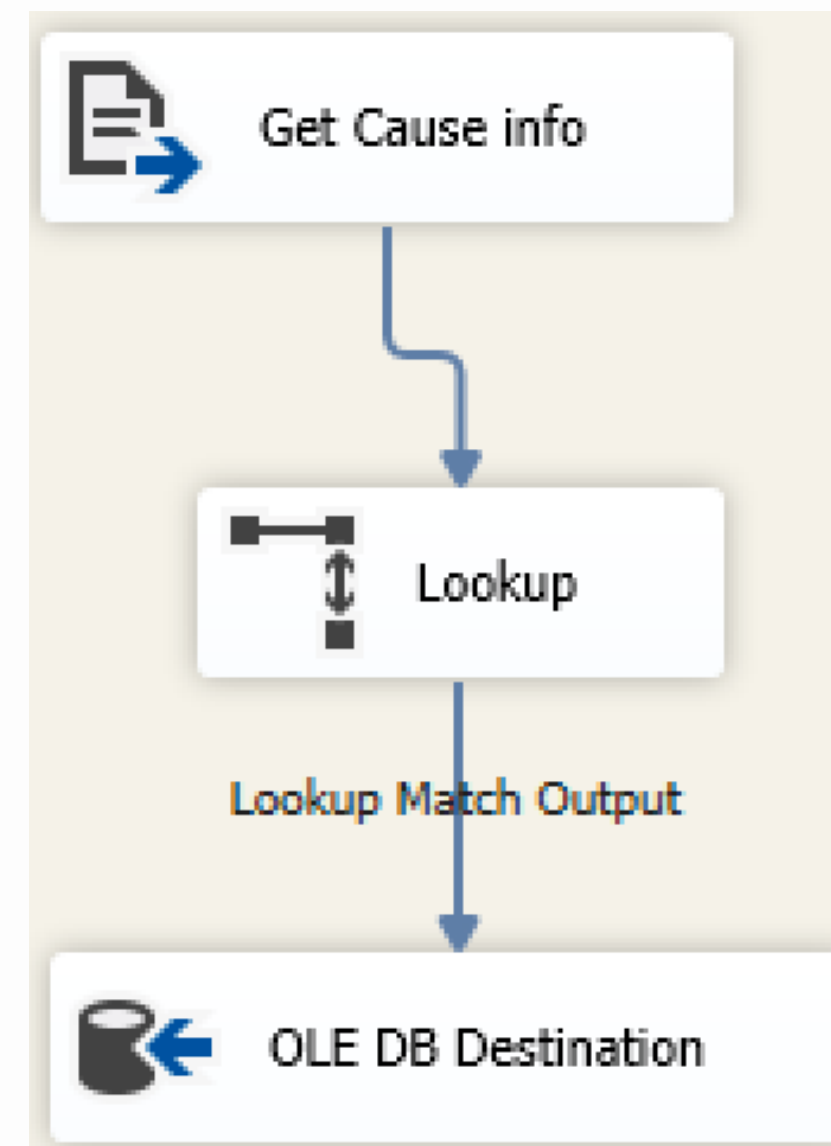


TABLES POPULATION WITH SSIS

Step 1: Select 10% from fact table

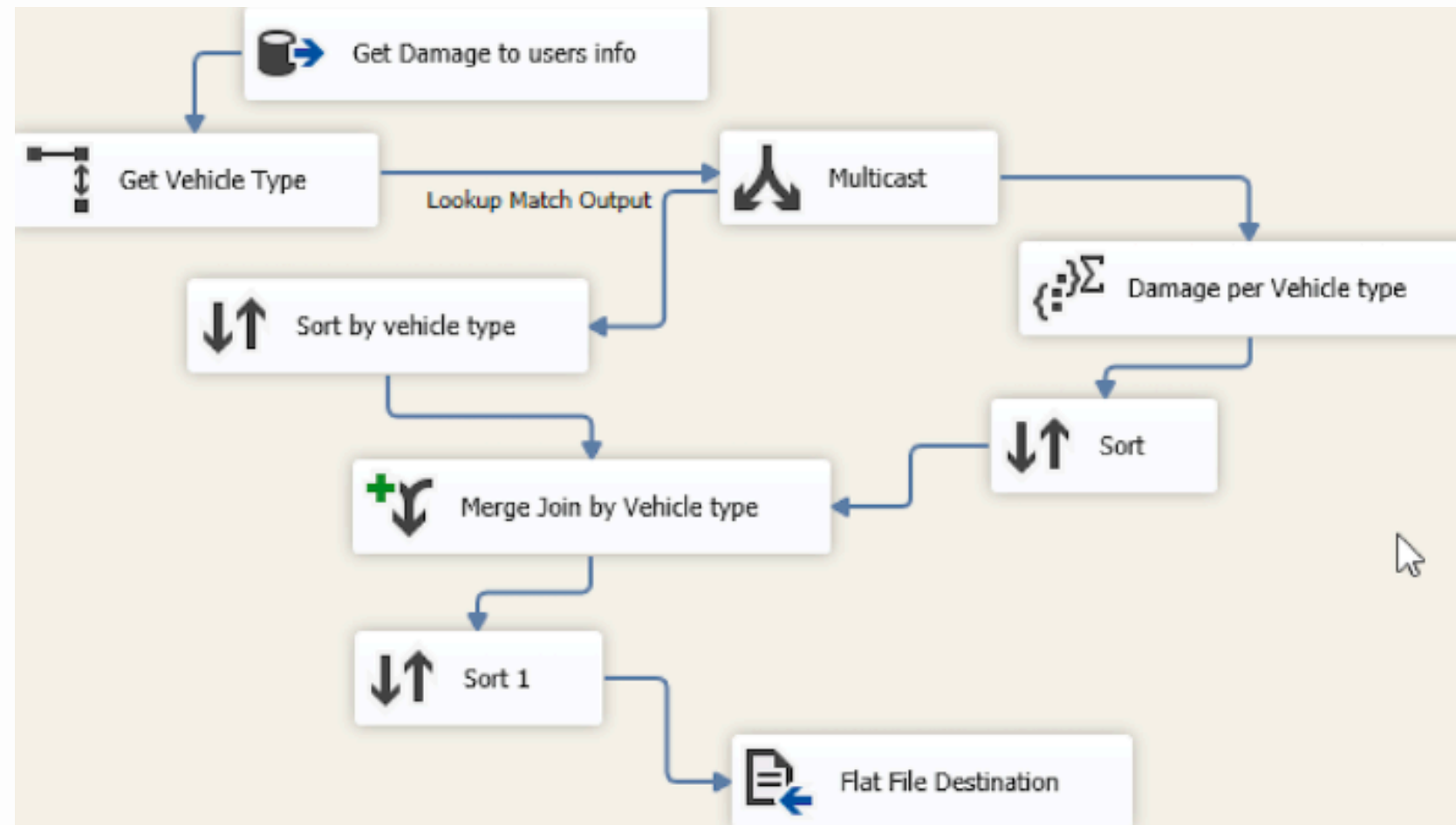


Step 2: This procedure has been carried out for all dimensions



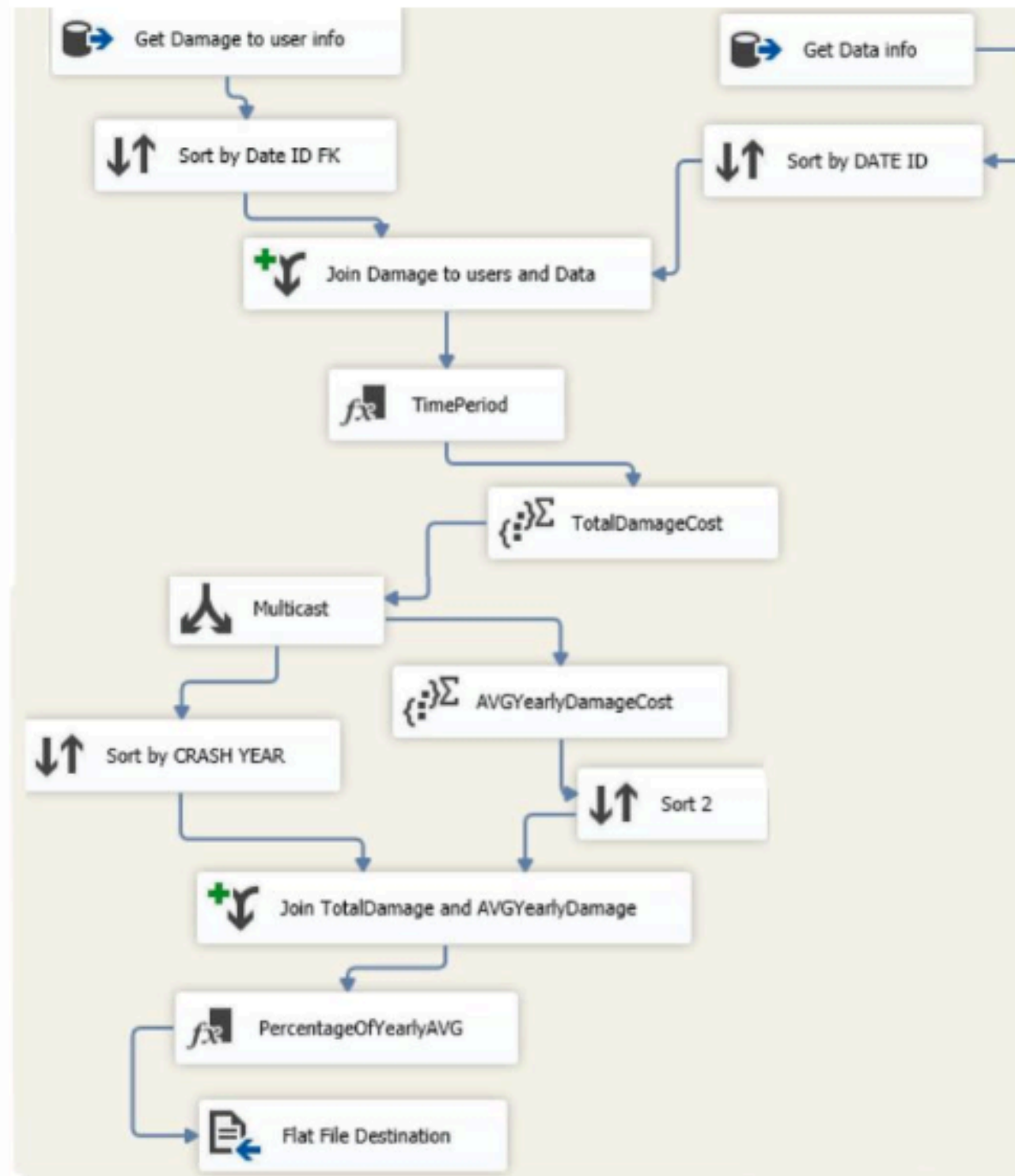
SSIS QUERY

Show all participants ordered by the total damage costs for every vehicle type.



- **Objective:** Calculate the total damage for each vehicle type and sort them in descending order.
- **Steps:**
 - Lookup for Vehicle Type: Perform a lookup on the vehicle type using the CRASH:UNIT_ID field to connect the "vehicle" dimension to the fact table.
 - Use of Multicast: Split the output data to work in parallel and ensure the results are organized by vehicle type.
 - Calculate Total Damage: Aggregate the data using GROUP BY and calculate the total damage for each vehicle type with the SUM function.
 - Merge Join: Combine the data on the "vehicle type" field while also including the associated damage and person_id.
 - Final Sorting: Sort the results in descending order based on the total damage.

SSIS QUERY

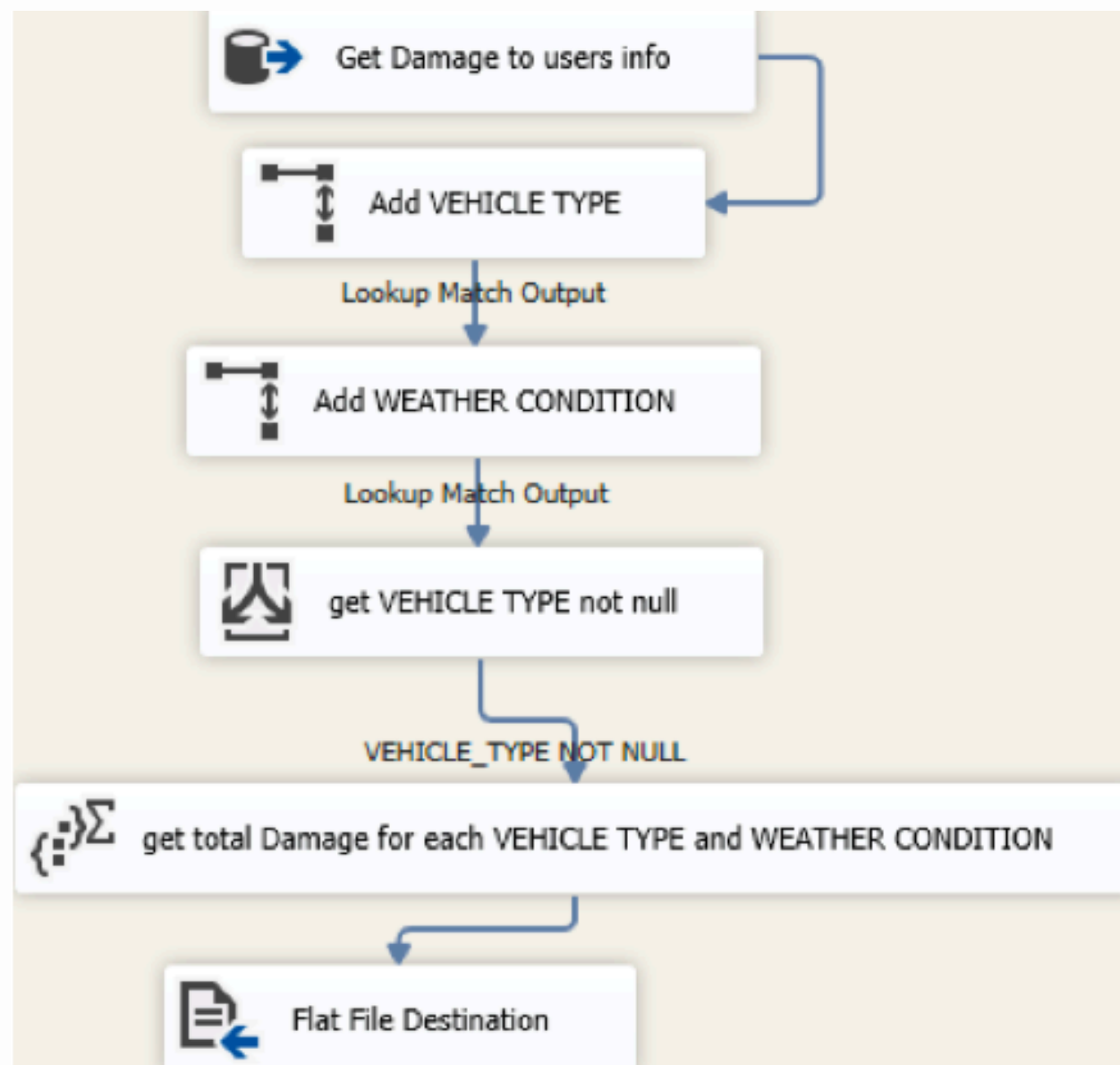


For each month, calculate the percentage of the total damage costs caused by incidents occurring between 9 pm and 8 am and incidents occurring between 8 am and 9 pm, with respect to the average total damage costs for all months within the same year

Steps:

- Merge Join su DATE_ID: Unione dei dati per ottenere informazioni come danno, mese e anno del crash.
- Divisione Giorno/Notte: Utilizzo di un'espressione per identificare gli incidenti notturni (CRASH_HOUR >= 21 è notte, altrimenti giorno).
- Calcolo Total_Damage_Cost: Somma dei danni totali per incidenti notturni, raggruppati per anno, mese e periodo della giornata.
- Uso di Multicast: Dividere i dati per elaborazioni parallele.
- Calcolo della Media Annuo: Costo danno medio annuale raggruppato per CRASH_YEAR.
- Merge Join su CRASH_YEAR: Unione dei dati includendo mese, periodo notturno, danni medi e totali.
- Calcolo della Percentuale: Percentuale del Total_Damage_Cost rispetto al Avg_Damage_Cost.
- Esportazione Finale: Risultati ordinati e salvati in un file piatto.

SSIS QUERY



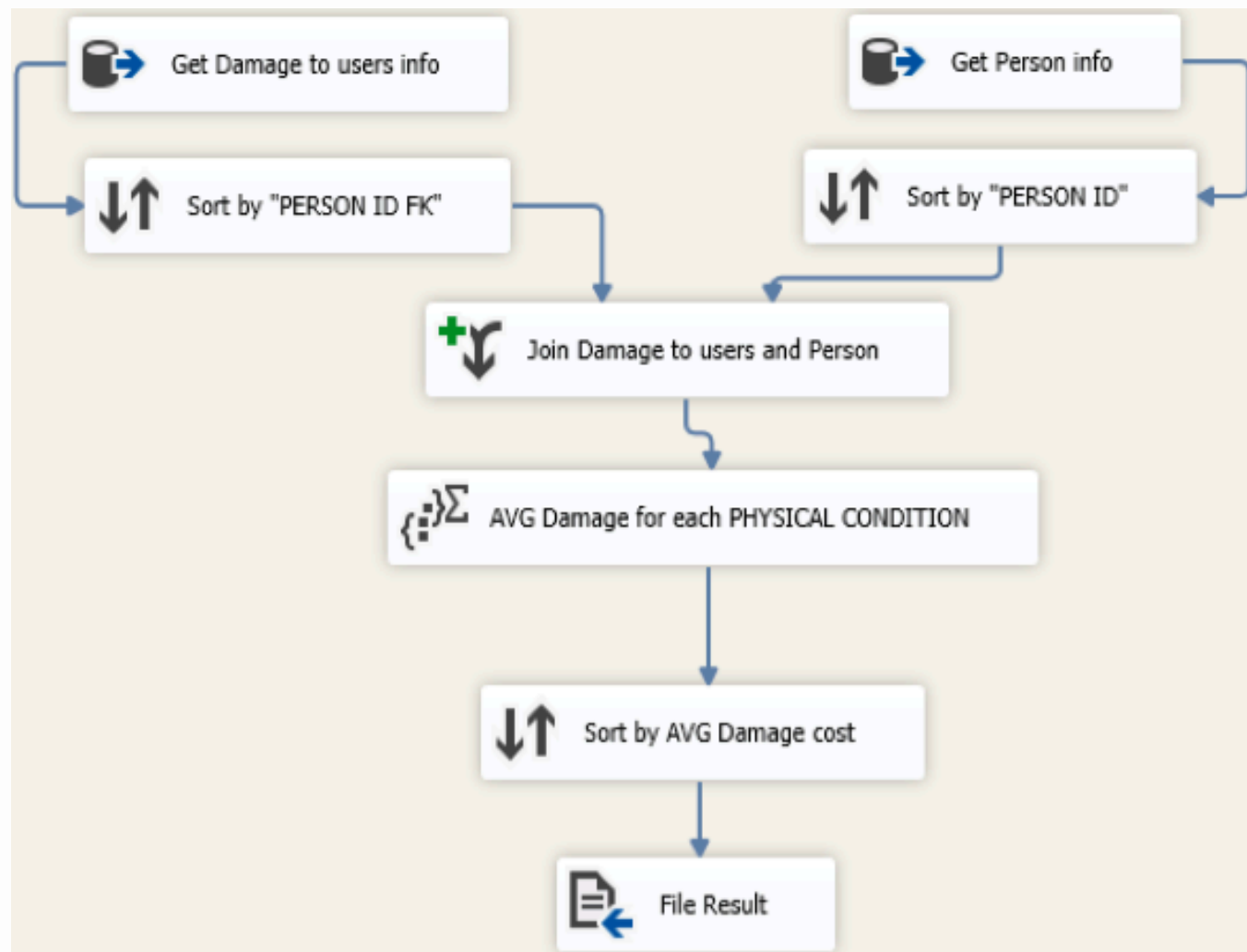
Show the total crash damage costs for each vehicle type and weather condition.

Steps:

- Two Lookups:
 - Retrieve VEHICLE TYPE by joining "Damage to User" with the Vehicle Dimension table.
 - Retrieve WEATHER CONDITION by joining with the Weather Dimension table.
- Filter Non-Null Vehicle Types:
 - Use a Lookup Match Output with a Conditional Split Transformer to include only rows where VEHICLE TYPE is not NULL.
- Aggregation:
 - Group data by VEHICLE TYPE and WEATHER CONDITION.
 - Calculate the Total Damage using the SUM function.
- Output:
 - Save the aggregated results to a Flat File Destination for further analysis.

SSIS QUERY

It calculates the average damage based on the driver's condition, sorting the averages in ascending order.



Steps:

- Merge Join:
 - Combine Person data and the fact table using PERSON ID.
 - Extract relevant information such as Damage and Physical Conditions of individuals.
- Aggregation:
 - Calculate the Average Damage grouped by each Physical Condition.
- Sorting:
 - Sort the results in descending order based on the average damage cost.
- Output:
 - Save the final results to a File Destination for reporting and analysis.

OLAP CUBE

To create the data cube, we chose to retain all attributes for each dimension rather than combining only those required to address the specific business question. This approach ensures the system can accommodate potential future business inquiries providing flexibility for unforeseen analytical needs.



Hierarchies:

- Date: Year, Quarter, Month, Day, Day of the Week
- Cause: Primary and Secondary Contributory Causes
- Geography: Region, City, Location
- Person: State, City

Dimension:

- Crash
- Vehicle
- Person
- Geography
- Injuries
- Date
- Cause
- Weather

Measures:

- DAMAGE: Damage costs reimbursement each client incurs for each crash event.
- NUM_UNITS: Represents the number of units involved in each incident.

MDX QUERY

2 - For each month, show the total damage costs for each location and the grand total with respect to the location.

		MonthlyDamage	TotalDamageByLocation
1	POINT (-87.56194603 41.760710194)	10961.9513015376	92815.8820690904
1	POINT (-87.63875619 41.885609917)	8529.68701790615	9529.68701790615
2	POINT (-87.701142758 41.884016475)	2625.35909502131	82367.1718864812
8	POINT (-87.607090628 41.751143218)	8599.16187106087	8599.16187106087
11	POINT (-87.662996825 41.998269694)	566.665907445764	566.665907445764
1	POINT (-87.53663321 41.713787169)	1325.65902772341	15464.352981877
1	POINT (-87.640628921 41.86897905)	1000	26294.8537699183
2	POINT (-87.667570048 41.925268887)	2479.68130812546	21859.3355945357

Steps:

- Definition of the Monthly Total of Damages
 - New member [Measures]. [MonthlyDamage], management of null values with the IIF function, Avoid errors in calculations due to missing data:
 - If the value is NULL replace with 0
 - Otherwise use the original value [Measures]. [DAMAGE]
- Creating the member [Measures]. [TotalDamageByLocation]
 - Use the SUM function to add up the damage over all months
 - Aggregation based on hierarchy [Dimension Date]. [Hierarchy]. [CRASH MONTH]
- The NONEMPTY filter ensures relevant and readable data
 - valid combinations of months and locations with data

```

WITH
-- Totale dei danni per ciascun mese
MEMBER [Measures].[MonthlyDamage] AS
    IIF(
        ISNULL([Measures].[DAMAGE]),
        0,
        [Measures].[DAMAGE]
    )

-- Totale complessivo dei danni per ciascuna località
MEMBER [Measures].[GrandTotal] AS
    SUM(
        [Dimension Date].[Hierarchy].[CRASH MONTH].Members,
        [Measures].[DAMAGE]
    )

SELECT
{
    [Measures].[MonthlyDamage],
    [Measures].[GrandTotal]
} ON COLUMNS,
NONEMPTY(
    [Dimension Date].[Hierarchy].[CRASH MONTH].Members *
    [Dimension Geography].[Hierarchy].[LOCATION].Members
) ON ROWS
FROM [Group ID 21]

```


MDX QUERY

4 - For each location, show the damage costs increase or decrease, in percentage, with respect to the previous year

Once the total values for the current and previous years are obtained, the difference between the two is calculated to highlight the magnitude of the change. Additionally, a percentage is computed to represent how significant this change is relative to the total of the previous year. This analysis is repeated for each location, taking into account all the years available in the dataset.

POINT (-87.620133867 41.707325932)	2015	9900.08169097808	(Null)	9900.08169097808	(Null)
POINT (-87.620133867 41.707325932)	2016	1000	9900.08169097808	-8900.08169097808	-89.90%
POINT (-87.620133867 41.707325932)	2017	3134.36628689832	1000	2134.36628689832	213.44%
POINT (-87.620133867 41.707325932)	2018	7359.53056822798	3134.36628689832	4225.16428132966	134.80%
POINT (-87.598913808 41.707678016)	2017	1000	(Null)	1000	(Null)
POINT (-87.598913808 41.707678016)	2018	4320.30024206541	1000	3320.30024206541	332.03%
POINT (-87.598891418 41.707458804)	2017	1508.9778141286	(Null)	1508.9778141286	(Null)

```

WITH
-- Valore cumulativo del danno per l'anno corrente, con controllo per NULL e 0
MEMBER [Measures].[YearDamage] AS
    IIF(
        ISNULL([Measures].[DAMAGE]) OR [Measures].[DAMAGE] = 0,
        0,
        SUM(
            PERIODSTODATE(
                [Dimension Date].[Hierarchy].[CRASH YEAR],
                [Dimension Date].[Hierarchy].CurrentMember
            ),
            [Measures].[DAMAGE]
        )
    )

-- Valore cumulativo del danno per l'anno precedente, con controllo per NULL e 0
MEMBER [Measures].[YearDamagePrev] AS
    IIF(
        ISNULL([Measures].[DAMAGE]) OR [Measures].[DAMAGE] = 0,
        0,
        SUM(
            PERIODSTODATE(
                [Dimension Date].[Hierarchy].[CRASH YEAR],
                [Dimension Date].[Hierarchy].PrevMember
            ),
            [Measures].[DAMAGE]
        )
    )

-- Differenza assoluta del danno tra l'anno corrente e quello precedente
MEMBER [Measures].[DiffDamage] AS
    [Measures].[YearDamage] - [Measures].[YearDamagePrev]

-- Percentuale di variazione del danno, con controllo per valori nulli o 0 nel denominatore
MEMBER [Measures].[DiffPercDamage] AS
    IIF(
        [Measures].[YearDamagePrev] = 0 OR ISNULL([Measures].[YearDamagePrev]),
        NULL,
        ([Measures].[DiffDamage] / [Measures].[YearDamagePrev])
    ), FORMAT_STRING = "Percent"

-- Selezione dei dati
SELECT
{
    [Measures].[YearDamage],
    [Measures].[YearDamagePrev],
    [Measures].[DiffDamage],
    [Measures].[DiffPercDamage]
} ON COLUMNS,
NONEMPTY(
    [Dimension Geography].[Hierarchy].[LOCATION].Members *
    [Dimension Date].[Hierarchy].[CRASH YEAR].Members
) ON ROWS
FROM [Group ID 21]

```

MDX QUERY

5 - For each quarter, show all the locations where the number of vehicles involved exceeds the average number of vehicles involved in the corresponding quarter of the previous year. Also, report the increase in both percentages.

		NUM UNITS	AvgVehiclesPrevYear	PercIncreaseVehicles
1	POINT (-87.724499209 41.836972624)	8	4	10000.00%
2	POINT (-87.61444973 41.736613967)	14	4	25000.00%
2	POINT (-87.62139584 41.894063163)	10	4	15000.00%
2	POINT (-87.577547876 41.780503236)	38	2	180000.00%
2	POINT (-87.624424687 41.888432092)	18	4	35000.00%
2	POINT (-87.626328333 41.89239608)	8	6	3333.33%
2	POINT (-87.620324126 41.867403983)	12	6	10000.00%
2	POINT (-87.623174208 41.895863834)	18	4	35000.00%
2	POINT (-87.765268498 41.887804621)	14	4	25000.00%
2	POINT (-87.74565352 41.886851263)	16	4	30000.00%
2	POINT (-87.634293268 41.915249987)	18	4	35000.00%
2	POINT (-87.620803561 41.887910042)	14	2	60000.00%

We analyze changes in vehicle involvement in crashes across specific locations and quarters, emphasizing cases where there is a notable increase compared to the previous year. By identifying locations and quarters with above-average growth, this analysis reveals potential patterns or anomalies in crash trends.

```

WITH
-- Calcolo della media dei veicoli coinvolti per il trimestre dell'anno precedente
MEMBER [Measures].[AvgVehiclesPrevYear] AS
    AVG(
        ParallelPeriod(
            [Dimension Date].[Hierarchy].[CRASH YEAR], 1, [Dimension Date].[Hierarchy].CurrentMember
        ),
        CoalesceEmpty([Measures].[NUM UNITS], 0)
    )

-- Percentuale di incremento dei veicoli rispetto alla media dell'anno precedente
MEMBER [Measures].[PercIncreaseVehicles] AS
    IIF(
        [Measures].[AvgVehiclesPrevYear] = 0,
        NULL,
        (([Measures].[NUM UNITS] - [Measures].[AvgVehiclesPrevYear]) / [Measures].[AvgVehiclesPrevYear]) * 100
    ), FORMAT_STRING = "Percent"

SELECT
{
    [Measures].[NUM UNITS],
    [Measures].[AvgVehiclesPrevYear],
    [Measures].[PercIncreaseVehicles]
} ON COLUMNS,

NONEMPTY(
    FILTER(
        [Dimension Date].[Hierarchy].[QUARTER].Members *
        [Dimension Geography].[Hierarchy].[LOCATION].Members,
        CoalesceEmpty([Measures].[NUM UNITS], 0) > 0 AND
        CoalesceEmpty([Measures].[AvgVehiclesPrevYear], 0) > 0 AND
        [Measures].[NUM UNITS] > [Measures].[AvgVehiclesPrevYear]
    )
) ON ROWS

FROM [Group ID 21]

```

MDX QUERY

7 - For each location, it calculates the total damage, the damage caused by bad weather as the primary cause, and displays the percentage of weather-related damage relative to the total.

The goal is to analyze the impact of weather as a primary cause of damage across various locations, comparing the damage associated with weather to the total recorded damage. The objective is to provide a clear understanding of how much weather has contributed to overall damages in each geographic area, evaluating both the absolute and relative impact.

		WeatherDamage	WeatherDamagePercentage
POINT (-87.592695408 41.799608974)	WEATHER	4966.27754450055	1663.77%
POINT (-87.607472951 41.73467445)	WEATHER	7984.85149920139	2675.03%
POINT (-87.764666198 41.953320153)	WEATHER	11783.3138717689	3947.57%
POINT (-87.654040394 41.936160611)	WEATHER	5115.11919488863	1713.63%

```

WITH
-- Calcolo del danno totale per ciascuna località
MEMBER [Measures].[TotalDamageByLocation] AS
    SUM(
        [Dimension Geography].[Hierarchy].[LOCATION].MEMBERS,
        IIF(
            ISEMPTY([Measures].[DAMAGE]),
            0,
            [Measures].[DAMAGE]
        )
    )

-- Calcolo del danno totale associato al maltempo come causa primaria
MEMBER [Measures].[WeatherDamage] AS
    SUM(
        {[Dimension Cause].[Hierarchy].[PRIM CONTRIBUTORY CAUSE].[WEATHER]},
        IIF(
            ISEMPTY([Measures].[DAMAGE]),
            0,
            [Measures].[DAMAGE]
        )
    )

-- Percentuale del danno legato al maltempo rispetto al danno totale
MEMBER [Measures].[WeatherDamagePercentage] AS
    IIF(
        [Measures].[TotalDamageByLocation] = 0,
        NULL,
        ([Measures].[WeatherDamage] / [Measures].[TotalDamageByLocation]) * 100
    ), FORMAT_STRING = "Percent"

SELECT
    {
        [Measures].[WeatherDamage],
        [Measures].[WeatherDamagePercentage]
    } ON COLUMNS,
    NONEMPTY(
        [Dimension Geography].[Hierarchy].[LOCATION].MEMBERS *
        {[Dimension Cause].[Hierarchy].[PRIM CONTRIBUTORY CAUSE].[WEATHER]}
    ) ON ROWS
FROM [Group ID 21]

```

MDX QUERY

8 - For each year, show the most frequent cause of crashes and the corresponding total damage costs. The primary crash contributing factor is given twice the weight of the secondary factor in the analysis. Additionally, show the overall most frequent crash cause across all years.

	Most Frequent Cause Per Year	TOT DMG for Most Freq per Year	Most Frequent Cause Overall
2014	UNABLE TO DETERMINE	20320.8738914341	UNABLE TO DETERMINE
2015	UNABLE TO DETERMINE	15819111.3048716	UNABLE TO DETERMINE
2016	UNABLE TO DETERMINE	76092040.9310556	UNABLE TO DETERMINE
2017	UNABLE TO DETERMINE	141370041.004631	UNABLE TO DETERMINE
2018	UNABLE TO DETERMINE	195426130.837452	UNABLE TO DETERMINE
2019	UNABLE TO DETERMINE	4370729.93410934	UNABLE TO DETERMINE

Analysis of Primary Accident Causes

- Most Frequent Cause Per Year: The primary contributory cause with the highest frequency each year.
- Total Damage for Most Frequent Cause Per Year: The total damage associated with the most frequent cause for each year.
- Most Frequent Cause Overall: The single most recurring primary contributory cause across all years.

```
WITH
-- Conteggio della frequenza delle cause primarie (considera tutte le righe con questa causa)
MEMBER [Measures].[PRIM CAUSE Frequency] AS
    SUM(
        NONEMPTY(
            [Dimension Date].[CRASH DATE].Members,
            ([Measures].[DAMAGE], [Dimension Cause].[PRIM CONTRIBUTORY CAUSE].CurrentMember,
            [Dimension Date].[CRASH YEAR].CurrentMember)
        ),
        1 -- Conta tutte le righe associate
    )

-- Conteggio della frequenza delle cause secondarie (considera tutte le righe con questa causa)
MEMBER [Measures].[SEC CAUSE Frequency] AS
    SUM(
        NONEMPTY(
            [Dimension Date].[CRASH DATE].Members,
            ([Measures].[DAMAGE], [Dimension Cause].[SEC CONTRIBUTORY CAUSE].CurrentMember,
            [Dimension Date].[CRASH YEAR].CurrentMember)
        ),
        1 -- Conta tutte le righe associate
    )

-- Calcolo della frequenza pesata
MEMBER [Measures].[WEIGHTED Freq] AS
    ([Measures].[PRIM CAUSE Frequency] * 2) +
    ([Measures].[SEC CAUSE Frequency] * 1)

-- Set per trovare la causa più frequente per anno
SET [Most Frequent Cause Per Year Set] AS
    TopCount(
        Filter(
            [Dimension Cause].[PRIM CONTRIBUTORY CAUSE].Members,
            [Dimension Cause].[PRIM CONTRIBUTORY CAUSE].CurrentMember.Name <> "All"
        ),
        1, -- Prendi la causa più frequente
        [Measures].[WEIGHTED Freq]
    )

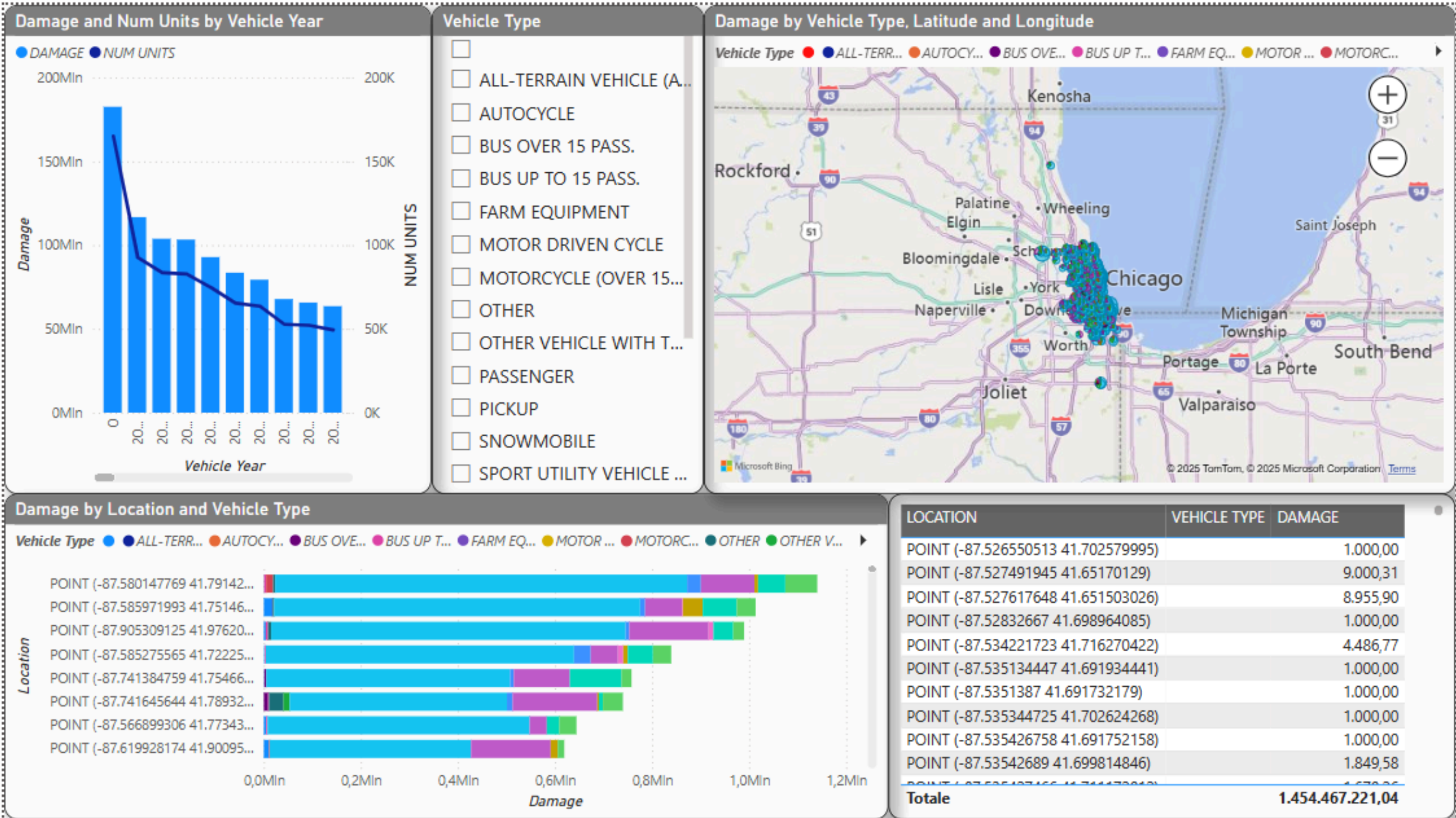
-- Nome della causa più frequente per anno
MEMBER [Measures].[Most Frequent Cause Per Year] AS
    Generate(
        [Most Frequent Cause Per Year Set],
        [Dimension Cause].[PRIM CONTRIBUTORY CAUSE].CurrentMember.Name
    )
```

```
-- Totale danni associati alla causa più frequente per anno
MEMBER [Measures].[TOT DMG for Most Freq per Year] AS
    Sum(
        [Most Frequent Cause Per Year Set],
        [Measures].[DAMAGE]
    )

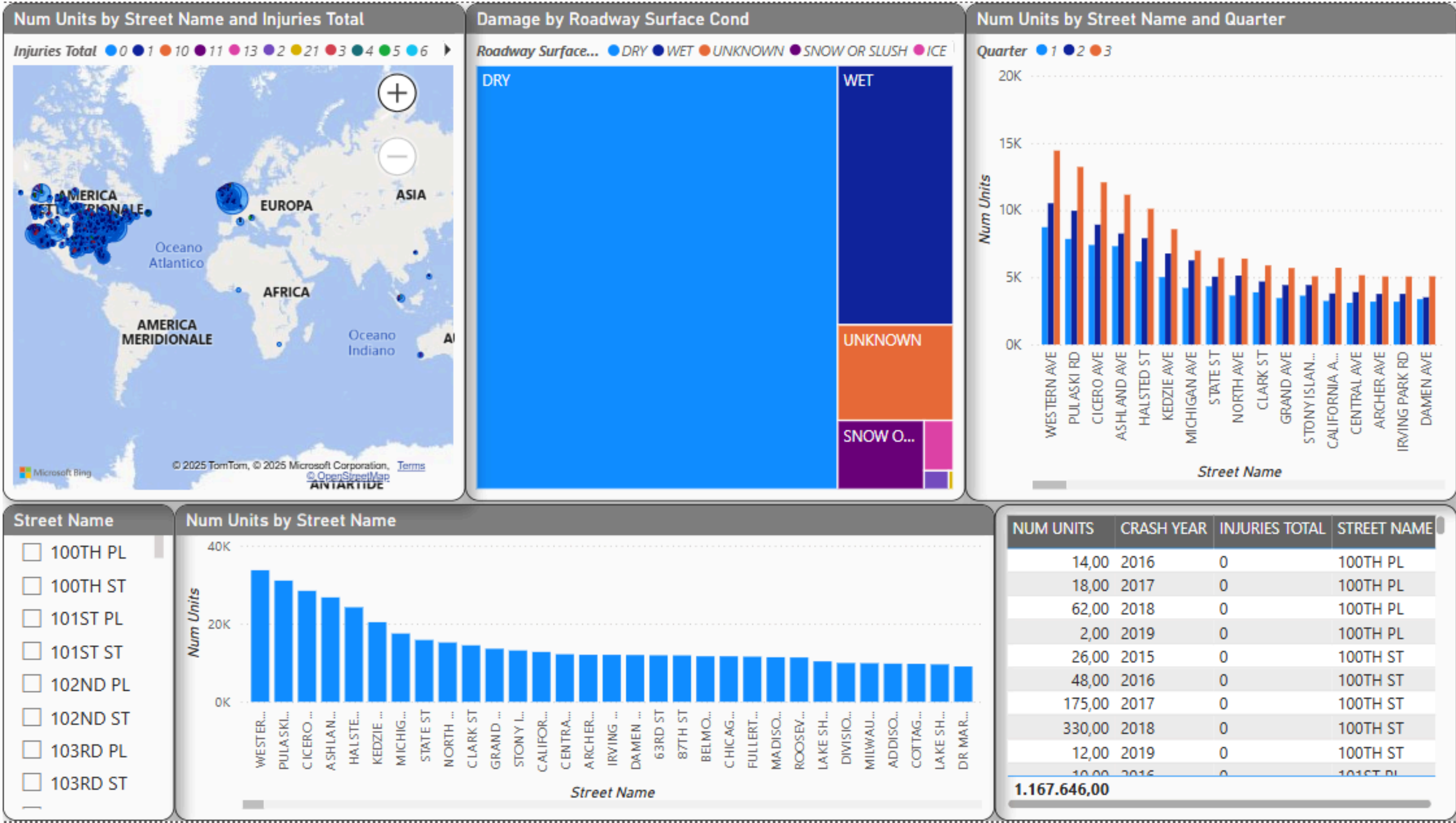
-- Causa più frequente generale (considerando tutti gli anni)
MEMBER [Measures].[Most Frequent Cause Overall] AS
    Generate(
        TopCount(
            Filter(
                [Dimension Cause].[PRIM CONTRIBUTORY CAUSE].Members,
                [Dimension Cause].[PRIM CONTRIBUTORY CAUSE].CurrentMember.Name <> "All"
            ),
            1, -- Prendi la causa più frequente
            Sum(
                [Dimension Date].[CRASH YEAR].Members,
                [Measures].[WEIGHTED Freq]
            )
        ),
        [Dimension Cause].[PRIM CONTRIBUTORY CAUSE].CurrentMember.Name
    )

-- Risultati finali
SELECT
{
    [Measures].[Most Frequent Cause Per Year],
    [Measures].[TOT DMG for Most Freq per Year],
    [Measures].[Most Frequent Cause Overall]
} ON COLUMNS,
{
    Filter(
        [Dimension Date].[CRASH YEAR].Members,
        [Dimension Date].[CRASH YEAR].CurrentMember.Name <> "All"
    )
} ON ROWS
FROM [Group ID 21]
```


CREATE A DASHBOARD THAT SHOWS THE GEOGRAPHICAL DISTRIBUTION OF THE TOTAL DAMAGE COSTS FOR EACH VEHICLE CATEGORY.



CREATE A PLOT/DASHBOARD THAT YOU DEEM INTERESTING W.R.T. THE DATA AVAILABLE IN YOUR CUBE, FOCUSSING ON DATA ABOUT THE STREET.



CREATE A PLOT/DASHBOARD THAT YOU DEEM INTERESTING W.R.T. THE DATA AVAILABLE IN YOUR CUBE, FOCUSSING ON DATA ABOUT THE PEOPLE INVOLVED IN A CRASH.

