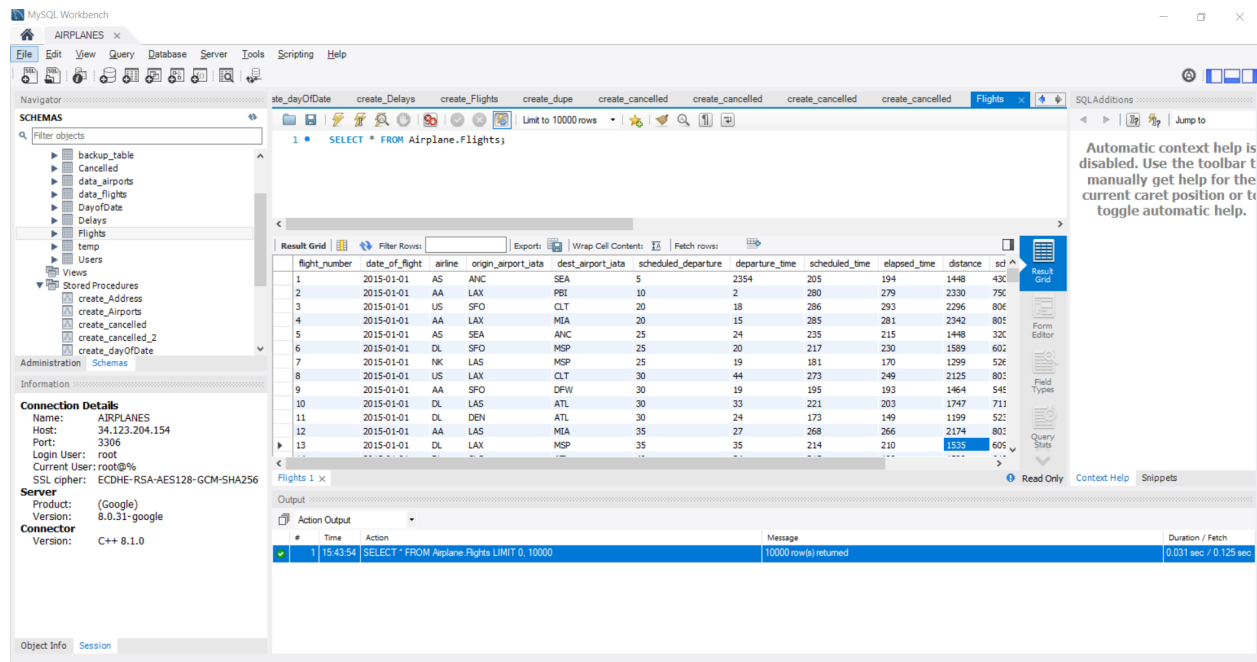


Group 112: Orange
NetIDs: zxchoo2, pnp4, hansenp2

Part 1: Database Implementation

1. Implemented the database tables locally or on GCP



2. Provide the Data Definition Language (DDL) commands you all used to create each of these tables in the database.

```
CREATE TABLE Users (  
    user_id INT PRIMARY KEY,  
    username VARCHAR(255),  
    password VARCHAR(255),  
    email VARCHAR(255)  
);
```

```
CREATE TABLE Searches (  
    search_id INT PRIMARY KEY,  
    date_searched DATETIME,
```

```
search_start_date DATE,  
search_end_date DATE,  
search_airline VARCHAR(255),  
user_id INT,  
FOREIGN KEY (user_id) REFERENCES Users(user_id),  
FOREIGN KEY (search_airline) REFERENCES Airlines(airline_iata_code)  
);
```

```
CREATE TABLE Flights (  
    flight_number INT PRIMARY KEY,  
    date_of_flight DATE,  
    airline VARCHAR(255),  
    origin_airport_iata VARCHAR(255),  
    dest_airport_iata VARCHAR(255),  
    scheduled_departure INT,  
    departure_time INT,  
    scheduled_time INT,  
    elapsed_time REAL,  
    distance REAL,  
    scheduled_arrival INT,  
    arrival_time INT,  
    FOREIGN KEY (origin_airport_iata) REFERENCES Airports(airport_iata_code),  
    FOREIGN KEY (dest_airport_iata) REFERENCES Airports(airport_iata_code),  
    FOREIGN KEY (date_of_flight) REFERENCES DayofDate(date_of_flight)  
);
```

```
CREATE TABLE DayofDate (  
    date_of_flight DATE PRIMARY KEY,  
    day_of_week VARCHAR(255)  
);
```

```
CREATE TABLE Cancelled (  
    flight_number INT PRIMARY KEY,  
    canceled BOOLEAN,  
    cancellation_reason VARCHAR(255),  
    FOREIGN KEY (flight_number) REFERENCES Flights(flight_number)  
);
```

```
CREATE TABLE Delays (  
    flight_number INT PRIMARY KEY,
```

```
departure_delay INT,  
arrival_delay INT,  
FOREIGN KEY (flight_number) REFERENCES Flights(flight_number)  
);
```

```
CREATE TABLE Airports (  
    airport_iata_code VARCHAR(255) PRIMARY KEY,  
    airport VARCHAR(255),  
);
```

```
CREATE TABLE Address (  
    airport_code VARCHAR(255) PRIMARY KEY,  
    city VARCHAR(255),  
    state VARCHAR(255),  
    country VARCHAR(255),  
    latitude REAL,  
    longitude REAL,  
    FOREIGN KEY (airport_code) REFERENCES Airports(airport_iata_code)  
);
```

```
CREATE TABLE Airlines (  
    airline_iata_code VARCHAR(255) PRIMARY KEY,  
    airline_name VARCHAR(255)  
);
```

3. Insert data into these tables. You should insert at least 1000 rows each in three of the tables. Try to use real data, but if you cannot find a good dataset for a particular table, you may use auto-generated data

ate_dayOfDate create_Delays create_Flights create_dupe create_cancelled create_c

Limit to 10000 rows

```
1 • SELECT Count(*) FROM Airplane.Flights;
```

<

Result Grid Filter Rows: Export: Wrap Cell Content:

	Count(*)
▶	809116

Limit to 10000 rows

```
1 • SELECT Count(*) FROM Airplane.Cancelled;
```

Result Grid

Count(*)
809116

Limit to 10000 rows

```
1 • SELECT Count(*) FROM Airplane.Delays;
```

Result Grid

Count(*)
809116

To note for the other tables:

- Airlines matches IATA_Code to airline name, which has count < 1000
- Airports matches airport_IATA_Code to airport, which has count < 1000
- DayofDate matches date_of_flight to day of week, which has count < 1000

Advanced Queries (show top 15 results)

Advanced Query 1

This query gives aggregate information about the airlines between the dates the user selects. These dates are represented below as start_date and end_date, and will be replaced by the frontend of the website when it gets the actual dates from the user. When generating the screenshot of the top 15 rows, we used 2015-01-27 as the start_date and 2015-02-01 as the end_date. The information will be displayed from highest to lowest delay number by default. Right now the information shown to the user is the number of delays, average severity of a delay, and number of cancellations. More information could be added in the future. There are only 14 airlines, so the query returns less than 15 entries.

```
SELECT *
FROM (
    SELECT airline, AVG(arrival_delay) AS delay_severity, COUNT(flight_number) AS
    delay_number
    FROM Delays NATURAL JOIN Flights
    WHERE date_of_flight >= start_date AND date_of_flight <= end_date
    GROUP BY airline
) AS DelayData NATURAL JOIN (
    SELECT airline, SUM(Cancelled) AS cancellation_number
    FROM Cancelled NATURAL JOIN Flights
    WHERE date_of_flight >= start_date AND date_of_flight <= end_date
    GROUP BY airline
) AS CancelData
ORDER BY delay_number DESC
LIMIT 15;
```

	airline	delay_severity	delay_number	cancellation_number
▶	WN	-1.6697	18927	909
	DL	-4.6599	12060	458
	EV	-0.3635	9068	735
	OO	3.6939	8911	286
	AA	-1.0662	8256	597
	UA	-1.1961	7100	726
	US	4.9669	6350	618
	MQ	2.7969	5707	686
	B6	3.3314	4083	710
	AS	1.6416	2536	31
	NK	4.4483	1684	60
	F9	12.6900	1258	40
	HA	2.5529	1239	10
	VX	-3.2769	892	66

Advanced Query 2

This second query shows the aggregate information about airlines on a specific day of the week. This is represented by week_day in the query, which will be replaced by what the user selects on the website. To test the query and generate the top 15 rows for the query, we used the value 'Friday' for week_day. The query currently sorts things by the number of delays by default, but could potentially be changed by the user. Right now we show the number of delays, average time of each delay, and the number of cancellations. More information could be added. There are only 14 airlines, so the query returns less than 15 entries.

```
SELECT *
FROM (
    SELECT airline, AVG(arrival_delay) AS delay_severity, COUNT(flight_number) AS
    delay_number
    FROM Delays NATURAL JOIN Flights NATURAL JOIN DayofDate
    WHERE day_of_week = week_day
    GROUP BY airline
) AS DelayData NATURAL JOIN (
    SELECT airline, SUM(Cancelled) AS cancellation_number
    FROM Cancelled NATURAL JOIN Flights NATURAL JOIN DayofDate
    WHERE day_of_week = week_day
    GROUP BY airline
) AS CancelData
ORDER BY delay_number DESC
LIMIT 15;
```

	airline	delay_severity	delay_number	cancellation_number
▶	WN	2.3473	27201	311
	DL	0.5489	18485	33
	EV	5.7957	14102	170
	OO	8.0546	13253	383
	AA	4.8002	11707	86
	UA	7.6203	11039	68
	US	3.6486	9339	97
	MQ	13.7975	8108	363
	B6	6.5977	5807	22
	AS	2.8898	3603	29
	NK	9.1865	2311	13
	F9	13.8379	1844	3
	HA	9.2829	1778	1
	VX	13.6102	1347	12

Part 2: Indexing Analysis

Default Index Design

First we needed to use EXPLAIN ANALYZE on both of our advanced queries. This is what the results were for advanced query 1:

```
| -> Limit: 15 row(s) (actual time=2924.239..2924.242 rows=14 loops=1)
      -> Sort: Delaydata.delay number DESC, limit input to 15 row(s) per chunk (actual time=2924.239..2924.240 rows=14 loops=1)
            -> Stream results (cost=74721205740226.30 rows=0) (actual time=2924.202..2924.216 rows=14 loops=1)
                  -> Filter: (Cancelled.airline = Delaydata.airline) (cost=74721205740226.30 rows=0) (actual time=2924.198..2924.208 rows=14 loops=1)
                        -> Inner hash join (<hash>(Cancelled.airline)<hash>(Delaydata.airline)) (cost=74721205740226.30 rows=0) (actual time=2924.195..2924.203 rows=14 loops=1)
                              -> Table scan on CancelData (cost=2.50..2.50 rows=0) (actual time=1444.468..1444.470 rows=14 loops=1)
                                    -> Materialize (cost=0.00..0.00 rows=0) (actual time=1444.468..1444.468 rows=14 loops=1)
                                          -> Table scan on <temporary> (actual time=1444.437..1444.440 rows=14 loops=1)
                                                -> Aggregate using temporary table (actual time=1444.432..1444.432 rows=14 loops=1)
                                                      -> Inner hash join (Cancelled.flight number = Flights.flight number) (cost=7228507136.60 rows=7228417342) (actual time=922.821..1355.893 rows=88233 loops=1)
                                                            -> Table scan on Cancelled (cost=0.18 rows=807820) (actual time=0.039..456.215 rows=809117 loops=1)
                                                                  -> Hash
                                                                        -> Filter: ((Flights.date_of_flight >= DATE'2015-01-27') and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=81518.10 rows=89481) (actual time=278.473..590.288 rows=88233 loops=1)
                                                                              -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.035..505.259 rows=809116 loops=1)
                                                                                                            -> Hash
                                                                                              -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=1479.679..1479.681 rows=14 loops=1)
                                                                                                    -> Materialize (cost=0.00..0.00 rows=0) (actual time=1479.678..1479.678 rows=14 loops=1)
                                                                                                          -> Table scan on <temporary> (actual time=1479.637..1479.641 rows=14 loops=1)
                                                                                                                -> Aggregate using temporary table (actual time=1479.634..1479.634 rows=14 loops=1)
                                                                                                                      -> Inner hash join (Delays.flight number = Flights.flight number) (cost=7226270690.24 rows=7226180328) (actual time=943.819..1397.001 rows=88233 loops=1)
                                                                                                                            -> Table scan on Delays (cost=0.19 rows=807570) (actual time=0.052..479.562 rows=809117 loops=1)
                                                                                                                                  -> Hash
                                                                                                                                        -> Filter: ((Flights.date_of_flight >= DATE'2015-01-27') and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=81518.10 rows=89481) (actual time=280.625..597.842 rows=88233 loops=1)
                                                                                                                                              -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.061..515.624 rows=809116 loops=1)
                                                                                                                                                                            |
```

The results come back as a tree with two main branches. This is because the query joins Flights with Delays and Flights with Cancelled separately, collects aggregate data about them, and then joins the results. The highest cost in both amount and placement up the tree is for Stream Results, with a cost 74721205740226.30.

Here are the results for query 2:

```
| -> Limit: 15 row(s) (actual time=4123.080..4123.084 rows=14 loops=1)
      -> Sort: DelayData.delay number DESC, limit input to 15 row(s) per chunk (actual time=4123.079..4123.081 rows=14 loops=1)
            -> Stream results (cost=60527793882894.46 rows=0) (actual time=4123.033..4123.056 rows=14 loops=1)
                  -> Filter: (Cancelled.airline = Delaydata.airline) (cost=60527793882894.46 rows=0) (actual time=4123.029..4123.046 rows=14 loops=1)
                        -> Inner hash join (<hash>(Cancelled.airline)<hash>(Delaydata.airline)) (cost=60527793882894.46 rows=0) (actual time=4123.027..4123.040 rows=14 loops=1)
                              -> Table scan on CancelData (cost=2.50..2.50 rows=0) (actual time=2093.776..2093.779 rows=14 loops=1)
                                    -> Materialize (cost=0.00..0.00 rows=0) (actual time=2093.775..2093.775 rows=14 loops=1)
                                          -> Table scan on <temporary> (actual time=2093.738..2093.741 rows=14 loops=1)
                                                -> Aggregate using temporary table (actual time=2093.734..2093.734 rows=14 loops=1)
                                                      -> Inner hash join (Cancelled.flight number = Flights.flight number) (cost=6505786279.69 rows=6504662113) (actual time=1078.597..1947.833 rows=129924 loops=1)
                                                            -> Table scan on Cancelled (cost=2.05 rows=807545) (actual time=0.043..496.715 rows=809117 loops=1)
                                                                  -> Hash
                                                                        -> Nested loop inner join (cost=967552.70 rows=80549) (actual time=20.343..1009.169 rows=129924 loops=1)
                                                                              -> Filter: (Flights.date_of_flight is not null) (cost=81518.10 rows=805486) (actual time=0.042..590.369 rows=809116 loops=1)
                                                                                    -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.040..526.360 rows=809116 loops=1)
                                                                                          -> Filter: (DayofDate.day_of_week = 'Friday') (cost=1.00 rows=0.1) (actual time=0.000..0.000 rows=0 loops=809116)
                                                                                              -> Single-row index lookup on DayofDate using PRIMARY (date_of_flight=Flights.date_of_flight) (cost=1.00 rows=1) (actual time=0.000..0.000 rows=1 loops=809116)
                                                                                                      -> Hash
                                                                                                            -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=2029.214..2029.216 rows=14 loops=1)
                                                                                                                  -> Materialize (cost=0.00..0.00 rows=0) (actual time=2029.213..2029.213 rows=14 loops=1)
                                                                                                                        -> Table scan on <temporary> (actual time=2029.080..2029.083 rows=14 loops=1)
                                                                                                                              -> Aggregate using temporary table (actual time=2029.076..2029.076 rows=14 loops=1)
                                                                                                                                    -> Inner hash join (Delays.flight number = Flights.flight number) (cost=6505992253.23 rows=6504863484) (actual time=1069.850..1903.607 rows=129924 loops=1)
                                                                                                                                          -> Table scan on Delays (cost=2.18 rows=807570) (actual time=0.031..460.979 rows=809117 loops=1)
                                                                                                                                                  -> Hash
                                                                                                                                                        -> Nested loop inner join (cost=967552.70 rows=80549) (actual time=24.988..1016.131 rows=129924 loops=1)
                                                                                                                                                              -> Filter: (Flights.date_of_flight is not null) (cost=81518.10 rows=805486) (actual time=0.046..585.953 rows=809116 loops=1)
                                                                                                                                                                    -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.044..522.832 rows=809116 loops=1)
                                                                                                                                                                          -> Filter: (DayofDate.day_of_week = 'Friday') (cost=1.00 rows=0.1) (actual time=0.000..0.000 rows=0 loops=809116)
                                                                                                                                                                              -> Single-row index lookup on DayofDate using PRIMARY (date_of_flight=Flights.date_of_flight) (cost=1.00 rows=1) (actual time=0.000..0.000 rows=1 loops=809116)
                                                                                                                                                                                  |
```

The query is still a tree with two branches, but the overall result is bigger, probably due to the extra join with the DayofDate table required to get the Day of the week for each flight. The cost for Stream Result is actually slightly lower, at 60527793882894.46.

Index Design 1

The idea behind our first design was that we would create an index over Flights.date_of_flight called date_idx. The hope was that since we were only looking at flights on certain dates, an index on the Flights.date_of_flight would allow the system to skip just to the tuples that were within the date range. After implementing the index, the EXPLAIN ANALYZE results for the first advanced query were:


```

1 -> Limit: 15 row(s) (actual time=2083.947..2083.950 rows=14 loops=1)
    -> Sort: DelayData.delay number DESC, limit input to 15 row(s) per chunk (actual time=2083.946..2083.948 rows=14 loops=1)
        -> Stream results (cost=293322305385123.10 rows=0) (actual time=2083.910..2083.924 rows=14 loops=1)
            -> Filter: (CancelData.airline = DelayData.airline) (cost=20322205385123.10 rows=0) (actual time=2083.906..2083.917 rows=14 loops=1)
                -> Inner hash join (chash>CancelData.airline)=chash>(DelayData.airline) (cost=293322305385123.10 rows=0) (actual time=2083.904..2083.912 rows=14 loops=1)
                    -> Table scan on CancelData (cost=2.50..2.50 rows=0) (actual time=1013.851..1013.853 rows=14 loops=1)
                        -> Materialize (cost=0.00..0.00 rows=0) (actual time=1013.850..1013.850 rows=14 loops=1)
                            -> Table scan on <temporary> (actual time=1013.820..1013.823 rows=14 loops=1)
                                -> Aggregate using temporary table (actual time=1013.817..1013.817 rows=14 loops=1)
                                    -> Inner hash join (Cancelled.flight number = Flights.flight number) (cost=14321775143.47 rows=14321679429) (actual time=470.305..918.008 rows=88233 loops=1)
                                        -> Table scan on Cancelled (cost=0.14 rows=807520) (actual time=0.036..464.659 rows=809117 loops=1)
                                            -> Hash
                                                -> Index range scan on Flights using date_idx over ('2015-01-27' <= date of flight <= '2015-02-01'), with index condition: ((Flights.date_of_flight >= DATE'2015-01-27') and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=79779.86 rows=177288) (actual time=0.056..130.370 rows=88233 loops=1)
                                                    -> Hash
                                                        -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=1070.019..1070.022 rows=14 loops=1)
                                                            -> Materialize (cost=0.00..0.00 rows=0) (actual time=1070.019..1070.019 rows=14 loops=1)
                                                                -> Table scan on <temporary> (actual time=1069.985..1069.988 rows=14 loops=1)
                                                                    -> Aggregate using temporary table (actual time=1069.982..1069.982 rows=14 loops=1)
                                                                        -> Inner hash join (Delays.flight number = Flights.flight number) (cost=14317344035.90 rows=14317247229) (actual time=545.245..981.628 rows=88233 loops=1)
                                                                            -> Table scan on Delays (cost=0.14 rows=807570) (actual time=0.034..453.471 rows=809117 loops=1)
                                                                                -> Hash
                                                                                    -> Index range scan on Flights using date_idx over ('2015-01-27' <= date of flight <= '2015-02-01'), with index condition: ((Flights.date_of_flight >= DATE'2015-01-27') and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=79779.86 rows=177288) (actual time=0.045..218.018 rows=88233 loops=1)

```

Looking at this result, we can see that date_idx was used at the bottom of the two legs of the tree, replacing two Table Scan on Flights with Index range scans on Flights. The table scans cost 81510.1 each, while the new index scans cost 79779.86 each. This actually leads to a higher cost further up the tree. The highest up cost in the tree is the Stream Results. With the default design, the cost of Stream Results with the original database design was 74721205740226.30. The cost of Stream Results using date_idx is 293322305385123.10.

The second advanced query looks at flights on certain days of the week. Because of that we switched out date_idx for day_idx, an index on DayofDate.day_of_week. After implementing this new index, the EXPLAIN ANALYZE results for the second advanced query were:

```

1 -> Limit: 15 row(s) (actual time=4010.535..4010.538 rows=14 loops=1)
    -> Sort: DelayData.delay number DESC, limit input to 15 row(s) per chunk (actual time=4010.534..4010.536 rows=14 loops=1)
        -> Stream results (cost=122850138424029.50 rows=0) (actual time=4010.502..4010.516 rows=14 loops=1)
            -> Filter: (CancelData.airline = DelayData.airline) (cost=122850138424029.50 rows=0) (actual time=4010.496..4010.506 rows=14 loops=1)
                -> Inner hash join (chash>CancelData.airline)=chash>(DelayData.airline) (cost=122850138424029.50 rows=0) (actual time=4010.494..4010.502 rows=14 loops=1)
                    -> Table scan on CancelData (cost=2.50..2.50 rows=0) (actual time=2013.261..2013.263 rows=14 loops=1)
                        -> Materialize (cost=0.00..0.00 rows=0) (actual time=2013.260..2013.260 rows=14 loops=1)
                            -> Table scan on <temporary> (actual time=2013.231..2013.234 rows=14 loops=1)
                                -> Aggregate using temporary table (actual time=2013.228..2013.228 rows=14 loops=1)
                                    -> Inner hash join (Cancelled.flight number = Flights.flight number) (cost=9267502232.38 rows=9266915873) (actual time=1087.327..1872.296 rows=129924 loops=1)
                                        -> Table scan on Cancelled (cost=2.01 rows=807545) (actual time=0.036..461.732 rows=809117 loops=1)
                                            -> Hash
                                                -> Nested loop inner join (cost=363438.20 rows=114754) (actual time=22.864..1033.793 rows=129924 loops=1)
                                                    -> Filter: (Flights.date of flight is not null) (cost=81518.10 rows=805486) (actual time=0.032..604.798 rows=809116 loops=1)
                                                        -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.031..540.067 rows=809116 loops=1)
                                                            -> Filter: (DayofDate.day of week = 'Friday') (cost=0.25 rows=0.1) (actual time=0.000..0.000 rows=0 loops=809116)
                                                                -> Single-row index lookup on DayofDate using PRIMARY (date of flight=Flights.date of flight) (cost=0.25 rows=1) (actual time=0.000..0.000 rows=1 loops=809116)
                                                                    -> Hash
                                                                        -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=1997.197..1997.199 rows=14 loops=1)
                                                                            -> Materialize (cost=0.00..0.00 rows=0) (actual time=1997.196..1997.196 rows=14 loops=1)
                                                                                -> Table scan on <temporary> (actual time=1997.155..1997.159 rows=14 loops=1)
                                                                                    -> Aggregate using temporary table (actual time=1997.152..1997.153 rows=14 loops=1)
                                                                                        -> Inner hash join (Delays.flight number = Flights.flight number) (cost=9267804208.49 rows=9267202759) (actual time=1074.324..1871.686 rows=129924 loops=1)
                                                                                            -> Table scan on Delays (cost=2.15 rows=807570) (actual time=0.035..460.656 rows=809117 loops=1)
                                                                                                -> Hash
                                                                                                    -> Nested loop inner join (cost=363438.20 rows=114754) (actual time=21.526..1021.095 rows=129924 loops=1)
                                                                                                        -> Filter: (Flights.date of flight is not null) (cost=81518.10 rows=805486) (actual time=0.030..597.330 rows=809116 loops=1)
                                                                                                            -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.028..534.982 rows=809116 loops=1)
                                                                                                                -> Filter: (DayofDate.day of week = 'Friday') (cost=0.25 rows=0.1) (actual time=0.000..0.000 rows=0 loops=809116)
                                                                                                                    -> Single-row index lookup on DayofDate using PRIMARY (date of flight=Flights.date of flight) (cost=0.25 rows=1) (actual time=0.000..0.000 rows=1 loops=809116)

```

Interestingly the results never reference day_idx, which would make one think that day_idx was never used. However, it had a noticeable effect on performance, with the cost of Stream Results now at 122850138424029.50, which was worse than the original cost. We are unsure why this design makes performance worse.

Index Design 2

The idea behind our second design was that we would create an index over Flights.airline called airline_idx. The hope was that it would improve the efficiency of grouping and aggregating the information. Here is the EXPLAIN ANALYZE results of the first query with airline_idx:

```

1 -> Limit: 15 row(s) (actual time=3634.856..3634.859 rows=14 loops=1)
   -> Sort: DelayData.delay number DESC, limit input to 15 row(s) per chunk (actual time=3634.855..3634.856 rows=14 loops=1)
       -> Stream results (cost=74721205740226.30 rows=0) (actual time=3634.819..3634.833 rows=14 loops=1)
           -> Filter: (CancelData.airline = DelayData.airline) (cost=74721205740226.30 rows=0) (actual time=3634.815..3634.826 rows=14 loops=1)
           -> Inner hash join (chash>(CancelData.airline)=chash>(DelayData.airline)) (cost=74721205740226.30 rows=0) (actual time=3634.814..3634.822 rows=14 loops=1)
               -> Table scan on CancelData (cost=2.50..2.50 rows=0) (actual time=1896.086..1896.088 rows=14 loops=1)
               -> Materialize (cost=0.00..0.00 rows=0) (actual time=1896.085..1896.085 rows=14 loops=1)
               -> Table scan on <temporary> (actual time=1896.097..1896.099 rows=14 loops=1)
                   -> Aggregate using temporary table (actual time=1896.053..1896.053 rows=14 loops=1)
                       -> Inner hash join (Cancelled.flight_number = Flights.flight_number) (cost=7228507136.60 rows=7228417342) (actual time=1392.962..1812.282 rows=88233 loops=1)
                           -> Table scan on Cancelled (cost=0.18 rows=807820) (actual time=0.045..472.377 rows=809117 loops=1)
                           -> Hash
                               -> Filter: ((Flights.date_of_flight >= DATE'2015-01-27') and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=81518.10 rows=89481) (actual time=485.309..1036.159 rows=88233 loops=1)
                                   -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.041..867.776 rows=809116 loops=1)
                                       -> Hash
                                           -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=1738.666..1738.669 rows=14 loops=1)
                                           -> Materialize (cost=0.00..0.00 rows=0) (actual time=1738.665..1738.665 rows=14 loops=1)
                                           -> Table scan on <temporary> (actual time=1738.614..1738.618 rows=14 loops=1)
                                               -> Aggregate using temporary table (actual time=1738.609..1738.609 rows=14 loops=1)
                                                   -> Inner hash join (Delays.flight_number = Flights.flight_number) (cost=722670690.24 rows=7226180328) (actual time=927.666..1603.455 rows=88233 loops=1)
                                                       -> Table scan on Delays (cost=0.19 rows=807570) (actual time=0.035..579.434 rows=809117 loops=1)
                                                       -> Hash
                                                           -> Filter: ((Flights.date_of_flight >= DATE'2015-01-27') and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=81518.10 rows=89481) (actual time=291.524..605.525 rows=88233 loops=1)
                                                               -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.038..517.387 rows=809116 loops=1)
1

```

Unfortunately, it seems like the database did not use airline_idx for anything, and the costs are the exact same as with the default design. This also held true for the other advanced query:

```

1 -> Limit: 15 row(s) (actual time=4007.335..4007.338 rows=14 loops=1)
   -> Sort: DelayData.delay number DESC, limit input to 15 row(s) per chunk (actual time=4007.334..4007.336 rows=14 loops=1)
       -> Stream results (cost=60527793882894.46 rows=0) (actual time=4007.298..4007.312 rows=14 loops=1)
           -> Filter: (CancelData.airline = DelayData.airline) (cost=60527793882894.46 rows=0) (actual time=4007.293..4007.303 rows=14 loops=1)
           -> Inner hash join (chash>(CancelData.airline)=chash>(DelayData.airline)) (cost=60527793882894.46 rows=0) (actual time=4007.291..4007.299 rows=14 loops=1)
               -> Table scan on CancelData (cost=2.50..2.50 rows=0) (actual time=2006.903..2006.905 rows=14 loops=1)
               -> Materialize (cost=0.00..0.00 rows=0) (actual time=2006.902..2006.902 rows=14 loops=1)
               -> Table scan on <temporary> (actual time=2006.872..2006.875 rows=14 loops=1)
                   -> Aggregate using temporary table (actual time=2006.863..2006.869 rows=14 loops=1)
                       -> Inner hash join (Cancelled.flight_number = Flights.flight_number) (cost=6505192165.19 rows=6504662113) (actual time=1064.246..1879.228 rows=129924 loops=1)
                           -> Table scan on Cancelled (cost=2.05 rows=807545) (actual time=0.031..455.344 rows=809117 loops=1)
                           -> Hash
                               -> Nested loop inner join (cost=363438.20 rows=80549) (actual time=21.965..1019.771 rows=129924 loops=1)
                                   -> Filter: (Flights.date_of_flight is not null) (cost=81518.10 rows=805486) (actual time=0.049..591.495 rows=809116 loops=1)
                                   -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.047..529.358 rows=809116 loops=1)
                                   -> Filter: (DayofDate.day_of_week = 'Friday') (cost=0.25 rows=0.1) (actual time=0.000..0.000 rows=0 loops=809116)
                                       -> Single-row index lookup on DayofDate using PRIMARY (date_of_flight=Flights.date_of_flight) (cost=0.25 rows=1) (actual time=0.000..0.000 rows=1 loops=809116)
                                           -> Hash
                                               -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=2000.346..2000.348 rows=14 loops=1)
                                               -> Materialize (cost=0.00..0.00 rows=0) (actual time=2000.344..2000.344 rows=14 loops=1)
                                               -> Table scan on <temporary> (actual time=2000.294..2000.297 rows=14 loops=1)
                                                   -> Aggregate using temporary table (actual time=2000.290..2000.290 rows=14 loops=1)
                                                       -> Inner hash join (Delays.flight_number = Flights.flight_number) (cost=6504863484.73 rows=6504863484) (actual time=1061.050..1877.619 rows=129924 loops=1)
                                                           -> Table scan on Delays (cost=2.18 rows=807570) (actual time=0.023..443.312 rows=809117 loops=1)
                                                           -> Hash
                                                               -> Nested loop inner join (cost=363438.20 rows=80549) (actual time=21.456..1037.245 rows=129924 loops=1)
                                                                   -> Filter: (Flights.date_of_flight is not null) (cost=81518.10 rows=805486) (actual time=0.043..603.875 rows=809116 loops=1)
                                                                   -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.041..539.202 rows=809116 loops=1)
                                                                   -> Filter: (DayofDate.day_of_week = 'Friday') (cost=0.25 rows=0.1) (actual time=0.000..0.000 rows=0 loops=809116)
                                                                       -> Single-row index lookup on DayofDate using PRIMARY (date_of_flight=Flights.date_of_flight) (cost=0.25 rows=1) (actual time=0.000..0.000 rows=1 loops=809116)
1

```

Index Design 3

For our final design, we decided to create indexes on the attributes that we were aggregating. These attributes are Delays.flight_number, Delays.arrival_delay, and Cancelled.cancelled, with indexes num_idx, delay_idx, and cancelled_idx respectively. The hope was that indexing these functions would help speed up the aggregation process in some way. Here is the EXPLAIN ANALYZE results of using these indexes on advanced query 1:

```

1 -> Limit: 15 row(s) (actual time=2704.948..2704.951 rows=14 loops=1)
   -> Sort: DelayData.delay number DESC, limit input to 15 row(s) per chunk (actual time=2704.947..2704.949 rows=14 loops=1)
       -> Stream results (cost=1015608329.17 rows=0) (actual time=2704.913..2704.927 rows=14 loops=1)
           -> Filter: (CancelData.airline = DelayData.airline) (cost=1015608329.17 rows=0) (actual time=2704.910..2704.920 rows=14 loops=1)
           -> Inner hash join (chash>(CancelData.airline)=chash>(DelayData.airline)) (cost=1015608329.17 rows=0) (actual time=2704.908..2704.916 rows=14 loops=1)
               -> Table scan on CancelData (cost=2.50..2.50 rows=0) (actual time=1421.184..1421.186 rows=14 loops=1)
               -> Materialize (cost=0.00..0.00 rows=0) (actual time=1421.183..1421.183 rows=14 loops=1)
               -> Table scan on <temporary> (actual time=1421.152..1421.156 rows=14 loops=1)
                   -> Aggregate using temporary table (actual time=1421.149..1421.149 rows=14 loops=1)
                       -> Inner hash join (Cancelled.flight_number = Flights.flight_number) (cost=7228507136.60 rows=7228417342) (actual time=917.174..1340.709 rows=88233 loops=1)
                           -> Table scan on Cancelled (cost=0.18 rows=807820) (actual time=0.027..453.216 rows=809117 loops=1)
                           -> Hash
                               -> Filter: ((Flights.date_of_flight >= DATE'2015-01-27') and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=81518.10 rows=89481) (actual time=289.294..599.521 rows=88233 loops=1)
                                   -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.016..514.666 rows=809116 loops=1)
                                       -> Hash
                                           -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=1283.681..1283.683 rows=14 loops=1)
                                           -> Materialize (cost=0.00..0.00 rows=0) (actual time=1283.680..1283.680 rows=14 loops=1)
                                           -> Table scan on <temporary> (actual time=1283.631..1283.636 rows=14 loops=1)
                                               -> Aggregate using temporary table (actual time=1283.627..1283.627 rows=14 loops=1)
                                                   -> Nested loop inner join (cost=112836.29 rows=89481) (actual time=297.608..1171.684 rows=88233 loops=1)
                                                       -> Filter: ((Flights.date_of_flight >= DATE'2015-01-27') and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=81518.10 rows=89481) (actual time=296.599..618.202 rows=88233 loops=1)
                                                           -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.029..526.048 rows=809116 loops=1)
                                                           -> Index lookup on Delays using num_idx (flight_number=Flights.flight_number) (cost=0.25 rows=1) (actual time=0.004..0.006 rows=1 loops=88233)
1

```

num_index was the only index used by the query. The index was used in a new step added at the end of one of the tree branches, titled Index Lookup on Delays. The cost of this new step was 0.25. Rather than help with the aggregation process as expected, we believe the index helped with the process of joining Flights with Delays, as flight_number was what they were joining on. This apparently resulted in a net time save, as the highest cost up the tree, Stream Results, dropped from 74695768970595.86 to 1015608329.17.

The results for the second advanced query were:

```
| -> Limit: 15 row(s) (actual time=3853.857..3853.860 rows=14 loops=1)
|   -> Sort: DelayData.delay_number DESC, limit input to 15 row(s) per chunk (actual time=3853.856..3853.858 rows=14 loops=1)
|     -> Stream results (cost=830806944.22 rows=0) (actual time=3853.820..3853.834 rows=14 loops=1)
|       -> Filter: (CancelData.airline = DelayData.airline) (cost=830806944.22 rows=0) (actual time=3853.816..3853.826 rows=14 loops=1)
|         -> Inner hash join (hash=>CancelData.airline)=hash=>(DelayData.airline) (cost=830806944.22 rows=0) (actual time=3853.814..3853.822 rows=14 loops=1)
|           -> Table scan on CancelData (cost=2.50..2.50 rows=0) (actual time=2005.433..2005.435 rows=14 loops=1)
|             -> Materialize (cost=0.00..0.00 rows=0) (actual time=2005.432..2005.432 rows=14 loops=1)
|           -> Table scan on <temporary> (actual time=2005.405..2005.408 rows=14 loops=1)
|             -> Aggregate using temporary table (actual time=2005.402..2005.402 rows=14 loops=1)
|               -> Inner hash join (Cancelled.flight_number = Flights.flight_number) (cost=6505182165.19 rows=6504662113) (actual time=1075.935..1881.263 rows=129924 loops=1)
|                 -> Table scan on Cancelled (cost=2.05 rows=807545) (actual time=0.036..453.685 rows=809117 loops=1)
|                   -> Hash
|                     -> Nested loop inner join (cost=363438.20 rows=80549) (actual time=21.971..1031.540 rows=129924 loops=1)
|                       -> Filter: (Flights.date_of_flight is not null) (cost=81518.10 rows=805486) (actual time=0.034..596.838 rows=809116 loops=1)
|                         -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.019..532.792 rows=809116 loops=1)
|                           -> Filter: (DayofDate.day_of_week = 'Friday') (cost=0.25 rows=0.1) (actual time=0.000..0.000 rows=0 loops=809116)
|                             -> Single-row index lookup on DayofDate using PRIMARY (date_of_flight=Flights.date_of_flight) (cost=0.25 rows=1) (actual time=0.000..0.000 rows=1 loops=809116)
|                   -> Hash
|                     -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=1848.339..1848.341 rows=14 loops=1)
|                       -> Materialize (cost=0.00..0.00 rows=0) (actual time=1848.338..1848.338 rows=14 loops=1)
|                     -> Table scan on <temporary> (actual time=1848.274..1848.279 rows=14 loops=1)
|                       -> Aggregate using temporary table (actual time=1848.269..1848.269 rows=14 loops=1)
|                         -> Nested loop inner join (cost=391630.21 rows=80549) (actual time=24.723..1678.674 rows=129924 loops=1)
|                           -> Nested loop inner join (cost=363438.20 rows=80549) (actual time=23.292..1094.812 rows=129924 loops=1)
|                             -> Filter: ((Flights.date_of_flight is not null) and (Flights.flight_number is not null)) (cost=81518.10 rows=805486) (actual time=0.040..649.438 rows=809116 loops=1)
|                               -> Table scan on Flights (cost=81518.10 rows=805486) (actual time=0.035..557.922 rows=809116 loops=1)
|                                 -> Filter: (DayofDate.day_of_week = 'Friday') (cost=0.25 rows=0.1) (actual time=0.000..0.000 rows=0 loops=809116)
|                                   -> Single-row index lookup on DayofDate using PRIMARY (date_of_flight=Flights.date_of_flight) (cost=0.25 rows=1) (actual time=0.000..0.000 rows=1 loops=809116)
|                                     -> Index lookup on Delays using num_idx (flight_number=Flights.flight_number) (cost=0.25 rows=1) (actual time=0.003..0.004 rows=1 loops=129924)
```

A very similar result to what happened with the first query also happened here. A new step was added to the end of one of the tree branches, with the name Index Lookup on Delays and with a cost of 0.25. The Stream Results cost dropped to 830806944.22.

Final Index Analysis and Design:

After looking at all 3 indexes, we came to some conclusions. First off, the second design did not do anything, so we did not consider it. It also looked like the indexes were never used to speed up aggregation, so we decided not to try to do that in the final design. We decided to add indexes for all of the points that tables were joined with, to hopefully capitalize on the speed up from design 3. In the end, this created the following indexes: `f_date_idx` on `Flights(date_of_week)`, `d_date_idx` on `DayofDate(date_of_week)`, `f_num_idx` on `Flights(flight_number)`, `d_num_idx` on `Delays(flight_number)`, and `c_num_idx` on `Cancelled(flights_number)`. This is the result for our first advanced query:

```
| -> Limit: 15 row(s) (actual time=2411.184..2411.188 rows=14 loops=1)
|   -> Sort: DelayData.delay number DESC, limit input to 15 row(s) per chunk (actual time=2411.183..2411.186 rows=14 loops=1)
|     -> Stream results (cost=47234.28 rows=0) (actual time=2411.134..2411.159 rows=14 loops=1)
|       -> Filter: (CancelData.airline = DelayData.airline) (cost=47234.28 rows=0) (actual time=2411.130..2411.149 rows=14 loops=1)
|         -> Inner hash join (chash)(CancelData.airline)=(chash)(DelayData.airline) (cost=47234.28 rows=0) (actual time=2411.129..2411.143 rows=14 loops=1)
|           -> Table scan on CancelData (cost=2.50..2.50 rows=0) (actual time=1202.961..1202.964 rows=14 loops=1)
|             -> Materialize (cost=0.00..0.00 rows=0) (actual time=1202.960..1202.960 rows=14 loops=1)
|           -> Aggregate using temporary table (actual time=1202.955..1202.925 rows=14 loops=1)
|             -> Nested loop inner join (cost=141830.66 rows=177288) (actual time=1.080..982.049 rows=88233 loops=1)
|               -> Filter: (Flights.flight_number is not null) (cost=79779.86 rows=177288) (actual time=0.028..327.006 rows=88233 loops=1)
|                 -> Index range scan on Flights using f_date_idx over ('2015-01-27' <= date of flight <= '2015-02-01'), with index condition: ((Flights.date_of_flight >= DATE'2015-01-27')
|                   and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=79779.86 rows=177288) (actual time=0.027..311.914 rows=88233 loops=1)
|                 -> Index lookup on Cancelled using c_num_idx (flight_number=Flights.flight_number) (cost=0.25 rows=1) (actual time=0.005..0.007 rows=1 loops=88233)
|               -> Hash
|                 -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=1208.089..1208.092 rows=14 loops=1)
|                   -> Materialize (cost=0.00..0.00 rows=0) (actual time=1208.088..1208.088 rows=14 loops=1)
|                 -> Table scan on <temporary> (actual time=1208.040..1208.044 rows=14 loops=1)
|                   -> Aggregate using temporary table (actual time=1208.037..1208.037 rows=14 loops=1)
|                 -> Nested loop inner join (cost=141830.66 rows=177288) (actual time=0.939..983.952 rows=88233 loops=1)
|                   -> Filter: (Flights.flight_number is not null) (cost=79779.86 rows=177288) (actual time=0.048..372.096 rows=88233 loops=1)
|                     -> Index range scan on Flights using f_date_idx over ('2015-01-27' <= date of flight <= '2015-02-01'), with index condition: ((Flights.date_of_flight >= DATE'2015-01-
|                       27') and (Flights.date_of_flight <= DATE'2015-02-01')) (cost=79779.86 rows=177288) (actual time=0.046..356.788 rows=88233 loops=1)
|                     -> Index lookup on Delays using d_num_idx (flight_number=Flights.flight_number) (cost=0.25 rows=1) (actual time=0.005..0.007 rows=1 loops=88233)
```

It creates the extra steps Index Lookup on Delays using `d_num_idx` with a cost of 0.25 and Index Lookup on Cancelled using `c_num_idx` with a cost of 0.25. It also replaces the table scans on Flights in each branch of the tree with Index Range Scans on Flights using `f_date_idx` with a cost of 79779.86 each. The cost of Stream Results went down to 47234.28.

The EXPLAIN ANALYZE results of the second advanced query were:

```
| -> Limit: 15 row(s) (actual time=0.085..0.085 rows=0 loops=1)
|   -> Sort: DelayData.delay number DESC, limit input to 15 row(s) per chunk (actual time=0.085..0.085 rows=0 loops=1)
|     -> Stream results (cost=5.03 rows=0) (actual time=0.075..0.075 rows=0 loops=1)
|       -> Filter: (CancelData.airline = DelayData.airline) (cost=5.03 rows=0) (actual time=0.073..0.073 rows=0 loops=1)
|         -> Inner hash join (chash)(CancelData.airline)=(chash)(DelayData.airline) (cost=5.03 rows=0) (actual time=0.072..0.072 rows=0 loops=1)
|           -> Table scan on CancelData (cost=2.50..2.50 rows=0) (never executed)
|             -> Materialize (cost=0.00..0.00 rows=0) (never executed)
|           -> Table scan on <temporary> (never executed)
|             -> Aggregate using temporary table (never executed)
|               -> Nested loop inner join (cost=1.07 rows=1) (never executed)
|                 -> Filter: (Flights.flight_number is not null) (cost=0.71 rows=1) (never executed)
|                   -> Index range scan on Flights using f_date_idx over ('2015-11-01' <= date of flight <= '2015-11-01'), with index condition: ((Flights.date_of_flight >= DATE'2015-10-27')
|                     and (Flights.date_of_flight <= DATE'2015-11-01')) (cost=0.71 rows=1) (never executed)
|                   -> Index lookup on Cancelled using c_num_idx (flight_number=Flights.flight_number) (cost=0.36 rows=1) (never executed)
|                 -> Hash
|                   -> Table scan on DelayData (cost=2.50..2.50 rows=0) (actual time=0.062..0.062 rows=0 loops=1)
|                     -> Materialize (cost=0.00..0.00 rows=0) (actual time=0.061..0.061 rows=0 loops=1)
|                   -> Table scan on <temporary> (actual time=0.046..0.046 rows=0 loops=1)
|                     -> Aggregate using temporary table (actual time=0.044..0.044 rows=0 loops=1)
|                   -> Nested loop inner join (cost=1.06 rows=1) (actual time=0.033..0.033 rows=0 loops=1)
|                     -> Filter: (Flights.flight_number is not null) (cost=0.71 rows=1) (actual time=0.032..0.032 rows=0 loops=1)
|                       -> Index range scan on Flights using f_date_idx over ('2015-10-27' <= date of flight <= '2015-11-01'), with index condition: ((Flights.date_of_flight >= DATE'2015-10-
|                         27') and (Flights.date_of_flight <= DATE'2015-11-01')) (cost=0.71 rows=1) (actual time=0.030..0.030 rows=0 loops=1)
|                       -> Index lookup on Delays using d_num_idx (flight_number=Flights.flight_number) (cost=0.35 rows=1) (never executed)
```

We once again created the extra Index Lookup steps, but this time the cost for the Index Lookup on Cancelled was only 0.26. The cost of Stream Results was 10685.02. Considering that both queries had smaller costs under this final design than any of the previous designs, this is the design that we will be using in our final database.