



Dr. Vishwanath Karad

**MIT WORLD PEACE
UNIVERSITY | PUNE**

TECHNOLOGY, RESEARCH, SOCIAL INNOVATION & PARTNERSHIPS

Mini Project Report

on

“YOUTUBE DATA ANALYSIS”

Submitted by

Atharv Nikam 1032180297 PA10

Ketaki Patil 1032180421 PA17

Saisrijan Gupta 1032180626 PA27

Anjali Nair 1032180751 PA31

Shivraj Patil 1032181709 PA51

Under the Guidance of

Prof. Vasundhara Ghate

At

School of Computer Engineering and Technology

ABSTRACT

There is a tremendous growth and popularity of YouTube. It has the potential to touch billions of lives globally as the no. of YouTube users is growing day by day. YouTube, owned by Google, a video streaming site which has billions of users and 400 hours of videos are being uploaded every minute. Almost billions of videos are watched on YouTube every single day, generating a mammoth amount of data daily. Since YouTube data is generally in unstructured form, there is an increased demand to store, process and analyze such real time Big Data. YouTubers can analyze their own channel performance with YouTube Analytics. But one can not analyze other channels. The proposed system uses the MapReduce framework of Hadoop for processing and analyzing real time YouTube datasets. It will help in discovering how competitors are performing on YouTube. One can easily identify what content works best on YouTube. These analytical data can be represented in demographic form which can be used by individuals and organizations for making immediate actionable decisions so as to gain competitive advantages.

CONTENTS

1. Introduction.....	5
1.1 Motivation	5
2. Problem Statement.....	6
2.1 Problem Definition.....	6
2.2 Objectives	6
3. Tools.....	7
3.1 Hadoop.....	7
3.1.1 Modules of Hadoop.....	7
3.2 Map Reduce.....	8
3.2.1 Steps in Map Reduce.....	8
3.3 Node Js.....	9
3.4 Cassandra.....	10
3.4.1 Data Model of Cassandra.....	11
3.4.2 Cassandra Query Language.....	11
3.5 Power BI.....	12
4. Dataset.....	13
5. System Architecture.....	14
6. Output.....	15
7. Visualization screenshots.....	22
8. Conclusion.....	24
9. References in IEEE Format.....	25

LIST OF FIGURES

Sr.No.	Figure	Page No
1	HDFS Architecture	7
2.	Map Reduce	8
3.	Features of Power BI	12
4	CSV Dataset	13
5	System Architecture of Youtube Data Analysis	14
6	Front-End UI	15
7	Data Extraction using Node js	15
8	Categories-wise Videos	16
9	Videos Uploaded by Channel	16
10	Views received by each video	17
11	Table Creation and Dataset Import	17
12	Displaying all records	18
13	Video id and view count for gaming	18
14	Delete one record from gaming	19
15	Maximum likes obtained by any video	19
16	Minimum likes obtained by any video	19
17	Homepage	20
18	Visualizations embedded in a webpage	20
19	Hadoop MapReduce Output embedded in a webpage	21
20	Visualizations filtered by education category	22
21	PoweBI Dashboard	22
22	Key Influencers	23

1. INTRODUCTION

YouTube, owned by Google, is the platform where people can upload, watch and share videos with others. They can also give their feedback accordingly through comments. YouTube has attracted users across the globe such as advertisers, various media, politicians and other interest groups. 400 hours of videos are uploaded on YouTube every minute. Billions of videos are watched on YouTube every day. This indicates that YouTube has the great potential to reach out to many people. Thus generating massive amounts of data daily. Such a mammoth amount of data is usually in unstructured form. Proposed system presents analytics of YouTube datasets in demographic form using the Apache Hadoop MapReduce framework. It will help its users make targeted real time and informed decisions.

1.1 MOTIVATION

YouTube is one of the most popular and engaging social media tools for uploading, viewing videos and an amazing platform that reveals the users response through comments for published videos, number of likes, dislikes, number of subscribers for a particular channel. YouTube collects a wide variety of traditional data points including View Counts, Likes, and Comments. The analysis of the above listed data points makes a very interesting data source to extract implicit knowledge about users, videos, categories and community interests. In Today's world as the 4 V's of Big data (Volume, Variety, Velocity & Veracity) are very rapidly increasing it has become a must to come up with meaningful insights in order to project and derive meanings from the data to help organizations grow business rapidly. To process this large amount of data, a high-level parallel and distributed system such as Hadoop is necessary. The data obtained from YouTube can be used for analysis and various visualizations. The information generated by this analysis can be used by multiple organizations to make key decisions about their services, targets to improve the business value.

2. PROBLEM STATEMENT

2.1 PROBLEM DEFINITION

"YouTube has over a billion users and every day people watch hundreds of millions of hours on YouTube and generate billions of views". Every day, people across the world are uploading 1.2 million videos to YouTube, or over 100 hours per minute and this number is ever increasing. To analyze and understand the activity occurring on such a massive scale, a relational SQL database is not enough. Such kind of data is well suited to a massively parallel and distributed system like Hadoop.

2.2 OBJECTIVES

The main objective of this project is to focus on how data generated from YouTube can be mined and utilized by different companies to make targeted, real time and informed decisions about their product that can increase their market share.

3.TOOLS

3.1 Hadoop

Hadoop is an Apache open source framework written in java that allows distributed processing of large datasets across clusters of computers using simple programming models. The Hadoop framework application works in an environment that provides distributed storage and computation across clusters of computers. Hadoop is designed to scale up from a single server to thousands of machines, each offering local computation and storage.

3.1.1 Modules of Hadoop

1. **HDFS:** Hadoop Distributed File System. Google published its paper GFS and on the basis of that HDFS was developed. It states that the files will be broken into blocks and stored in nodes over the distributed architecture.
2. **Yarn:** Yet another Resource Negotiator is used for job scheduling and managing the cluster.
3. **Map Reduce:** This is a framework which helps Java programs to do the parallel computation on data using key value pairs. The Map task takes input data and converts it into a data set which can be computed in Key value pairs. The output of Map task is consumed by reduce task and then the output of reducer gives the desired result.
4. **Hadoop Common:** These Java libraries are used to start Hadoop and are used by other Hadoop modules.

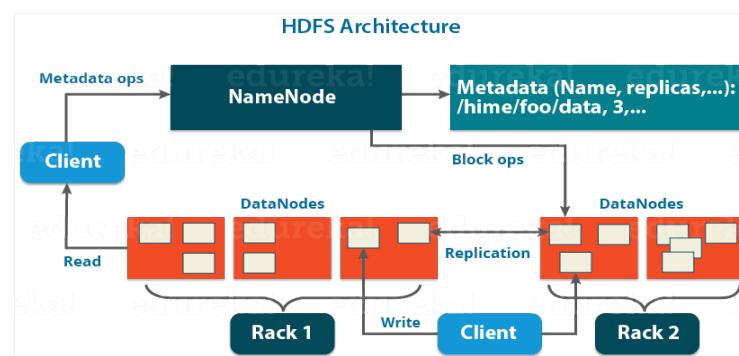


Fig 1.:HDFS Architecture

3.2 Map Reduce

The MapReduce is a paradigm which has two phases, the mapper phase, and the reducer phase. In the Mapper, the input is given in the form of a key-value pair. The output of the Mapper is fed to the reducer as input. The reducer runs only after the Mapper is over. The reducer too takes input in key-value format, and the output of the reducer is the final output.

3.2.1 Steps in Map Reduce

1. The map takes data in the form of pairs and returns a list of <key, value> pairs. The keys will not be unique in this case.
2. Using the output of Map, sort and shuffle are applied by the Hadoop architecture. This sort and shuffle acts on these lists of <key, value> pairs and sends out unique keys and a list of values associated with this unique key <key, list(values)>.
3. An output of sort and shuffle sent to the reducer phase. The reducer performs a defined function on a list of values for unique keys, and Final output <key, value> will be stored/displayed.

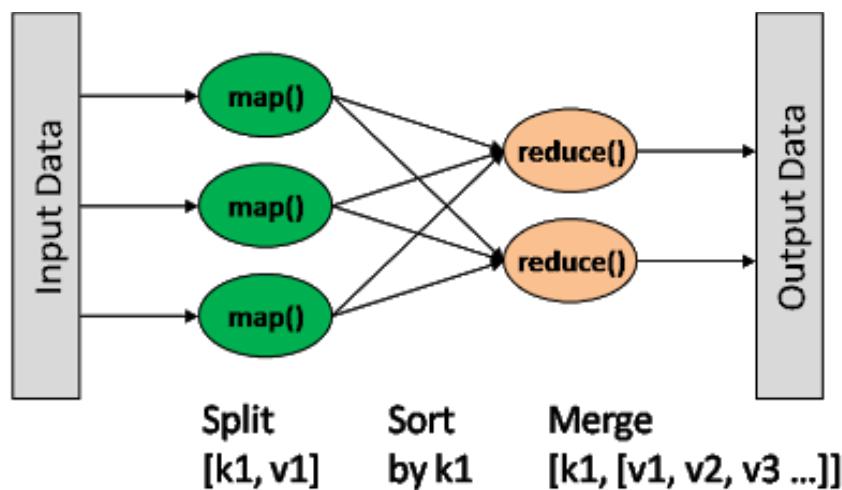


Fig 2.: Map Reduce

3.3 Node.js

Node.js is a cross-platform runtime environment and library for running JavaScript applications outside the browser. It is used for creating server-side and networking web applications. It is open source and free to use. Many of the basic modules of Node.js are written in JavaScript. Node.js is mostly used to run real-time server applications.

Features of Node.js

Following is a list of some important features of Node.js that makes it the first choice of software architects.

1. **Extremely fast:** Node.js is built on Google Chrome's V8 JavaScript Engine, so its library is very fast in code execution.
2. **I/O is Asynchronous and Event Driven:** All APIs of Node.js library are asynchronous i.e. non-blocking. So a Node.js based server never waits for an API to return data. The server moves to the next API after calling it and a notification mechanism of Events of Node.js helps the server to get a response from the previous API call. It is also a reason that it is very fast.
3. **Single threaded:** Node.js follows a single threaded model with event looping.
4. **Highly Scalable:** Node.js is highly scalable because the event mechanism helps the server to respond in a non-blocking way.
5. **No buffering:** Node.js cuts down the overall processing time while uploading audio and video files. Node.js applications never buffer any data. These applications simply output the data in chunks.
6. **Open source:** Node.js has an open source community which has produced many excellent modules to add additional capabilities to Node.js applications.

3.4 Cassandra

Apache Cassandra is an open source, distributed and decentralized/distributed storage system (database), for managing very large amounts of structured data spread out across the world. It provides highly available service with no single point of failure.

Features of Cassandra

Cassandra has become so popular because of its outstanding technical features. Given below are some of the features of Cassandra:

1. **Elastic scalability** – Cassandra is highly scalable; it allows to add more hardware to accommodate more customers and more data as per requirement.
2. **Always on architecture** – Cassandra has no single point of failure and it is continuously available for business-critical applications that cannot afford a failure.
3. **Fast linear-scale performance** – Cassandra is linearly scalable, i.e., it increases your throughput as you increase the number of nodes in the cluster. Therefore it maintains a quick response time.
4. **Flexible data storage** – Cassandra accommodates all possible data formats including: structured, semi-structured, and unstructured. It can dynamically accommodate changes to your data structures according to your needs.
5. **Easy data distribution** – Cassandra provides the flexibility to distribute data where you need by replicating data across multiple data centers.
6. **Transaction support** – Cassandra supports properties like Atomicity, Consistency, Isolation, and Durability (ACID).
7. **Fast writes** – Cassandra was designed to run on cheap commodity hardware. It performs blazingly fast writes and can store hundreds of terabytes of data, without sacrificing the read efficiency.

3.4.1 Data Model of Cassandra

1. Cluster

Cassandra database is distributed over several machines that operate together. The outermost container is known as the Cluster. For failure handling, every node contains a replica, and in case of a failure, the replica takes charge. Cassandra arranges the nodes in a cluster, in a ring format, and assigns data to them.

2. Keyspaces

Cassandra data model consists of keyspaces at the highest level. Keyspaces are the containers of data, similar to the schema or database in a relational database. Typically, keyspaces contain many tables.

3. Column Family

A column family is a container for an ordered collection of rows. Each row, in turn, is an ordered collection of columns.

4. Table

Within the keyspaces, the tables are defined. Tables are also referred to as Column Families in the earlier versions of Cassandra. Tables contain a set of columns and a primary key, and they store data in a set of rows.

3.4.2 Cassandra Query Language

Cassandra Query Language (CQL) is used to access Cassandra through its nodes. CQL treats the database (Keyspace) as a container of tables. Programmers use cqlsh: a prompt to work with CQL or separate application language drivers. The client can approach any of the nodes for their read-write operations. That node (coordinator) plays a proxy between the client and the nodes holding the data.

3.5 Power BI

Microsoft Power BI is a business intelligence platform that provides nontechnical business users with tools for aggregating, analyzing, visualizing and sharing data. Power BI's user interface is fairly intuitive for users familiar with Excel and its deep integration with other Microsoft products makes it a very versatile self-service tool that requires little upfront training.

Microsoft Power BI is used to find insights within data. Power BI can help connect disparate data sets, transform and clean the data into a data model and create charts or graphs to provide visuals of the data. All of this can be shared with other Power BI users.

The data models created from Power BI can be used in several ways for organizations, including telling stories through charts and data visualizations and examining "what if" scenarios within the data. Power BI reports can also answer questions in real time and help with forecasting to make sure departments meet business metrics.

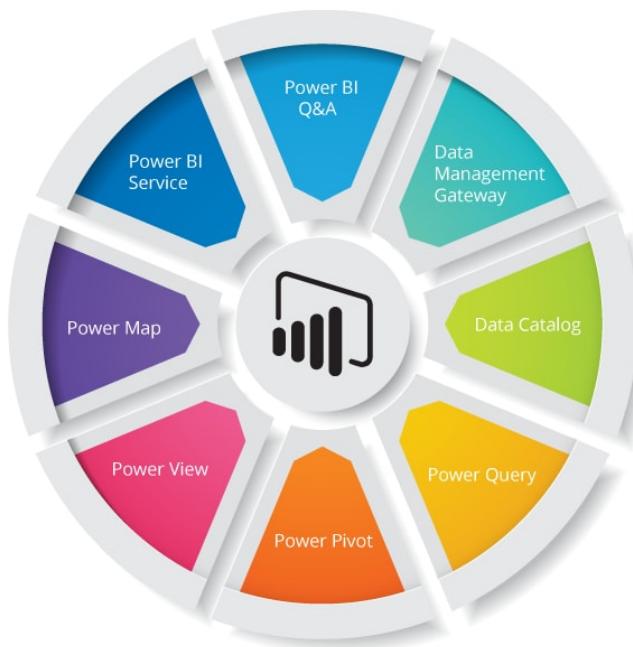


Fig 3.: Features of Power BI

4. DATASET

A	B	C	D	E	F	G	H	I	J	K	L	M
1 title	videoId	categoryId	Category	Duration	viewCount	commentCount	channelId	likeCount	dislikeCount	language	publishedAt	
2 Lordi - Lebt denn der alte Holzmichl noch?	cYe12DijydA	10	Music	PT3M22S	155694	68	UCPp1Vv3fMb8aoZzjvAVRslg	236	112	en-GB	2008-11-02T21:02:28Z	
3 De Randfichten- Steig E!Mir Fahrn in der Tschechei GnhydAGvkE	VfnbdBt3nGc	10	Music	PT3M38S	907	0	UCPp1Vv3fMb8aoZzjvAVRslg	15	0	en-GB	2020-09-12T13:45:16Z	
4 De Spackfettbemm		10	Music	PT3M41S	128320	6	UCPp1Vv3fMb8aoZzjvAVRslg	691	47	en-GB	2014-12-25T17:29:35Z	
5 De Randfichten - Du Kleine Fliege	4IS3qWAHxN4	10	Music	PT1M35S	313220	45	UCPp1Vv3fMb8aoZzjvAVRslg	870	107	en-GB	2012-01-09T12:49:24Z	
6 De Randfichten - Das kommt vom Rudern 2006	dI68_6PDLoFw	10	Music	PT3M2S	148745	14	UCPp1Vv3fMb8aoZzjvAVRslg	195	9	en-GB	2007-12-07T22:25:27Z	
7 GrÄlliss Gott Mei Arzgeberg	7INpjcrm2M	10	Music	PT3M4S	48335	13	UCPp1Vv3fMb8aoZzjvAVRslg	213	20	en-GB	2014-11-09T01:45:56Z	
8 De Randfichten - s Alte Kanapee	eQIPtQjOH6c	10	Music	PT2M38S	111672	4	UCPp1Vv3fMb8aoZzjvAVRslg	350	27	en-GB	2012-10-05T20:31:11Z	
9 Heja Heja Ho - De Randfichten Sel Do (Live From GY8HA3WcB4		10	Music	PT5M18S	28532	1	UCPp1Vv3fMb8aoZzjvAVRslg	154	9	en-GB	2015-07-08T02:46:34Z	
10 Dr Holzmichl (Live)	diiZ5IMhjyK	10	Music	PT7M57S	33756	5	UCPp1Vv3fMb8aoZzjvAVRslg	180	17	en-GB	2014-11-08T18:43:06Z	
11 Jetzt geht die Party richtig los (Remix)	MD-mLYCK_kU	10	Music	PT3M32S	3459	1	UCPp1Vv3fMb8aoZzjvAVRslg	27	7	en-GB	2015-07-09T13:42:24Z	
12 De Lustigen Holzhackerlein	TVxuzjjT21	10	Music	PT2M32S	773	0	UCPp1Vv3fMb8aoZzjvAVRslg	5	1	en-GB	2018-08-16T14:10:18Z	
13 De Randfichten - Rups am Grill 2010	DWBBOUlkHls	10	Music	PT3M8S	156983	54	UCPp1Vv3fMb8aoZzjvAVRslg	687	123	en-GB	2010-08-09T19:23:18Z	
14 Do pfeift dr Fuchs	3x3f6fQ9QZ4o	10	Music	PT2M57S	2799	1	UCPp1Vv3fMb8aoZzjvAVRslg	13	2	en-GB	2018-07-06T02:11:07Z	
15 Mei Schiene Laadergack	TQNriidgTdg	10	Music	PT3M30S	22849	5	UCPp1Vv3fMb8aoZzjvAVRslg	134	15	en-GB	2014-12-25T17:23:45Z	
16 De Randfichten - Wir feiern heut ne Party 2012	8DXE8SBWhJU	10	Music	PT2M53S	12548	2	UCPp1Vv3fMb8aoZzjvAVRslg	45	3	en-GB	2012-07-04T04:14:52Z	
17 De Randfichten - Ich bin gut drauf Offizielles Musik_NBK_ci0ds8	PT3M6S	10	Music	PT3M6S	61063	22	UCPp1Vv3fMb8aoZzjvAVRslg	236	28	en-GB	2015-09-11T16:00:01Z	
18 De Randfichten - Die alten Zeiten sind vorbei 2011	BbvW7JK0To	10	Music	PT3M37S	65769	13	UCPp1Vv3fMb8aoZzjvAVRslg	321	13	en-GB	2011-06-05T14:06:27Z	
19 Ich hatt heit Nacht n Traam	9djKnaad_c	10	Music	PT3M27S	4583	0	UCPp1Vv3fMb8aoZzjvAVRslg	31	1	en-GB	2018-07-04T23:34:06Z	
20 Sauerkraut	m58vK6Giv8A	10	Music	PT3M40S	5339	0	UCPp1Vv3fMb8aoZzjvAVRslg	30	5	en-GB	2018-05-12T03:33:45Z	
21 De Randfichten - Kaan Staapilz find ich net 2008	VALUF3h8VPok	10	Music	PT2M40S	51541	4	UCPp1Vv3fMb8aoZzjvAVRslg	145	8	en-GB	2009-09-24T04:53:21Z	
22 Dann ist Weihnacht	20alom6FWO	10	Music	PT3M9S	10206	0	UCPp1Vv3fMb8aoZzjvAVRslg	70	3	en-GB	2015-01-29T12:53:32Z	
23 De Randfichten - Lebt den alte Holzmichel noch - XERpMOG0C03A	10	Music	PT3M48S	3228	0	UCPp1Vv3fMb8aoZzjvAVRslg	11	5	en-GB	2016-08-16T11:55:06Z		
24 De Randfichten - Lebt denn der alte Winter-Holzn E-usnOPVOKw	10	Music	PT3M44S	78746	2	UCPp1Vv3fMb8aoZzjvAVRslg	99	8	en-GB	2008-12-01T06:46:27Z		
25 Klitscherlied	8Qf68ssIK5oc	10	Music	PT2M58S	6613	0	UCPp1Vv3fMb8aoZzjvAVRslg	46	1	en-GB	2014-11-08T23:27:15Z	
26 De Randfichten - Alles was ich brauch 2011	qFiaI5loZ5w	10	Music	PT2M40S	21973	2	UCPp1Vv3fMb8aoZzjvAVRslg	73	1	en-GB	2011-09-24T06:41:19Z	
27 Drham Is Orham	coWpWNUIhno	10	Music	PT3M48S	6191	0	UCPp1Vv3fMb8aoZzjvAVRslg	44	4	en-GB	2014-11-08T08:12:16Z	
28 De Randfichten - Hey Du lass doch bitte mein klein Ox_Op_JFM	10	Music	PT2M24S	5070	2	UCPp1Vv3fMb8aoZzjvAVRslg	29	2	en-GB	2019-02-21T15:15:13Z		
29 Nur noch mal heim - De Randfichten	4KMHY0if8s	10	Music	PT3M33S	66628	31	UCPp1Vv3fMb8aoZzjvAVRslg	209	10	en-GB	2011-12-22T19:59:40Z	
30 Hab Dank DatÅ'r	hVRxhfskNH8	10	Music	PT3M2S	3458	0	UCPp1Vv3fMb8aoZzjvAVRslg	24	2	en-GB	2014-12-25T17:30:12Z	
31 Heilig-Obnd-Lied	rpCyXQAvlu0	10	Music	PT3M10S	4922	0	UCPp1Vv3fMb8aoZzjvAVRslg	42	1	en-GB	2014-11-08T20:36:53Z	
32 Randfichten - Rups am Grill der Grill Hit	cpDhWcVPNg8	10	Music	PT3M11S	27676	0	UCPp1Vv3fMb8aoZzjvAVRslg	62	17	en-GB	2010-08-02T13:30:33Z	
33 De Randfichten - Rups am Grill der Grill Hit	6Lw18C74A4	10	Music	PT3M20S	4004	0	UCPp1Vv3fMb8aoZzjvAVRslg	47	0	en-GB	2010-08-02T13:30:33Z	

Fig.4:CSV Dataset

5. SYSTEM ARCHITECTURE

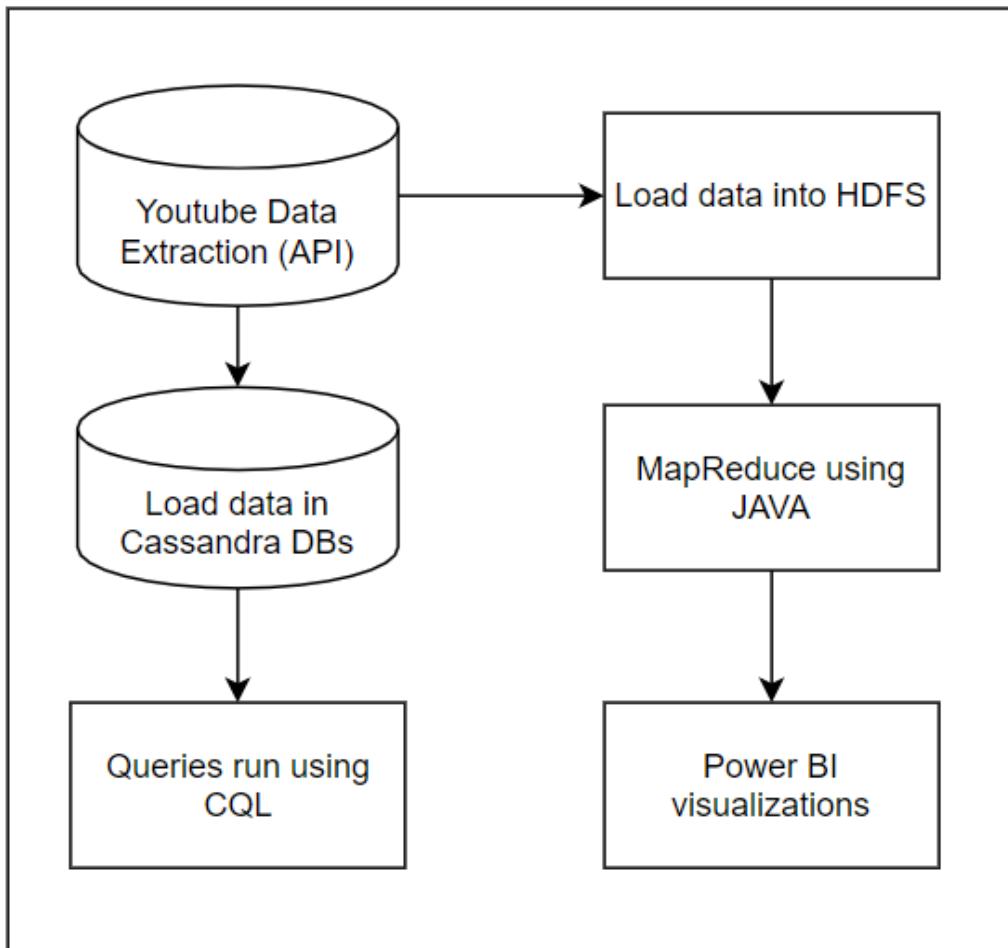


Fig.5: System Architecture for Youtube Data Analysis

6. OUTPUT ANALYSIS

6.1 Fetch Data using Node JS

- Selection of duration from front-end UI:

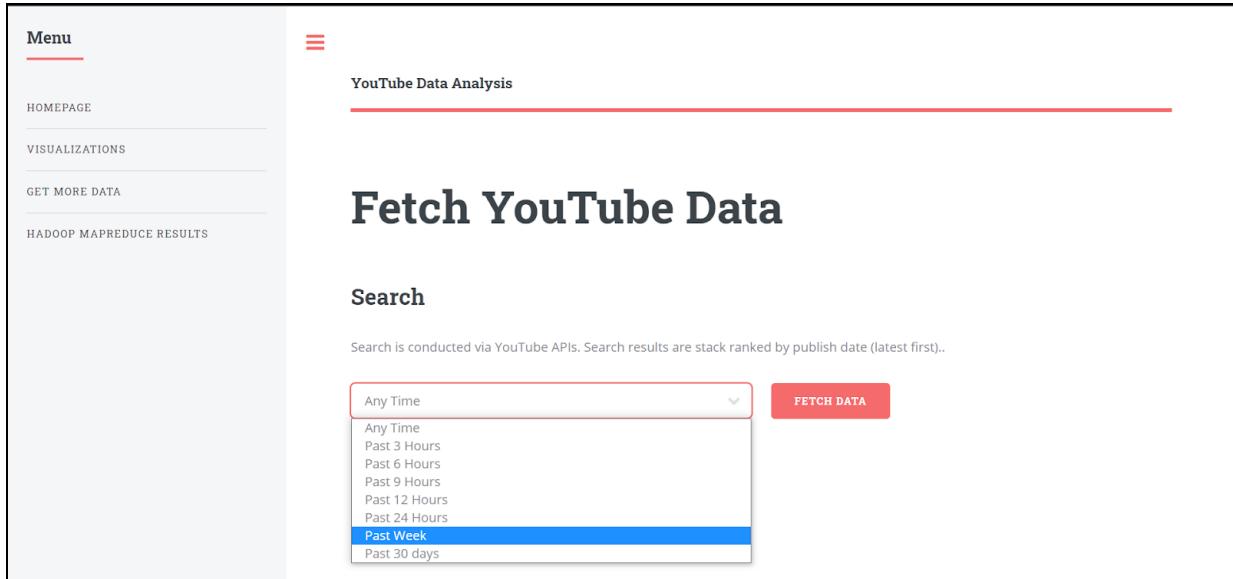


Fig.6:Front-End UI

- Command prompt output:

```

D:\Downloads\YouTube-Data-Analysis\YouTube-Data-Analysis\WebContent>node app.js
Server is running at port 8080
Connected using sockets
Next Playlist call number: 0 Total Videos Added :14
Next Playlist call number: 1 Total Videos Added :26
Next Playlist call number: 2 Total Videos Added :64
Next Playlist call number: 3 Total Videos Added :94
Next Playlist call number: 4 Total Videos Added :113
Next Playlist call number: 5 Total Videos Added :129
Next Playlist call number: 6 Total Videos Added :159
Next Playlist call number: 7 Total Videos Added :189
Next Playlist call number: 8 Total Videos Added :201
Next Playlist call number: 9 Total Videos Added :218
Next Playlist call number: 10 Total Videos Added :243
Next Playlist call number: 11 Total Videos Added :257
Next Playlist call number: 12 Total Videos Added :317
Next Playlist call number: 13 Total Videos Added :330
Next Playlist call number: 14 Total Videos Added :358
Next Playlist call number: 15 Total Videos Added :416
Next Playlist call number: 16 Total Videos Added :456
Next Playlist call number: 17 Total Videos Added :476
Next Playlist call number: 18 Total Videos Added :497
Next Playlist call number: 19 Total Videos Added :506
Next Playlist call number: 20 Total Videos Added :513

```

Fig.7:Data Extraction using Node js

6.2 Hadoop MapReduce

- Category

part-r-00000 - Notepad	
File	Edit
"Autos & Vehicles"	1923
"Comedy"	318
"Education"	1885
"Entertainment"	3592
"Film & Animation"	4025
"Gaming"	6596
"Howto & Style"	729
"Music"	8705
"News & Politics"	2741
"Nonprofits & Activism"	1096
"People & Blogs"	2856
"Pets & Animals"	30
"Science & Technology"	2121
"Sports"	1966
"Travel & Events"	301

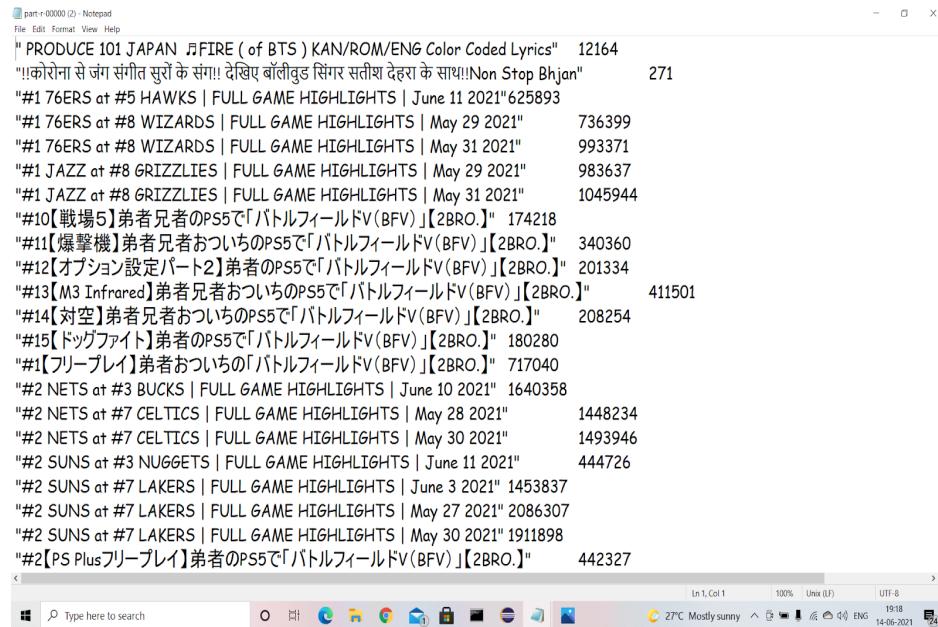
Fig.8:Categories-wise Videos

- Uploaders

part-r-00000 (3) - Notepad	
File	Edit
"UC-2Y8dQb0S6DttxNgAKoJKA"	18
"UC-BaA0SD5-4acdInXWIDRoUg"	78
"UC-CSyyi47VX1D9zyeABW3w"	29
"UC-XWpctw55Q6b_AHo8rkJgw"	54
"UC_-S&f5Zhq-MtppuOJ0Nw"	21
"UC0076UMUgEng8HORUw_MYHA"	124
"UC07-d0wgza1IguKA86jpxNA"	16
"UC0DM_mHV2u6dj8ig51GkQwg"	54
"UC0DVwAbADMlPcpVc3wiYbQ"	25
"UC0qrcKuayvbMTpj-WXNFylA"	18
"UC0yw39DOcCOwyUj2uY8vbuA"	8
"UC11JFCosIUUMTuAsUPH8RdQ"	8
"UC1BdpFCsvKmqzetThZ09M6A"	14
"UC1L2JoMpcY6MRLhf3gg5Xg"	16
"UC1XW5EjmkkW46Oic8FVncPQ"	17
"UC1mDL6TXcTrZpIG1wWLqsSQ"	25
"UC1oPBWUifc0QOOY8DEKhLuQ"	133
"UC1pfsmDBnMQB8sOuQvmTvRQ"	11
"UC24_Z2L-8K183AI9zJJzNQ"	40
"UC2D0dmLHKbIe9aACnclTLPg"	60
"UC2GuoutVye6gPUK88ILpjw"	15
"UC2L4gT_zPTtgzCkBdGdi7eNg"	19

Fig.9:Videos Uploaded by Channel

● Views



```

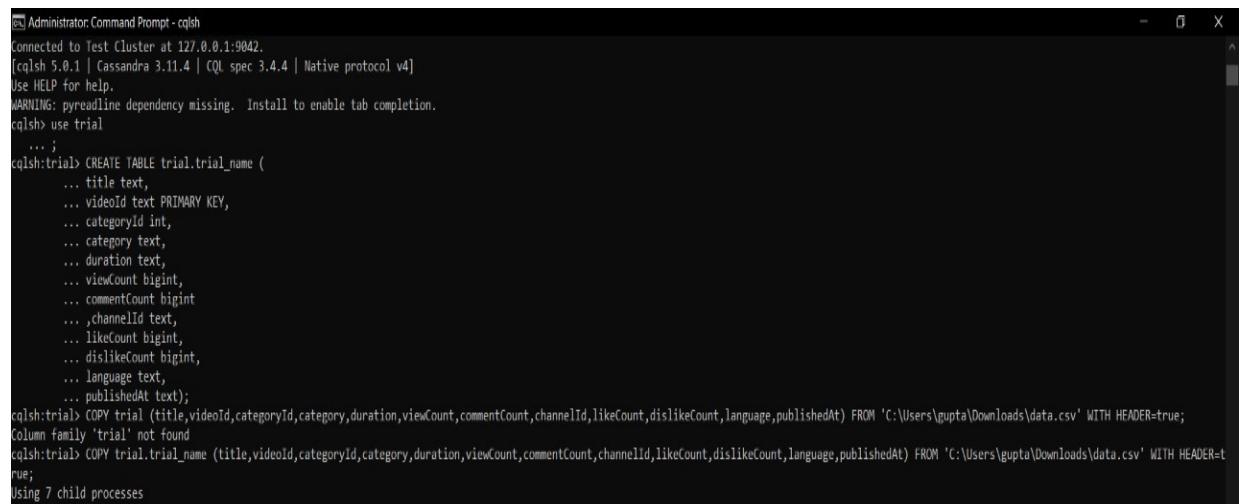
part-r-00000 (2) - Notepad
File Edit Format View Help
" PRODUCE 101 JAPAN ピュア・ラブ (of BTS) KAN/ROM/ENG Color Coded Lyrics" 12164
"!कोरोना से जंग सीति चुमो के संग!! देखिए बॉलीवुड सिंगर सरीष देहरा के साथ!Non Stop Bhajan" 271
"#1 76ERS at #5 HAWKS | FULL GAME HIGHLIGHTS | June 11 2021" 625893
"#1 76ERS at #8 WIZARDS | FULL GAME HIGHLIGHTS | May 29 2021" 736399
"#1 76ERS at #8 WIZARDS | FULL GAME HIGHLIGHTS | May 31 2021" 993371
"#1 JAZZ at #8 GRIZZLIES | FULL GAME HIGHLIGHTS | May 29 2021" 983637
"#1 JAZZ at #8 GRIZZLIES | FULL GAME HIGHLIGHTS | May 31 2021" 1045944
"#10【戦場5】弟者兄者のPS5で「バトルフィールドV(BFV)」【2BRO.】" 174218
"#11【爆撃機】弟者兄者おついちのPS5で「バトルフィールドV(BFV)」【2BRO.】" 340360
"#12【オプション設定パート2】弟者のPS5で「バトルフィールドV(BFV)」【2BRO.】" 201334
"#13【M3 Infrared】弟者兄者おついちのPS5で「バトルフィールドV(BFV)」【2BRO.】" 411501
"#14【対空】弟者兄者おついちのPS5で「バトルフィールドV(BFV)」【2BRO.】" 208254
"#15【ドッグファイト】弟者のPS5で「バトルフィールドV(BFV)」【2BRO.】" 180280
"#1【フリープレイ】弟者おついちの「バトルフィールドV(BFV)」【2BRO.】" 717040
"#2 NETS at #3 BUCKS | FULL GAME HIGHLIGHTS | June 10 2021" 1640358
"#2 NETS at #7 CELTICS | FULL GAME HIGHLIGHTS | May 28 2021" 1448234
"#2 NETS at #7 CELTICS | FULL GAME HIGHLIGHTS | May 30 2021" 1493946
"#2 SUNS at #3 NUGGETS | FULL GAME HIGHLIGHTS | June 11 2021" 444726
"#2 SUNS at #7 LAKERS | FULL GAME HIGHLIGHTS | June 3 2021" 1453837
"#2 SUNS at #7 LAKERS | FULL GAME HIGHLIGHTS | May 27 2021" 2086307
"#2 SUNS at #7 LAKERS | FULL GAME HIGHLIGHTS | May 30 2021" 1911898
"#2【PS Plusフリープレイ】弟者のPS5で「バトルフィールドV(BFV)」【2BRO.】" 442327

```

Fig.9:Views received by each video

6.3 Cassandra Queries

● Table Creation and Import of Dataset



```

Administrator:Command Prompt - cqsh
Connected to Test Cluster at 127.0.0.1:9042.
[cqlsh 5.0.1 | Cassandra 3.11.4 | CQL spec 3.4.4 | Native protocol v4]
Use HELP for help.
WARNING: pyreadline dependency missing. Install to enable tab completion.
cqlsh> use trial
...
cqlsh:trial> CREATE TABLE trial.trial_name (
...   title text,
...   videoId text PRIMARY KEY,
...   categoryId int,
...   category text,
...   duration text,
...   viewCount bigint,
...   commentCount bigint,
...   channelId text,
...   likeCount bigint,
...   dislikeCount bigint,
...   language text,
...   publishedAt text);
cqlsh:trial> COPY trial (title,videoId,categoryId,category,duration,viewCount,commentCount,channelId,likeCount,dislikeCount,language,publishedAt) FROM 'C:\Users\gupta\Downloads\data.csv' WITH HEADER=true;
Column family 'trial' not found
cqlsh:trial> COPY trial.trial_name (title,videoId,categoryId,category,duration,viewCount,commentCount,channelId,likeCount,dislikeCount,language,publishedAt) FROM 'C:\Users\gupta\Downloads\data.csv' WITH HEADER=true;
Using 7 child processes

```

Fig.10:Table Creation and Dataset Import

- Dataset imported in Cassandra DBs

Fig.11:Displaying all records

- Aggregation queries using CQL

Video id and view count from gaming category

```
cqlsh:trial> select videoId, viewCount from trial_name where category='Gaming' LIMIT 10 ALLOW FILTERING;
+-----+-----+
| videoId | viewcount |
+-----+-----+
| nYe0Nwvvrrg | 87743 |
| oyEph5fxBy4 | 589997 |
| _BG98e_w6d0 | 20200 |
| ZKGt23WHUU0 | 1204 |
| WLLUKLhVl0E | 184444 |
| c0uf82rJnI8 | 2149855 |
| WQDyvTUVaEY | 10167 |
| XIyf-HfIak8 | 1973995 |
| vEndFKJM-fo | 5093959 |
| g0-QnBH5Jwo | 82347 |
+-----+-----+
(10 rows)
cqlsh:trial>
```

Fig.12:Video id and view count for gaming

Delete a video from gaming category

```
cqlsh:trial> DELETE FROM trial_name where videoid='ZKGt23WHUU0';
cqlsh:trial> select videoId, viewCount from trial_name where category='Gaming' LIMIT 10 ALLOW FILTERING;
+-----+-----+
| videoid | viewcount |
+-----+-----+
| nYe0Nwvvrrg | 87743 |
| oyEph5fxBy4 | 589997 |
| _BG98e_w6d0 | 20200 |
| WLLUKLhVl0E | 184444 |
| c0uf82rJnI8 | 2149855 |
| WQDyvTUVaEY | 10167 |
| XIyf-HfIak8 | 1973995 |
| vEndFKJM-fo | 5093959 |
| g0-QnBH5Jwo | 82347 |
| qeOooiDR27g | 15721 |
+-----+-----+
(10 rows)
cqlsh:trial>
```

Fig.13:Delete one record from gaming

Maximum likes obtained by any video

```
cqlsh:trial> select min(likeCount) from trial_name;
+-----+
| system.min(likecount) |
+-----+
| 0 |
+-----+
(1 rows)

Warnings :
Aggregation query used without partition key

cqlsh:trial>
```

Fig.14:Maximum likes obtained by any video

Maximum likes obtained by any video

```
cqlsh:trial> select max(likeCount) from trial_name;
+-----+
| system.max(likecount) |
+-----+
| 28793939 |
+-----+
(1 rows)

Warnings :
Aggregation query used without partition key

cqlsh:trial>
```

Fig.15: Minimum likes obtained by any video

6.4 Frontend

- Homepage:

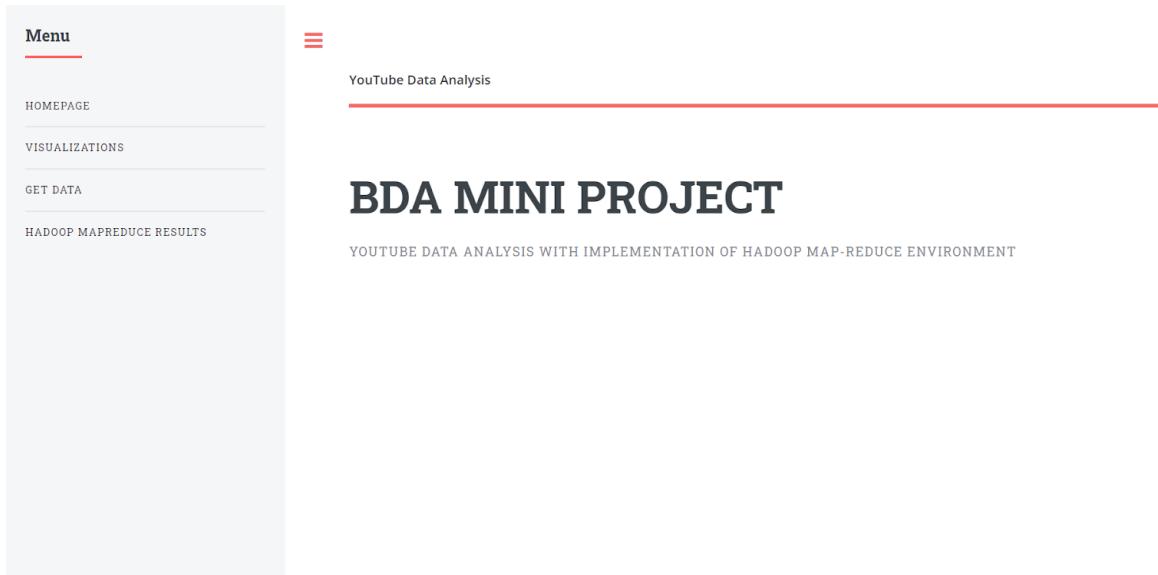


Fig.16: Homepage

- PowerBI Visualizations:

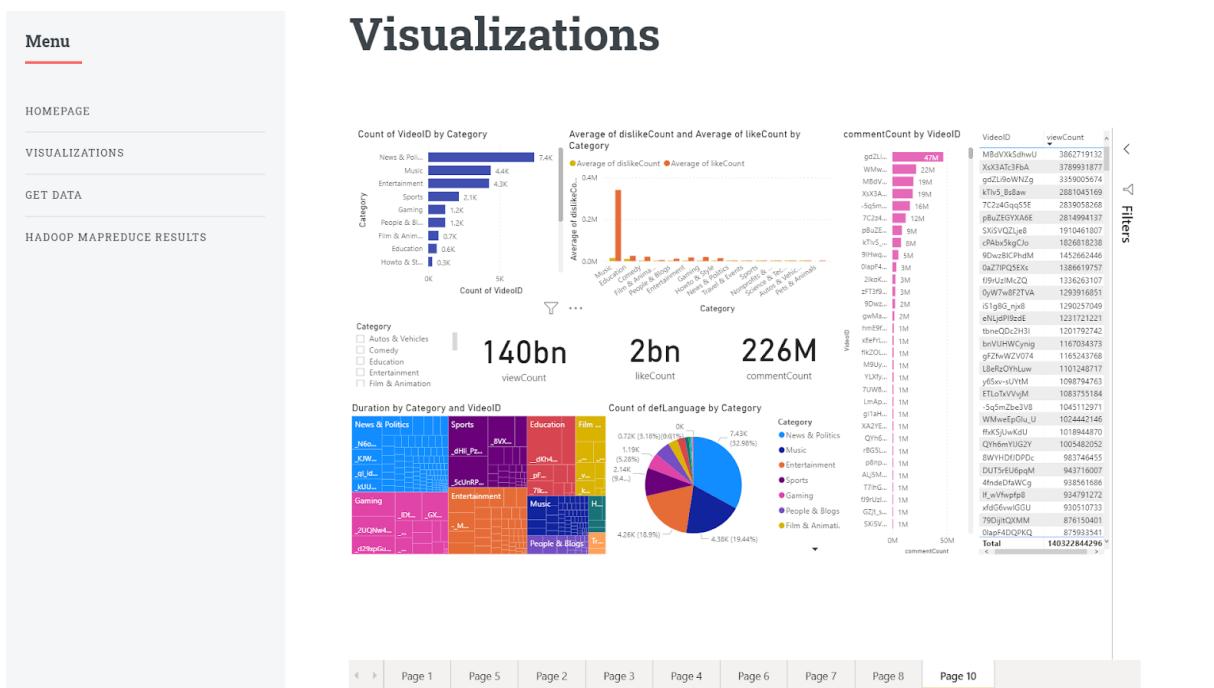


Fig.17: Visualizations embedded in a webpage

• Hadoop Map-Reduce Outputs:

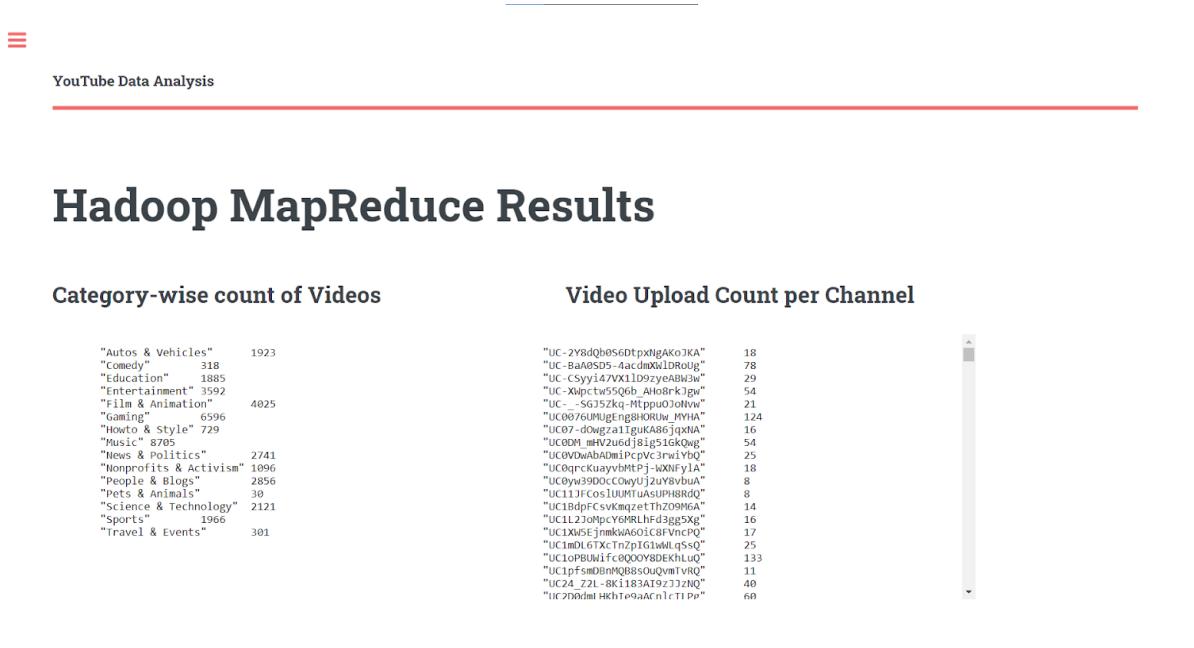


Fig.18: Hadoop MapReduce Output embedded in a webpage

7. VISUALIZATION SCREENSHOTS

● Visualizations filtered by Education Category

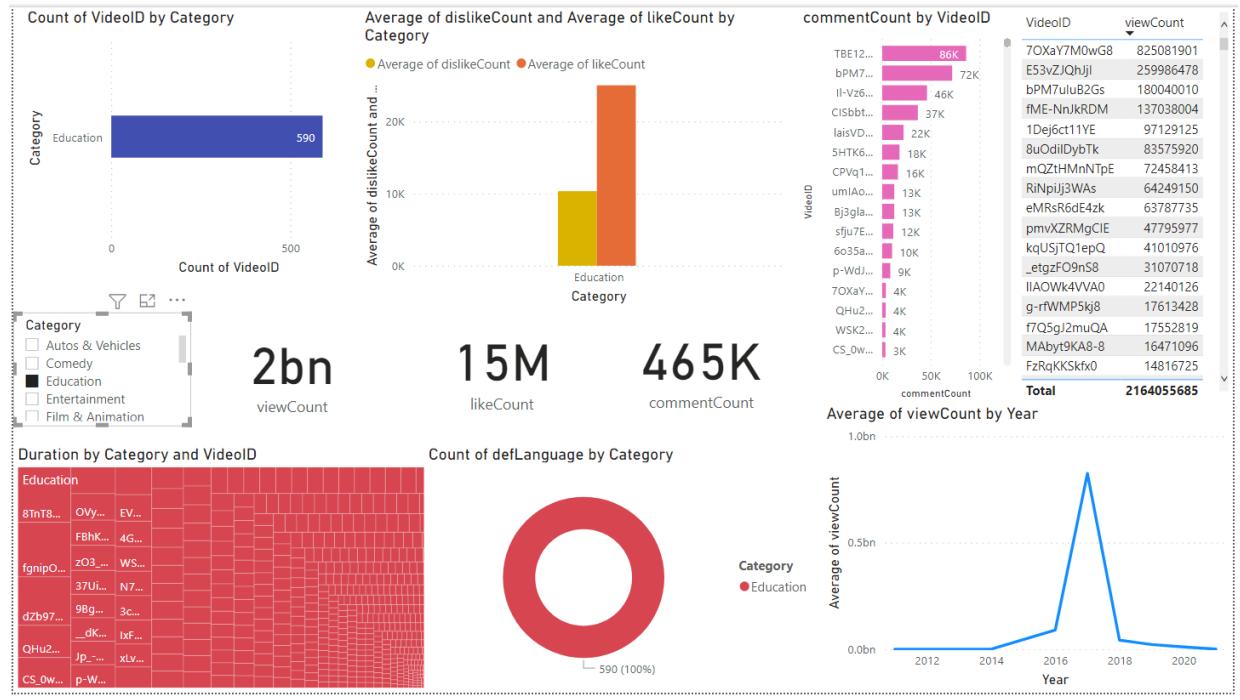


Fig.19: Visualizations filtered by education category

● Dashboard

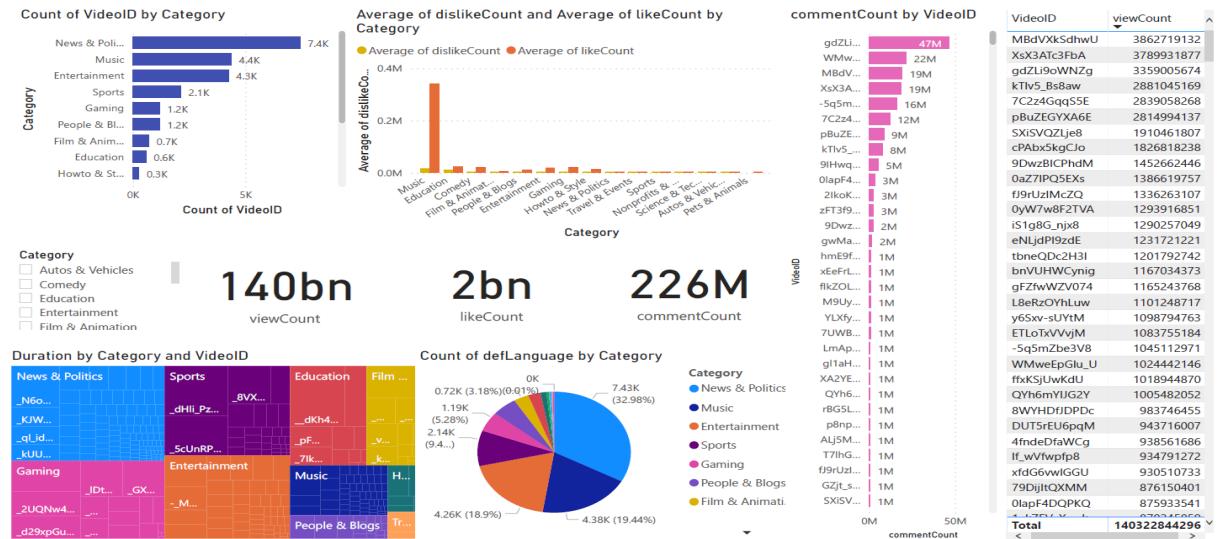


Fig.20: PoweBI Dashboard

● Key Influencers

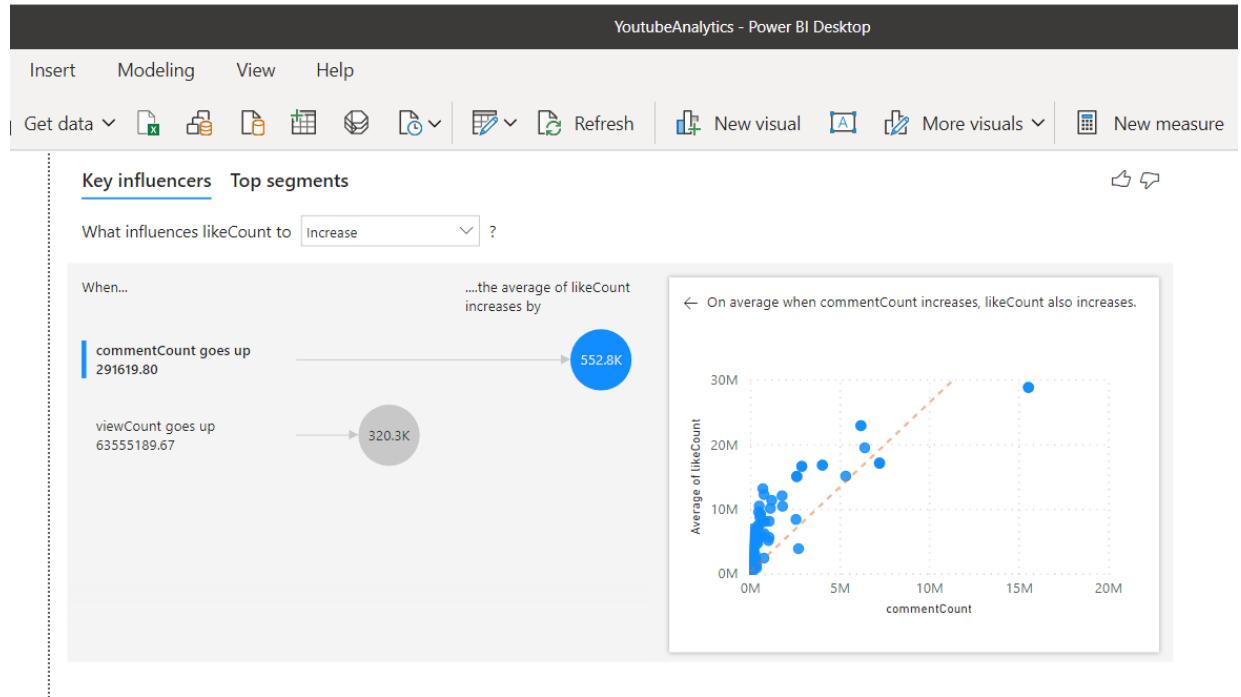


Fig.21:Key Influencers

8. CONCLUSION AND FUTURE WORK

8.1 Conclusion

In the digitalized era today, companies use YouTube for marketing and promoting their products and brand by uploading their product advertisement video to YouTube and movie makers promote their movies by uploading songs and movie trailers to YouTube. The measure of how well the product and movie is received by the public are determined by the number of views, likes (ratings) and comments on the video. This project intends to hit on those key areas which companies and organizations use or can use to measure their product's/movie's success against their competitors. As seen from the methodology, the basic algorithm retrieves reports to better understand and view statistics and trends for users' channels depending on the number of views and likes not only on their respective videos but also check if their competitors' are at the top. Another output result gives us insights on what categories of videos interest the public more. This can be done by analyzing the top video categories. This also helps budding YouTubers who upload YouTube videos to earn money. Thus, we used MapReduce based on JAVA, Cassandra CQL queries and Power BI tool to gain insights from the extracted Youtube data.

8.2 Future Scope

This project can be further advanced by designing a MapReduce algorithm to perform sentiment analysis on YouTube video comments. Also, an algorithm on comment analysis can help analyze and identify the numbers of trolls harassing authentic users and spam users. Designing a MapReduce algorithm for analysing how many users click on the ads or how many users enable or disable Adblock for their website can give a very interesting and important insight to the company. These algorithms can also be implemented on any social networking site promoting ads.

9. REFERENCES

1. F. Shaikh, D. Pawaskar, A. Siddiqui and U. Khan, "YouTube Data Analysis using MapReduce on Hadoop," 2018 3rd IEEE International Conference on Recent Trends in Electronics, Information & Communication Technology (RTEICT), 2018, pp. 2037-2041, doi: 10.1109/RTEICT42901.2018.9012635.
2. Hota, S. (2018). Big Data Analysis on YouTube Using Hadoop And Mapreduce. International Journal of Computer Engineering in Research Trends, 5, 98-104.
3. Shelke, M.B., 2017. YDA: Youtube Data Analysis Using Hadoop and Mapreduce. Open Access International Journal of Science and Engineering, 2(11).