

SJSU ChatGPT Tweets Analysis

DATA 225: LAB 2

Anusmriti Sikdar

Data Analytics

SJSU

San Jose, USA

anusmriti.sikdar@sjsu.edu

Akanksha Tyagi

Data Analytics

SJSU

San Jose, USA

akanksha.tyagi@sjsu.edu

Anjali Mishra

Data Analytics

SJSU

San Jose, USA

anjalihimanshu.ojha@sjsu.edu

Priyanka Bhyregowda

Data Analytics

SJSU

San Jose, USA

priyanka.bhyregowda@sjsu.edu

Sahana Thoravalli Prabhuswamy

Data Analytics

SJSU

San Jose, USA

sahana.thoravalliprabhuswamy@sjsu.edu

Abstract—Recently, ChatGPT has been one of the most discussed topics. Twitter is a popular social media platform for people to connect and share their thoughts. Analyzing Twitter data provides deep insights into the way people think. By analyzing Twitter data on ChatGPT we can obtain valuable information. Hence, there is a need for a ChatGPT tweet analyzer application. In this paper, we're going to analyse the ChatGPT dataset using MongoDB, which is a NoSQL database. NoSQL database provides many advantages like high scalability, flexibility, performance, high availability, cost-effectiveness, and drive towards big data. The data is stored in MongoDB in the form of documents, which are similar to JSON objects. Manipulating and storing data in its natural form makes this database a simple choice, without requiring complex normalization and mapping operations.

Index Terms—ChatGPT, Atlas, MongoDB

I. INTRODUCTION

One of the most discussed revolutionary innovations in recent times in the field of Artificial Intelligence and Machine learning is the introduction of ChatGPT which is Large Language Model (LLM). People have shared their ChatGPT usage experience on social media platforms like twitter, facebook etc. In this project, we analyze the conversations involving ChatGPT using Twitter Analysis dataset from Kaggle [1]. The main purpose of this analysis is to provide insights into SJSU community's feedback and experiences with the language model using popularly used hashtags, most discussed tweets, etc. This also gives insight into users' emotions and opinions. We build an application to perform the analysis. The application system will be built using MongoDB as the NoSQL database to store and retrieve the data. The ChatGPT analyzer app can be run on any machine with a command line interface, python, and access to the internet, which will be available in most of the systems. Twitter Database in MongoDB has a reduced number of Tables or collections when compared to RDBMS, because of normalization. Even though

there is data redundancy caused due to data denormalization, the performance is significantly improved.

In section II, we describe the functional requirement of the application.

II. FUNCTIONAL REQUIREMENT ANALYSIS

A. Functionalities

SJSU ChatGPT Analyzer application has a client-server architecture as shown in figure 1. The Analyzer app acts as a client to the MongoDB Atlas cluster which acts as a server. A cluster in MongoDB Atlas provides a cloud-based database environment that caters to the data need of our application. The raw data which is in the form of a .csv file is extracted using a Python program and then cleaned and transformed into the required format and then loaded into the twitter database in the MongoDB Atlas cluster. The application runs on a command line interface and the user of our application can select options to view a few of the interesting metrics that can be obtained from the application. The developer can connect to MongoDB Atlas cluster endpoints either through a Mongo shell client or MongoDB compass or through the Python application. Only a DB admin can connect to the MongoDB Atlas and is in control of access privileges and performance tuning. Our application provides information about the following -

- Users most discussing about ChatGPT.
- Most viral tweet.
- Most used hashtags.
- Most Discussed tweet with respect to replies and conversation.
- Most used source.
- Most active users.
- Number of tweets each day.
- Most used language.
- Most used media.

- Most mentioned users.
- Most liked tweets about ChatGPT.
- Users with most number of retweets.
- Most commonly used words in tweets about ChatGPT.

Functional Components of SJSU ChatGPT Analyzer using MongoDB

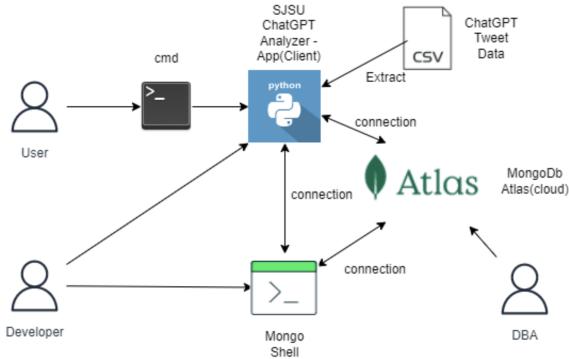


Fig. 1. Application design.

B. Limitations

Our application has following limitations -

- It provides a fixed number of metrics.
- The metrics are based on a small set of data. So they may not provide full insight into the real metrics.

III. DATABASE DESIGN

MongoDB is a NoSQL document-oriented database that stores information in flexible BSON documents. To achieve a well-structured and efficient database it is important to identify and model the relationship between data collections, define their attributes, and identify the appropriate schema design for storage in the database.

To design the schema for the database in MongoDB, it is important to understand the structure of data and requirements of the application. The dataset consists of a collection of tweets that have been scraped from Twitter consisting of the hashtag #chatgpt. People have been tweeting about this language model, sharing their experiences with it, and talking about articles or resources related to ChatGPT. Based on the requirements defined in section II-A, the dataset is split into various collections. The collections are chosen such that they can exist and be accessed independently. The following sections describe each collection in detail.

A. Collections

1) *Users*: The collection users stores Twitter users related information. Each document in this collection represents a single user and has a `username`, `userUrl` and `userInfo` fields. An example document in the collection is shown below.

```
{
  "_id": {
```

```
    "$oid": "64437059ddd851158e3777c7"
  },
  "username": "ciffi",
  "userUrl": "https://twitter.com/ciffi",
  "userInfo": "username=ciffi,id=19856336,
  displayname=christianfüller..."}
```

2) *User Mentions*: The collection `user_mentions` stores data about other Twitter users who are mentioned in the tweets. Each document in the collection represents a single mention and has `tweetId`, `username`, `userid` and `displayname`. An example document in the collection is shown below.

```
{
  "_id": {
    "$oid": "64437090ddd851158e39baee"
  },
  "tweetId": {
    "$numberLong": "1617156291046133761"
  },
  "username": "AlexandrovnaIng",
  "userid": "2827059006",
  "displayname": "Alexandrovna"
}
```

3) *Tweets*: The collection `tweets` is the main collection of the entire database design. It stores details of all the information related to a tweet. Each document in this collection represents a single tweet posted by a user and has `tweetId`, `tweetText`, `username`, `createdTime`, `lang`, `conversationId`, `replyCount`, `retweetCount`, `likeCount`, `quoteCount`, `hashtagCounts`, `Used_source` and an embedded document `urls` containing `permalinkUrl`, `outlink`, `countLink`, and `quotedTweet`. An example document in the collection is shown below.

```
{
  "_id": {
    "$oid": "64437068ddd851158e3832f7"
  },
  "tweetId": {
    "$numberLong": "1617156291046133761"
  },
  "tweetText": "@AlexandrovnaIng Prohibition
of ChatGPT has been added to the honor
code of my daughter's school",
  "username": "Caput_LupinumSG",
  "createdTime": {
    "$date": "2023-01-22T13:44:39.000Z"
  },
  "lang": "en",
  "conversationId": {
    "$numberLong": "1617148639993806848"
  },
  "replyCount": 1,
  "retweetCount": 0,
  "likeCount": 5,
  "quoteCount": 0,
  "hashtagCounts": 0,
  "Used_source": "Twitter for iPhone",
  "urls": {
    "permalinkUrl":
      "https://twitter.com/Caput_LupinumSG..",
    "outLink": "",
    "countLink": "",
    "quotedTweet": ""
  }
}
```

}

4) *Hastags*: The collection hastags stores data about the used hashtags by the user in each tweet. Each document in this collection has a tweetid and hashtag. An example document of this collection is shown below.

```
{
  "_id": {
    "$oid": "64437084ddd851158e38f647"
  },
  "tweetId": {
    "$numberLong": "1617156308926349312"
  },
  "hashtag": "#ChatGPT"
}
```

5) *Media*: The collection media stores data about other Twitter users who are mentioned in the tweets. Each document in the collection represents a single mention and has corresponding tweetId, mediaType, mediaDetails.

```
{
  "_id": {
    "$oid": "6443708cddd851158e398f54"
  },
  "tweetId": {
    "$numberLong": "1617156308926349312"
  },
  "mediaType": "Photo",
  "mediaDetails": "Photo(
    previewUrl='https://pbs.twimg.com/media/...')"
}
```

B. Denormalization

It is always important to consider how often the data changes versus how often you read data from the database. Normalization makes inserting efficient and it is easier to maintain data integrity whereas denormalization makes data reading efficient. A fully de-normalized database might achieve higher read performance but takes longer to write because every field should be updated in a document if one of the fields must be changed. Fields that do not change too frequently can be embedded and inserted into documents. If the document is too small, does not need to change regularly, grows by a small amount, needs faster reads, and does not require immediate consistency, we use the concept of embedding the document. On the other hand, if there are big subdocuments, data changes frequently, needs up-to-date information, grows by a large amount, and faster writes then, we must use the concept of referencing. In our design, the de-normalized Data Structure of the collection tweets has an embedded document field urls which is shown below.

```
"urls": {
  "permalinkUrl": "https://twitter.com/mochico0123/status/...",
  "outLink": "",
  "countLink": "",
  "quotedTweet": ""
}
```

We can also have an array of objects for media_info, mentioned_users and hashtags in tweets which uses the concept of referencing. An example is shown in figure 2.

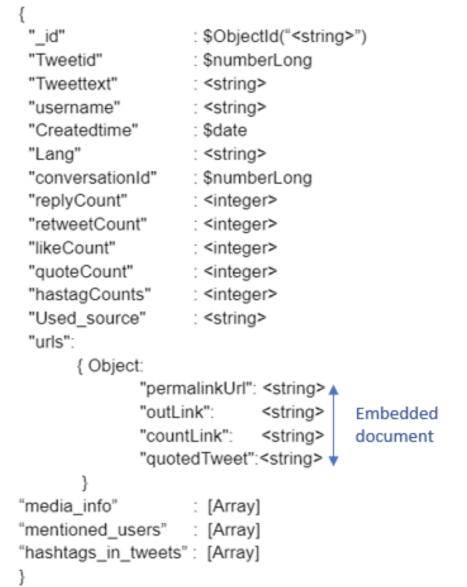


Fig. 2. Denormalization example.

IV. CONNECTION TO MONGODB ATLAS.

We use pymongo to connect to MongoDB atlas and upload all the collections there. The uploaded collections are shown in figure 3.

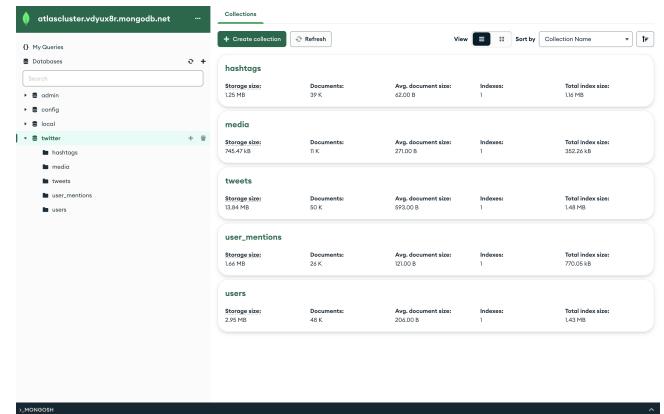


Fig. 3. Connection to MongoDB atlas.

V. QUERIES

In this section, we run queries for each requirement described in section II-A.

A. Number of tweets each day

The following query retrieves the number of tweets per day from tweets collection. The query uses the aggregation framework to group the tweets by day, using the \$dateToString operator to extract the date from the createdTime field of each tweet and format it as "YYYY-MM-DD". Then, it calculates the number of tweets per day using the \$sum operator.

```

db.tweets.aggregate([
  {
    "$group": {
      "_id": {
        "$dateToString": {
          "format": "%Y-%m-%d",
          "date": "$createdTime"
        }
      },
      "tweet_count": {
        "$sum": 1
      }
    }
  },
  {
    "$project": {
      "_id": 0,
      "tweeted_date": "$_id",
      "tweet_count": 1
    }
  },
  {
    "$sort": {
      "tweeted_date": 1
    }
  }
])

```

The results of the query are shown in figure 4.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
...   {
...     $group: {
...       _id: {
...         $dateToString: {
...           format: "%Y-%m-%d",
...           date: "$createdTime"
...         }
...       },
...       tweet_count: { $sum: 1 }
...     }
...   },
...   {
...     $project: {
...       _id: 0,
...       tweeted_date: "$_id",
...       tweet_count: 1
...     }
...   },
...   {
...     $sort: {
...       tweeted_date: 1
...     }
...   }
... ])
[ { tweet_count: 10068, tweeted_date: '2023-01-22' },
  { tweet_count: 31700, tweeted_date: '2023-01-23' },
  { tweet_count: 8233, tweeted_date: '2023-01-24' } ]

```

Fig. 4. Number of tweets each day.

B. Most active users

The following query returns a list of active users along with their total tweet count, sorted by the number of tweets in descending order. It can help identify the most active users and help to better understand the user behavior and engagement on the platform.

```

db.tweets.aggregate([
  {
    "$group": {
      "_id": "$username",
      "tweet_count": {
        "$sum": 1
      }
    }
  },
  {
    "$project": {
      "_id": 0,
      "active_user": "$_id",
      "total_tweets": {
        "$sum": 1
      }
    }
  }
])

```

```

    "tweet_count": 1
  }
},
{
  "$sort": {
    "tweet_count": -1
  }
},
{
  "$limit": 10
}
])

```

The results of the query are shown in figure 5.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
...   {
...     $group: {
...       _id: "$username",
...       tweet_count: { $sum: 1 }
...     }
...   },
...   {
...     $project: {
...       _id: 0,
...       active_user: "$_id",
...       tweet_count: 1
...     }
...   },
...   {
...     $sort: {
...       tweet_count: -1
...     }
...   }
... ])
[ { tweet_count: 60, active_user: 'translation_ja' },
  { tweet_count: 47, active_user: 'SaveToLotion' },
  { tweet_count: 44, active_user: 'brandanimo' },
  { tweet_count: 43, active_user: 'richardkimph' },
  { tweet_count: 38, active_user: 'VelleCyber3' },
  { tweet_count: 33, active_user: 'ChatGPTSpecial' },
  { tweet_count: 33, active_user: 'mitsuke' },
  { tweet_count: 31, active_user: 'jimaskade' },
  { tweet_count: 31, active_user: 'infotecnica' },
  { tweet_count: 30, active_user: 'ArtInNews' },
  { tweet_count: 30, active_user: 'AronMarcelino' },
  { tweet_count: 27, active_user: 'ba_zie' },
  { tweet_count: 27, active_user: 'RubenKaravelho' },
  { tweet_count: 27, active_user: 'itswhitchouhan' },
  { tweet_count: 26, active_user: 'genericgranda' },
  { tweet_count: 26, active_user: 'ArtIntel_BI' },
  { tweet_count: 26, active_user: 'halfteammind' },
  { tweet_count: 25, active_user: 'Alan_Nishihara' },
  { tweet_count: 25, active_user: 'MidJourneyAI' },
  { tweet_count: 25, active_user: 'yosikananoru' } ]

```

Fig. 5. Most active users.

C. Most used hashtags

The following query retrieves the 10 most used hashtags from the hashtags collection by grouping the tweets by hashtag field and using the \$sum operator to calculate the number of times each hashtag appears in the collection. The \$sort stage sorts the output in descending order based on the hashtag_count field.

```

db.hashtags.aggregate([
  {
    "$group": {
      "_id": "$hashtag",
      "hashtag_count": {
        "$sum": 1
      }
    }
  },
  {
    "$sort": {
      "hashtag_count": -1
    }
  },
  {
    "$limit": 10
  }
])

```

The results of the query are shown in figure 6.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.hashtags.aggregate([
...   {
...     $group: {
...       _id: "$hashtag",
...       hashtag_count: { $sum: 1 }
...     }
...   },
...   {
...     $sort: {
...       hashtag_count: -1
...     }
...   },
...   {
...     $limit: 10
...   }
... ])
[ { _id: '#ChatGPT', hashtag_count: 7821 },
{ _id: '#AI', hashtag_count: 1913 },
{ _id: '#chatgpt', hashtag_count: 1255 },
{ _id: '#openAI', hashtag_count: 744 },
{ _id: '#ai', hashtag_count: 617 },
{ _id: '#chatGPT', hashtag_count: 560 },
{ _id: '#ArtificialIntelligence', hashtag_count: 538 },
{ _id: '#microsoft', hashtag_count: 441 },
{ _id: '#openai', hashtag_count: 248 },
{ _id: '#IA', hashtag_count: 230 } ]

```

Fig. 6. Most used hashtags.

D. Most used source

The following query retrieves the top 10 most used sources from the tweet collection by grouping the tweets by the `Used_source` field and using the `$sum` operator to calculate the number of times each source appears in the collection. The `$sort` stage sorts the output in descending order based on the `count` field.

```

db.tweets.aggregate([
{
  "$group": {
    "_id": "$Used_source",
    "count": {
      "$sum": 1
    }
  }
},
{
  "$sort": {
    "count": -1
  }
},
{
  "$limit": 10
}
])

```

The results of the query are shown in figure 7.

E. Most used language

The following query returns a list of languages used in the tweets collection along with the total number of tweets in each language, sorted by the number of tweets in descending order. It can help to identify the most common languages used on the platform and better understand the language preferences of the users.

```

db.tweets.aggregate([
{
  "$group": {
    "_id": "$lang",
    "totalNumberOfTweets": {
      "$sum": 1
    }
  }
},
{
  ...
}
])

```

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
...   {
...     $group: {
...       _id: "$Used_source",
...       count: { $sum: 1 }
...     }
...   },
...   {
...     $sort: { count: -1 }
...   },
...   {
...     $limit: 20
...   }
... ])
[ { _id: 'Twitter Web App', count: 17814 },
{ _id: 'Twitter for iPhone', count: 12281 },
{ _id: 'Twitter for Android', count: 8972 },
{ _id: 'Twitter', count: 103 },
{ _id: 'DuckDuckGo', count: 859 },
{ _id: 'TweetDeck', count: 928 },
{ _id: 'Twitter for iPad', count: 756 },
{ _id: 'Jetpack.com', count: 629 },
{ _id: 'Buffer', count: 542 },
{ _id: 'LinkedIn', count: 423 },
{ _id: 'Hootsuite Inc.', count: 384 },
{ _id: 'Twitter for Mac', count: 218 },
{ _id: 'Microsoft Power Platform', count: 139 },
{ _id: 'SocialFlow', count: 131 },
{ _id: 'Search for Twitter', count: 131 },
{ _id: 'Echofon', count: 126 },
{ _id: 'Typefully', count: 122 },
{ _id: 'The Tweeted Times', count: 119 },
{ _id: 'Zapier.com', count: 114 },
{ _id: 'Hypefury', count: 103 } ]

```

Fig. 7. Most used source.

```

{
  "$$sort": {
    "totalNumberOfTweets": -1
  }
},
{
  "$limit": 10
}
])

```

The results of the query are shown in figure 8.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
...   {
...     $group: {
...       _id: "$lang",
...       totalNumberOfTweets: { $sum: 1 }
...     }
...   },
...   {
...     $sort: { totalNumberOfTweets: -1 }
...   }
... ])
[ { _id: 'en', totalNumberOfTweets: 32076 },
{ _id: 'ja', totalNumberOfTweets: 5046 },
{ _id: 'es', totalNumberOfTweets: 3315 },
{ _id: 'fr', totalNumberOfTweets: 2492 },
{ _id: 'de', totalNumberOfTweets: 1287 },
{ _id: 'pt', totalNumberOfTweets: 1175 },
{ _id: 'it', totalNumberOfTweets: 443 },
{ _id: 'tr', totalNumberOfTweets: 436 },
{ _id: 'und', totalNumberOfTweets: 423 },
{ _id: 'qmc', totalNumberOfTweets: 395 },
{ _id: 'ar', totalNumberOfTweets: 392 },
{ _id: 'nl', totalNumberOfTweets: 319 },
{ _id: 'in', totalNumberOfTweets: 251 },
{ _id: 'th', totalNumberOfTweets: 193 },
{ _id: 'fa', totalNumberOfTweets: 187 },
{ _id: 'ru', totalNumberOfTweets: 181 },
{ _id: 'zh', totalNumberOfTweets: 149 },
{ _id: 'ko', totalNumberOfTweets: 141 },
{ _id: 'iw', totalNumberOfTweets: 113 },
{ _id: 'ca', totalNumberOfTweets: 113 } ]

```

Fig. 8. Most used language.

F. Most used media

The following query returns a list of media types used in the collection along with the number of times each type was used, sorted by the number of times used in descending order. It can help to identify the most popular media types on the platform and to better understand the preferences of the users.

```

db.media.aggregate([
{
  "$group": {
    "_id": "$mediaType",
    "NumberofTimesUsed": {
      "$sum": 1
    }
  }
},
{
  ...
}
])

```

```

        }
    },
    {
      "$sort": {
        "NumberOfTimesUsed": -1
      }
    },
    {
      "$limit": 10
    }
  ]
)

```

The results of the query are shown in figure 9.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.media.aggregate([
...   {
...     $group :{
...       _id : "$mediaType",
...       NumberOfTimesUsed: { $sum: 1}
...     }
...   },
...   {
...     $sort: { NumberOfTimesUsed: -1 }
...   }
... ])
[
{ _id: 'Photo', NumberOfTimesUsed: 10233 },
{ _id: 'Video', NumberOfTimesUsed: 522 },
{ _id: 'Gif', NumberOfTimesUsed: 407 }
]

```

Fig. 9. Most used media.

G. Most liked tweets about ChatGPT

The following query will first filter the tweets to only include those containing "chatgpt" Then, it will sort the tweets in descending order by likeCount and limit the results to the first document.

```

db.tweets.aggregate([
{
  $match: {
    tweetText: /chatgpt/i
  }
},
{
  "$sort": {
    "likeCount": -1
  }
},
{
  "$limit": 1
}
])

```

The results of the query are shown in figure 10.

H. Top 10 users who have tweeted the most about #chatgpt.

This query retrieves the top 10 users who have used the #chatgpt hashtag in their tweets the most number of times.

```

db.tweets.aggregate([
{
  "$match": {
    tweetText: { "$regex": /\#chatgpt/i }
  }
},
{
  "$group": {
    "_id": "$username",
    "count": {
      "$sum": 1
    }
  }
})

```

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
...   {
...     $match: {
...       tweetText: /\#chatgpt/i
...     }
...   },
...   {
...     $sort: {
...       likeCount: -1
...     }
...   },
...   {
...     $limit: 1
...   }
... ])
{
  _id: ObjectId("64a37668dd851158e383485"),
  tweetId: Long("1617162355112124421"),
  tweetText: "ChatGPT passed a Wharton MBA exam. \n\nTime to overhaul education.",
  username: "GRDector",
  createdTime: ISODate("2023-01-22T14:08:45.000Z"),
  lang: "en",
  conversationId: Long("1617162355112124421"),
  replyCount: 1421,
  retweetCount: 6815,
  likeCount: 2397,
  quoteCount: 1947,
  hashtagCounts: 8,
  Used_source: "Twitter for iPhone",
  urls: [
    permanentUrl: "https://twitter.com/GRDector/status/1617162355112124421"
  ],
  outlinks: [],
  countLink: 1,
  quotedTweet: null
}

```

Fig. 10. Most liked tweets about ChatGPT.

```

  },
  {
    "$sort": {
      "count": -1
    }
  },
  {
    "$limit": 10
  }
])

```

The results of the query are shown in figure 11.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
...   {
...     $match: {
...       tweetText: { $regex: /\#chatgpt/i }
...     }
...   },
...   {
...     $group: {
...       _id: "$username",
...       count: { $sum: 1 }
...     }
...   },
...   {
...     $sort: { count: -1 }
...   },
...   {
...     $limit: 10
...   }
... ])
[
{ _id: 'translation_jp', count: 60 },
{ _id: 'richardkimhd', count: 43 },
{ _id: 'VeilleCyber3', count: 38 },
{ _id: 'Artilex', count: 38 },
{ _id: 'MidJourneyAI', count: 25 },
{ _id: 'Jhengster', count: 21 },
{ _id: 'ChatGPT_Tweeter', count: 21 },
{ _id: 'abushayeb', count: 20 },
{ _id: 'Bar_zie', count: 19 },
{ _id: 'PriorityDomains', count: 19 }
]

```

Fig. 11. Top 10 users who have tweeted the most about #chatgpt.

I. Top 10 users who have the highest number of retweets in their tweets containing #chatgpt

This query retrieves the top 10 users who have tweeted with the #chatgpt hashtag and have received the highest total number of retweets for their tweets.

```

db.tweets.aggregate([
{
  $match: { tweetText: { $regex: /\#chatgpt/i } }
},
{

```

```

        "$group": {
            "_id": "$username",
            "totalRetweets": {
                "$sum": "$retweetCount"
            }
        },
        {
            "$sort": {
                "totalRetweets": -1
            }
        },
        {
            "$limit": 10
        }
    ]
)

```

The results of the query are shown in figure 12.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
...   {
...     $match: {
...       tweetText: { $regex: /#chatgpt/i }
...     }
...   },
...   {
...     $group: {
...       _id: "$username",
...       totalRetweets: { $sum: "$retweetCount" }
...     }
...   },
...   {
...     $sort: { totalRetweets: -1 }
...   },
...   {
...     $limit: 10
...   }
... ])
[{"_id": "DataChar", "totalRetweets": 331}, {"_id": "kyle_chasse", "totalRetweets": 222}, {"_id": "Schuldenuehner", "totalRetweets": 138}, {"_id": "StevenXGe", "totalRetweets": 125}, {"_id": "CerfiaFR", "totalRetweets": 107}, {"_id": "Veillecyber3", "totalRetweets": 102}, {"_id": "healthy_pockets", "totalRetweets": 93}, {"_id": "rankurconic", "totalRetweets": 88}, {"_id": "lgoodterro", "totalRetweets": 83}, {"_id": "PriorityDomains", "totalRetweets": 82}]

```

Fig. 12. Top 10 users who have the highest number of retweets in their tweets containing #chatgpt.

J. Most viral tweet

This query finds tweets that contain the hashtag #chatgpt and then calculates the "viral count" by adding up the number of replies, retweets, likes, and quotes for each tweet. Then, it sorts the tweets in descending order based on their viral count and returns the tweet with the highest viral count, which is considered the most viral tweet about chatGPT.

```

db.tweets.aggregate([
{
    $match: { tweetText: { $regex: /chatgpt/i } }
},
{
    "$addFields": {
        "viral_count": {
            "$add": [
                "$replyCount",
                "$retweetCount",
                "$likeCount",
                "$quoteCount"
            ]
        }
    }
},
{
    "$sort": {
        "viral_count": -1
    }
}
])

```

```

        },
        {
            "$limit": 1
        }
    ]
)

```

The results of the query are shown in figure 13.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
...   {
...     $match: {
...       tweetText: { $regex: /#chatgpt/i }
...     }
...   },
...   {
...     $addFields: {
...       viralcount: {
...         $add: ["$replyCount", "$retweetCount", "$likeCount", "$quoteCount"]
...       }
...     }
...   },
...   {
...     $sort: { viralcount: -1 }
...   },
...   {
...     $limit: 1
...   }
... ])
[{"_id": ObjectId("64a47988dd8d81158e3aaac0"), "username": "CerfiaFR", "createdAt": "2023-01-23T16:58:38.888Z", "lang": "fr", "conversationId": Long("1617565947676424047"), "replyCount": 107, "retweetCount": 107, "likeCount": 1199, "quoteCount": 119, "hashtagCount": 1, "Used_source": "TweetDeck", "urls": [{"permalinkUrl": "https://twitter.com/CerfiaFR/status/1617565947676424047"}], "outlink": "", "counted": 1, "quotedCount": 1, "viralcount": 1356}]

```

Fig. 13. Most viral tweet.

K. Most commonly used words in tweets about ChatGPT

The following query searches for tweets containing the word "chatgpt" (case-insensitive) in their text. Then it splits each tweet's text into individual words and groups them by word, counting the number of occurrences of each word. The resulting list is sorted in descending order of frequency and limited to the top 10 most used words. This can provide insights into the topics and themes associated with the use of the hashtag #chatgpt on Twitter.

```

db.tweets.aggregate([
{
    $match: { tweetText: { $regex: /chatgpt/i } }
},
{
    "$project": {
        "words": {
            "$split": [
                "$tweetText",
                " "
            ]
        }
    }
},
{
    "$unwind": "$words"
},
{
    "$group": {
        "_id": "$words",
        "count": {
            "$sum": 1
        }
    }
},
{
    "$sort": {
        "count": -1
    }
},
{
    "$limit": 10
}
])

```

```

    {
      "$limit": 10
    }
  ])

```

The results of the query are shown in figure 14.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
  ...
  {
    $match: {
      tweetText: { $regex: /chatgpt/i }
    }
  },
  {
    $project: {
      words: { $split: ["$tweetText", " "] }
    }
  },
  {
    $unwind: "$words"
  },
  {
    $group: {
      _id: "$words",
      count: { $sum: 1 }
    }
  },
  {
    $sort: { count: -1 }
  },
  {
    $limit: 10
  }
])
[{"_id": "CHATGPT", "count": 21120}, {"_id": "the", "count": 18801}, {"_id": "to", "count": 17800}, {"_id": "a", "count": 15373}, {"_id": "and", "count": 11882}, {"_id": "you", "count": 11733}, {"_id": "is", "count": 10298}, {"_id": "it", "count": 8397}, {"_id": "for", "count": 7831}, {"_id": "it", "count": 6752}]

```

Fig. 14. Most commonly used words in tweets about ChatGPT.

L. Most Discussed tweet with respect to replies and conversation

The following query helps to identify the most active conversations in the tweets collection. The \$match filters the tweets to include only those with a non-null conversationId. The \$group groups the tweets by conversationId and calculates the total number of replies for each conversation using \$sum operator. Then, the conversations are sorted by the total number of replies in descending order using -1.

```

db.tweets.aggregate([
  {
    "$match": {
      "conversationId": {
        "$ne": null
      }
    }
  },
  {
    "$group": {
      "_id": "$conversationId",
      "count": {
        "$sum": "$replyCount"
      }
    }
  },
  {
    "$sort": {
      "count": -1
    }
  },
  {
    "$limit": 1
  }
])

```

The results of the query are shown in figure 15.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.tweets.aggregate([
  ...
  {
    $match: {
      conversationId: { $ne: null }
    }
  },
  {
    $group: {
      _id: "$conversationId",
      count: { $sum: "$replyCount" }
    }
  },
  {
    $sort: { count: -1 }
  },
  {
    $limit: 1
  }
])
[{"_id": Long("1617161446202245120"), count: 3099}]

```

Fig. 15. Most Discussed tweet with respect to replies and conversation.

M. Most mentioned users

The following query will return the top 10 users who have been mentioned the most in the tweets, along with the count of how many times each user has been mentioned.

```

db.user_mentions.aggregate([
  {
    "$group": {
      "_id": "$username",
      "count": {
        "$sum": 1
      }
    }
  },
  {
    "$sort": {
      "count": -1
    }
  },
  {
    "$limit": 10
  }
])

```

The results of the query are shown in figure 16.

```

Atlas atlas-9ps6sb-shard-0 [primary] twitter> db.user_mentions.aggregate([
  ...
  {
    $group: {
      _id: "$username",
      count: { $sum: 1 }
    }
  },
  {
    $sort: { count: -1 }
  },
  {
    $limit: 10
  }
])
[{"_id": "OpenAI", "count": 744}, {"_id": "GRDeeter", "count": 495}, {"_id": "elonmusk", "count": 368}, {"_id": "YouTube", "count": 308}, {"_id": "chatgpt_isaac", "count": 196}, {"_id": "Microsoft", "count": 173}, {"_id": "noor_siddiqui_", "count": 158}, {"_id": "ChatGPT", "count": 153}, {"_id": "sama", "count": 141}, {"_id": "sejournal", "count": 131}]

```

Fig. 16. Most mentioned users.

VI. CHATGPT APPLICATION

We have also implemented an interactive python app. Users can select from the available options which are as discussed in the functional requirements to get the necessary metrics as shown in figure 17, 18, 19 and 20.

```
(base) C:\Users\amith\PycharmProjects\Test>mongoqueries.py -p
Password:
['testdb', 'tips', 'twitter', 'twitterdb', 'admin', 'local']
created index on tweets collection - Tweetid_1
created index on media collection - Tweetid_1
created index on hashtags collection - Tweetid_1
created index on users collection - username_1
created index on user_mentions collection - Tweetid_1

-----
method_id | method_type
-----|-----
1 | Get most viral tweet
2 | Get most used hashtags
3 | Get most used source
4 | Get most mentioned user
5 | Get most active users in descending order
6 | Get most liked tweet about chatgpt
7 | Get most used media type
8 | Get Users with more No of retweets
9 | Get most used language to tweet about chat GPT
10 | number of Tweets on Chat GPT
11 | Get most retweeted tweet in each language
12 | Get number of tweets in a day on ChatGpt
13 | Get most retweeted tweet with the comments
exit | Exits program

-----
Enter the method_name: 1
```

Fig. 17. App example 1.

```
method_id | method_type
-----|-----
1 | Get most viral tweet
2 | Get most used hashtags
3 | Get most used source
4 | Get most mentioned user
5 | Get most active users in descending order
6 | Get most liked tweet about chatgpt
7 | Get most used media type
8 | Get Users with more No of retweets
9 | Get most used language to tweet about chat GPT
10 | number of Tweets on Chat GPT
11 | Get most retweeted tweet in each language
12 | Get number of tweets in a day on ChatGpt
13 | Get most retweeted tweet with the comments
exit | Exits program

-----
Enter the method_name: 6
#####
##### Get most liked tweet about chatgpt #####
#####

{'Lang': 'en',
'Tweetid': '1617162355112124421',
'Tweettext': "ChatGPT passed a Wharton MBA exam.\xa0\n\n",
'Used_source': 'Twitter for iPhone',
'_id': ObjectId('64431942599c74d7583a940'),
'conversationId': '1617162355112124421',
'createdAt': datetime.datetime(2023, 1, 22, 14, 8, 45),
'hashtags_in_tweet': [],
'likeCount': 5073,
'media_info': [],
'mentioned_users': [],
'quoteCount': 1947,
'replyCount': 1421,
'retweetCount': 6815,
"urls': {'countlink': '',
'outlink': '',
'permalinkUrl': 'https://twitter.com/GRDector/status/1617162355112124421',
'quotedTweet': ''},
'username': 'GRDector',
'verticalCount': 6626}
```

Fig. 19. App example 3.

```
method_id | method_type
-----|-----
1 | Get most viral tweet
2 | Get most used hashtags
3 | Get most used source
4 | Get most mentioned user
5 | Get most active users in descending order
6 | Get most liked tweet about chatgpt
7 | Get most used media type
8 | Get Users with more No of retweets
9 | Get most used language to tweet about chat GPT
10 | number of Tweets on Chat GPT
11 | Get most retweeted tweet in each language
12 | Get number of tweets in a day on ChatGpt
13 | Get most retweeted tweet with the comments
exit | Exits program

-----
Enter the method_name: 1
#####
##### Get most viral tweet #####
#####

{'Lang': 'en',
'Tweetid': '1617162355112124421',
'Tweettext': "ChatGPT passed a Wharton MBA exam.\xa0\n\n",
'Used_source': 'Twitter for iPhone',
'_id': ObjectId('64431942599c74d7583a940'),
'conversationId': '1617162355112124421',
'createdAt': datetime.datetime(2023, 1, 22, 14, 8, 45),
'hashtags_in_tweet': [],
'likeCount': 5073,
'media_info': [],
'mentioned_users': [],
'quoteCount': 1947,
'replyCount': 1421,
'retwtweetCount': 6815,
"urls': {'countlink': '',
'outlink': '',
'permalinkUrl': 'https://twitter.com/GRDector/status/1617162355112124421',
'quotedTweet': ''},
'username': 'GRDector',
'verticalCount': 6626}
```

Fig. 18. App example 2.

```
method_id | method_type
-----|-----
1 | Get most viral tweet
2 | Get most used hashtags
3 | Get most used source
4 | Get most mentioned user
5 | Get most active users in descending order
6 | Get most liked tweet about chatgpt
7 | Get most used media type
8 | Get Users with more No of retweets
9 | Get most used language to tweet about chat GPT
10 | number of Tweets on Chat GPT
11 | Get most retweeted tweet in each language
12 | Get number of tweets in a day on ChatGpt
13 | Get most retweeted tweet with the comments
exit | Exits program

-----
Enter the method_name: 2
#####
##### Get most used hashtags #####
#####

{'NumberOfTimesUsed': 7821, '_id': '#ChatGPT'}
{'NumberOfTimesUsed': 1913, '_id': '#AI'}
{'NumberOfTimesUsed': 1255, '_id': '#chatgpt'}
{'NumberOfTimesUsed': 744, '_id': '#OpenAI'}
{'NumberOfTimesUsed': 617, '_id': '#ai'}
{'NumberOfTimesUsed': 560, '_id': '#chatGPT'}
{'NumberOfTimesUsed': 538, '_id': '#ArtificialIntelligence'}
{'NumberOfTimesUsed': 441, '_id': '#Microsoft'}
{'NumberOfTimesUsed': 248, '_id': '#openai'}
{'NumberOfTimesUsed': 230, '_id': '#IA'}
```

Fig. 20. App example 4.

VII. DATA VISUALIZATION

We have visualized the tweet data in MongoDB atlas. The chart is shown in Figure 21. The chart can also be accessed at the following public url: <https://charts.mongodb.com/charts-lab-2---chatgpt-tweet-ana-ifsvx/public/dashboards/64433f00-8295-4353-892c-04104f0e8333>.

VIII. PERFORMANCE ANALYSIS

We use `explain("executionStats")` method to analyse the performance of each query. The `explain` method results present the query plans as a tree of stages. figure

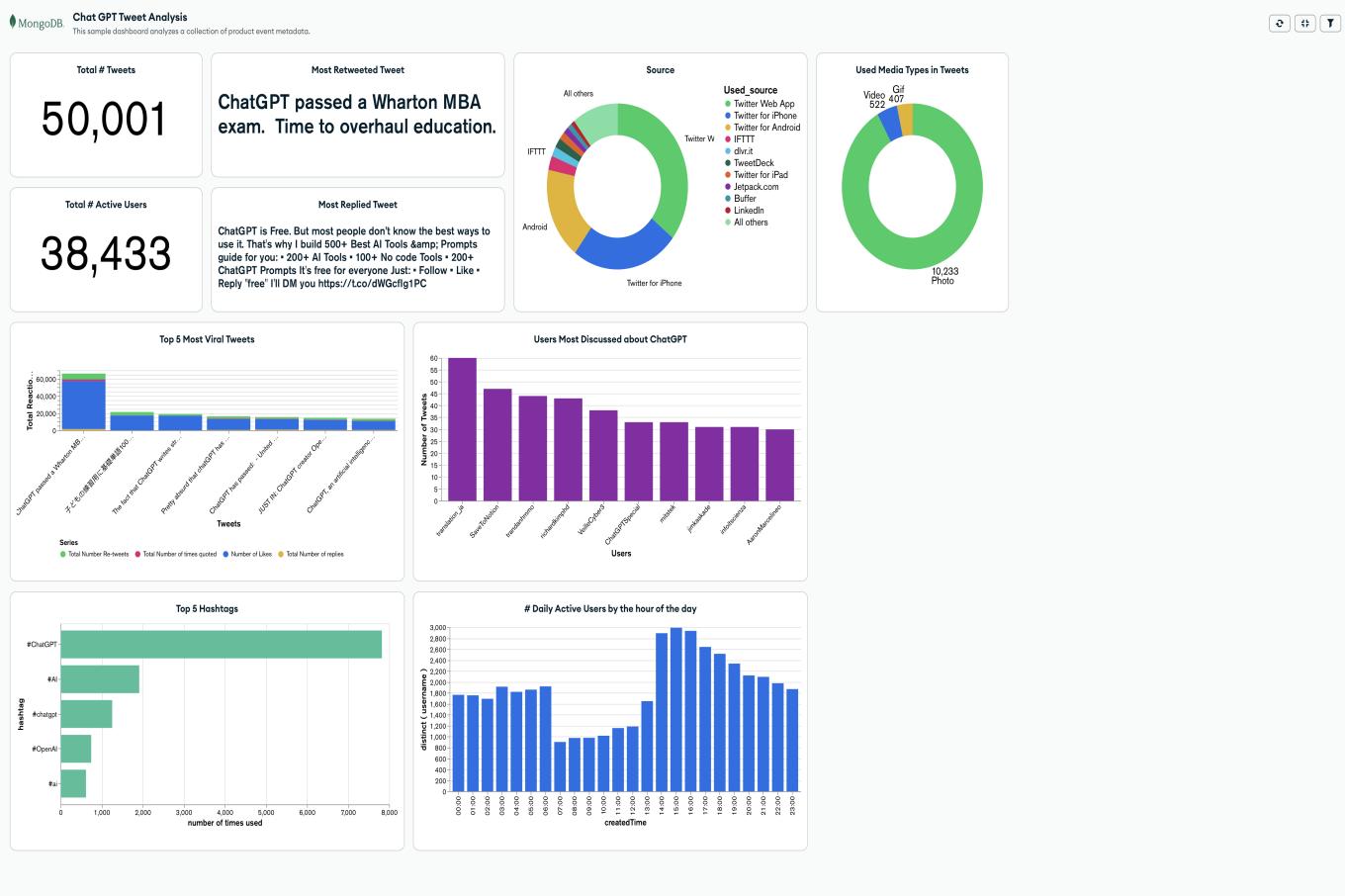


Fig. 21. Data visualization using MongoDb atlas.

23 shows the comparison of the execution times of all the queries. The detailed analysis of each stage is provided in the jupyter notebook. We see that the queries to get the most commonly used words takes the most amount of time. Detailed analysis shows that MongoDB uses `COLLSCAN` to scan all the documents in the collection, which results in high execution time. The execution stats are shown below.

```
{  
  "executionStats": {  
    "executionStages": {  
      "advanced": 46546,  
      "executionTimeMillisEstimate": 65,  
      "inputStage": {  
        "advanced": 46546,  
        "direction": "forward",  
        "docsExamined": 50001,  
        "executionTimeMillisEstimate": 35  
        "filter": {  
          "tweetText": {  
            "$options": "i",  
            "$regex": "chatgpt"  
          }  
        },  
        "isEOF": 1,  
        "nReturned": 46546,  
        "needTime": 3456,  
        "needYield": 0.  
      }  
    }  
  }  
}
```

```
        "restoreState": 60,
        "saveState": 60,
        "stage": "COLLSCAN",
        "works": 50003
    },
    "isEOF": 1,
    "nReturned": 46546,
    "needTime": 3456,
    "needYield": 0,
    "restoreState": 60,
    "saveState": 60,
    "stage": "PROJECTION_DEFAULT",
    "transformBy": {
        "_id": true,
        "words": {
            "$$split": [
                "$tweetText",
                {
                    "$const": " "
                }
            ]
        }
    },
    "works": 50003
},
"executionSuccess": true,
"executionTimeMillis": 878,
"nReturned": 46546,
"totalDocsExamined": 50001,
"totalKeysExamined": 0
```

```
},
"queryPlanner": {
    "indexFilterSet": false,
    "maxIndexedAndSolutionsReached": false,
    "maxIndexedOrSolutionsReached": false,
    "maxScansToExplodeReached": false,
    "namespace": "twitter.tweets",
    "parsedQuery": {
        "tweetText": {
            "$options": "i",
            "$regex": "chatgpt"
        }
    },
    "planCacheKey": "5B027F7E",
    "queryHash": "5B027F7E",
    "rejectedPlans": [],
    "winningPlan": {
        "inputStage": {
            "direction": "forward",
            "filter": {
                "tweetText": {
                    "$options": "i",
                    "$regex": "chatgpt"
                }
            },
            "stage": "COLLSCAN"
        },
        "stage": "PROJECTION_DEFAULT",
        "transformBy": {
            "_id": true,
            "words": {
                "$split": [
                    "$tweetText",
                    {
                        "$const": " "
                    }
                ]
            }
        }
    }
}
```

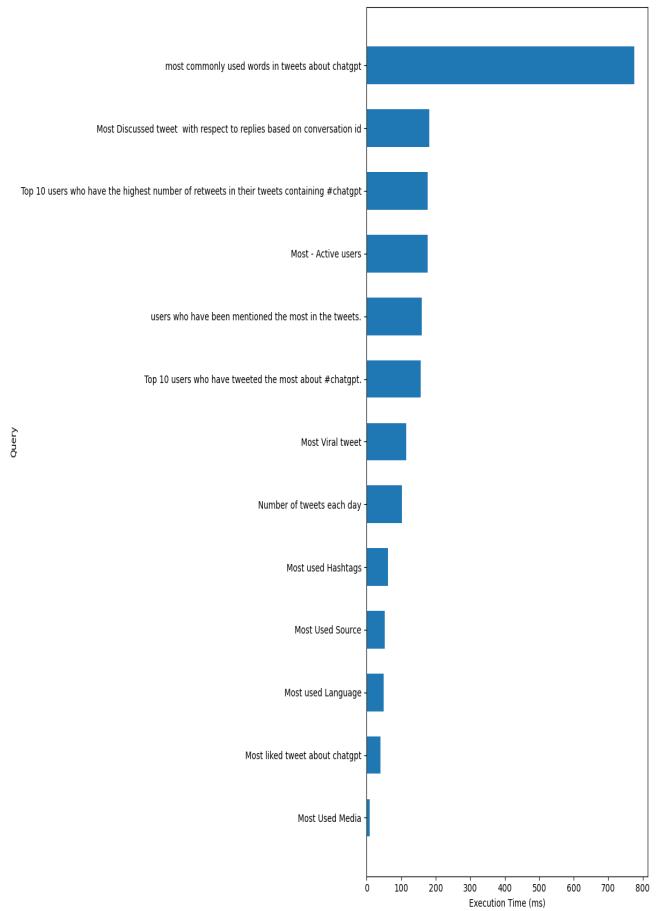


Fig. 23. Timing analysis of each query.

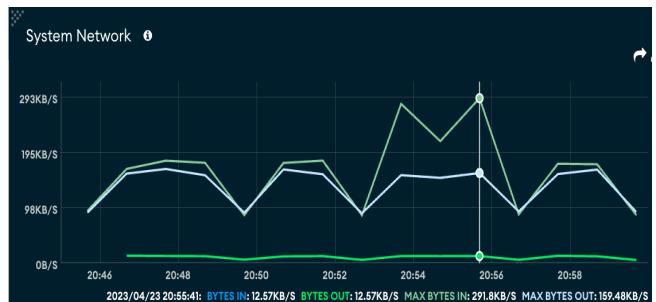


Fig. 24. System network usage.



Fig. 22. Execution times.

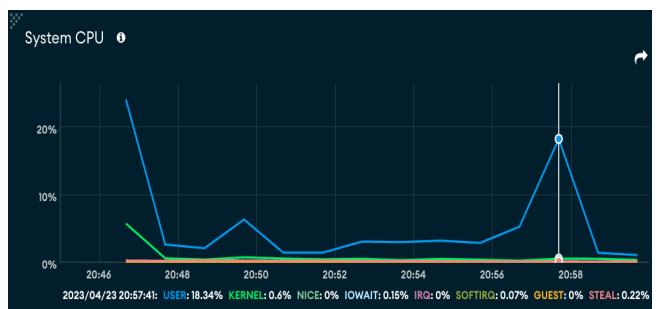


Fig. 25. CPU usage.

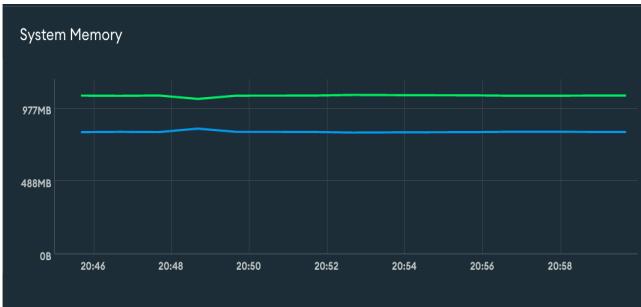


Fig. 26. Memory usage.



Fig. 27. Network usage.

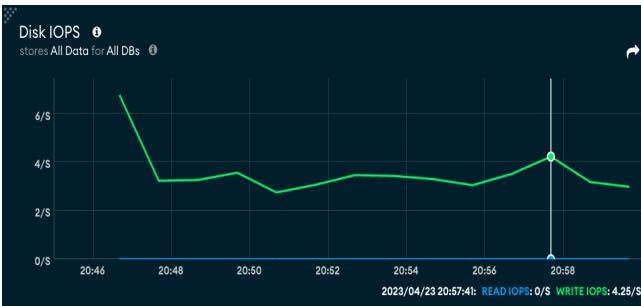


Fig. 28. Disk IO usage.



Fig. 29. Disk latency.

A. Comparison with MySQL

In figure 30, we compare the run times of each query on both MySQL and MongoDB. We find that most of the queries

run faster on MongoDB. However, we must take into account the factor that these queries were run on different clusters.

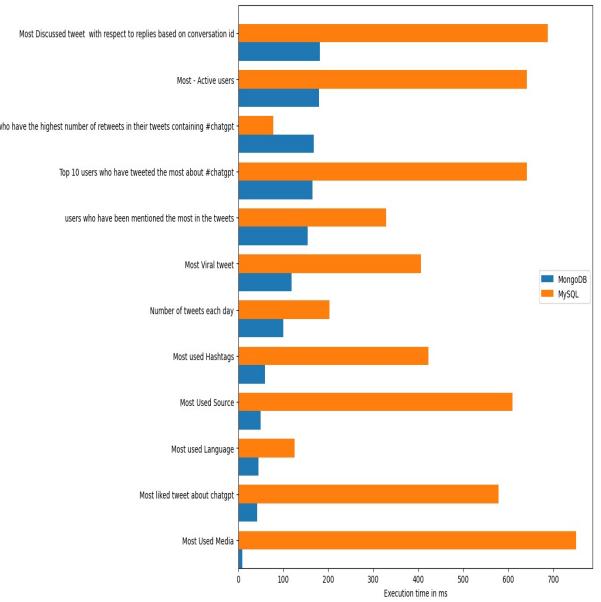


Fig. 30. Comparison of each queries on MySQL and MongoDB.

REFERENCES

- [1] ChatGPT Twitter Dataset: <https://www.kaggle.com/datasets/tariqsays/chatgpt-twitter-dataset>.
- [2] MongoDB: <https://www.mongodb.com/docs/manual/introduction/>