

**CSI4142 Fundamentals of Data Science**  
**Project Phase 3 - OLAP Queries and BI Dashboard**  
**April 5, 2024**

**Shannon Noah - 300163898**  
**Serena Iyoha - 300187757**  
**Tina Trinh - 300175427**

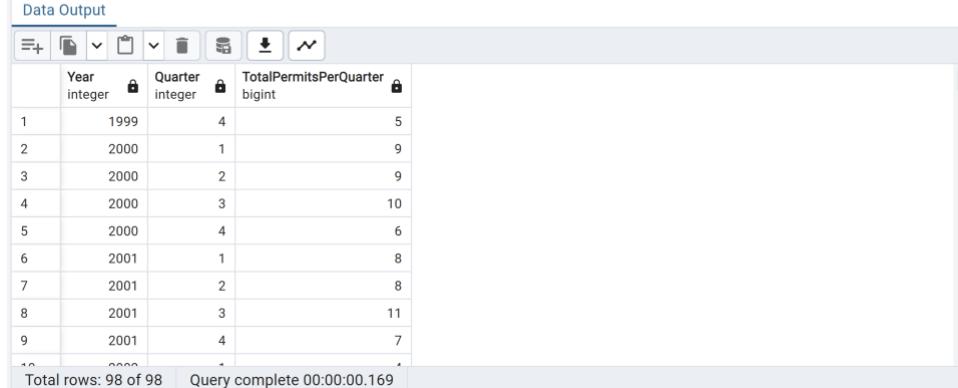
**GitHub Repo Link to Project:** <https://github.com/tinatrinh8/Data-Science-Gentrification-Study>  
 (Scripts to execute the SQL queries are found in “Phase 3 OLAP Queries” Folder)

## Part A. OLAP queries

### Part 1. Standard OLAP operations

#### a. Drill down and roll up

*Query 1: Drilling down the date dimension to get the number of permits by quarter*

Query History		
<pre>1 SELECT d."Year", d."Quarter", COUNT(*) AS "TotalPermitsPerQuarter" 2 FROM "BuildingPermitsDimension" p 3 JOIN "DateDimension" d ON p."Application_Date_Key" = d."Date_Key" 4 GROUP BY d."Year", d."Quarter" 5 ORDER BY d."Year", d."Quarter";</pre>		
Data Output		
		
Total rows: 98 of 98    Query complete 00:00:00.169		

*Query 2: Drilling down to get the number of permits by street name, rather than per Ward*

Query History		
<pre>1 SELECT 2     (p."Street_Name"    ' '    p."Street_Type"    ' '    3      COALESCE(p."Street_Direction", '')) AS "FullStreetName", 4     COUNT(*) AS "TotalPermits" 5 FROM 6     "BuildingPermitsDimension" p 7 GROUP BY 8     "FullStreetName" 9 ORDER BY 10    "FullStreetName";</pre>		
Data Output		
		
Total rows: 1000 of 2351    Query complete 00:00:00.192		

## b. Slice

*Query 3: Slicing to get the total number of building permits in the 13th Ward*

The screenshot shows a data analysis interface with two main sections: a query editor and a data output viewer.

**Query Editor:**

```
1 SELECT
2     f."Ward_ID",
3     w."Ward_Name",
4     COUNT(*) AS "TotalPermitsWard13"
5 FROM
6     "development_fact_table" f
7 JOIN
8     "BuildingPermitsDimension" p
9     ON f."Permit_Key" = p."Permit_Key"
10 JOIN
11     "WardDimension" w ON f."Ward_ID" = w."Ward_ID"
12 WHERE
13     f."Ward_ID" = 13
14 GROUP BY
15     f."Ward_ID", w."Ward_Name";|
```

**Data Output:**

	Ward_ID	Ward_Name	TotalPermitsWard13
1	13	Toronto Centre	182

## c. Dice

*Query 4: Dicing to form a sub-cube, which has dimensions 'Ward\_ID' and 'Year'*

Query History

```

1 SELECT
2     f."Ward_ID",
3     d."Year",
4     COUNT(*) AS "TotalPermits"
5 FROM
6     "development_fact_table" f
7 JOIN
8     "BuildingPermitsDimension" p ON f."Permit_Key" = p."Permit_Key"
9 JOIN
10    "DateDimension" d ON p."Application_Date_Key" = d."Date_Key"
11 GROUP BY
12    f."Ward_ID",
13    d."Year"

```

Data Output

	Ward_ID	Year	TotalPermits
1	0	1999	1
2	0	2004	11
3	0	2008	1
4	0	2010	3

*Query 5: Dices the data to show annual totals for demolitions and replacements of affordable, mid-range, and high-end rental homes.*

Query History

```

1 SELECT
2     d."Year",
3     COUNT(*) AS "TotalDemolitions",
4     SUM(dd."Affordable Rental Homes for Demolition") AS "TotalAffordableDemolitions",
5     SUM(dd."Mid-Range Rental Homes for Demolition") AS "TotalMidRangeDemolitions",
6     SUM(dd."High-End Rental Homes for Demolition") AS "TotalHighEndDemolitions",
7     SUM(dd."Affordable Rental Homes Replaced") AS "TotalAffordableReplacements",
8     SUM(dd."Mid-Range Rental Homes Replaced") AS "TotalMidRangeReplacements",
9     SUM(dd."High-End Rental Homes Replaced") AS "TotalHighEndReplacements"
10    FROM
11        "DemolitionDimension" dd
12    JOIN
13        "DateDimension" d ON dd."Approval_Date_Key" = d."Date_Key"
14    GROUP BY
15        d."Year"
16    ORDER BY
17        d."Year";
18

```

Data Output

	Year	TotalDemolitions	TotalAffordableDemolitions	TotalMidRangeDemolitions	TotalHighEndDemolitions	TotalAffordableReplacements
1	2017	11	184	101	13	13
2	2018	6	150	174	8	8
3	2019	17	323	397	23	23
4	2020	12	152	39	39	39
5	2021	9	295	31	26	26
6	2022	23	519	268	80	80
7	2023	24	1047	511	49	49

\*Note that "TotalDemolitions" represents the count of demolition events, while the other columns, like "TotalAffordableDemolitions" represent the count of individual housing units affected by those demolitions.

#### d. Combining OLAP operations

*Query 6: Dices the data by year and occupation type to calculate the total population and the average ratio of individuals without certifications for each occupation.*

```

SELECT
    e."Employment" AS "OccupationType",
    w."Year",
    SUM(w."Population") AS "TotalPopulationForOccupation",
    AVG(w."No_Cert_Ratio") AS "AverageNoCertRatio"
FROM
    "EmploymentDimension" e
JOIN
    "ward_profile_fact_table" w ON e."Employment_Key" = w."Dimension_Key"
GROUP BY
    e."Employment",
    w."Year"
ORDER BY
    w."Year",
    e."Employment";
  
```

Data Output					
	OccupationType text	Year bigint	TotalPopulationForOccupation double precision	AverageNoCertR double precision	
1	All occupations	2016	4590990	0.1668316100C	
2	Business, finance and administration occupations	2016	784215	0.16653528937	
3	Health occupations	2016	602570	0.166369254E	
4	Management occupations	2016	769655	0.16461633587	
5	Natural and applied sciences and related occupations	2016	563000	0.1649047392	
6	Natural resources, agriculture and related production occupations	2016	431335	0.16469697073	
7	Occupations in art, culture, recreation and sport	2016	488380	0.1673460498C	
8	Occupations in education, law and social, community and government services	2016	857275	0.1643605888E	
9	Occupations in manufacturing and utilities	2016	532800	0.16637349357	
10	Sales and service occupations	2016	1181390	0.1664354035E	

From 2016 to 2021, there seems to be a greater percentage of the population who have a diploma or certificate (16% to 14% no certificate ratio) indicating an increase in the number of individuals with diplomas or certifications over time in these employment groups.

*Query 7: Analyzing how shelter costs and the percentage of income spent on shelter vary over time and by household type (Owned vs. Rented)*

Query History

```

1 SELECT
2     sd."Household_Type",
3     wp."Year",
4     AVG(sd."Average_Monthly_Shelter_Costs") AS "AvgShelterCosts",
5     AVG(sd."Percent_Spending_30_Percent_Or_More_On_Shelter")
6     AS "AvgPercentSpendingMoreThan30Pct",
7     SUM(wp."Population") AS "TotalPopulation"
8 FROM
9     "ShelterDimension" sd
10 JOIN
11     "ward_profile_fact_table" wp ON sd."Shelter_Key" = wp."Dimension_Key"
12 GROUP BY
13     sd."Household_Type",
14     wp."Year"
15 ORDER BY
16     wp."Year",
17     sd."Household_Type";
18
19 AS "AvgPercentSpendingMoreThan30Pct"

```

Data Output

	Household_Type	Year	AvgShelterCosts	AvgPercentSpendingMoreThan30Pct	TotalPopulation
	text	bigint	numeric	double precision	double precision
1	Owner	2016	1677.7307692307692308	27.35	1174155
2	Tenant	2016	1209.4615384615384615	46.64615384615385	1051645
3	Owner	2021	2027.4615384615384615	25.423076923076923	1205830
4	Tenant	2021	1513.0769230769230769	38.84615384615385	1115940

Query 8: Compares the change in population for each occupation type between 2016 and 2021 by summing up the populations across all wards

Query History

```

1 WITH Population2016 AS (
2     SELECT
3         e."Employment" AS "OccupationType",
4         SUM(w."Population") AS "Population2016"
5     FROM
6         "EmploymentDimension" e
7     JOIN
8         "ward_profile_fact_table" w ON e."Employment_Key" = w."Dimension_Key"
9     WHERE
10        w."Year" = 2016
11     GROUP BY
12        e."Employment"
13 ),
14 Population2021 AS (
15     SELECT
16         e."Employment" AS "OccupationType",
17         SUM(w."Population") AS "Population2021"
18     FROM
19         "EmploymentDimension" e
20     JOIN
21         "ward_profile_fact_table" w ON e."Employment_Key" = w."Dimension_Key"
22     WHERE

```

```

    WHERE
        w."Year" = 2021
    GROUP BY
        e."Employment"
)
SELECT
    p2016."OccupationType",
    p2016."Population2016",
    p2021."Population2021",
    (p2021."Population2021" - p2016."Population2016") AS "ChangeInPopulation"
FROM
    Population2016 p2016
JOIN
    Population2021 p2021 ON p2016."OccupationType" = p2021."OccupationType"
ORDER BY
    "ChangeInPopulation" DESC;

```

Data Output

	OccupationType text	Population2016 double precision	Population2021 double precision	ChangeInPopulation double precision
1	Sales and service occupations	1181390	2205210	1023820
2	Natural and applied sciences and related occupations	563000	543270	-19730
3	Business, finance and administration occupations	784215	709350	-74865
4	Trades, transport and equipment operators and related occupations	698845	559460	-139385
5	Occupations in art, culture, recreation and sport	488380	344435	-143945
6	Health occupations	602570	339270	-263300
7	Occupations in manufacturing and utilities	532800	203735	-329065
8	Occupations in education, law and social, community and government services...	857275	479710	-377565
9	Natural resources, agriculture and related production occupations	431335	28020	-403315
10	Management occupations	769655	310585	-459070
11	All occupations	4590990	3117675	-1473315

Total rows: 11 of 11    Query complete 00:00:00.096

*Query 9: Slice and dice to get the permits for new houses where the estimated construction cost is over 10000000*

Query    Query History

```

1 SELECT
2     bp."Permit_Key",
3     bp."Permit_Type",
4     df."Est_Const_Cost",
5     df."Ward_ID"
6 FROM
7     "BuildingPermitsDimension" bp
8 JOIN
9     development_fact_table df ON bp."Permit_Key" = df."Permit_Key"
10 WHERE
11     df."Est_Const_Cost" > 10000000
12     AND bp."Permit_Type" = 'New Houses'
13

```

## Data Output

	Permit_Key bigint	Permit_Type text	Est_Const_Cost integer	Ward_ID integer
1	574	New Houses	23000000	23
2	3337	New Houses	20000000	11
3	4489	New Houses	12000000	11
4	4512	New Houses	11000000	11
5	4646	New Houses	15000000	11
6	5922	New Houses	61013000	2

## Part 2. Explorative operation

### a. Iceberg

*Query 10: Returns the top five wards with the highest total estimated construction costs above \$50,000, suggesting areas with the most high-investment construction activity*

Query

Query History

```
1 SELECT
2     dft."Ward_ID",
3     SUM(dft."Est_Const_Cost") AS "TotalHighEndConstructionCost"
4 FROM
5     "development_fact_table" dft
6 WHERE
7     dft."Est_Const_Cost" > 50000
8 GROUP BY
9     dft."Ward_ID"
10 ORDER BY
11     "TotalHighEndConstructionCost" DESC
12 LIMIT 5;
```

Data Output

	Ward_ID integer	TotalHighEndConstructionCost bigint
1	24	11445323016
2	21	4272518596
3	12	4005277104
4	3	3429929003
5	15	3341595758

b. Widowing

*Query 11: Ranks wards by the average project duration for construction projects, displaying the ranks per ward for each of the last five years.*

```
Query  Query History

1  SELECT
2    dft."Ward_ID",
3    dd."Year",
4    AVG(dft."Application_to_Issuance_Duration") AS "AverageProjectDuration",
5    RANK() OVER (PARTITION BY dd."Year" ORDER BY AVG(dft."Application_to_Issuance_Duration")) DESC
6  FROM
7    "development_fact_table" dft
8  JOIN
9    "DateDimension" dd ON dft."Application_Date_Key" = dd."Date_Key"
10 WHERE
11   dd."Year" BETWEEN EXTRACT(YEAR FROM CURRENT_DATE) - 5 AND EXTRACT(YEAR FROM CURRENT_DATE)
12 GROUP BY
13   dft."Ward_ID",
14   dd."Year"
15 ORDER BY
16   dd."Year" DESC, "AverageProjectDuration" DESC;
17
```

## Data Output

	Ward_ID integer	Year integer	AverageProjectDuration numeric	DurationRank bigint
1	13	2024	56.000000000000000000	1
2	3	2024	45.000000000000000000	2
3	4	2024	40.000000000000000000	3
4	20	2024	40.000000000000000000	3
5	1	2024	36.600000000000000000	5
6	12	2024	35.000000000000000000	6
7	17	2024	33.500000000000000000	7
8	9	2024	32.000000000000000000	8
9	0	2024	30.500000000000000000	9
10	15	2024	29.000000000000000000	10

### c. Using the window clause

*Query 12: Compares the number of building permit issued per ward to that of the previous and next years*

```
Query Query History

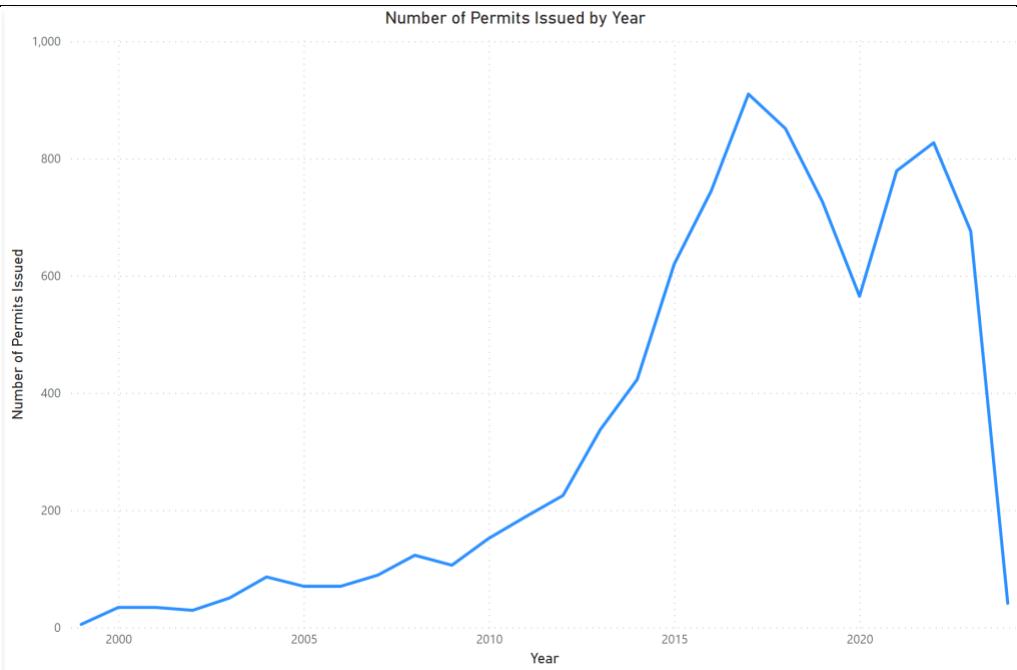
1 SELECT
2     dft."Ward_ID",
3     dd."Year",
4     COUNT(dft."Permit_Key") AS PermitsIssued,
5     LAG(COUNT(dft."Permit_Key"), 1) OVER (PARTITION BY dft."Ward_ID" ORDER BY dd."Year") AS Previous,
6     LEAD(COUNT(dft."Permit_Key"), 1) OVER (PARTITION BY dft."Ward_ID" ORDER BY dd."Year") AS Next,
7 FROM
8     "development_fact_table" dft
9 JOIN
10    "DateDimension" dd ON dft."Issued_Date_Key" = dd."Date_Key"
11 GROUP BY
12     dft."Ward_ID",
13     dd."Year";
```

Data Output

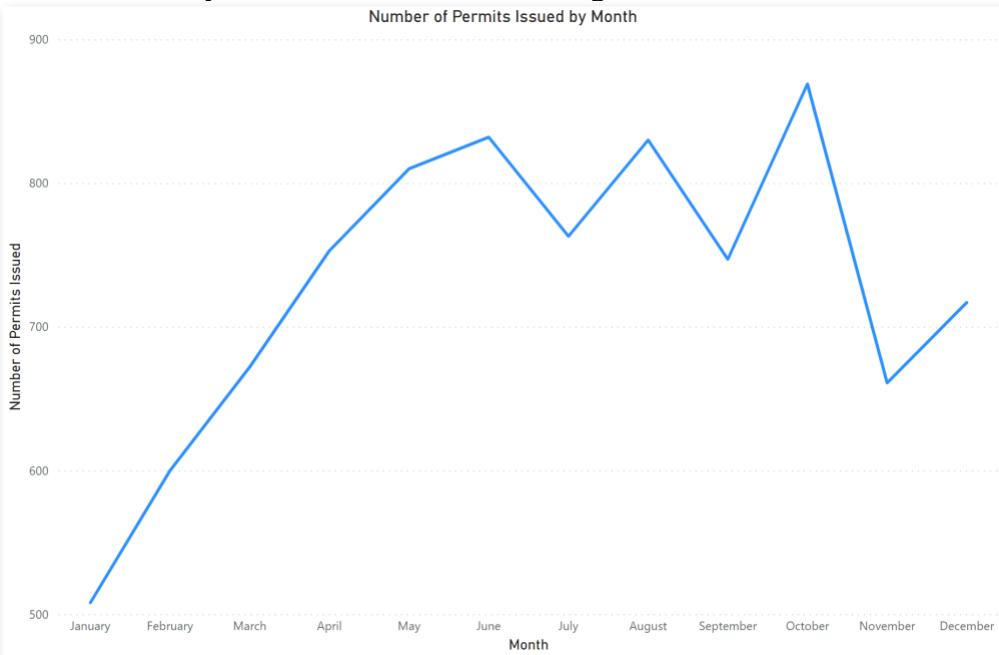
	Ward_ID integer	Year integer	permitsissued bigint	previousyearpermits bigint	nextyearpermits bigint
1	0	2000	1	[null]	11
2	0	2007	11	1	1
3	0	2008	1	11	2
4	0	2010	2	1	2
5	0	2011	2	2	2
6	0	2012	2	2	2
7	0	2013	2	2	2
8	0	2014	2	2	7
9	0	2015	7	2	3
10	0	2016	3	7	10
11	0	2017	10	3	9
12	0	2018	9	10	5
13	0	2019	5	9	7

## Part B. BI dashboard and Information Visualization

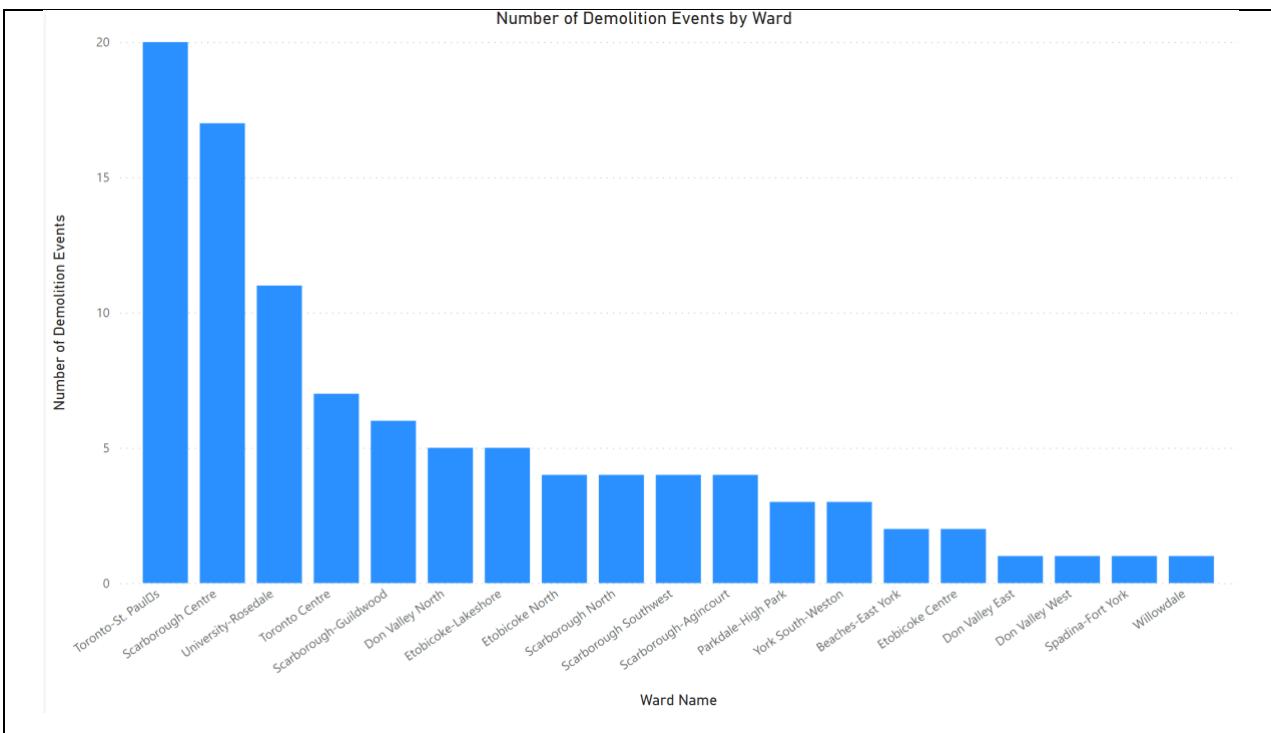
*Visualization I: Compares the number of permits issued over the years*



Ability to drill down to months using the date dimension:



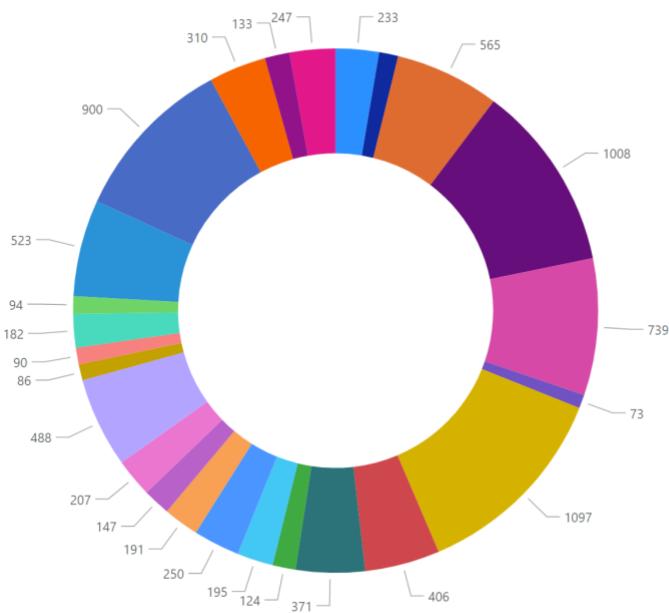
*Visualization 2: Examines the number of demolition events occurring within each ward*



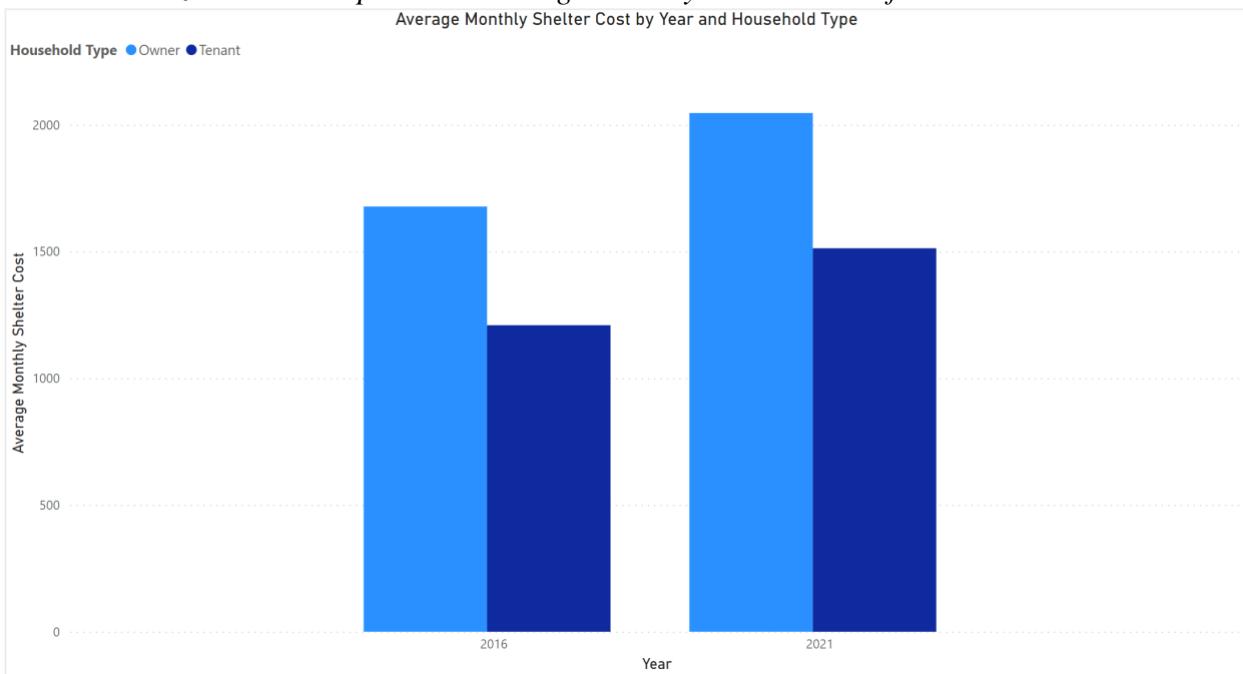
*Visualization 3: Displays the number of permits issued within each ward in a pie chart*

Number of Permits by Ward

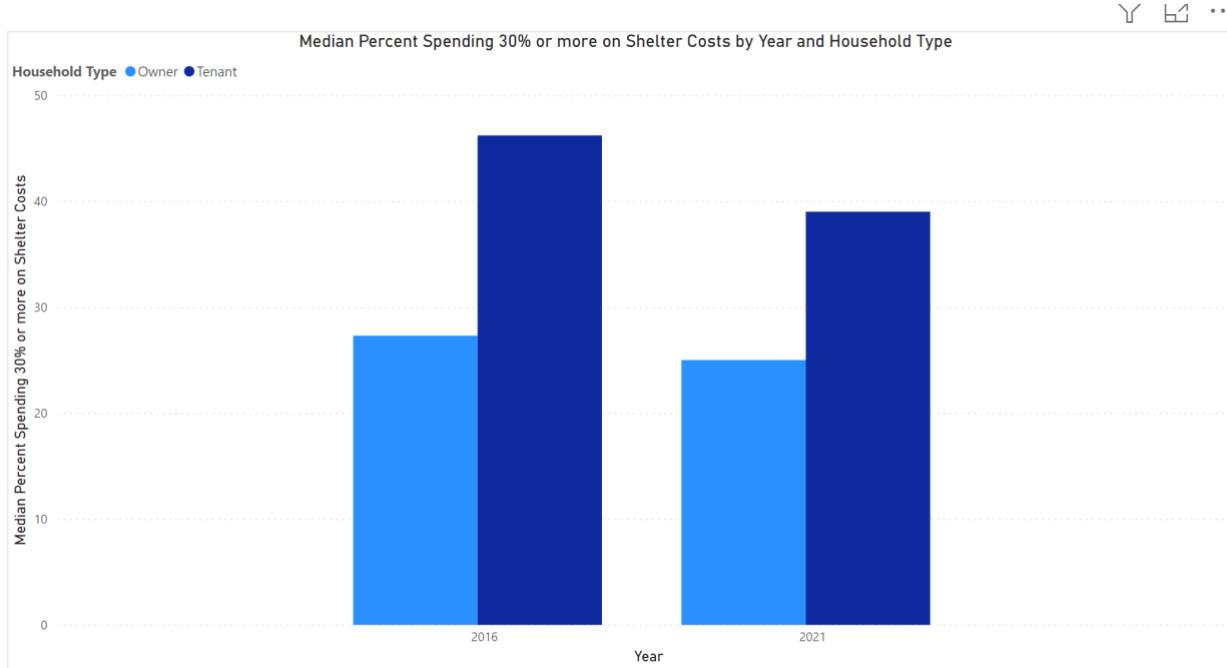
Ward Name
Beaches-East York
Davenport
Don Valley East
Don Valley North
Don Valley West
Eglinton-Lawrence
Etobicoke Centre
Etobicoke North
Etobicoke-Lakeshore
Humber River-Black Creek
Parkdale-High Park
Scarborough Centre
Scarborough North
Scarborough Southwest
Scarborough-Agincourt
Scarborough-Guildwood
Spadina-Fort York
Toronto
Toronto Centre
Toronto-Danforth
Toronto-St. Pauls
University-Rosedale
Willowdale
York Centre
York South-Weston



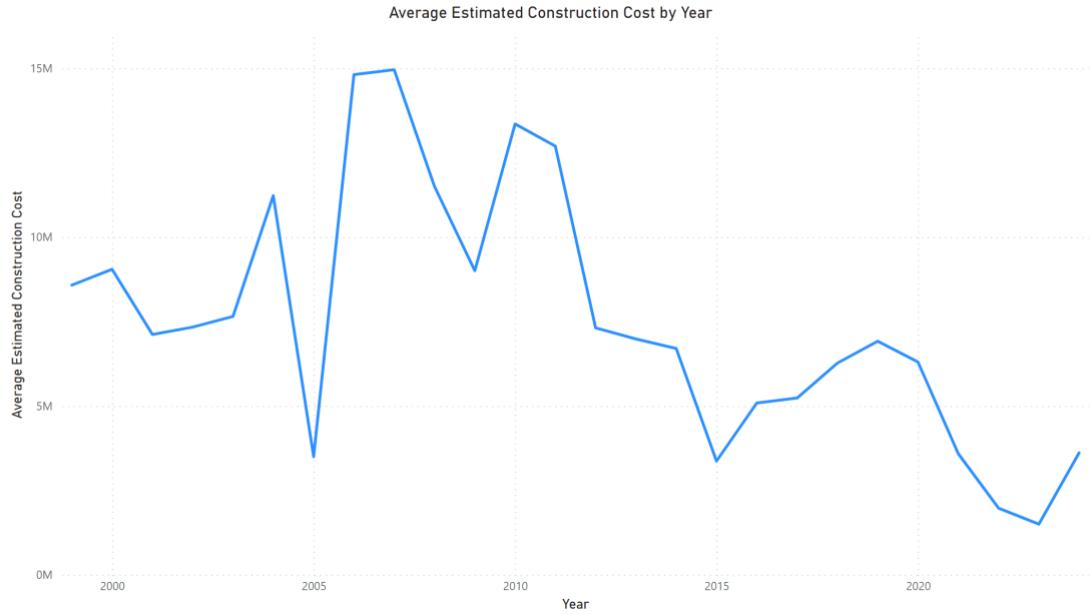
*Visualization 4: Compares the average monthly shelter costs of Owners vs. Tenants*



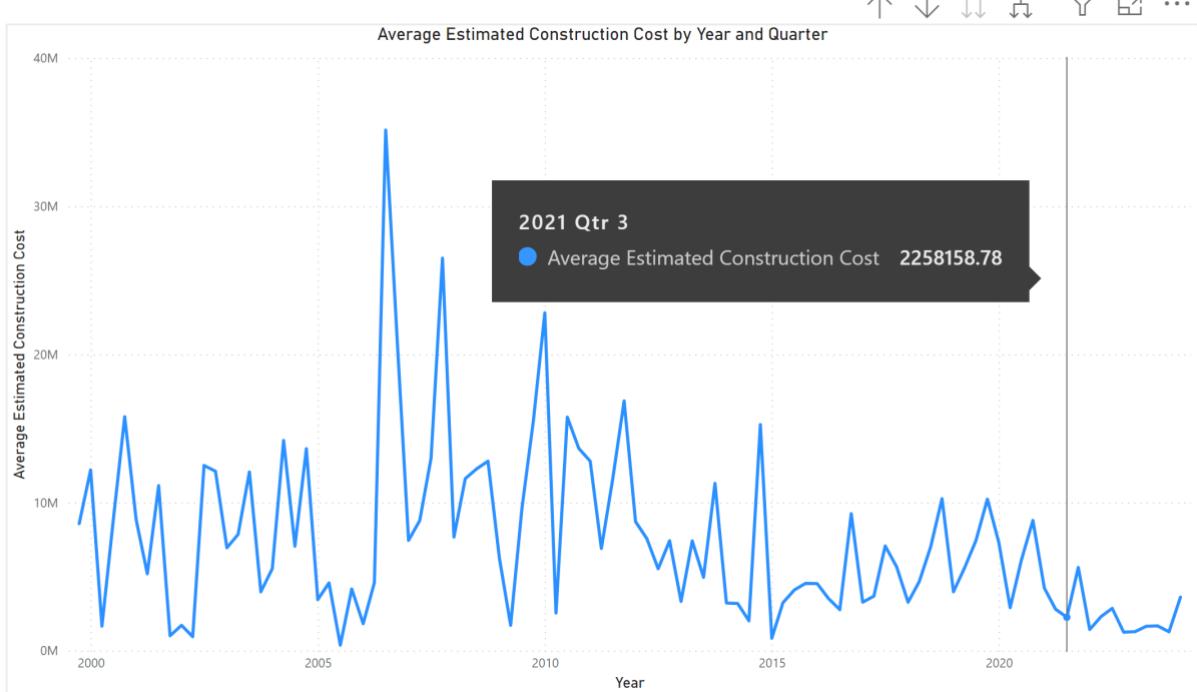
*Visualization 5: Compares the median percentage of population spending 30% or more of income on shelter costs*



*Visualization 6: Displays the trends in the average estimated construction cost for all wards over the years*



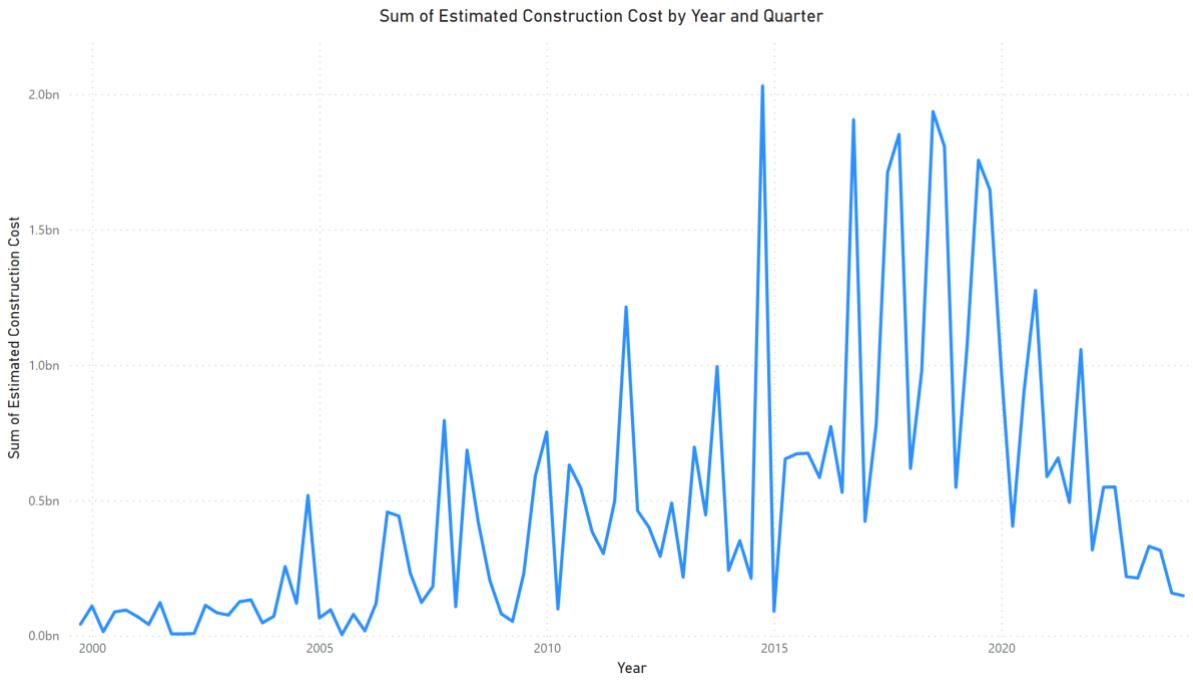
Ability to drill-down to show quarter with date dimension:



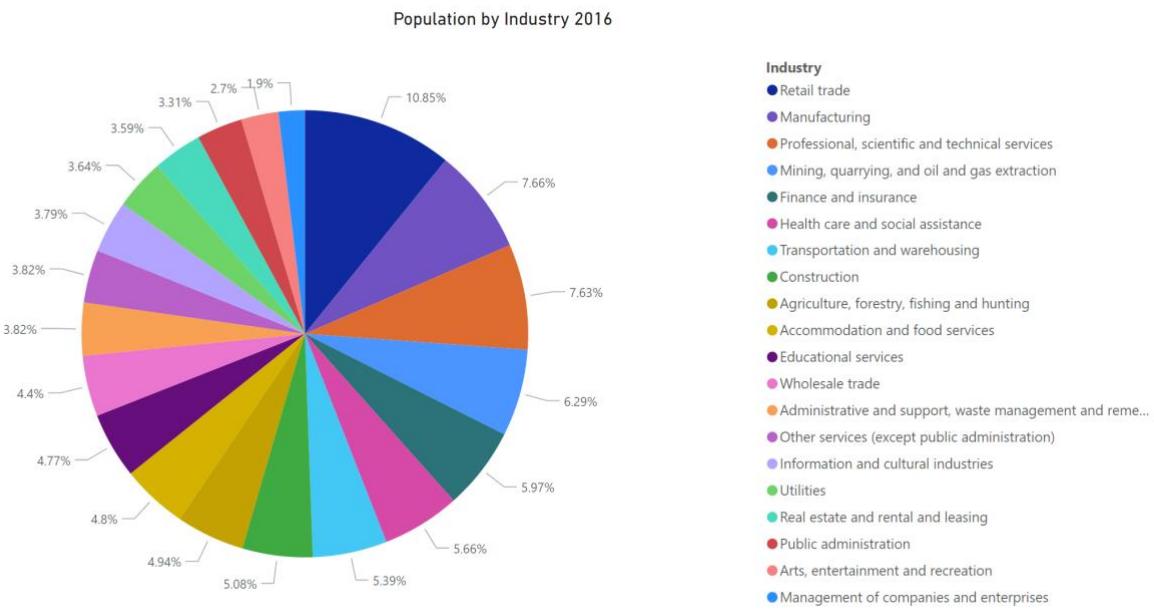
Visualization 7: Displays the trends in the sum of estimated construction cost for all wards over the years



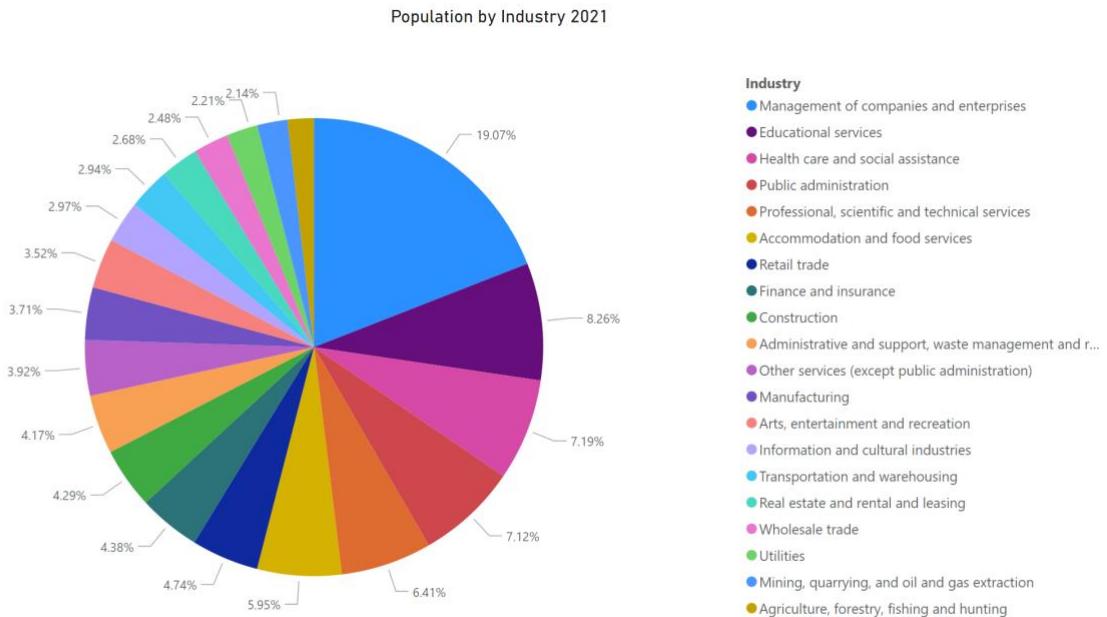
**Ability to drill-down to show quarter with date dimension:**



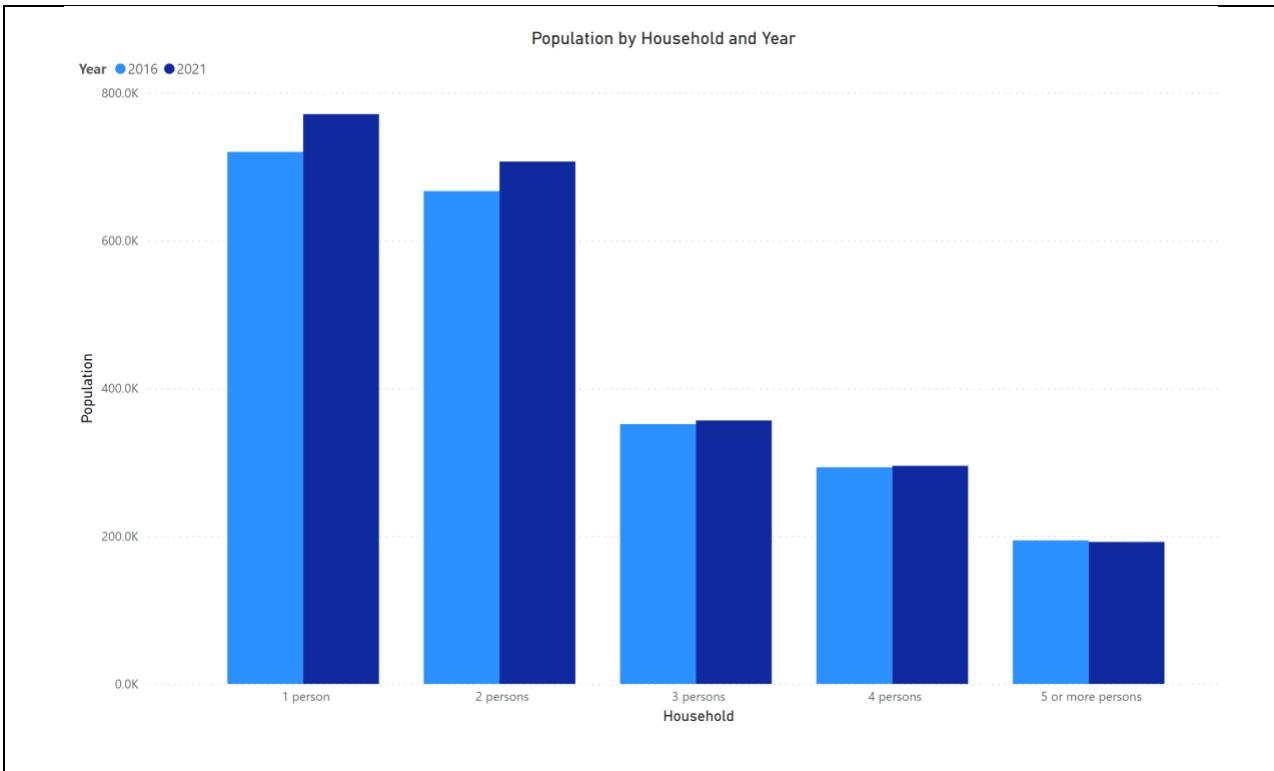
*Visualization 9: Displays the percentage of the population participating in each employment industry type for the year 2016*



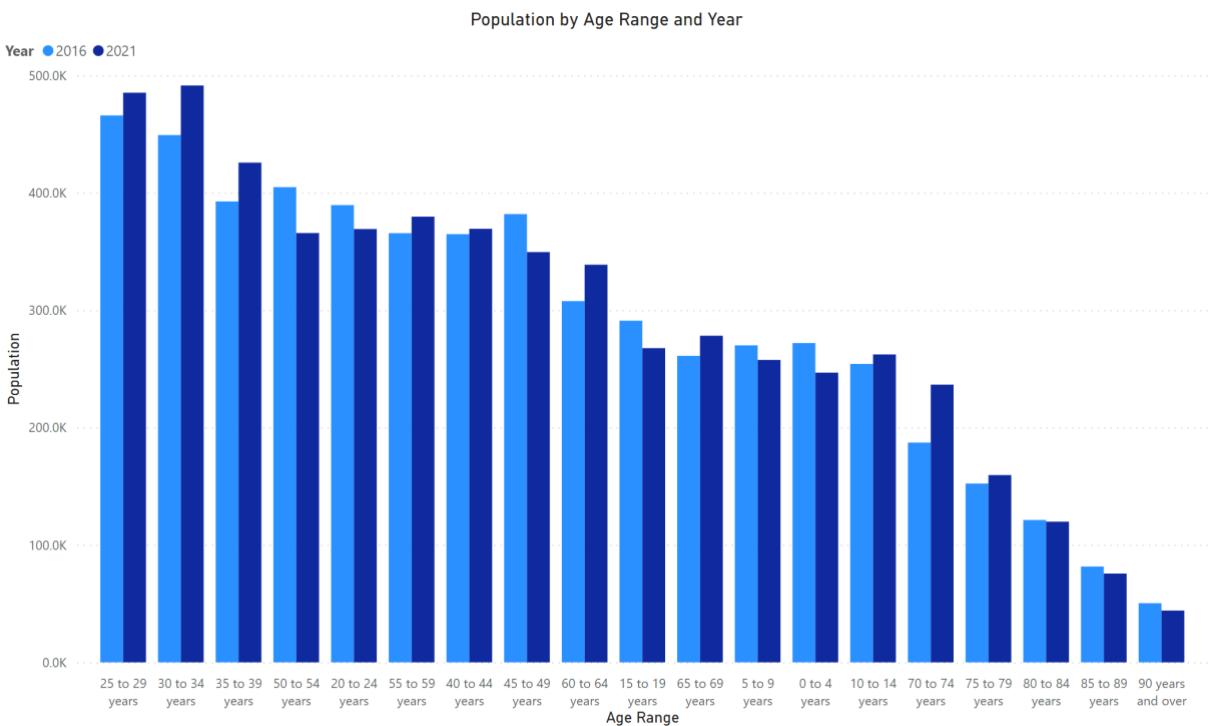
*Visualization 10: Displays the percentage of the population participating in each employment industry type for the year 2021*



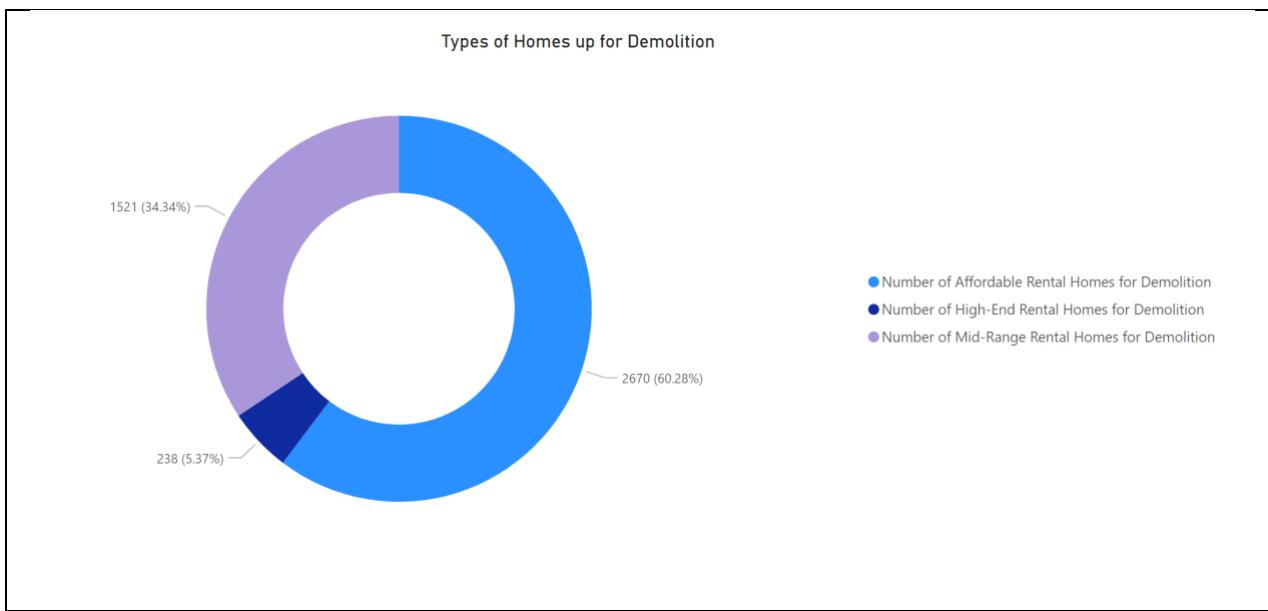
*Visualization 11: Comparing the change in population for household size from 2016 to 2021*



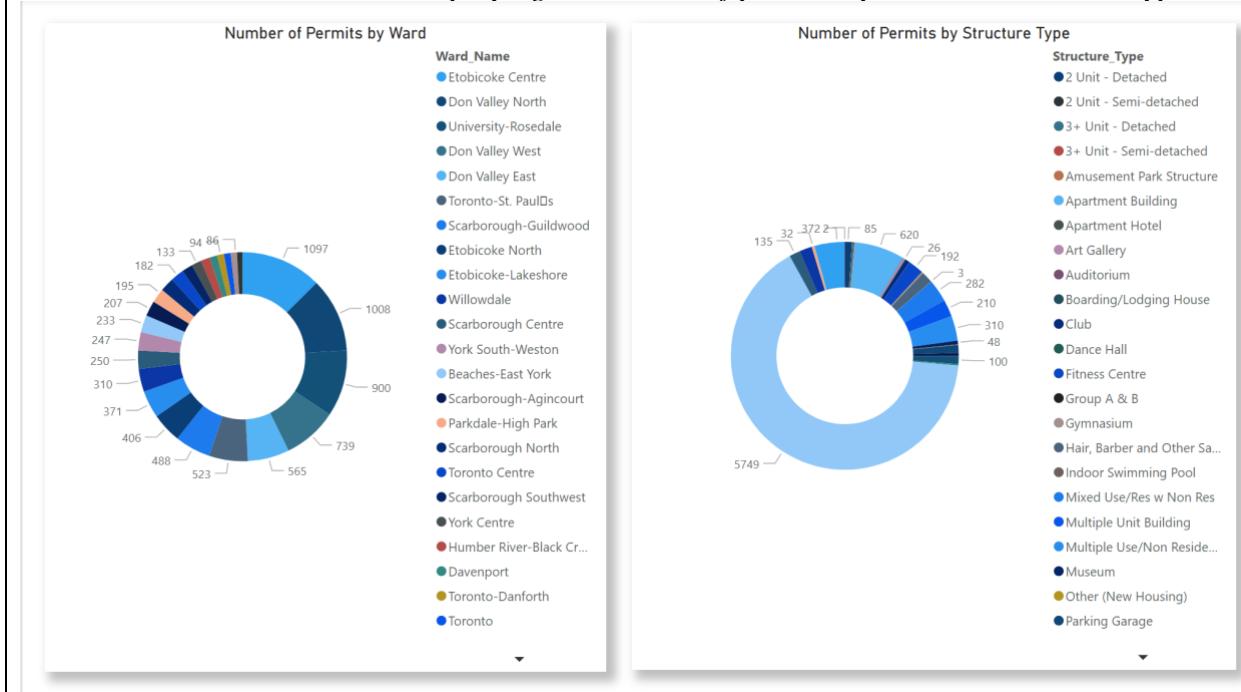
*Visualization 12: Comparing the change in population for age brackets from 2016 to 2021*



*Visualization 13: Displaying the number of total demolitions for each type of home*

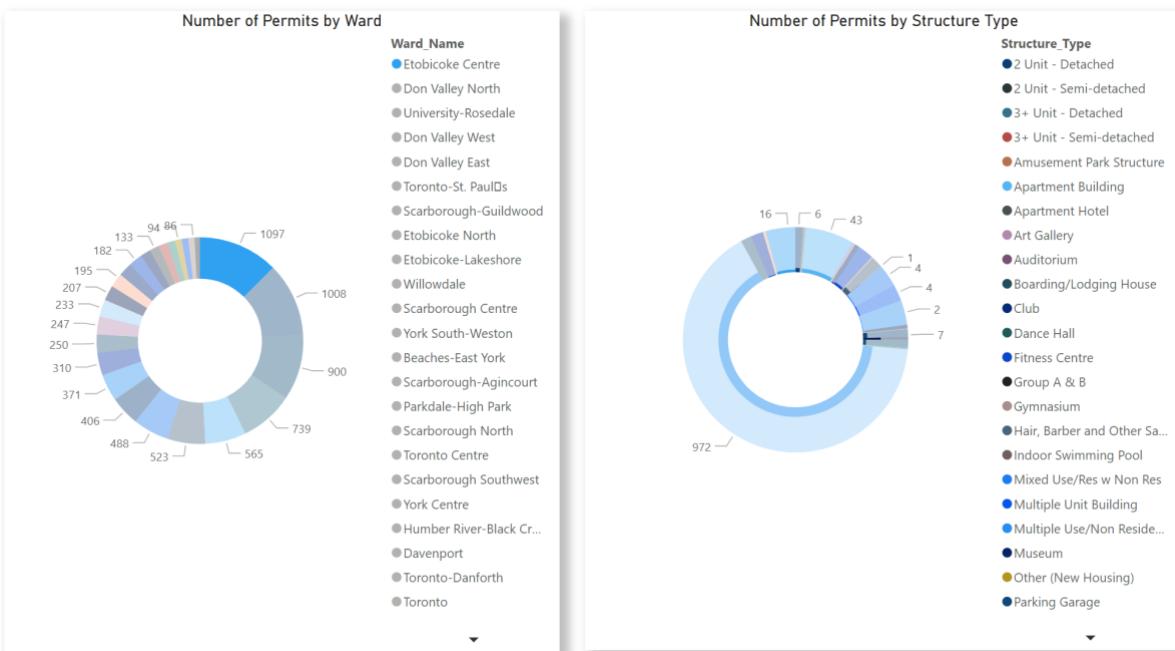


*Visualizations 14 & 15: Displaying the number of permits by ward and structure type*

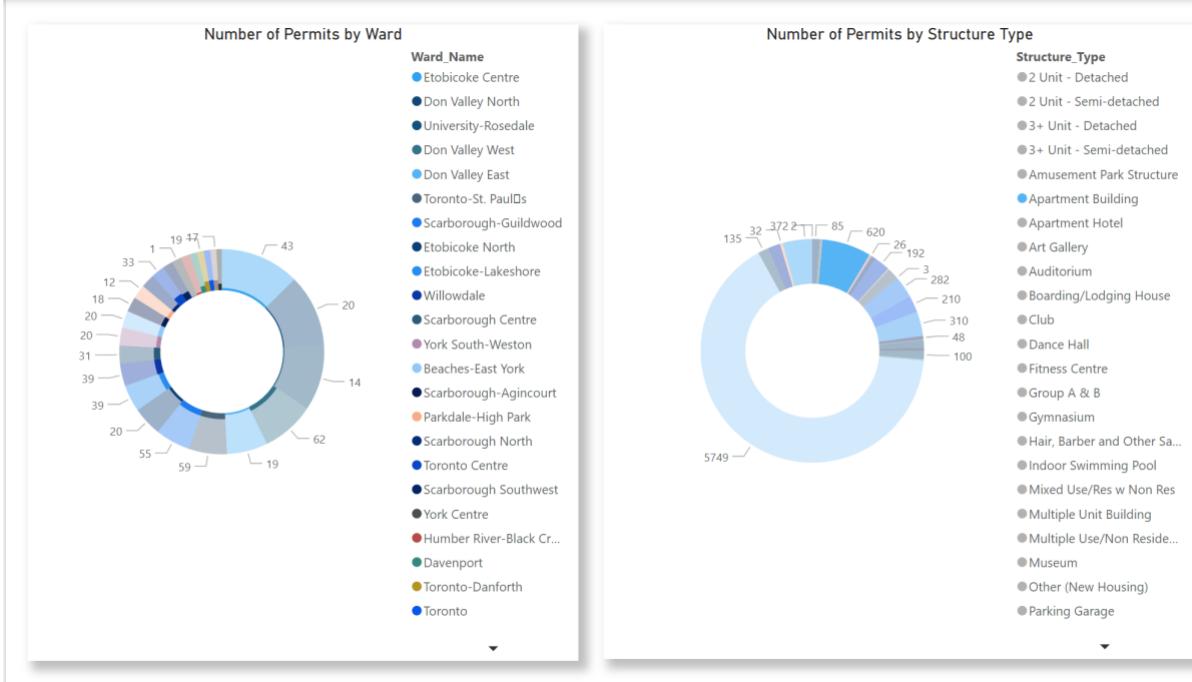


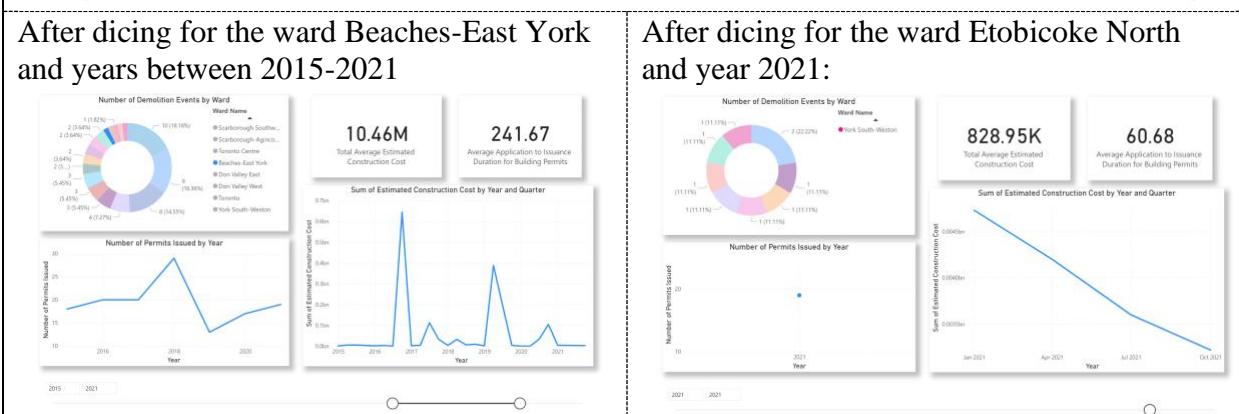
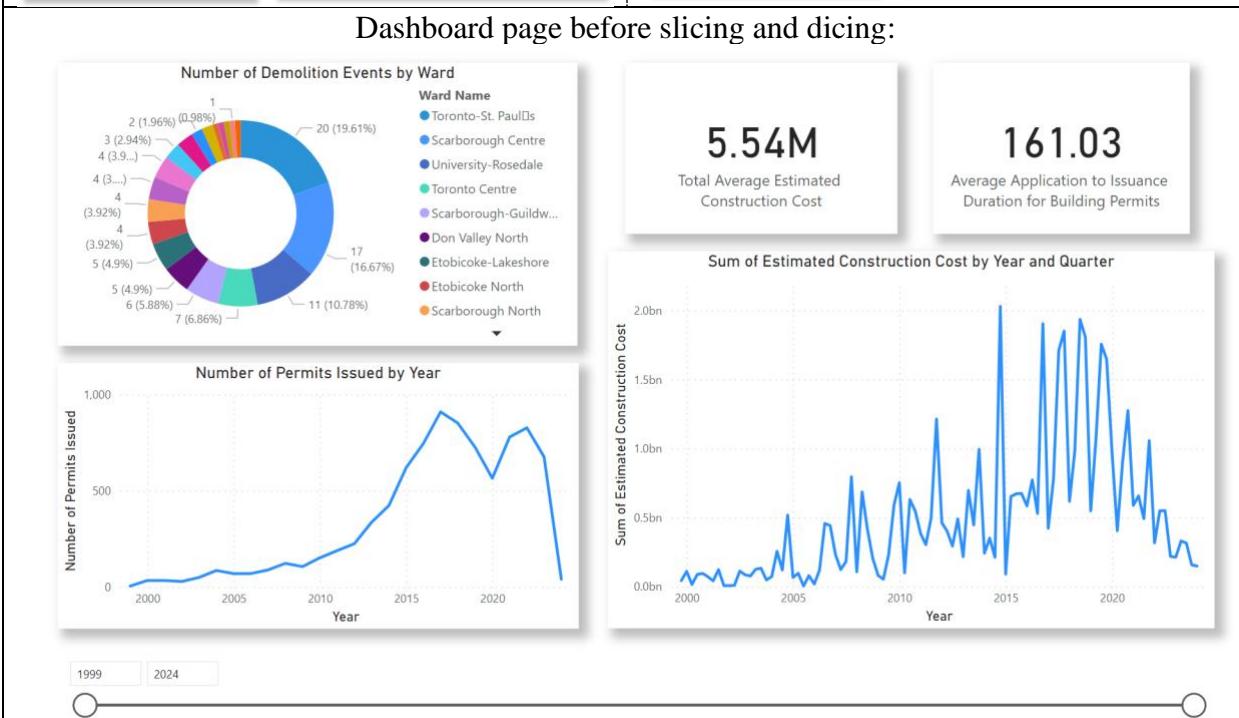
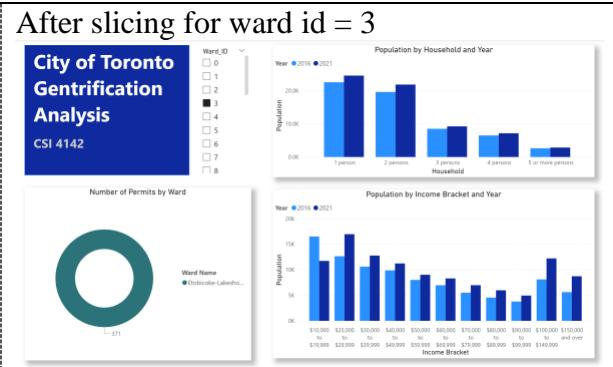
# Slicing and Dicing

By selecting Etobicoke Center as the ward, the user can see the breakdown of the structure types of the permits associated with that ward:



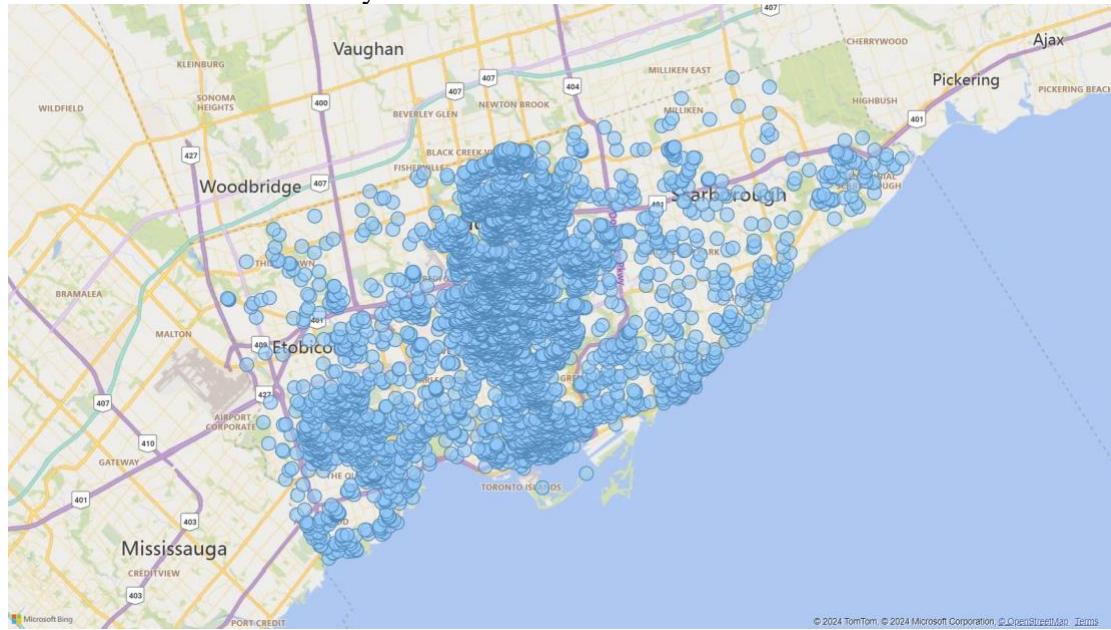
By selecting structure type “Apartment Building”, the user can see in which wards the permits of this structure type originated from:



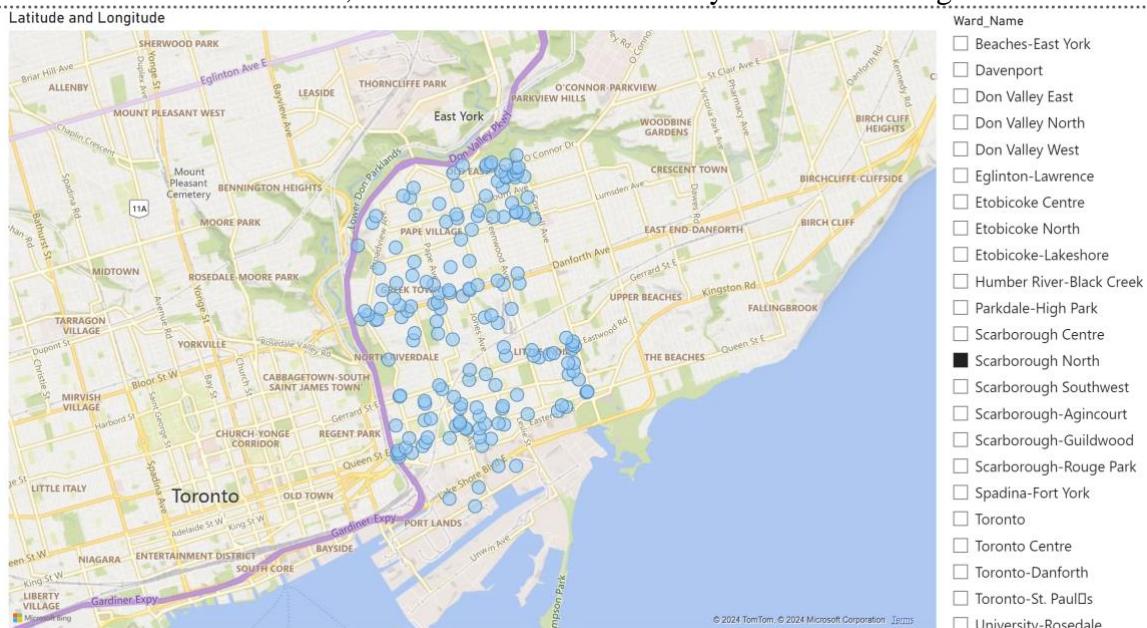


## Executing Top N and Bottom N Queries

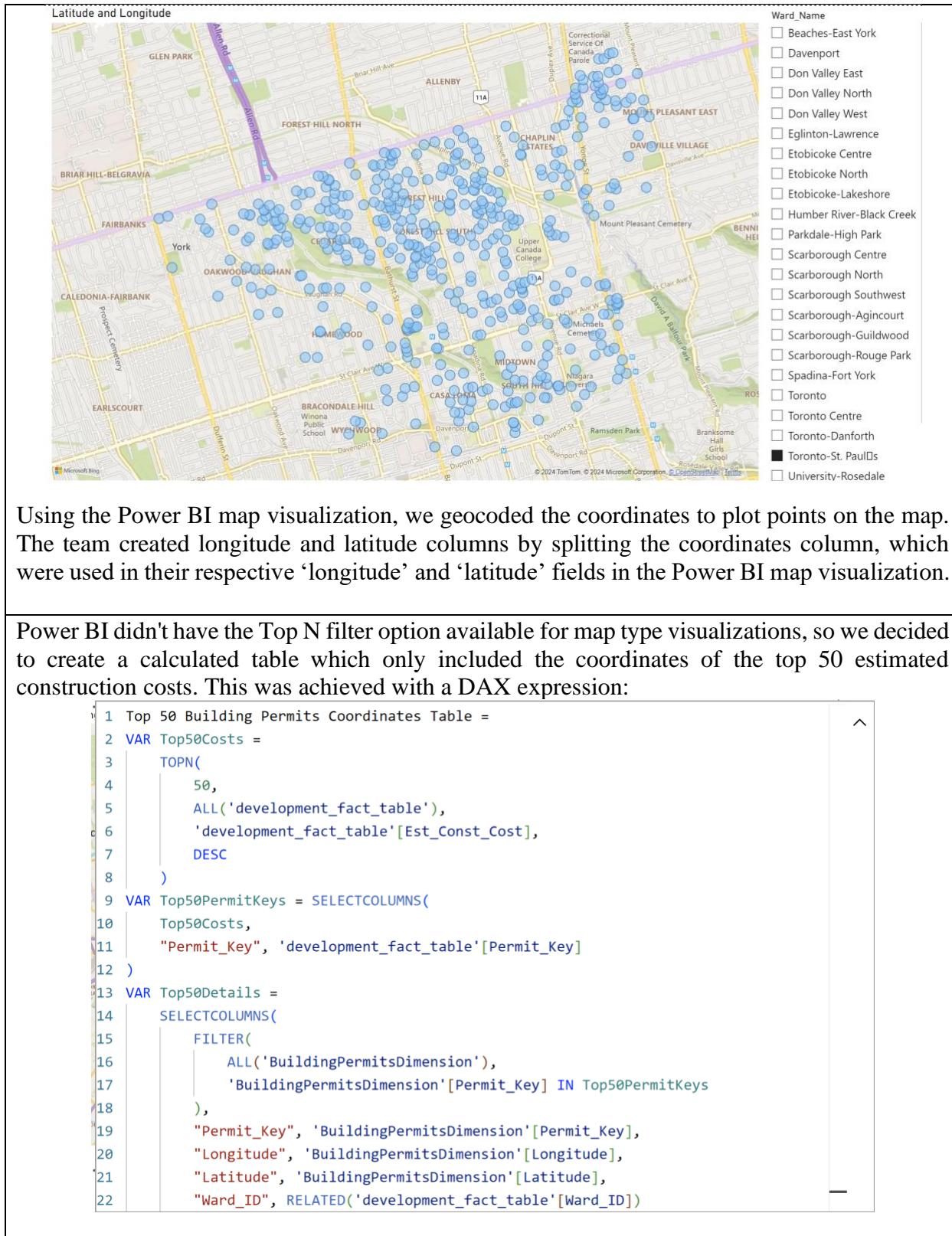
To visualize areas of increased construction activity, we plotted coordinates to map visualizations in Power BI. Below, all 8000+ construction projects are plotted, showing a trend of increased construction activity in central Toronto.



These projects can be divided by Toronto wards to explore trends in construction at the ward level. In the screenshot below, the coordinates are sliced by ward Scarborough North:



Here, they are sliced by Toronto St. Paul's:



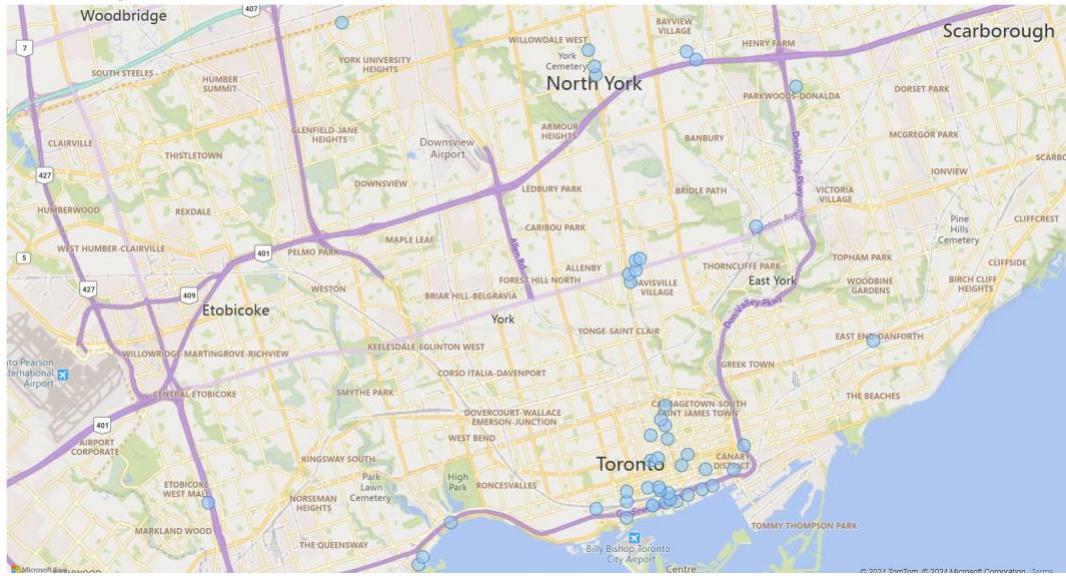
```

23 )
24 RETURN
25 Top50Details
26

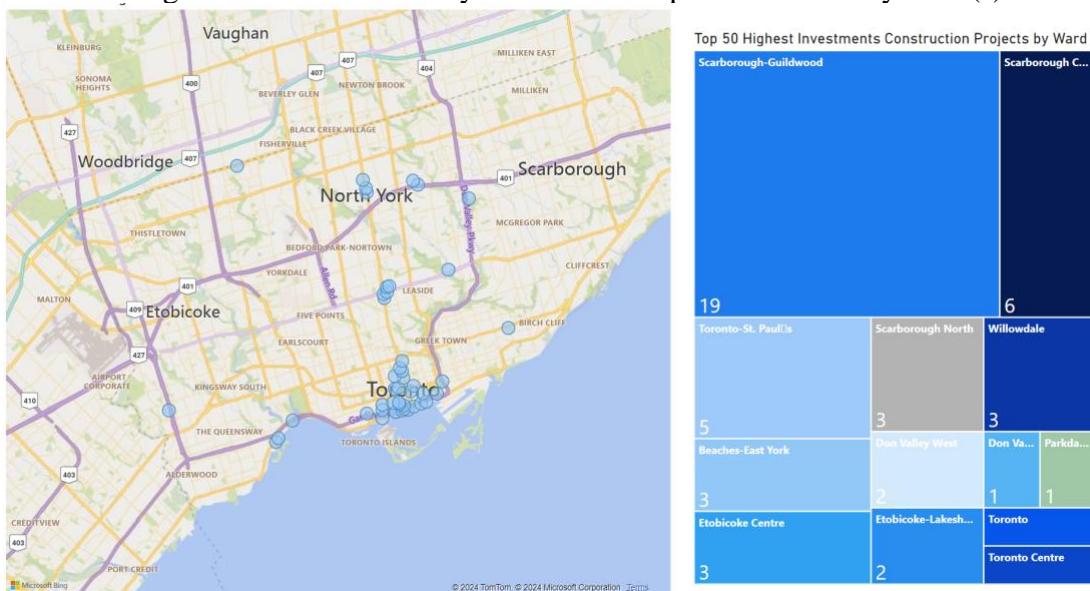
```

A relationship between the Ward Dimension table and the Top 50 Building Permit Coordinates table was created, with the linking column being Ward ID. This was to ensure that visuals were able to be filtered by Ward Name for user understandability and clarity.

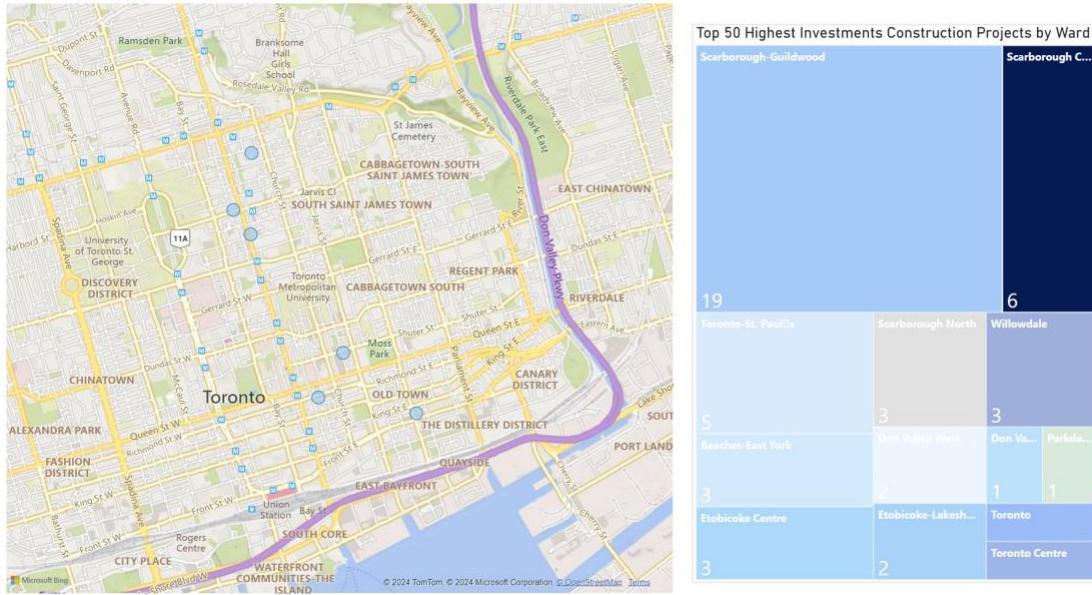
The results of the top 50 query were a much cleaner map visual, allowing the user to get a clearer picture of where the highest investment construction activity is.



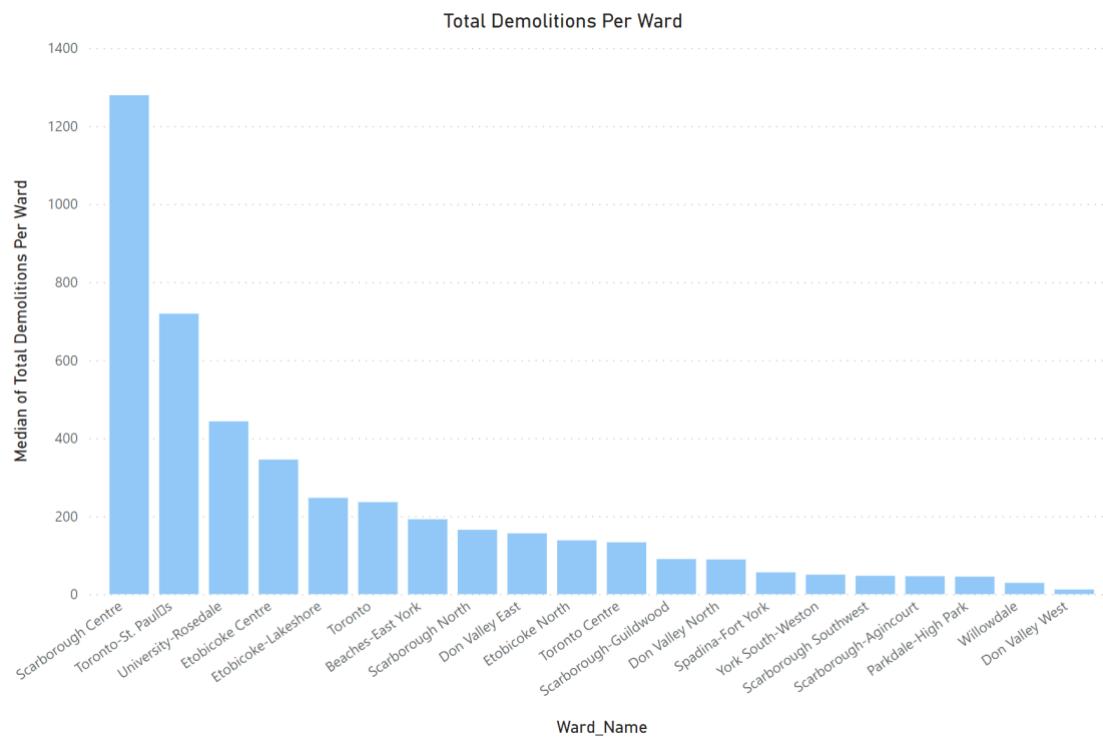
To provide a clearer picture of the breakdown of the top 50 highest construction cost projects, the team included a tree map showing the count of permits organized by each ward on the dashboard. This gives the user the ability to slice the map visualization by ward(s) seamlessly.



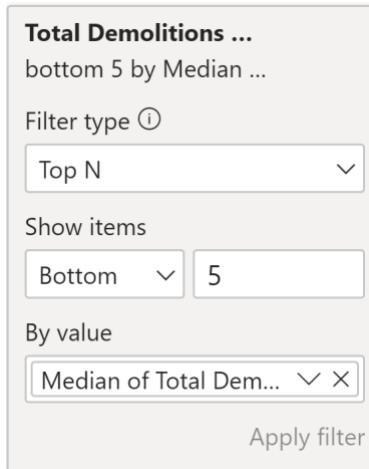
By clicking on one of the wards in the tree map, for example, Scarborough Center, the user can view the high-investment projects mapped to only this ward. In the below screenshot, only the 6 construction projects of Scarborough Center can be seen on the map.



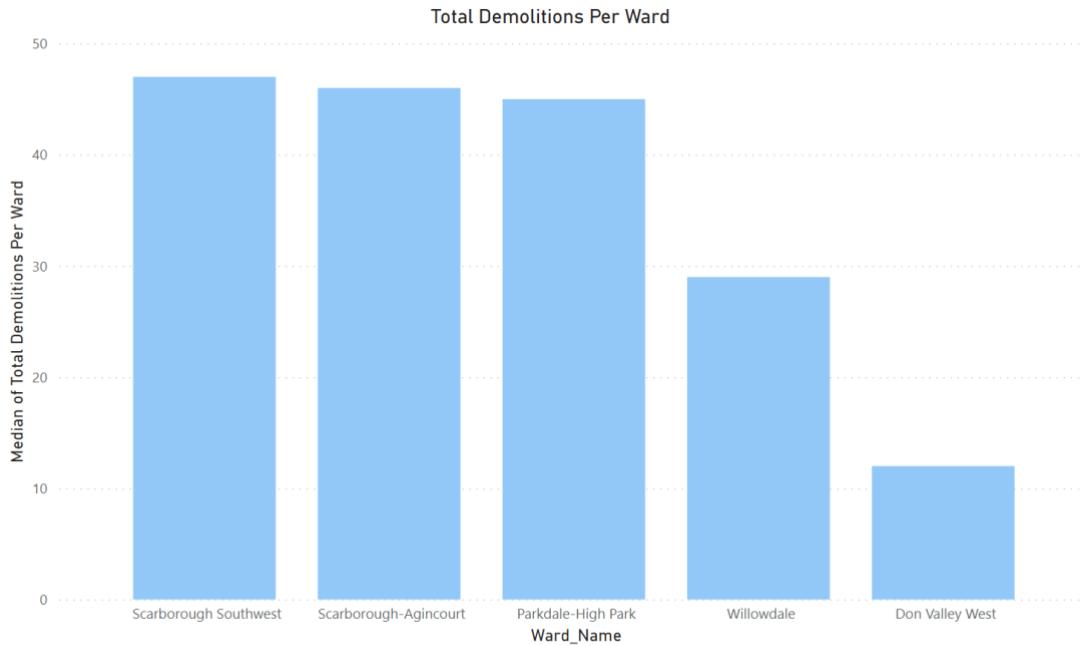
Using the “total demolition” fact in the demolition fact table, a visual displaying the total demolitions per ward was created.



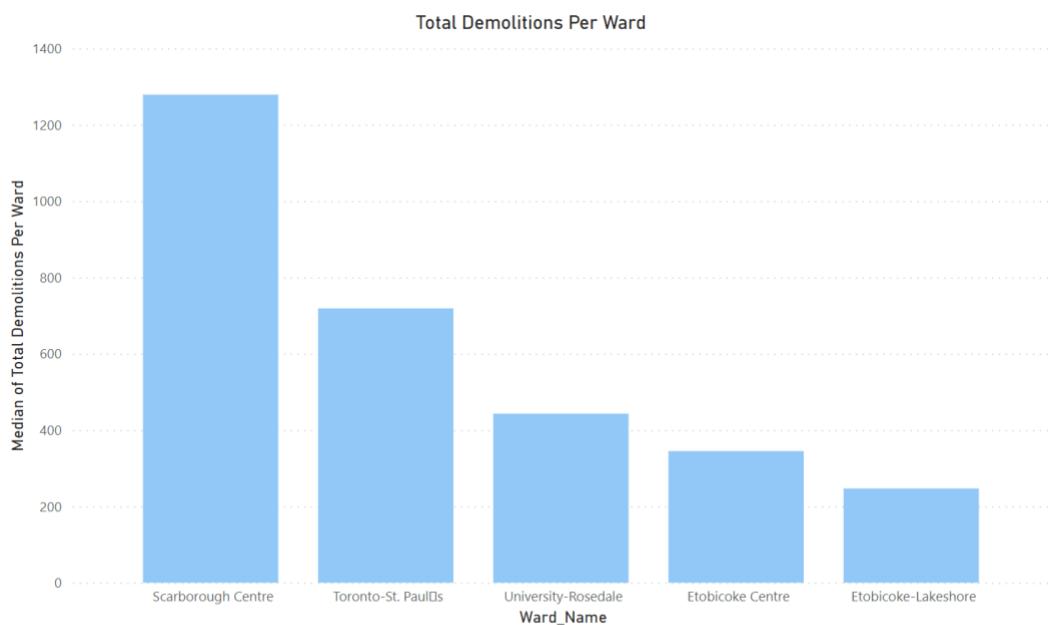
To identify the wards with the lowest demolition activity, we can use a bottom N query to find the bottom 5 wards in terms of demolition rates. This was achieved using Power BI's Top N filter type.



The resulting visualization filtered for the bottom 5 wards is shown below.



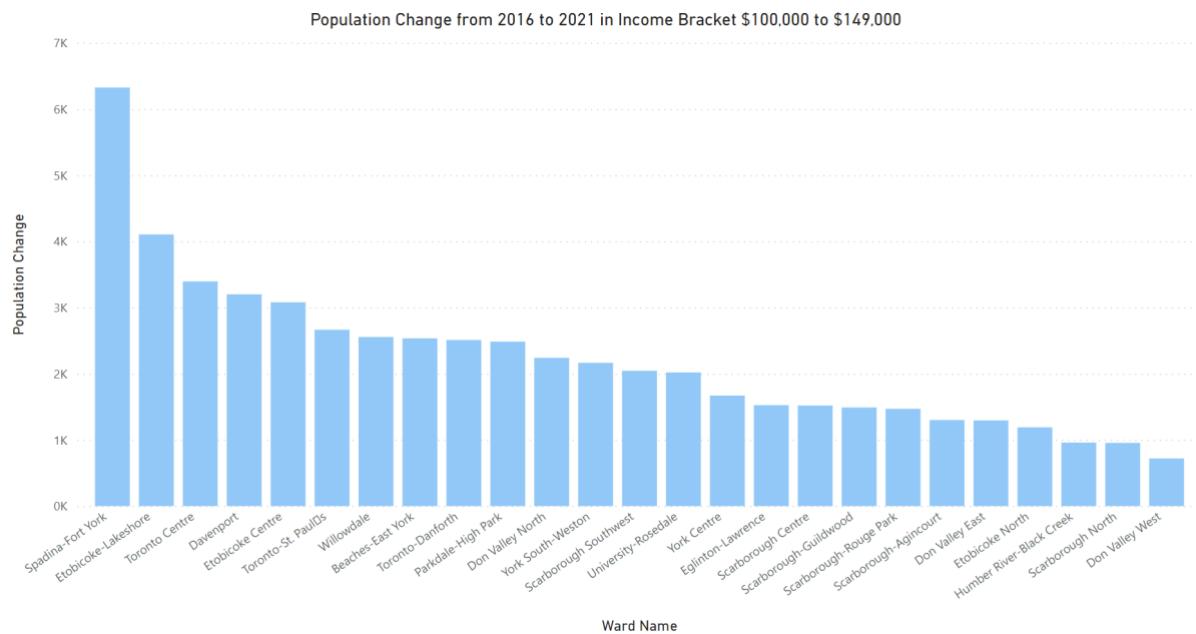
The visual was also filtered to display the Top 5 wards with the highest demolition rates, as seen in the screenshot below. This is particularly relevant in the context of our analysis, as high rates of demolition can be a sign of gentrification making way for new developments.



After observing the results from the top N query for wards with the highest demolition rates, as well as the top N query for wards with the highest construction investment activity, the intersection of wards for both query results is the Scarborough Center, Etobicoke Center, and Etobicoke Lakeshore.

Our team wanted to explore the change in the population in terms of income. Specifically, we wanted to examine the change in population from 2016 to 2021 in terms of the \$100,000 to \$149,000 income bracket to see in which wards the greatest increase in wealth was taking place.

The visualization below shows the population increase for each ward for the \$100,000 to \$149,000 income bracket.



This was achieved with the following DAX expression:

```

1 Population_Change =
2 CALCULATE (
3     SUM (ward_profile_fact_table[Population]),
4     ward_profile_fact_table[Year] = 2021,
5     IncomeDimension[Income] = "$100,000 to $149,999"
6 ) - CALCULATE (
7     SUM (ward_profile_fact_table[Population]),
8     ward_profile_fact_table[Year] = 2016,
9     IncomeDimension[Income] = "$100,000 to $149,999"
10 )
11

```

To filter the above visualization for the top 5 wards, we created a DAX expression to rank the wards by the increase in population:

```

1 Rank_Wards_by_Population_Change =
2 RANKX (
3     ALL (ward_profile_fact_table[Ward_ID]),
4     [Population_Change],
5     ,DESC,
6     Dense
7 )
8

```

Our Top N query is shown below to filter to get the top 5 ranks.

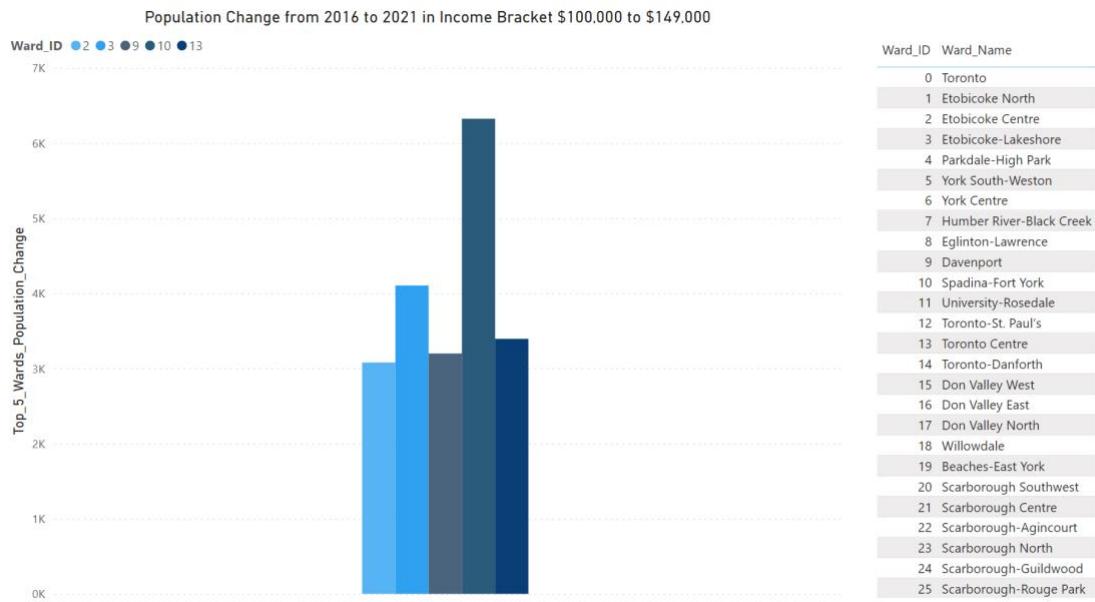
```

1 Top_5_Wards_Population_Change =
2 IF (
3     [Rank_Wards_by_Population_Change] <= 6,
4     [Population_Change],
5     BLANK()
6 )

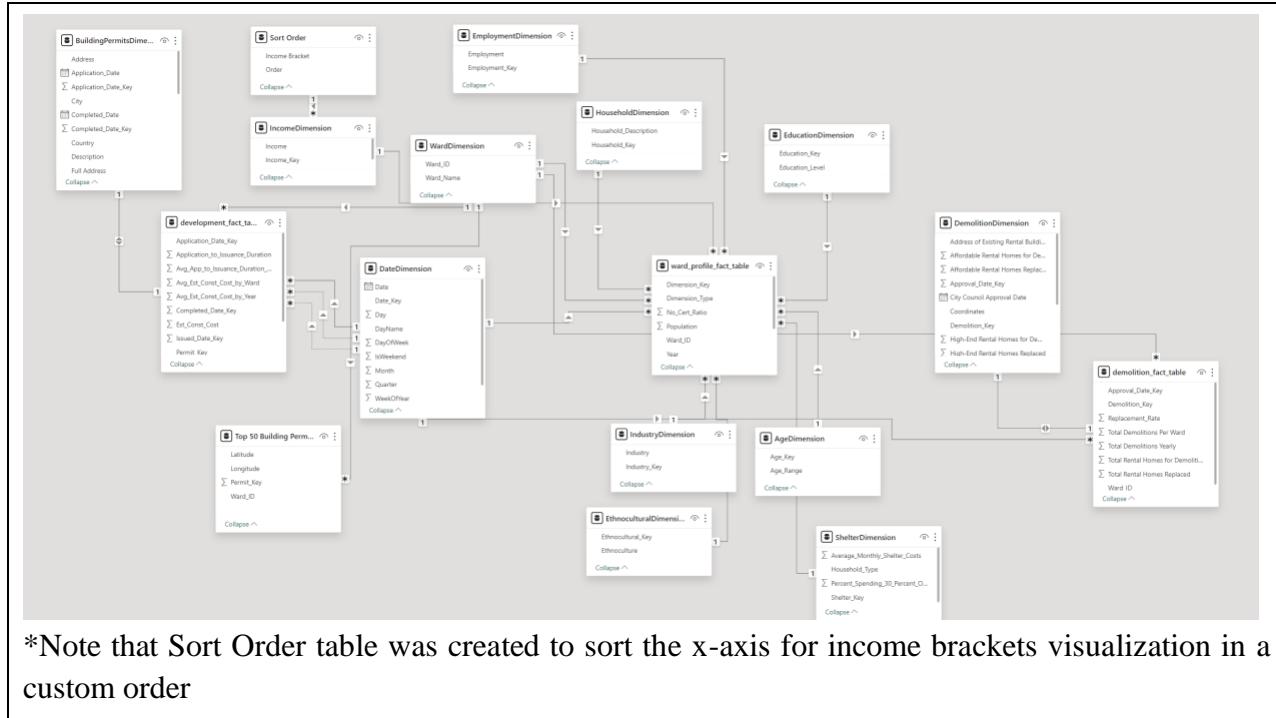
```

\*Note that “<=6” is due to Toronto having a Ward\_ID of 0, which we excluded from the analysis to focus on the wards themselves.

Our visualization below displays the top 5 wards with the highest increase in the population in income bracket \$100,000 to \$149,000.



## Power BI Datamart Design



\*Note that Sort Order table was created to sort the x-axis for income brackets visualization in a custom order