

DAV6100 FINAL PROJECT

NYC311 Complaints Analysis

Group 2

Agenda

Data Profile

Conceptual Architecture

Demo

Project Milestones & Timeline

Team Responsibilities

Challenges

Lessons Learned

Next Steps

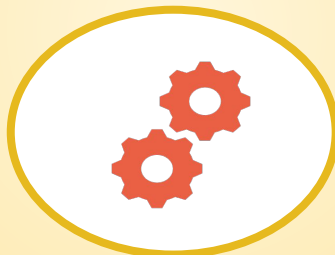
Project Overview

*311 Requests
NYC Covid19
NYC Median Income*



Many data sources

*Combines data
sources into DW
that is easy to
maintain*



Data Warehouse

*Tableau
dashboard
Python*



Data Visualization

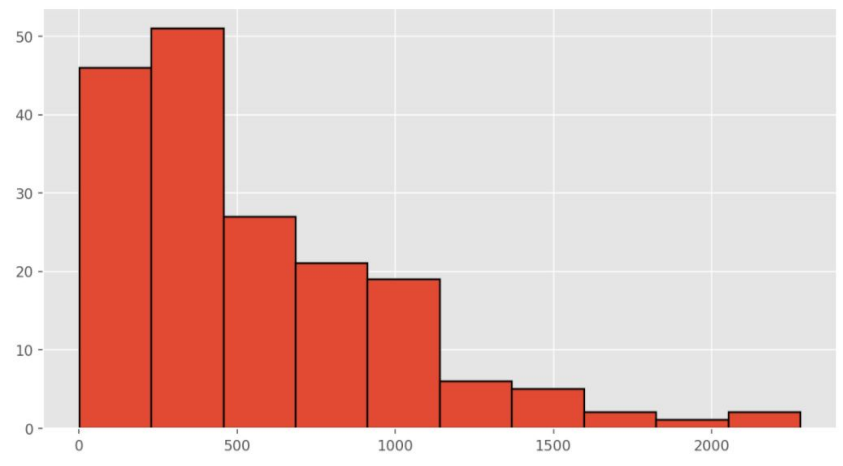
Data Profile (1): NYC311 Dataset



Dataset Summary

Source of Information	https://portal.311.nyc.gov
Number of Records	1.3M (Original 25M)
Frequency of updates	Daily
Data type and structure	Text, Integers, Float, Date, Time API, CSV
Number of columns	16
Granularity	Each 311 complaint

Frequency of Complaints

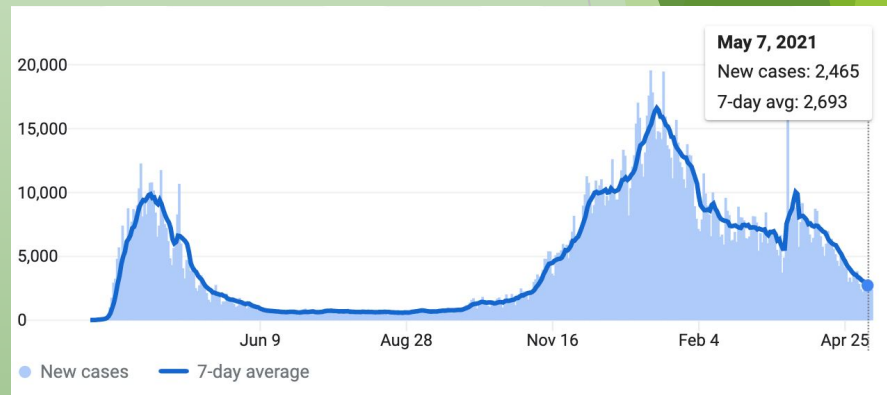


Data Profile (2): NYC COVID-19 2019



Dataset Summary

Source of Information	https://github.com/nychealth/coronavirus-data
Number of Records	6764
Frequency of updates	Daily
Data type and structure	Text, Float API, CSV
Number of columns	3
Granularity	Caserate grouping by ZIP Code



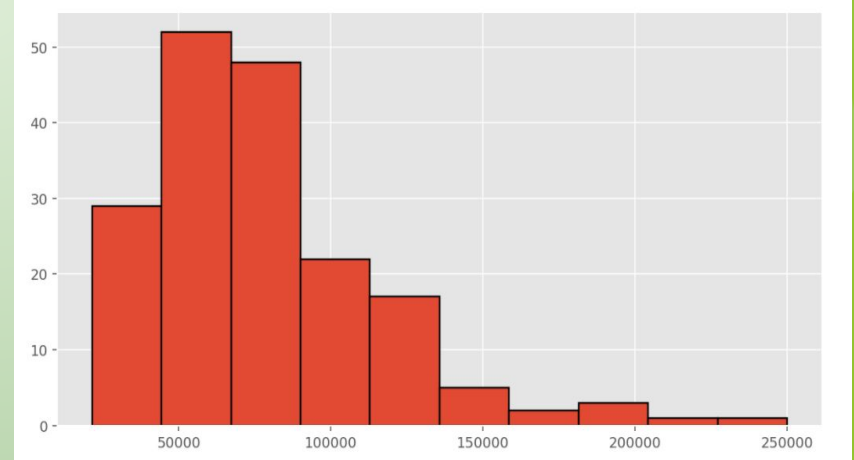
Data Profile (3): NYC Median Income



Dataset Summary

Source of Information	https://data.cccnewyork.org
Number of Records	181
Frequency of updates	Yearly
Data type and structure	Integers CSV
Number of columns	3
Granularity	Income by ZIP code

Average Income Per Household



Schema Selection

The first thing to notice about the dimensional schema is its simplicity and symmetry. The simplicity of a dimensional model also has performance benefits. Dimensional models are gracefully extensible to accommodate

--- Kimball The Data Warehouse Toolkit by Ralph Kimball, Margy Ross

1

Millions of Records

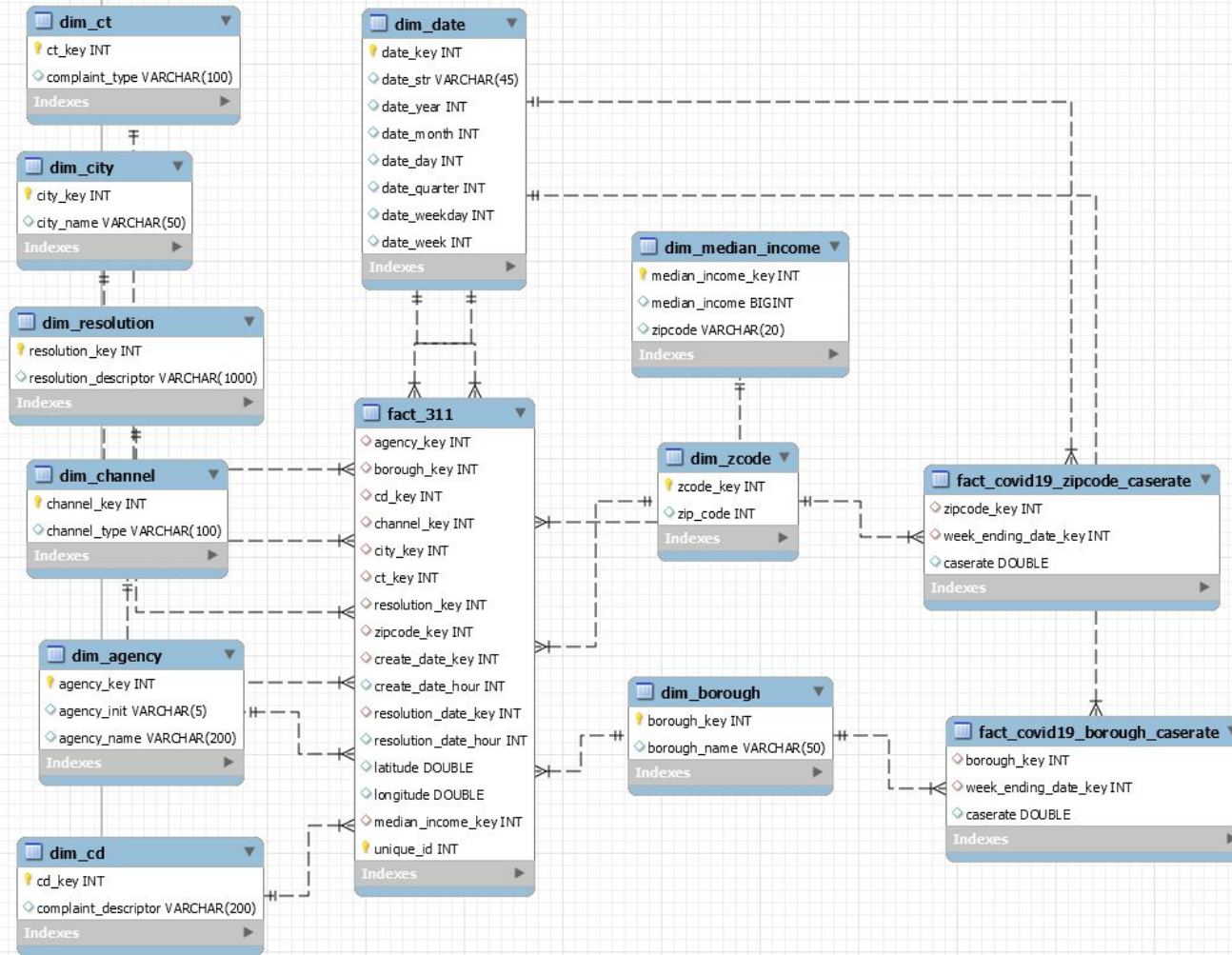
2

Aggregation

3

Not Only for Us

ER Diagram



Assumptions

1

The variables we are integrating to our warehouse affected the 311 complaint rate

2

Stakeholders would want daily updates

3

Two Choices:

1. Combine API to csv & create a new file.
 2. Apply star schema to 311 data
- For practicing and data updating purpose, the second choice was decided.

4

Easy to work asynchronously

5

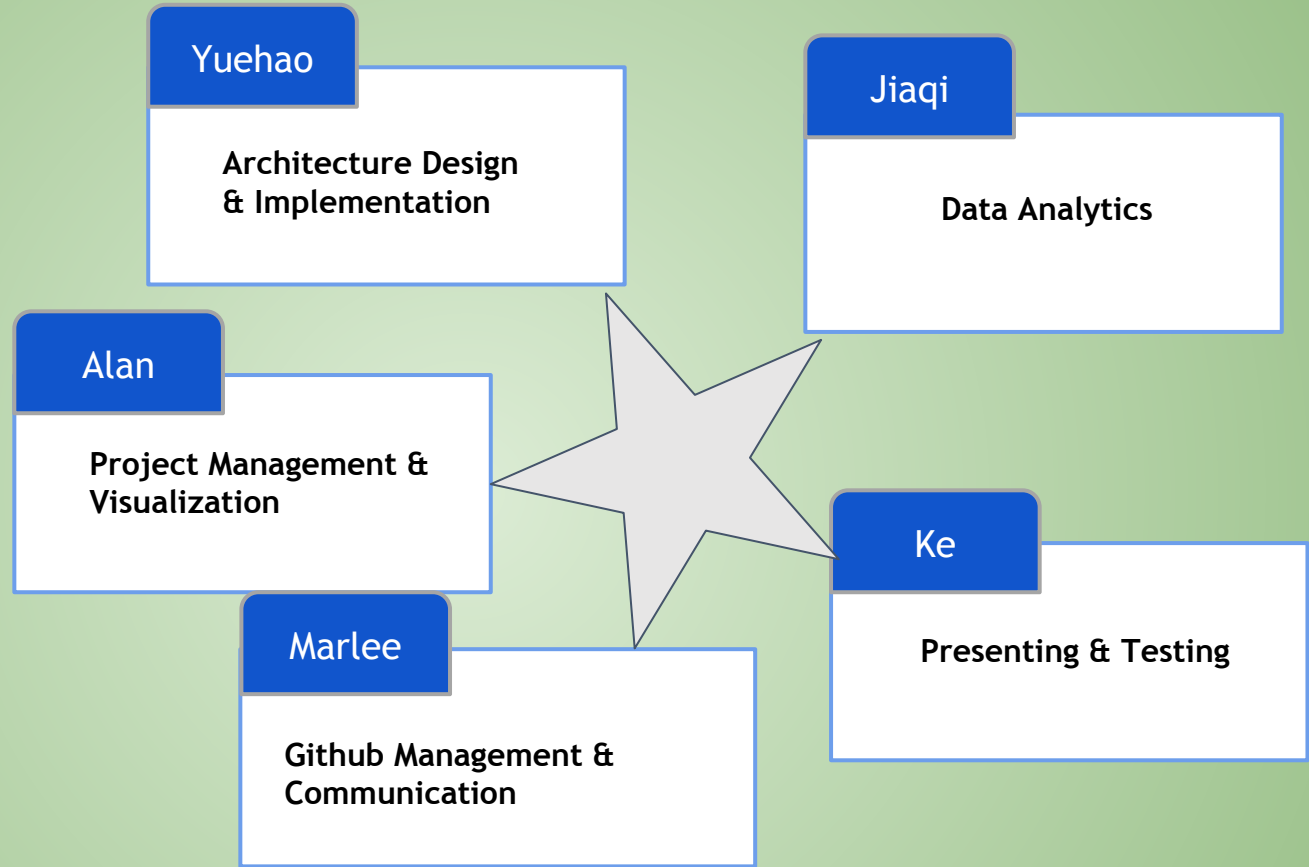
No problem with doing class while also learning how to do the project 'on the run'

6

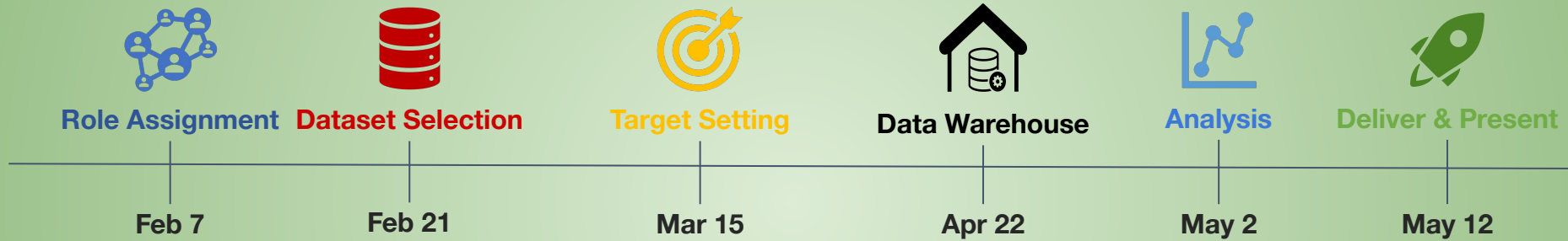
Data amount not overwhelming



Team Leads

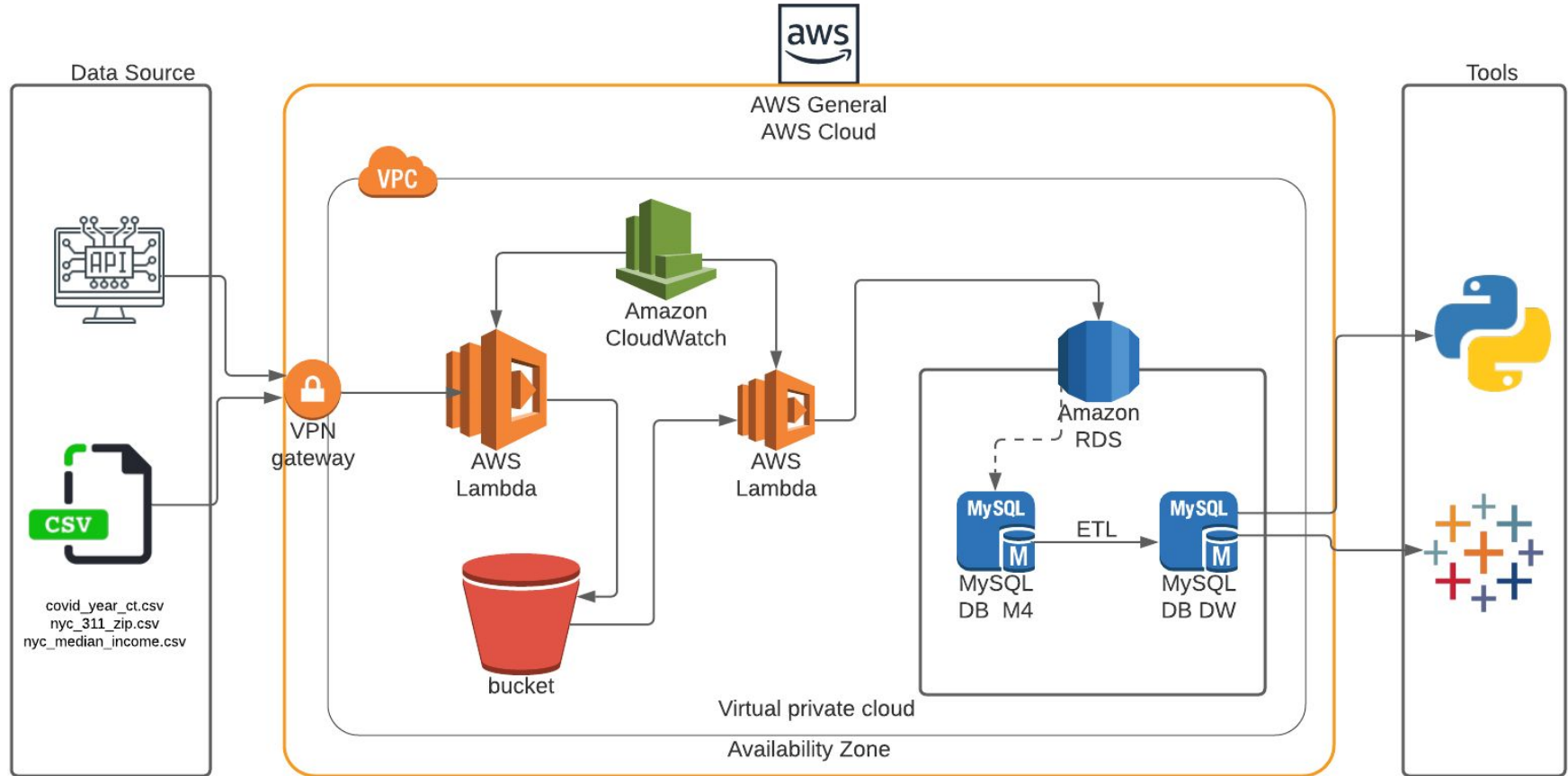


Milestones



Timeline

Conceptual Architecture



DEMO

Yuehao: Architecture

Alan: Tableau

Marlee: Statistics &
hypothesis testing



Challenges

Blank Slate
Asynchronous Team
Big Data

Lessons Learned

**ENSURE CLARITY OF
PROJECT PURPOSE,
ROLE EXPECTATIONS,
AND DELIVERABLES AS
EARLY AS POSSIBLE**

**DIAGRAMS BEFORE
IMPLEMENTATION IS
GOOD, BUT NOT AS
GOOD AS STARTING
SMALL THEN TESTING**

**THERE IS A REAL
TRADEOFF TO BUILDING
AN ARCHITECTURE FOR
READING VS WRITING**

**RESEARCH COSTS AND
SET A BUDGET**

Next Steps

Incorporate new datasets to measure new variables [race, twitter sentiments etc]

Build a more robust architecture

Leverage Glue and other AWS products to create a more robust architecture instead of custom scripts. ETL and transferring is done with scripts now, instead of aws tools. All hail Amazon

Expand the dataset to NY state

Thank you!



Any Questions?