

# Title: - NTD\_Public Transit Improvement

## Introduction:

## Dataset Overview:

The dataset, "2022 NTD Annual Data – Vehicle Age Distribution & Vehicles Type Count by Agency", likely contains comprehensive information about public transit systems, vehicle types and ages including types of vehicles used by various transit agencies in a particular year. Such datasets typically include:

- Vehicle type (buses & it's types etc.)
- Quantities or capacities
- Operating agencies
- Geographical information
- Operational data (like usage statistics)

## Data Life Cycle:

1. **Data Processing Using Open Refine tool:** It is used to enhance the quality of the data like cleaning and transforming. We used our both static and streaming data source in this tool for better refinement and finally ensured that the data is accurate, consistent, and suitable for the intended analysis.

The screenshot displays the OpenRefine interface for clustering data. The left panel, titled "Cluster and edit column 'City'", shows a table of clusters with columns for "Cluster size", "Row count", "Values in cluster", "Merge?", and "New cell value". The "Merge?" column contains checkboxes, and the "New cell value" column contains the merged values. The "Merge" button is circled in red. The central data table shows the results of the clustering operation, with columns for "Average Age of Fleet (in Years)" and "Average Lifetime Miles per Vehicle". The right panel shows the column filters. The "Merge selected & re-cluster" button is also circled in red at the bottom of the interface.

Cluster size	Row count	Values in cluster	Merge?	New cell value
2	3	• Lagrange (2 rows) • La Grange	<input checked="" type="checkbox"/>	Lagrange
2	5	• Lafayette (4 rows) • La Fayette	<input checked="" type="checkbox"/>	Lafayette
2	2	• Cœur D'Alene (2 rows) • Cœur D'Alene	<input checked="" type="checkbox"/>	Cœur D'Alene
2	3	• Southgate (2 rows) • South Gate	<input checked="" type="checkbox"/>	Southgate
2	5	• New Castle (4 rows) • Newcastle	<input checked="" type="checkbox"/>	New Castle

icles	Average Age of Fleet (in Years)	Average Lifetime Miles per Vehicle	
4.12	118,610	4.12	118,610
4.12	118,610	4.12	118,610
8.18	184,696	8.18	184,696
6.93	231,206	6.93	231,206
26.07	1,391,691	26.07	1,391,691
4.59	146,004	4.59	146,004
5.89	181,352	5.89	181,352
5.89	181,352	5.89	181,352
6.96	175,524	6.96	175,524

Permalink

Fill down 23 cells in column City [Undo](#)

2771 rows

Show as: **rows** records Show: 5 10 25 50 100 500 1000 rows

All	Agency	City	State	NTD ID	Organization Type	Reporter Type	UA
1.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217
2.	New Jersey Transit	Newark	NJ	20080	Publicly-Owned or Privately Chartered	Full Reporter	63217

OpenRefine NTD Vehicles Age Distribution [Facets](#)

Facet / Filter Undo / Redo 7909 rows

Using facets and filters

Use facets and filters to select subsets of your data to act on. Choose facet and filter methods from the menu at the top of each data column.

Not sure how to get started? Watch these screencasts

All	Agency	City	State	NTD ID	Organization Type	Reporter Type	UACE Code	UZA Name
1.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
2.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
3.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
4.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
5.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
6.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
7.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
8.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
9.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
10.	MTA New York City Transit	Brooklyn	NY	20008	Unit of a Transit Agency, Reporting Separately	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ
11.	New Jersey Transit	Newark	NJ	20080	Publicly-Owned or Privately Chartered	Full Reporter	63217	New York - Jersey City - Newark, NY - NJ

OpenRefine NTD VehiclesCount by Agency [Facets](#)

Facet / Filter Undo / Redo 2771 rows

Using facets and filters

Use facets and filters to select subsets of your data to act on. Choose facet and filter methods from the menu at the top of each data column.

Not sure how to get started? Watch these screencasts

All	Agency	City	State	NTD ID	Organization Type	Reporter Type	UACE Code	UZA Name
101.	Blue Water Area	Port Huron	MI	50145	Public Agency or Authority of Transit Service	Full Reporter	71155	Port Huron, MI
102.	City of Detroit	Detroit	MI	50119	County or Local Government Unit or Department of Transportation	Full Reporter	23824	Detroit, MI
103.	Southern California Regional Rail Authority - Blue Water Area	Los Angeles	CA	90151	Public Agency or Authority of Transit Service	Full Reporter	51445	Los Angeles - Long Beach - Inglewood, CA
104.	Central Arkansas Development Council (aka South Central Arkansas Transit)	Benton	AR	60248	Corporation	Rural Reporter	50392	Little Rock, AR
105.	Tampa Bay Area Regional Transit Agency	Tampa	FL	40200	Public Agency or Authority of Transit Service	Full Reporter	80099	Tampa-St. Petersburg, FL
106.	Riverside	Riverside	CA	90101	Public Agency or Authority of Transit Service	Full Reporter	75240	Riverside-San Bernardino, CA
107.	Interurban	Grand Rapids	MI	50103	Public Agency or Authority of Transit Service	Full Reporter	54300	Grand Rapids, MI
108.	Denton County	Leander	TX	60101	Public Agency or Authority of Transit Service	Full Reporter	20000	Denton-Leander, TX
109.	METRO	Albany	OH	50110	Public Agency or Authority of Transit Service	Full Reporter	766	Albany, OH

2. **Set up Storage:** Secure and scalable storage is essential for managing large datasets. Having this data in Cloud Storage allows for easy accessibility and manipulation as needed, crucial for transit agencies looking to analyze their fleet and operational efficiency.

← Bucket details

**ntdvehicle-finalproj**

Location	Storage class	Public access	Protection
us-south1 (Dallas)	Standard	Not public	None

OBJECTS	CONFIGURATION	PERMISSIONS	PROTECTION	LIFECYCLE	OBSERVABILITY	INVENTORY REPORTS
---------	---------------	-------------	------------	-----------	---------------	-------------------

Buckets > ntdvehicle-finalproj

UPLOAD FILES UPLOAD FOLDER CREATE FOLDER TRANSFER DATA MANAGE HOLDS DOWNLOAD DELETE

Filter by name prefix only Filter objects and folders

Name	Size	Type	Created	Storage class	Last modified	Public access
<input type="checkbox"/> NTD_VehiclesAge.csv	1.5 MB	text/csv	Dec 3, 2023, 9:39:06 AM	Standard	Dec 3, 2023, 9:39:06 AM	Not public
<input type="checkbox"/> NTD_VehiclesCount.csv	834 KB	text/csv	Dec 3, 2023, 9:39:06 AM	Standard	Dec 3, 2023, 9:39:06 AM	Not public

3. **Set Up Hadoop Ecosystem (Clusters) using Dataproc:** In this phase, we need to enable the GCP API's called Compute Engine API & Cloud Dataproc API. We created a Cluster in Dataproc and the nodes with 1 Manager & 2 Worker Nodes in Compute Engine→VM Instances. This gives transit agencies real-time insights into their operations, enabling quick decision-making for issues like vehicle deployment and route adjustments.

Google Cloud

DataMasters Group FianlProj

Search (/) for resources, docs, products, and more

Search

6

Clusters

CREATE CLUSTER

REFRESH

START

STOP

DELETE

REGIONS

+ 5 RECOMMENDED ALERTS

SHOW INFO PANEL

Filter

Search clusters, press Enter

	Name	Status	Region	Zone	Total worker nodes	Flexible VMs?	Scheduled deletion	Cloud Storage staging bucket	Created
	<a href="#">ned-hadoospark-cluster-finalproj</a>	Running	us-central1	us-central1-a	2	No	Off	<a href="#">ntdvehicle-finalproj</a>	Dec 3, 2023, 9:46:57 AM

Google Cloud

DataMasters Group FianlProj

Search (/) for resources, docs, products, and more

Search

6

?

Y

VM instances

CREATE INSTANCE

IMPORT VM

REFRESH

LEARN

INSTANCES

OBSERVABILITY

INSTANCE SCHEDULES

VM instances

Filter Enter property name or value

	Status	Name	Zone	Recommendations	In use by	Internal IP	External IP	Connect
		ned-hadoospark-cluster-finalproj-m	us-central1-a			10.128.0.3 (nic0)	35.192.106.59 (nic0)	SSH
		ned-hadoospark-cluster-finalproj-w-0	us-central1-a			10.128.0.4 (nic0)	34.31.248.0 (nic0)	SSH
		ned-hadoospark-cluster-finalproj-w-1	us-central1-a			10.128.0.2 (nic0)	34.66.13.201 (nic0)	SSH

#### 4. Creating Sample Queries: Performing moderate to complex queries and analytics using BigQuery, Hive, and Spark.

##### a. Big Query Studio (Running fast, ad-hoc queries on large datasets for operational reports.)

- For this we enabled Big query API and created a table from both the datasets directly from the cloud storage bucket we created.
- Performed some queries in Big Query Studio
- Below are the examples for both the datasets.

Big Queries for Data set 1:

Q	Total Vehicles by state	RUN	SAVE QUERY	JOB INFORMATION	RESULTS	CHART	PREVIEW
1	SELECT State, SUM(`Total_Vehicles`) AS TotalVehicles			Row	State	TotalVehicles	
2	FROM `VehiclesAge_2.VehiclesAge`			1	CA	26187	
3	GROUP BY State			2	NY	24551	
4	ORDER BY TotalVehicles DESC			3	IL	11408	
5	LIMIT 5;			4	TX	10030	
				5	WA	9222	

Q	Top Cities with more than 1...les	RUN	JOB INFORMATION	RESULTS	CHART	PREVIEW
1	SELECT City, SUM(`Total_Vehicles`) AS TotalVehicles			Row	City	TotalVehicles
2	FROM `VehiclesAge_2.VehiclesAge`			1	Brooklyn	14358
3	GROUP BY City			2	Chicago	6137
4	HAVING SUM(`Total_Vehicles`) > 1000			3	Los Angeles	5615
5	ORDER BY SUM(`Total_Vehicles`) DESC			4	Newark	5487
6	LIMIT 5;			5	Washington	4764

Big Studio Queries for Dataset 2:



### List of Agenci...



```

1 SELECT Agency, City
2 FROM `VehiclesCount_1.VehiclesCount`
3 WHERE City = 'Los Angeles';
4

```

Row	Agency	City
1	City of Los Angeles, dba: City of Los Angeles Department of Transportation	Los Angeles
2	Los Angeles County Metropolitan Transportation Authority . dba: Metro	Los Angeles
3	Southern California Regional R...	Los Angeles



### Count of Agencies by State



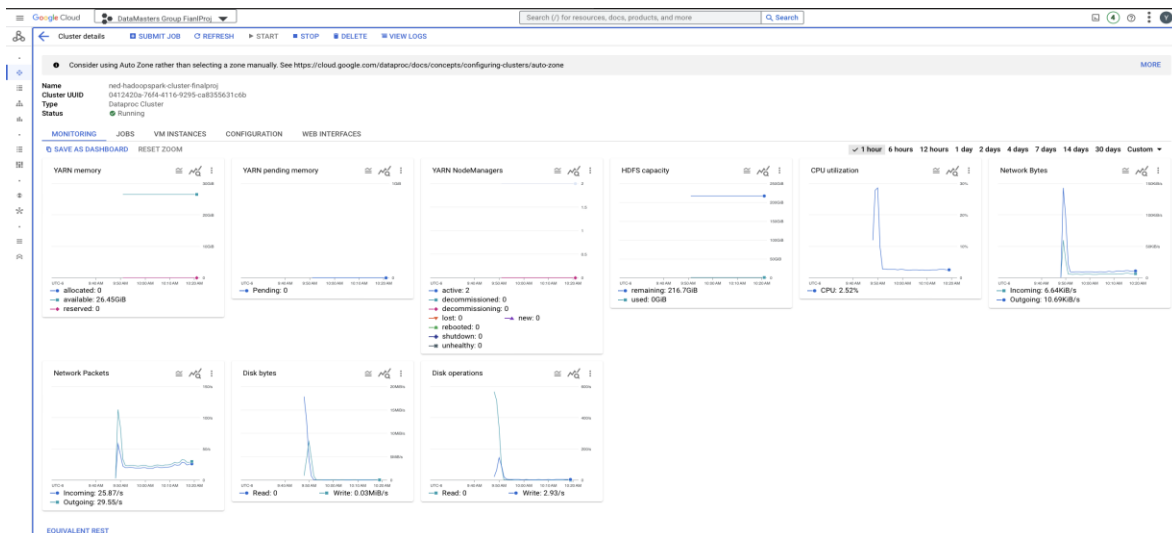
```

1 SELECT State, COUNT(*) as AgencyCount
2 FROM `VehiclesCount_1.VehiclesCount`
3 GROUP BY State
4 ORDER BY AgencyCount DESC
5 LIMIT 5;
6

```

Row	State	AgencyCount
1	CA	218
2	MI	197
3	WI	163
4	OH	153
5	NC	101

After this we can monitor the cluster and it's performance if needed.



## b. Hive (Managing and querying structured data for detailed analysis, like vehicle utilization rates.)

- Accessing Hive through SSH on the created Dataproc cluster.
- Use the Hive command line to execute Hive queries.
- Below are the 3 example's for both the datasets.

**We use the following Hive command:** `beeline -u jdbc:hive2://localhost:10000`

```

- t@ntd-hadoopsark-cluster-finalproj-m:~$ beeline -u jdbc:hive2://localhost:10000
Connecting to jdbc:hive2://localhost:10000
Connected to: Apache Hive (version 3.1.3)
Driver: Hive JDBC (version 3.1.3)
Transaction isolation: TRANSACTION_REPEATABLE_READ
Beeline version 3.1.3 by Apache Hive
0: jdbc:hive2://localhost:10000>

```

## Queries & Output in Hive: -

### 3 Hive Queries for Dataset-1 :

```

0: jdbc:hive2://localhost:10000> CREATE EXTERNAL TABLE IF NOT EXISTS ntd_vehicles
.....> (ntdvehicles string)
.....> ROW FORMAT DELIMITED
.....> STORED AS TEXTFILE
.....> LOCATION '/user/yashwanthjilla unt/data/NTD_VehiclesCount/';
INFO : Compiling command(queryId=hive_20231203194153_d45e8c8b-b373-4777-9ce6-e6aa8619ed66): CREATE EXTERNAL TABLE IF NOT EXISTS ntd_vehicles
(ntdvehicles string)
ROW FORMAT DELIMITED
STORED AS TEXTFILE
LOCATION '/user/yashwanthjilla unt/data/NTD_VehiclesCount/'
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldSchemas:null, properties:null)
INFO : Completed compiling command(queryId=hive_20231203194153_d45e8c8b-b373-4777-9ce6-e6aa8619ed66); Time taken: 0.216 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20231203194153_d45e8c8b-b373-4777-9ce6-e6aa8619ed66): CREATE EXTERNAL TABLE IF NOT EXISTS ntd_vehicles
(ntdvehicles string)
ROW FORMAT DELIMITED
STORED AS TEXTFILE
LOCATION '/user/yashwanthjilla unt/data/NTD_VehiclesCount/'
INFO : Starting task (Stage=0DDDL) in serial mode
INFO : Completed executing command(queryId=hive_20231203194153_d45e8c8b-b373-4777-9ce6-e6aa8619ed66); Time taken: 1.877 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager
No rows affected (2.787 seconds)
0: jdbc:hive2://localhost:10000>

```

```

0: jdbc:hive2://localhost:10000>
0: jdbc:hive2://localhost:10000> SELECT * FROM ntd_vehicles LIMIT 1;
INFO : Compiling command(queryId=hive_20231203194616_902bd8b8-a2fb-402c-a8cf-59aa96d0be04): SELECT * FROM ntd_vehicles LIMIT 1
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:ntd_vehicles.ntdvehicles, type:string, comment:null)], properties:null)
INFO : Completed compiling command(queryId=hive_20231203194616_902bd8b8-a2fb-402c-a8cf-59aa96d0be04); Time taken: 2.323 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20231203194616_902bd8b8-a2fb-402c-a8cf-59aa96d0be04): SELECT * FROM ntd_vehicles LIMIT 1
INFO : Query ID = hive_20231203194616_902bd8b8-a2fb-402c-a8cf-59aa96d0be04
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task (Stage=1MAPRED) in serial mode
INFO : Subscribed to counters: {} for queryId: hive_20231203194616_902bd8b8-a2fb-402c-a8cf-59aa96d0be04
INFO : Tex session hasn't been created yet. Opening session
INFO : Dag name: SELECT * FROM ntd_vehicles LIMIT 1 (Stage=
INFO : Status: Running (Executing on YARN cluster with App id application_1701629121542_0001)

INFO : Completed executing command(queryId=hive_20231203194616_902bd8b8-a2fb-402c-a8cf-59aa96d0be04); Time taken: 22.704 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager
-----
VERTICES      MODE      STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
-----
Map 1 ..... container SUCCEEDED      1      1      0      0      0      0
-----
VERTICES: 01/01 [=====>>>] 100% ELAPSED TIME: 12.94 s
-----
+-----+
|          ntd_vehicles.ntdvehicles          |
+-----+
| Agency, City, State, NTD ID, Organization Type, Reporter Type, UACE Code, UZA Name, Primary UZA Population, Agency VOMS, Bus, Bus with ULB Reported, Bus >= ULB, Articulated Bus, Articulated Bu
s with ULB Reported, Articulated Bus >= ULB, Over-the-Road Bus, Over-the-Road Bus with ULB Reported, Over-the-Road Bus >= ULB, Double Decker Bus, Double Decker Bus with ULB Reported, Doub
le Decker Bus >= ULB, School Bus, School Bus with ULB Reported, School Bus >= ULB, Van, Van with ULB Reported, Van >= ULB, Cutaway, Cutaway with ULB Reported, Cutaway >= ULB, Automobile, Autom
obile with ULB Reported, Automobile >= ULB, Minivan, Minivan with ULB Reported, Minivan >= ULB, Sport Utility Vehicle, Sport Utility Vehicle with ULB Reported, Sport Utility Vehicle >= UL
B, Trolleybus, Trolleybus with ULB Reported, Trolleybus >= ULB, Heavy Rail Passenger Car, Heavy Rail Passenger Car with ULB Reported, Heavy Rail Passenger Car >= ULB, Light Rail Vehicle, L
ight Rail Vehicle with ULB Reported, Light Rail Vehicle >= ULB, Commuter Rail Passenger Coach, Commuter Rail Passenger Coach with ULB Reported, Commuter Rail Passenger Coach >= ULB, Com
muter Rail Self-Propelled Passenger Car, Commuter Rail Self-Propelled Passenger Car with ULB Reported, Commuter Rail Self-Propelled Passenger Car >= ULB, Locomotive, Locomotive with UL
B Reported, Locomotive >= ULB, Automated Guideway Vehicle, Automated Guideway Vehicle with ULB Reported, Automated Guideway Vehicle >= ULB, Vintage/Historic Trolley, Vintage/Historic Tro
lley with ULB Reported, Vintage/Historic Trolley >= ULB, Streetcar, Streetcar with ULB Reported, Streetcar >= ULB, Aerial Tramway, Aerial Tramway with ULB Reported, Aerial Tramway >= ULB,
Monorail, Monorail >= ULB, Cable Car, Cable Car with ULB Reported, Cable Car >= ULB, Inclined Plane, Inclined Plane with ULB Reported, Inclined Plane >= ULB, Ferryboat, Ferryboat with ULB R
eported, Ferryboat >= ULB, Other, Other with ULB Reported, Other >= ULB, Total Revenue Vehicles, Total with ULB Reported, Total Revenue Vehicles >= ULB, Automobiles, Automobiles >= ULB, Truc
ks and Other Rubber Tire Vehicles, Trucks and Other Rubber Tire Vehicles >= ULB, Steel Wheel Vehicles, Steel Wheel Vehicles >= ULB, Total Service Vehicles, Total Service Vehicles >= ULB
|
+-----+
1 row selected (33.105 seconds)
0: jdbc:hive2://localhost:10000>

```





```

0: jdbc:hive2://localhost:10000> SELECT AVG(LENGTH(ntdvehiclesAge)) AS avg_length FROM ntd_vehiclesAge;
INFO : Compiling command(queryId=hive_20231203201239_d017061f-c6d6-4f26-a527-a2933781ff67): SELECT AVG(LENGTH(ntdvehiclesAge)) AS avg_length FROM ntd_vehiclesAge
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Semantic Analysis Completed (retry = false)
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:avg_length, type:double, comment:null)], properties:null)
INFO : Completed compiling command(queryId=hive_20231203201239_d017061f-c6d6-4f26-a527-a2933781ff67); Time taken: 0.214 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20231203201239_d017061f-c6d6-4f26-a527-a2933781ff67): SELECT AVG(LENGTH(ntdvehiclesAge)) AS avg_length FROM ntd_vehiclesAge
INFO : Query ID = hive_20231203201239_d017061f-c6d6-4f26-a527-a2933781ff67
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Subscribed to counters: [] for queryId: hive_20231203201239_d017061f-c6d6-4f26-a527-a2933781ff67
INFO : Session is already open
INFO : Dag name: SELECT AVG(LENGTH(ntdvehic...ntd_vehiclesAge (Stage-1)
INFO : Status: Running (Executing on YARN cluster with App id application_1701629121542_0002)

-----
VERTICES      MODE      STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    1          1          0          0          0          0
Reducer 2 ..... container  SUCCEEDED    1          1          0          0          0          0
-----
VERTICES: 02/02 [=====] 100% ELAPSED TIME: 5.41 s
-----
INFO : Completed executing command(queryId=hive_20231203201239_d017061f-c6d6-4f26-a527-a2933781ff67); Time taken: 6.2 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager
+-----+
|      avg_length      |
+-----+
| 201.6211125158028    |
+-----+
1 row selected (6.467 seconds)
0: jdbc:hive2://localhost:10000>

```

**c. Spark (Advanced analytics, possibly including predictive modeling for future fleet requirements.)**

- We used the Spark shell for data processing and analysis.

```

yashwanthjilla_unt@ned-hadoopsark-cluster-finalproj-m:~$ spark-sql
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
ivysettings.xml file not found in HIVE_HOME or HIVE_CONF_DIR,/etc/hive/conf.dist/ivysettings.xml will be used
23/12/03 20:17:24 INFO org.apache.spark.SparkEnv: Registering MapOutputTracker
23/12/03 20:17:24 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
23/12/03 20:17:25 INFO org.apache.spark.SparkEnv: Registering BlockManagerMasterHeartbeat
23/12/03 20:17:25 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
Spark master: yarn, Application Id: application_1701629121542_0003
spark-sql> show tables;
default ntd_vehicles      false
default ntd_vehiclesage   false
Time taken: 3.862 seconds, Fetched 2 row(s)
spark-sql>

```

For Dataset -1 (3 Queries & output's in Spark)

```

== SQL ==
CREATE EXTERNAL TABLE IF NOT EXISTS ntd_vehicles
(ntdvehicles string) /data/NTD_VehiclesCount/'
-----^^^
spark-sql> CREATE EXTERNAL TABLE IF NOT EXISTS ntd_vehicles
> (
>   ntdvehicles STRING
> )
> ROW FORMAT DELIMITED
> STORED AS TEXTFILE
> LOCATION '/user/_/NTD_VehiclesCount/';
23/12/03 20:22:10 WARN org.apache.hadoop.hive.q1.session.SessionState: M
thorizerFactory.
Time taken: 0.575 seconds
spark-sql>

```

```
spark-sql>
> SELECT * FROM ntd_vehicles LIMIT 1;
Agency, City, State, NTD ID, Organization Type, Reporter Type, UACE Code, UZA N
with ULB Reported, Articulated Bus >= ULB, Over-the-Road Bus, Over-the-Road
Decker Bus >= ULB, School Bus, SchoolBus with ULB Reported, School Bus >=
ile with ULB Reported, Automobile >= ULB, Minivan, Minivan with ULB Reporte
Trolleybus, TrolleyBus with ULB Reported, Trolleybus >= ULB, Heavy Rail Pas
ht Rail Vehicle with ULB Reported, Light Rail Vehicle >= ULB, Commuter Rail
ter Rail Self-Propelled Passenger Car, Commuter Rail Self-Propelled Passe
Reported, Locomotive >= ULB, Automated Guideway Vehicle, Automated Guideway
ey with ULB Reported, Vintage/Historic Trolley >= ULB, Streetcar, Streetcar
norrail, Monorail >= ULB, Cable Car, Cable Car with ULB Reported, Cable Car >
orted, Ferryboat >= ULB, Other, Other with ULB Reported, Other >= ULB, Total
and Other Rubber Tire Vehicles, Trucks and Other Rubber Tire Vehicles >=
Time taken: 5.717 seconds, Fetched 1 row(s)
spark-sql>
```

[illegible]

For Dataset-2 (3 Queries & outputs in Spark)

```
LOCATION '/user/yashwanthjilla_unt/data/NTD_VehiclesAge/'

spark-sql> CREATE EXTERNAL TABLE IF NOT EXISTS ntd_vehiclesAge
  (
    > ntdvehiclesAge STRING
    > )
    > ROW FORMAT DELIMITED
    > STORED AS TEXTFILE
    > LOCATION '/user/yashwanthjilla_unt/data/NTD_VehiclesAge/';
Time taken: 0.12 seconds
spark-sql> SELECT COUNT(*) FROM ntd_vehiclesAge;
23/12/03 20:27:09 WARN org.apache.hadoop.util.concurrent.ExecutorHelper: Thread (Thread[GetFileInfo #1,5,main]) interrupted:
java.lang.InterruptedException
    at com.google.common.util.concurrent.AbstractFuture.get(AbstractFuture.java:510)
    at com.google.common.util.concurrent.PluentFuture$TrustedFuture.get(PluentFuture.java:88)
    at org.apache.hadoop.util.concurrent.ExecutorHelper.logThrowableFromAfterExecute(ExecutorHelper.java:48)
    at org.apache.hadoop.util.concurrent.HadoopThreadPoolExecutor.afterExecute(HadoopThreadPoolExecutor.java:90)
    at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1157)
    at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:624)
    at java.lang.Thread.run(Thread.java:750)

7910
Time taken: 1.706 seconds, Fetched 1 row(s)
spark-sql>
```

```
spark-sql> CREATE EXTERNAL TABLE IF NOT EXISTS ntd_vehiclesAge
> (
>   ntdvehiclesAge STRING
> )
> ROW FORMAT DELIMITED
> STORED AS TEXTFILE
> LOCATION '/user/[REDACTED]data/NTD_VehiclesAge/';
Time taken: 0.12 seconds
spark-sql>
```



```

spark-sql> SELECT AVG(LENGTH(ntdvehiclesAge)) AS avg_length FROM ntd_vehiclesAge;
201.6211125158028
Time taken: 0.923 seconds, Fetched 1 row(s)
spark-sql>

```

### Comparison of Hive & Spark Execution Time: -

Dataset	Query Description	Hive Execution Time (sec)	Spark Execution Time (sec)
Dataset-1	CREATE EXTERNAL TABLE	2.787	0.575
Dataset-1	Selecting 1 row	33.105	5.717
Dataset-1	Selecting 5 rows	11.205	0.339
Dataset-2	CREATE EXTERNAL TABLE in VehiclesAge	0.145	0.12
Dataset-2	Select Count all from Vehicle Age	17.672	1.706
Dataset-2	Select Average length	6.467	0.923

### Key- Analysis: -

- Hive is suitable for traditional data warehousing and SQL queries.
- Spark's in-memory processing makes it a superior choice for tasks involving iterative algorithms and real-time analytics.
- In comparing Hive and Spark based on the provided data and queries, Spark consistently outperforms Hive with faster execution times, particularly for data processing and analytics.
- Developing robust data management through complex queries, deeper multidimensional analysis

### Conclusion:

Data experts can uncover patterns and relationships to better understand transit usage. This allows for data-driven forecasting and infrastructure planning to efficiently address current problems and proactively meet future needs.

- Predict optimal accessible stop locations via modeling.
- Forecast peak demand periods with predictive analytics.

- Collect and analyze emissions, mileage, fuel efficiency data.
- Model route optimization scenarios to increase sustainability.

**References:**

1. <https://spark.apache.org/>
2. <https://hive.apache.org/>
3. <https://cloud.google.com/>
4. <https://www.transit.dot.gov/ntd/data-product/2022-vehicles>