# Project Presentation

Mi Yon Kim, Yufei Tao

# Project Overview

- Dataset: Freeway data

- Tools & Languages:
    - MongoDB Compass
    - MongoDB Database
    - Python(Pymongo)

# Data model

**Original model(RDBMS)**

- Highways (2 rows)
- Stations (17 rows)
- Detectors (52 rows)
- Loopdata (large size)

**VS**

**Our model(NoSQL)**

- **Detectors**
  - Combined 3 tables(Highway, Stations, Detectors)
  - Embedded documents of Highways and Stations.
  - Contains duplicated data however it can be minor because data is small

- **Loopdata**
  - Added locationtext from the original model.

# Data model



```
_id: ObjectId("5fced998e5d0fb0e54a4e125")
detectorid: 1345
milepost: 14.32
locationtext: "Sunnyside NB"
detectorclass: 1
lanenumber: 1
schema: "1.0"
∨ station: Object
    stationid: 1045
    upstream: 0
    downstream: 1046
    stationclass: 1
    numberlanes: 4
    latlon: "45.43324,-122.565775"
    length: 0.94
∨ highway: Object
    highwayid: 3
    shortdirection: "N"
    direction: "NORTH"
    highwayname: "I-205"
```

**Detectors**

```
_id: ObjectId("5fced9b3e5d0fb0e54a4e159")
detectorid: 1345
starttime: "2011-09-15 00:01:40-07"
volume: 1
speed: 47
occupancy: 0
status: 3
dqflags: 0
schema: "1.0"
locationtext: "Sunnyside NB"
```

**Loopdata**

# ETL

- Implemented with python

- Drop loopdata which has the Null or zero value in the speed field

- Issues
- Some data in Detector's locationtext and Station's locationtext didn't match
*(e.g. Columbia to I-205 NB(Station table) vs I-205 NB at Columbia(Detectors))*

- Some detectorid in Loopdata doesn't included in Detectors table.
*(e.g. Loopdata with 1350 detectorid vs 1350 doesn't show up in Detectors table)*

- Solution 1: Modified the csv file

- Solution 2: Automate through script

# Query Examples

*Q1: Count low speeds and high speeds:*
Find the number of speeds < 5 mph or >
80 mph in the data set.

```
db.Loopdata.aggregate(
  [
    {
      $match: {
        $or: [
          {
            speed: {
              $lt: 5
            }
          },
          {
            speed: {
              $gt: 80
            }
          }
        ]
      }
    },
    {
      $count: "Number of speeds < 5 mph and > 80 mph: "
    }
  ]
)
```

< [ { 'Number of speeds < 5 mph and > 80 mph: ': 130182 } ]

# Query Examples

Q2: _Volume:_ Find the total volume for the station Foster NB on Sept 15, 2011.

```
[ { _id: 0, 'Total volume:': 49870 } ]
```

```
# use $options:'i' to make the query case-insensitive
query = db.Loopdata.aggregate(
  [
    {
      "$match": {
        "$and": [
          {"starttime" : {"$regex": "2011-09-15.*"}
          },
          {"locationtext" : "Foster NB"
          }
        ]
      }
    },
    {
      "$group": {
        "_id": 0,
        "Total volume:": {
          "$sum": "$volume"
        }
      }
    }
  ]
)
print(list(query))
```

# Query Examples

*Q5: Route Finding:* Find a route from Johnson Creek to Columbia Blvd on I-205 NB using the upstream and downstream fields.

```python
# Get all of the NB data
results = list(col_detectors.find({"locationtext": {"$regex": direction+"$", "$options":'i'_}}, {"station.stationid":1, "locationtext":

# save downstream of start point to find the next point
d_stream = start_doc["station"]["downstream"]
# Add start point
routes.append(start_doc["locationtext"])
for i in range(len(results)):
    depart = False
    for j in range(len(results)):
        # When finding previous point's downstream == stationid, add locationtext to routes
        # update the downstream
        # if updated downstream == end point's, ends of travel
        if d_stream == results[j]["station"]["stationid"]:
            routes.append(results[j]["locationtext"])
            d_stream = results[j]["station"]["downstream"]
            if d_stream == end_doc["station"]["stationid"]:
                depart = True
                break
        break
    if(depart == True):
        break

# Add end point
routes.append(end_doc["locationtext"])
```

```
(base) C:\gitRepo\CloudCluster\pdx-cs-cloud-cluster\queries>python3 query5.py
Johnson Cr NB
Foster NB
Powell to I-205 NB
Division NB
Glisan to I-205 NB
Columbia to I-205 NB
```

# Query Examples

_**Q6: Update**: Change the milepost of the Foster NB station to 22.6._

```
results = col_detectors.find({"locationtext": "Foster NB"})
milepost = "18.1"


print("\n\n---- Display the data before updating milepost ----")
for r in results:
    print("milepost:", r["milepost"])
    col_detectors.update_one({"milepost":r.get("milepost")}, {"$set":{"milepost": milepost}})


results = col_detectors.find({"locationtext": "Foster NB"})

print("\n\n---- Display the data after updating milepost to " + milepost + " --
for r in results:
    print("milepost", r["milepost"])
```

```
---- Display the data before updating milepost ----
milepost: 18.1
milepost: 18.1
milepost: 18.1



---- Display the data after updating milepost to 22.6 ----
milepost 22.6
milepost 22.6
milepost 22.6
```

Demo

# Critique

- Changed Loopdata model
  - Reason: MongoDB document size limit: 16mb

```
loopdata {
    _id
    schema
    detectorid
    data: [
        {
            starttime
            volume
            speed
            occupancy
            status
            dqflags
        }
    ]

    locationtext
    num_lowspeed
    num_highspeed
}
```
Previous loopdata

Pros
- Improve performance on query #1
- No duplicate data

Cons
- Complicated
- Doesn't fit in MongoDB System

```
_id: ObjectId("5fcdd22384a63b4ed4be430d")
detectorid: "1345"
starttime: "2011-09-15 00:01:40-07"
volume: 1
speed: 47
occupancy: 0
status: 3
dqflags: 0
schema: "1.0"
locationtext: "Sunnyside NB"
```
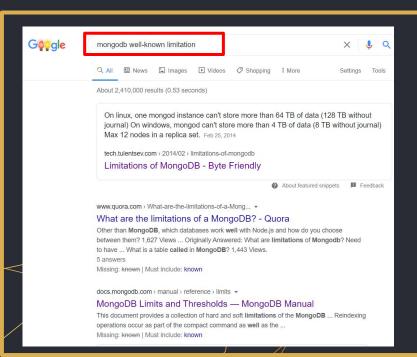current loopdata

Pros
- Simpler
- Simpler to make queries

Cons
- Duplicate data
- Bad performance on query #1 if the data size gets bigger

# Critique



- Advice
  - Spend more time on researching new system before you design

# Lessons learned

1. Declare each attributes with the proper data type.

_id: ObjectId("5fc731d4f1451a4748aad1b8")
detectorid: "1345"
milepost: "14.32"
locationtext: "Sunnyside NB"
detectorclass: "1"
lanenumber: 1
schema: "1.0"
station: Object
    stationid: "1045"
    upstream: "0"
    downstream: "1046"
    stationclass: "1"
    numberlanes: 4
    latlon: "45.43324,-122.565775"
    length: 0.94
highway: Object
    highwayid: "3"
    shortdirection: "N"
    direction: "NORTH"
    highwayname: "I-205"

_id: ObjectId("5fced998e5d0fb0e54a4e125")
detectorid: 1345
milepost: 14.32
locationtext: "Sunnyside NB"
detectorclass: 1
lanenumber: 1
schema: "1.0"
station: Object
    stationid: 1045
    upstream: 0
    downstream: 1046
    stationclass: 1
    numberlanes: 4
    latlon: "45.43324,-122.565775"
    length: 0.94
highway: Object
    highwayid: 3
    shortdirection: "N"
    direction: "NORTH"
    highwayname: "I-205"

# Lessons learned

2. Writing query purely using MongoDB could make things look

messy. But they are easy to understand.

```
cursor = db.Loopdata.aggregate([{"$match": {"$or": [{"speed": {"$lt": 5}}, {"speed": {"$gt": 80}}]}},{"$count": "Number of speeds < 5 mph and > 80 mph: "}])
```
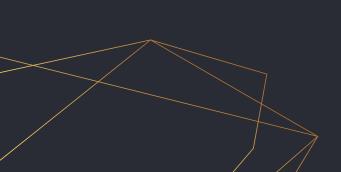
```
query = db.Loopdata.aggregate(
    [
        {
            "$match": {
                "$and": [
                    {"starttime" : {"$regex": "2011-09-15.*"}
                    },
                    {"locationtext" : "Foster NB"
                    }
                ]
            }
        },
        {
            "$group": {
                "_id": 0,
                "Total volume:": {
                    "$sum": "$volume"
                }
            }
        }
    ]
)
```

VS

```
...
import pymongo
import os
from pymongo import MongoClient
from dotenv import load_dotenv
load_dotenv()
MONGO_HOST = os.getenv("MONGO_HOST")
MONGO_DB = os.getenv("MONGO_DB")

# connect to cluster
cluster = MongoClient(MONGO_HOST, 27017)
db = cluster[MONGO_DB]

# Get collection for query
col_detectors = db["Detectors"]

results = col_detectors.find({"locationtext": "Foster NB"})
milepost = "18.1"

print("\n\n---- Display the data before updating milepost ----")
for r in results:
    print("milepost:", r["milepost"])
    col_detectors.update_one({"milepost":r.get("milepost")}, {"$set":{"milepost": milepost}})

results = col_detectors.find({"locationtext": "Foster NB"})

print("\n\n---- Display the data after updating milepost to " + milepost + " ----")
for r in results:
    print("milepost", r["milepost"])
```
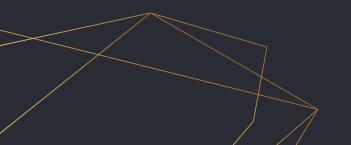
# Lessons learned

3. Solve complicated questions with combination of python and

queries.

```
...
import pymongo
import os
from pymongo import MongoClient
from dotenv import load_dotenv
load_dotenv()
MONGO_HOST = os.getenv("MONGO_HOST")
MONGO_DB = os.getenv("MONGO_DB")

# connect to cluster
cluster = MongoClient(MONGO_HOST, 27017)
db = cluster[MONGO_DB]

# Get collection for query
col_detectors = db["Detectors"]

results = col_detectors.find({"locationtext": "Foster NB"})
milepost = "18.1"

print("\n\n---- Display the data before updating milepost ----")
for r in results:
    print("milepost:", r["milepost"])
    col_detectors.update_one({"milepost":r.get("milepost")}, {"$set":{"milepost": milepost}})

results = col_detectors.find({"locationtext": "Foster NB"})

print("\n\n---- Display the data after updating milepost to " + milepost + " ----")
for r in results:
    print("milepost", r["milepost"])
```

Q & A

# Appendix

Project github link: <u>youn0125/pdx-cs-cloud-cluster (github.com)</u>

# Appendix

## Q1: Count low speeds and high speeds:
Find the number of speeds < 5 mph and > 80 mph in the data set.

```
db.Loopdata.aggregate(
  [
    {
      $match: {
        $or: [
          {
            speed: {
              $lt: 5
            }
          },
          {
            speed: {
              $gt: 80
            }
          }
        ]
      }
    },
    {
      $count: "Number of speeds < 5 mph and > 80 mph: "
    }
  ]
)
```

```
'
< [ { 'Number of speeds < 5 mph and > 80 mph: ': 130182 } ]
`
```

# Appendix

*Q2: Volume:* Find the total volume for the station Foster NB for Sept 15, 2011.

```
# use $options:'i' to make the query case-insensitive
query = db.Loopdata.aggregate(
  [
    {
      "$match": {
        "$and": [
          {"starttime" : {"$regex": "2011-09-15.*"}
          },
          {"locationtext" : "Foster NB"
          }
        ]
      }
    },
    {
      "$group": {
        "_id": 0,
        "Total volume:": {
          "$sum": "$volume"
        }
      }
    }
  ]
)
print(list(query))
```

```
  )
< [ { _id: 0, 'Total volume:': 49870 } ]
```

# Appendix

*Q5. Route Finding:* Find a route from Johnson Creek to Columbia Blvd on I-205 NB using the upstream and downstream fields.

```python
# Get all of the NB data
results = list(col_detectors.find({"locationtext": {"$regex": direction+"$", "$options":'i'_}}, {"station.stationid":1, "locationtext":

# save downstream of start point to find the next point
d_stream = start_doc["station"]["downstream"]
# Add start point
routes.append(start_doc["locationtext"])
for i in range(len(results)):
    depart = False
    for j in range(len(results)):
        # When finding previous point's downstream == stationid, add locationtext to routes
        # update the downstream
        # if updated downstream == end point's, ends of travel
        if d_stream == results[j]["station"]["stationid"]:
            routes.append(results[j]["locationtext"])
            d_stream = results[j]["station"]["downstream"]
            if d_stream == end_doc["station"]["stationid"]:
                depart = True
                break
        break
    if(depart == True):
        break

# Add end point
routes.append(end_doc["locationtext"])
```
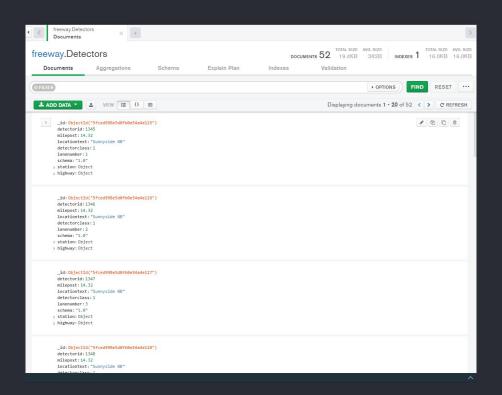
```
(base) C:\gitRepo\CloudCluster\pdx-cs-cloud-cluster\queries>python3 query5.py
Johnson Cr NB
Foster NB
Powell to I-205 NB
Division NB
Glisan to I-205 NB
Columbia to I-205 NB
```

# Appendix

*Q6: Update:* Change the milepost of the Foster NB station to 22.6.

```python
results = col_detectors.find({"locationtext": "Foster NB"})
milepost = "18.1"

print("\n\n---- Display the data before updating milepost ----")
for r in results:
    print("milepost:", r["milepost"])
    col_detectors.update_one({"milepost":r.get("milepost")}, {"$set":{"milepost": milepost}})

results = col_detectors.find({"locationtext": "Foster NB"})

print("\n\n---- Display the data after updating milepost to " + milepost + " ----")
for r in results:
    print("milepost", r["milepost"])
```

# Appendix

Data size:

# Appendix

Data size: