# Data Profiling, Cleaning and Ingestion using MapReduce

**Anirudh Garg**
Courant Institute of
Mathematical Sciences
New York University
`ag9563@nyu.edu`

## Objective

To perform data profiling, cleaning, and ingestion on a dataset containing incident-related information. Specifically, we aimed to:

- Filter and clean data based on specific criteria.
- Compute profiling statistics to understand data quality and characteristics.
- Save the cleaned dataset and profiling statistics for further analysis.

## Dataset Overview

### Fire Department and Emergency Medical Services Dispatched Calls for Service

- **Link**: Dataset Link
- **Ownership**: Publicly accessible data from the City of San Francisco.
- **Data Size**: 2.7 GB

These dataset is periodic, and being updated regularly.

The dataset consisted of the following columns:

- **Call Number:** A unique 9-digit number assigned by the 911 Dispatch Center (DEM) to this call. These number are used for both Police and Fire calls.
- **Unit ID:** Unit Identifier. For example E01 for Engine 1 or T01 for Truck 1.
- **Unit ID:** Unit Identifier. For example E01 for Engine 1 or T01 for Truck 1.
- **Incident Number:** A unique 8-digit number assigned by DEM to this Fire incident.
- **Call Type:** Type of call the incident falls into.
- **Call Date:** Date the call is received at the 911 Dispatch Center. Used for reporting purposes.
- **Watch Date:** Watch date when the call is received. Watch date starts at 0800 each morning and ends at 0800 the next day.
- **Received DtTm:** Date and time of call is received at the 911 Dispatch Center.
- **Entry DtTm:** Date and time the 911 operator submits the entry of the initial call information into the CAD system
- **Dispatch DtTm:** Date and time the 911 operator dispatches this unit to the call.
- **Response DtTm:** Date and time this unit acknowledges the dispatch and records that the unit is en route to the location of the call.
- **On Scene DtTm:** Date and time the unit records arriving to the location of the incident

- **Transport DtTm:** If this unit is an ambulance, date and time the unit begins the transport to the hospital
- **Hospital DtTm:** If this unit is an ambulance, date and time the unit arrives to the hospital.
- **Call Final Disposition:** Disposition of the call (Code). For example TH2: Transport to Hospital - Code 2, FIR: Resolved by Fire Department
- **Available DtTm:** Date and time this unit is not longer assigned to this call and it is available for another dispatch.
- **Address:** Address of intersection or call box point associated with incident (obfuscated address to protect caller privacy)
- **City:** City of incident
- **Zipcode of Incident:** Zip code of incident
- **Battalion:** Emergency Response District (There are 10 Fire Emergency Response Districts).
- **Station Area:** Fire Station First Response Area associated with the address of the incident.
- **Box:** Fire box associated with the address of the incident. A box is the smallest area used to divide the City. Each box is associated with a unique unit dispatch order. The City is divided into more than 2,400 boxes.
- **Original Priority:** Initial call priority (Code 2: Non-Emergency or Code 3:Emergency).
- **Priority:** Call priority once all information has been asssessed (Code 2: Non-Emergency or Code 3:Emergency).
- **Final Priority:** Final call priority (Code 2: Non-Emergency or Code 3:Emergency).
- **ALS Unit:** Does this unit includes ALS (Advance Life Support) resources? Is there a paramedic in this unit?
- **Call Type Group:** Call types are divided into four main groups: Fire, Alarm, Potential Life Threatening and Non Life Threatening.
- **Number of Alarms:** There are five levels of fire alarms (1-5). The number of alarms indicates the number of resources required in an incident. This number is a combination of engines, trucks, rescue squads, chiefs and EMS units.
- **Unit Type:** Type of unit responding
- **Unit sequence in call dispatch:** A number that indicates the order this unit was assigned to this call.
- **Fire Prevention District:** Bureau of Fire Prevention District associated with this address
- **Supervisor District:** Supervisor District number
- **Neighborhooods - Analysis Boundaries:** San Francisco Neighborhood associated with the incident address.
- **RowID:** Unique identifier used for managing data updates. It is the concatenation of Call Number and Unit ID separated by a dash.
- **case_location:** Latitude and Longitude for the call
- **data_as_of:** Timestamp when the record (row) was last updated in the source system.
- **data_loaded_at:** time the data was loaded into the Open Data Portal

## Data Acquisition and Access

- **Storage**: Data will be stored on distributed file system, HDFS, to accommodate storage needs and allow team-wide accessibility.
- **Permissions**: As these datasets are public, no additional permissions are required beyond accessing and downloading from the respective government data portals.
- **Access Time**: Instant access; no approval needed for these publicly hosted datasets.

# Data Profiling and Cleaning Approach

We focused on profiling and cleaning the dataset, ensuring that data is structured, consistent, and free of erroneous entries. MapReduce and other technolgies taught in the course were used to profile, clean, and transform the data.

### Profiling

- **Purpose**: Characterize each data column to understand its content, distribution, and structure. This includes identifying column data types, ranges, unique values, and possible data issues.
- **Tools**: MapReduce jobs will be written to process large datasets and output profiles on each column.

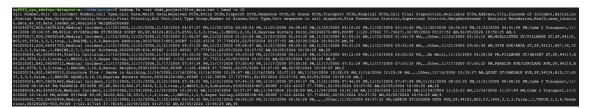### Steps Taken for Data Filtering, Transformation, and Analysis

- **Filter Columns:** Created the `FilterColumns.java` program, including a Mapper (`FilterColumnsMapper.java`) and a Reducer (`FilterColumnsReducer.java`) to filter specific columns from the input data. The columns filtered for further analysis are as follows:

  ```
  Call Number, Incident Number, Call Type, Call Date, Received DtTm,
  Response DtTm, On Scene DtTm, Transport DtTm, Hospital DtTm, Call Final
  Disposition, Address, City, Zipcode of Incident, Battalion, Station Area,
  Box, Final Priority, ALS Unit, Call Type Group, Number of Alarms, Unit
  Type, Neighborhooods - Analysis Boundaries
  ```

- **Filter and Transform:** Created the `FilterAndTransform.java` program with a Mapper (`FilterAndTransformMapper.java`) and a Reducer (`FilterAndTransformReducer.java`) to apply transformations to the filtered data. The transformation applied was to convert `mmddyyyy` to `yyyymmdd` and then taking only a subset of the dates, Jan 1st 2014 to October 31st 2024 (20140101 to 20241031). This was done using the `Call Date` column. Also, in the `city` column, we find there are multiple names referring to San Francisco, (such as SFO, SAN FRANCISCO, San Francisco), we combine them to refer to San Francisco.

- **Value Counts:** To get to know the value counts of each of the selected columns, I wrote the `ValueCounts.java` program along with the corresponding Mapper (`ValueCountsMapper.java`) and Reducer (`ValueCountsReducer.java`) to calculate the frequency of unique values in a specified column. The function was written such that apart from the 2 standard arguments of input path and output path, another argument was added which took in the parameter of column name for which the value count is required.

### Outputs

### Original Data:



### Filter Columns:

**Filter and Transform Columns:**

Call Number,Incident Number,Call Type,Call Date,Received DtTm,Response DtTm,On Scene DtTm,Transport DtTm,Hospital DtTm,Call Final Disposition,Address,City,Zipcode of Incident,Battalion,Station Area,Box,Final Priority,ALS Unit,Call Type Group,Number of Alarms,Unit Type,Neighborhooods - Analysis Boundaries
213132776,21137938,Medical Incident,11/09/2021,20211109,11/09/2021 06:23:34 PM,11/09/2021 06:49:41 PM,,,Patient Declined Transport,CLEMENT ST/09TH AVE,San Francisco,94118,B07,31,7135,2,false,Non Life-threatening,1,SUPPORT,Inner Richmond
213121482,21137277,Medical Incident,11/08/2021,20211108,11/08/2021 12:19:08 PM,11/08/2021 12:23:21 PM,11/08/2021 12:38:21 PM,11/08/2021 12:49:45 PM,Code 3 Transport,ANZA ST/ARGUELLO BLVD,San Francisco,94118,B07,31,7112,3,false,Potentially Life-Threatening,1,PRIVATE,Lone Mountain/USF
213151169,21138683,Medical Incident,11/11/2021,20211111,11/11/2021 10:47:17 AM,11/11/2021 10:50:22 AM,11/11/2021 11:05:33 AM,11/11/2021 11:39:27 AM,Code 2 Transport,VALENCIA ST/14TH ST,San Francisco,94103,B02,06,5126,3,true,Potentially Life-Threatening,1,MEDIC,Mission
213130539,21137656,Alarms,11/09/2021,20211109,11/09/2021 07:18:42 AM,11/09/2021 07:21:23 AM,,,Fire,EDDY ST/LEAVENWORTH ST,San Francisco,94102,B02,03,1545,3,false,Alarm,1,ENGINE,Tenderloin
213132443,21137899,Medical Incident,11/09/2021,20211109,11/09/2021 04:35:52 PM,11/09/2021 04:35:53 PM,,,Code 3 Transport,ELLIS ST/POWELL ST,San Francisco,94102,B03,13,1322,3,true,Potentially Life-Threatening,1,ENGINE,Financial District/South Beach
213151333,21138709,Medical Incident,11/11/2021,20211111,11/11/2021 11:38:18 AM,11/11/2021 12:00:19 PM,11/11/2021 12:04:57 PM,11/11/2021 12:36:51 PM,Code 2 Transport,JERROLD AVE/EARL ST,San Francisco,94124,B10,17,6713,2,false,Non Life-threatening,1,PRIVATE,Bayview Hunters Point
213121203,21137242,Medical Incident,11/08/2021,20211108,11/08/2021 11:51:56 AM,,,,Code 2 Transport,PLAZA ST/LAGUNA HONDA BLVD,San Francisco,94116,B08,20,8641,2,false,Non Life-threatening,1,SUPPORT,West of Twin Peaks
213151407,21138722,Medical Incident,11/11/2021,20211111,,,,,No Merit,MONTEREY BLVD/SAN RAFAEL WAY,San Francisco,94127,B09,19,8553,3,false,Non Life-threatening,1,TRUCK,West of Twin Peaks
213141329,21138224,Alarms,11/10/2021,20211110,11/10/2021 11:30:13 AM,11/10/2021 11:34:49 AM,,,Fire,FEDERAL ST/DELANCEY ST,San Francisco,94107,B03,35,2134,3,true,Alarm,1,ENGINE,Financial District/South Beach

**Value Counts: Call Type**

```
"Extrication / Entrapped (Machinery      636
Administrative   201
Aircraft Emergency        236
Alarms   407765
Assist Police    441
Citizen Assist / Service Call   49621
Confined Space / Structure Collapse      555
Electrical Hazard        14550
Elevator / Escalator Rescue       10676
Explosion        1062
Fuel Spill       3073
Gas Leak (Natural and LP Gases) 21610
HazMat   1176
High Angle Rescue        819
Industrial Accidents     796
Lightning Strike (Investigation)        15
Marine Fire      228
Medical Incident         2337276
Mutual Aid / Assist Outside Agency      470
Odor (Strange / Unknown)         3015
Oil Spill        5
Other    55174
Outside Fire     52995
Smoke Investigation (Outside)    7626
Structure Fire / Smoke in Building       262659
Suspicious Package       109
Traffic Collision        137564
Train / Rail Fire        129
Train / Rail Incident    896
Vehicle Fire     11310
Water Rescue     22121
Watercraft in Distress   640
```

**Value Counts: Call Final Disposition**

```
Against Medical Advice   62486
CHP       790
Cancelled         84212
Code 2 Transport         1555078
Code 3 Transport         157174
Duplicate        1038
Fire    868966
Gone on Arrival 13043
Medical Examiner         48194
Multi-casualty Incident 421
No Merit          160880
Other    181642
Patient Declined Transport      174690
SFPD     20674
Unable to Locate         75525
```

**Value Counts: Call Type Group**

```
Alarm    780915
Fire     129772
Non Life-threatening     815349
Potentially Life-Threatening    1649296
```

**Value Counts: Unit Type**

```
AIRPORT 1778
BLS      11490
CHIEF    243206
CP       35315
ENGINE   1171193
Fire     6
INVESTIGATION    3036
MEDIC    1025677
PRIVATE 329392
RESCUE CAPTAIN   108944
RESCUE SQUAD     48551
SUPPORT 87254
TRUCK    331527
```

**Value Counts: Neighborhoods - Analysis Boundaries**

```
Bayview Hunters Point     182161
Bernal Heights   58914
Castro/Upper Market       84680
Chinatown        63025
Excelsior        64271
Financial District/South Beach  244692
Glen Park        17286
Golden Gate Park          27139
Haight Ashbury   47562
Hayes Valley     81797
Inner Richmond   39130
Inner Sunset     47852
Japantown        34682
Lakeshore        47708
Lincoln Park     2225
Lone Mountain/USF         46250
Marina   66337
McLaren Park      2362
Mission 313973
Mission Bay      52148
Nob Hill         113069
Noe Valley       36150
None     2722
North Beach      64126
Oceanview/Merced/Ingleside        42629
Outer Mission    49717
Outer Richmond   82925
Pacific Heights  63650
Portola 32821
Potrero Hill     37521
Presidio         23304
Presidio Heights          27498
Russian Hill     55252
Seacliff         6092
South of Market 346105
Sunset/Parkside 126923
Tenderloin       499062
Treasure Island 16928
Twin Peaks       13862
Visitacion Valley         41495
West of Twin Peaks        73328
Western Addition          118398
```