# Storm Impacts on Aviation

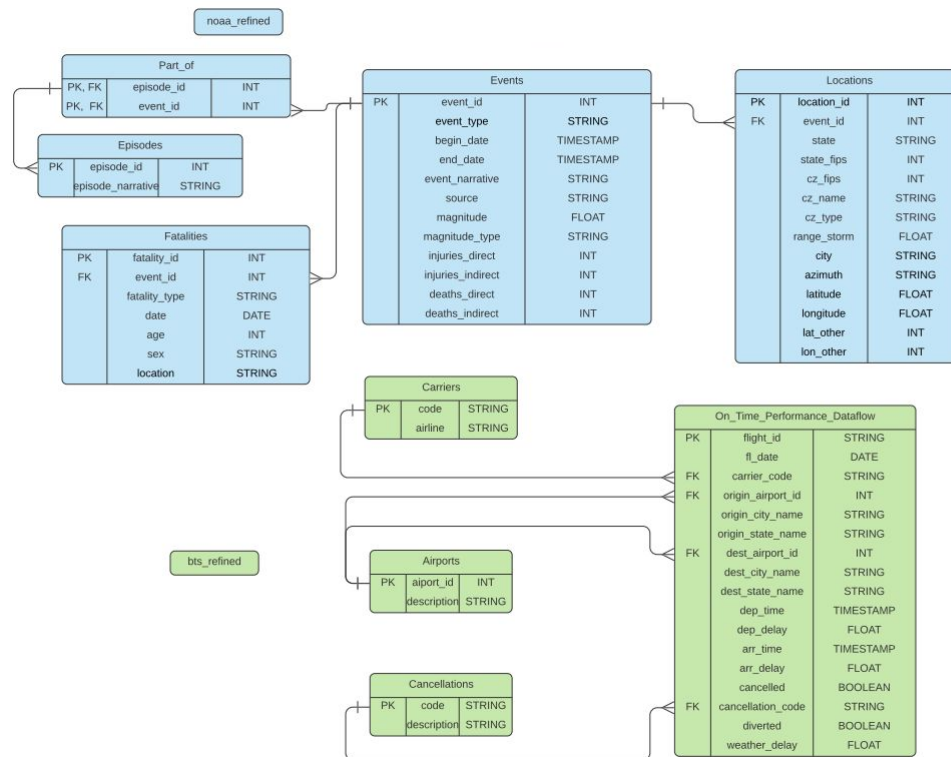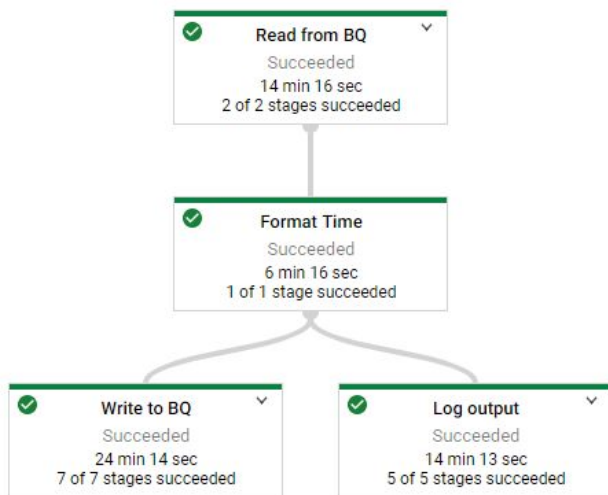Kathy Wang, Sierra Obermoeller-Gilmer

# Areas of Interest

- What storm events have the greatest effects on aviation?
- In what ways are different airlines affected by storm events?
- Which airports are most affected by storm related delays?

# The Datasets

- NOAA Storm Events Dataset
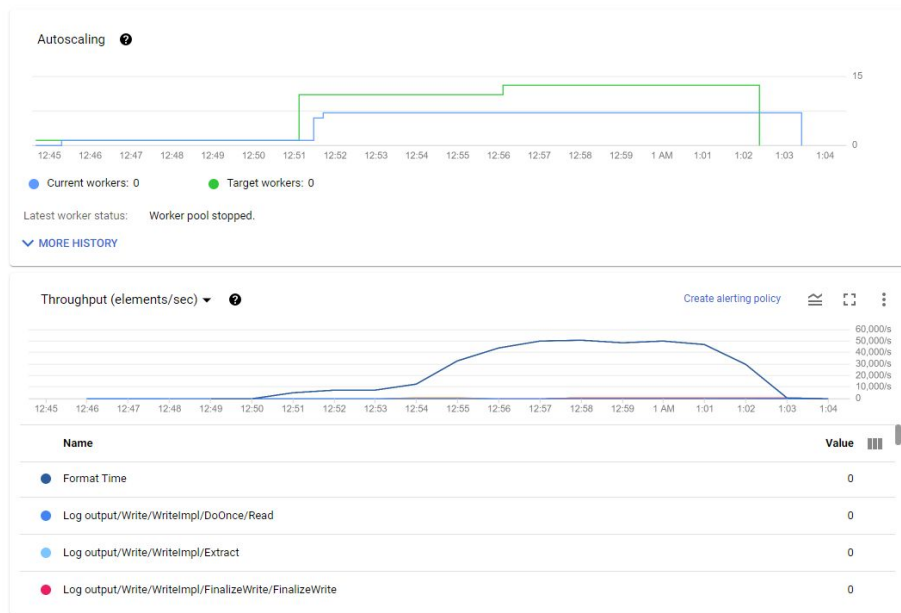- Bureau of Transportation Statistics Airline On-Time Performance Dataset

**noaa_refined**

**Part_of**

| PK, FK | episode_id | INT |
| PK, FK | event_id | INT |

**Episodes**

| PK | episode_id | INT |
| | episode_narrative | STRING |

**Events**

| PK | event_id | INT |
| | event_type | STRING |
| | begin_date | TIMESTAMP |
| | end_date | TIMESTAMP |
| | event_narrative | STRING |
| | source | STRING |
| | magnitude | FLOAT |
| | magnitude_type | STRING |
| | injuries_direct | INT |
| | injuries_indirect | INT |
| | deaths_direct | INT |
| | deaths_indirect | INT |

**Locations**

| PK | location_id | INT |
| FK | event_id | INT |
| | state | STRING |
| | state_fips | INT |
| | cz_fips | INT |
| | cz_name | STRING |
| | cz_type | STRING |
| | range_storm | FLOAT |
| | city | STRING |
| | azimuth | STRING |
| | latitude | FLOAT |
| | longitude | FLOAT |
| | lat_other | INT |
| | lon_other | INT |

**Fatalities**

| PK | fatality_id | INT |
| FK | event_id | INT |
| | fatality_type | STRING |
| | date | DATE |
| | age | INT |
| | sex | STRING |
| | location | STRING |

**Carriers**

| PK | code | STRING |
| | airline | STRING |

**bts_refined**

**Airports**

| PK | aiport_id | INT |
| | description | STRING |

**Cancellations**

| PK | code | STRING |
| | description | STRING |

**On_Time_Performance_Dataflow**

| PK | flight_id | STRING |
| | fl_date | DATE |
| FK | carrier_code | STRING |
| FK | origin_airport_id | INT |
| | origin_city_name | STRING |
| | origin_state_name | STRING |
| FK | dest_airport_id | INT |
| | dest_city_name | STRING |
| | dest_state_name | STRING |
| | dep_time | TIMESTAMP |
| | dep_delay | FLOAT |
| | arr_time | TIMESTAMP |
| | arr_delay | FLOAT |
| | cancelled | BOOLEAN |
| FK | cancellation_code | STRING |
| | diverted | BOOLEAN |
| | weather_delay | FLOAT |

# Beam Pipeline

# Conversion of Time

```python
class FormatTime(beam.DoFn):
 def process(self, element):
    from datetime import datetime, timedelta
    import datetime
#checks if arrival time is the next day, and if it is sets the arrival
date 1 day after flight
    arr_date= fl_date
    if dep_time is not None and arr_time is not None and dep_time!=''
and arr_time!='' and fl_date is not None and fl_date!='':
        if int(dep_time)> int(arr_time) or int(arr_time)==2400:
            temp_date=
datetime.datetime.strptime(str(fl_date),"%Y-%m-%d")
            arr_date=temp_date + timedelta(days=1)
            array_date= str(arr_date).split(" ")
            arr_date=array_date[0]
        else:
            arr_date=fl_date
```

```python
# Handle time conversion
    # Convert departure time from int to time
    # New format will be hh:mm:ss
    dep_time_new = dep_time
    dep_time = str(dep_time)
    if dep_time != 'None' and dep_time != '':
        if dep_time == '2400':
            dep_time_new = '00:00:00'
        elif len(dep_time) == 3:
            dep_time_new ='0' + dep_time[:1] + ':' +
dep_time[1:] + ':00'
        elif len(dep_time) == 4:
            dep_time_new = dep_time[:2] + ':' + dep_time[2:] +
':00'
        elif len(dep_time) == 2:
            dep_time_new ='00:' + dep_time + ':00'
        elif len(dep_time) == 1:
            dep_time_new ='00:0' + dep_time + ':00'
        else:
            dep_time_new = '00:00:00'
    else:
        dep_time_new = '00:00:00'
```
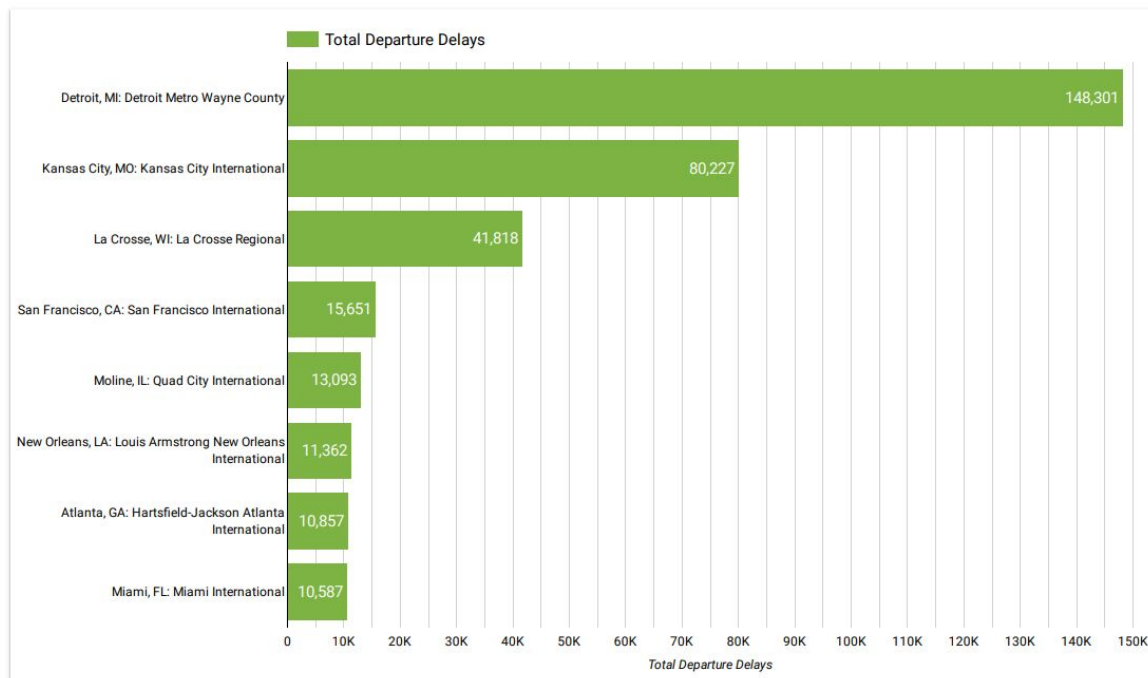
# Challenges

- \> 20,000,000 entities
- Formatting time fields for cross dataset joins
- Debugging data flow errors
  - 24:00:00 does not work in TIME
  - Decision to change to TIMESTAMP

# Query 1

```
SELECT a.description, SUM(a.dep_delay) AS
total_departure_delays
FROM (SELECT * FROM (SELECT * FROM
            bts_refined.On_Time_Performance_Dataflow
            WHERE EXTRACT(YEAR FROM fl_date) =
    2019) AS o
    INNER JOIN bts_refined.Airports c
    ON c.airport_id = o.origin_airport_id) AS a
INNER JOIN (SELECT *
    FROM noaa_refined.Events a
    INNER JOIN noaa_refined.Locations l
    ON l.event_id = a.event_id
    WHERE extract(year from a.begin_date) = 2019) AS b
ON UPPER(a.origin_state_name) = b.state
        AND UPPER(a.origin_city_name) = b.city
        AND a.dep_time BETWEEN b.begin_date AND
b.end_date
WHERE a.dep_delay > 0
GROUP BY a.description
HAVING SUM(a.dep_delay) > 10000
ORDER BY SUM(a.dep_delay) DESC LIMIT 8;
```
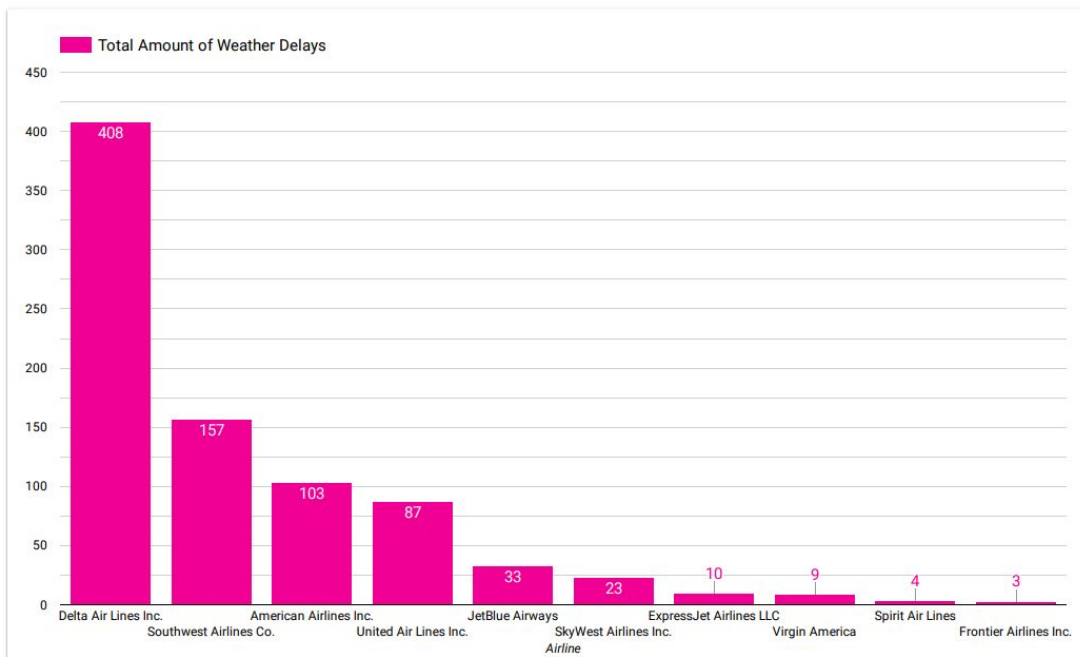


Airports with the Longest Storm Related Departure Delays in 2019

# Query 2

```
SELECT a.airline, COUNT(*) AS count
FROM (SELECT * FROM (SELECT * FROM
            bts_refined.On_Time_Performance_Dataflow
            WHERE EXTRACT(YEAR FROM fl_date) =
    2017) AS o
    INNER JOIN bts_refined.Carriers c
    ON c.code = o.carrier_code) AS a
INNER JOIN (SELECT *
    FROM noaa_refined.Events a
    INNER JOIN noaa_refined.Locations l
    ON l.event_id = a.event_id
    WHERE EXTRACT(YEAR FROM a.begin_date) = 2017)
AS b ON UPPER(a.origin_state_name) = b.state
    AND UPPER(a.origin_city_name) = b.city
    AND a.dep_time BETWEEN b.begin_date AND
b.end_date
WHERE a.weather_delay > 0
GROUP BY a.airline
ORDER BY COUNT(*) desc
LIMIT 10;
```
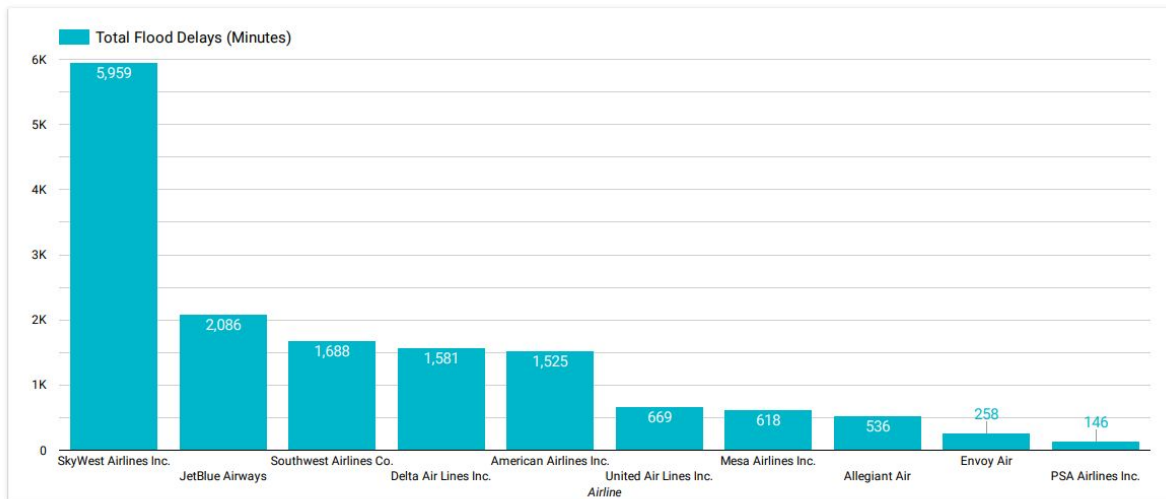
Airlines with the Most Storm Related Weather Delays in 2017

# Query 3

SELECT a.airline, SUM(a.dep_delay) AS total_flood_delays
FROM (SELECT * FROM (SELECT * from
                 bts_refined.On_Time_Performance_Dataflow
                 WHERE EXTRACT(YEAR FROM fl_date) =
       2018) AS o
       INNER JOIN bts_refined.Carriers c
       ON c.code = o.carrier_code) AS a
INNER JOIN (SELECT *
       FROM noaa_refined.Events a
       INNER JOIN noaa_refined.Locations l
       ON l.event_id = a.event_id
       WHERE EXTRACT(YEAR FROM a.begin_date) = 2018
       AND a.event_type = 'Flood') AS b
ON UPPER(a.origin_state_name) = b.state
       AND UPPER(a.origin_city_name) = b.city
       AND a.dep_time BETWEEN b.begin_date AND
b.end_date
WHERE a.weather_delay > 0
GROUP BY a.airline
ORDER BY SUM(a.dep_delay) desc
LIMIT 10;



Airlines with the Longest Total Flood Related Weather Delays in 2018

# Future Improvements

- Formatting airport location for more accurate cross-dataset joins
  - Get latitude and longitude for airport
- Flights that were cancelled by storms at the origin vs destination
- Use NTSB Aviation Accidents dataset