

Perbandingan Performa Apache Impala Dengan Apache Spark Dalam Mengeksekusi Kueri

Pratama Arjan Rangkuti, Muhammad Zihni Athallah, Tasya Harwani Barus*, Nadila Afisa Rani, Muhammad Ikhsan Setiawan
Sistem Komputer, Universitas Sriwijaya, Indonesia

*Korespondensi: tasyabarus87@gmail.com

ARTICLE INFO

Article History:

- Received 08 January 2023
- Received in revised form 23 March 2023
- Accepted 22 April 2023
- Available online 30 July 2023

ABSTRAK

Perkembangan teknologi komputasi modern membawa transformasi besar dalam pengelolaan dan analisis data. Saat ini, kita memiliki akses ke teknologi komputasi yang kuat, yang memudahkan pengolahan dan analisis data. Apache Impala dan Apache Spark adalah dua platform yang memainkan peran penting dalam ekosistem analisis data saat ini. Keduanya adalah bagian integral dari Hadoop, sebuah framework open-source yang efisien dalam penyimpanan dan pemrosesan data besar. Apache Impala adalah query engine yang memungkinkan eksekusi query SQL secara interaktif terhadap data yang tersimpan di HDFS atau HBase. Di sisi lain, Apache Spark adalah platform yang mengintegrasikan pemrosesan data streaming dan batch dalam satu kerangka kerja yang fleksibel. Kedua platform ini memiliki keunggulan dan perbedaan yang perlu dipahami dengan baik. Penelitian ini bertujuan untuk membandingkan performa eksekusi kueri antara Apache Impala dan Apache Spark. Eksperimen dilakukan dengan merancang kueri yang serupa dan menjalankannya pada kedua platform dengan dataset yang ada. Waktu eksekusi kueri dicatat dan dianalisis untuk memahami kinerja relatif dari masing-masing platform dalam menanggapi permintaan kueri. Hasil dari penelitian ini memberikan wawasan berharga tentang kecepatan dan efisiensi eksekusi kueri antara Apache Impala dan Apache Spark. Pemahaman ini dapat membantu organisasi dan praktisi data dalam memilih platform yang sesuai dengan kebutuhan dan konteks analisis data mereka.

Kata Kunci: Apache Impala, Apache Spark, Waktu Eksekusi, Kueri, Pemrosesan Data

ABSTRACT

The development of modern computing technology has brought significant transformation in data management and analysis. Today, we have access to powerful computing technology that facilitates data processing and analysis. Apache Impala and Apache Spark are two platforms that play a crucial role in the current data analysis ecosystem. Both are integral parts of Hadoop, an efficient open-source framework for storing and processing big data. Apache Impala is a query engine that enables the interactive execution of SQL queries on data stored in HDFS or HBase. On the other hand, Apache Spark is a platform that integrates streaming and batch data processing in a flexible framework. Both platforms have advantages and differences that need to be well-understood. This research aims to compare the query execution performance between Apache Impala and Apache Spark. Experiments were conducted by designing similar queries and running them on both platforms with existing datasets. Query execution times were recorded and analyzed to understand the relative performance of each platform in responding to query requests. The results of this research provide valuable insights into the speed and efficiency of query execution between Apache Impala and Apache Spark. This understanding can assist organizations and data practitioners in choosing a platform that suits their data analysis needs and context.

Keywords: Apache Impala, Apache Spark, Execution Time, Query, Data Processing

1. PENDAHULUAN

Dalam konteks analisis data yang semakin penting, performa dan efisiensi eksekusi kueri merupakan aspek utama dalam memilih platform yang tepat. Dua platform terkemuka yang sering digunakan untuk menganalisis data secara real-time adalah Apache Impala dan Apache Spark. Apache Impala memungkinkan eksekusi kueri SQL interaktif terhadap data yang disimpan, sementara Apache Spark memfasilitasi pemrosesan data yang terdistribusi. Penelitian ini bertujuan untuk membandingkan performa eksekusi kueri antara Apache Impala dan Apache Spark. Penelitian ini melibatkan uji coba eksekusi kueri pada data yang dihasilkan melalui pembuatan tabel dan pengisian konten. Hasil penelitian menunjukkan perbandingan performa antara kedua platform dalam merespon kueri. Penelitian ini memberikan wawasan yang berharga bagi para praktisi teknologi dan pengambil keputusan dalam memilih platform yang paling sesuai untuk kebutuhan analisis data yang efektif. Analisis perbandingan merupakan teknik analisis suatu aplikasi atau software yang dilakukan dengan cara memasukkan proses terhadap suatu aplikasi dan membandingkan antara satu menggunakan yg lain, serta memberikan keterangan atau data lain berupa kecepatan memproses data atau besar memori yang dibutuhkan untuk memproses sebuah data pada suatu aplikasi.

Apache Spark adalah suatu framework terdistribusi buat pengolahan data menggunakan standar yang kegunaannya mirip dengan Map-Reduce milik Apache Hadoop apalagi dipergunakan buat pengolahan sebagai parallel. Suatu komponen yang ada Apache Spark yang tak dipunyai sang platform yang lain ialah RDD (Resilient Distributed Dataset). RDD ialah memori tererovolusi yang bisa dilihat ketika kluster dijalankan sebagai parallel (Pan 2016). Pada saat pengoperasian di setiap RDD, nilai pembatas pada dalamnya akan menunjukkan kualitas paralelisme[1]. Apache Spark tersusun atas sebagian susunan yaitu (1) Spark SQL merupakan evolusi modern sesudah Shark SQL yang awalnya dipakai pada Spark. Spark SQL bisa mengerjakan pengarsipan memori atau pada Bahasa (In-Memory Columnar Storage) berasal dari Shark, Spark SQL bisa compatible menggunakan Hive suatu Query yang ada di Hadoop (Ivanov et al. 2014). Oleh sebab itu Spark SQL Ini mempunyai kemajuan besar pada hal pengumpulan data serta menumbuhkan kinerja (Han and Zhang 2016). (2) Spark Streaming ialah framework stream yang bergerak pada Apache Spark yang ada pada API serta spontan terhubung menggunakan streaming menjadi terhubung. Pada dasarnya Spark Streaming ini bisa handle kekurangan yang ada di Apache Hadoop pada bagian pemaparan pada streaming data secara realtime[1]. Apache Impala menyediakan kueri SQL berperforma tinggi dan berlatensi rendah pada data yang disimpan dalam format file Apache Hadoop yang populer. Respons cepat terhadap kueri memungkinkan eksplorasi interaktif dan penyesuaian kueri analitik, dibandingkan pekerjaan batch panjang yang biasanya dikaitkan dengan teknologi SQL-on-Hadoop. (Anda akan sering melihat istilah "interaktif" diterapkan pada kueri cepat semacam ini dengan waktu respons skala manusia[2].

Impala terintegrasi dengan database metastore Apache Hive, untuk berbagi database dan tabel antara kedua komponen. Integrasi tingkat tinggi dengan Hive, dan kompatibilitas dengan sintaksis HiveQL, memungkinkan Anda menggunakan Impala atau Hive untuk membuat tabel, mengeluarkan kueri, memuat data, dan sebagainya[2]. Berikut ini adalah beberapa keunggulan utama Impala:

- Impala terintegrasi dengan ekosistem CDH yang ada, artinya data dapat disimpan, dibagikan, dan diakses menggunakan berbagai solusi yang disertakan dengan CDH. Hal ini juga menghindari silo data dan meminimalkan perpindahan data yang mahal.

- Impala menyediakan akses ke data yang disimpan di CDH tanpa memerlukan keterampilan Java yang diperlukan untuk pekerjaan MapReduce. Impala dapat mengakses data langsung dari sistem file HDFS. Impala juga menyediakan front-end SQL untuk mengakses data di sistem database HBase, atau di Amazon Simple Storage System (S3) .
- Impala mengembalikan hasil biasanya dalam hitungan detik atau beberapa menit, bukan beberapa menit atau jam yang sering kali diperlukan untuk menyelesaikan kueri Hive.
- Impala memelopori penggunaan format file Parquet, tata letak penyimpanan berbentuk kolom yang dioptimalkan untuk kueri skala besar yang umum digunakan dalam skenario gudang data[2].

Perbedaan antara apache impala dan spark adalah impala merupakan *Kueri Waktu Nyata untuk Hadoop* . Impala adalah mesin kueri MPP SQL sumber terbuka dan modern untuk Apache Hadoop. Impala dikirimkan oleh Cloudera, MapR, dan Amazon. Dengan Impala, Anda dapat melakukan kueri data, baik yang disimpan di HDFS atau Apache HBase – termasuk fungsi SELECT, JOIN, dan agregat – secara real-time. Sedangkan apache spark merupakan Mesin cepat dan umum untuk pemrosesan data skala besar . Spark adalah mesin pemrosesan cepat dan umum yang kompatibel dengan data Hadoop. Itu dapat berjalan di kluster Hadoop melalui mode mandiri YARN atau Spark, dan dapat memproses data dalam HDFS, HBase, Cassandra, Hive, dan Hadoop InputFormat apa pun. Ini dirancang untuk melakukan pemrosesan batch (mirip dengan MapReduce) dan beban kerja baru seperti streaming, kueri interaktif, dan pembelajaran mesin[3][4].

2. TINJAUAN PUSTAKA

Kegagalan perangkat keras adalah hal yang biasa dan bukan pengecualian. Sebuah instansi HDFS bisa terbentuk dari ratusan atau ribuan mesin server, yang masing-masing menyimpan sebagian data sistem arsip. Berita bahwa adanya sejumlah besar komponen dan setiap komponen mempunyai kemungkinan kegagalan yg tak sepele berarti bahwa beberapa komponen HDFS selalu tidak berfungsi. Oleh karena itu, deteksi kesalahan dan pemulihan kesalahan secara cepat dan otomatis merupakan tujuan arsitektur inti HDFS. Pada project ini, kami melakukan perbandingan Apache Impala dengan Apache Spark dalam mengeksekusi kueri. Hadoop Distributed File System (HDFS) bisa dipaparkan berdasarkan harfiah dapat bermakna Sistem arsip tersalur pada Hadoop. meskipun disebut sistem arsip, tetapi HDFS tak sama dengan sistem arsip tipe sistem operasi, contohnya NTFS, FAT32, serta lain sebagainya. HDFS ialah sistem pengumpulan tersalur, yang melaksanakan metode memecah arsip besar jadi pecahan kecil serta setelah itu mendistribusikannya ke cluster komputer. Cluster ini umumnya terbentuk karena adanya banyak node ataupun computer atau server. Tiap node di cluster ini wajib menginstal Hadoop supaya dapat bekerja. Selaku sistem arsip terdistribusi, HDFS berfungsi buat mengatasi data besar yang ditaruh acak di segala cluster[3].

Apache Impala mengizinkan pengguna buat mengeksekusi kueri SQL berlatensi rendah pada data yang disimpan di HDFS serta Apache HBase tanpa membutuhkan perpindahan maupun transformasi data berkat pengenalan teknologi database paralel yang bisa diskalakan ke Hadoop. Impala berhubungan dengan Hadoop untuk menggunakan metadata, keamanan, manajemen sumber energi, dan format file serta informasi yang sama semacam Apache Hive, MapReduce, Apache Pig, serta aplikasi Hadoop lainnya. Impala dipasarkan untuk analis serta ilmuwan data untuk memakai SQL ataupun perlengkapan intelijen bisnis buat melaksanakan

analisis terhadap data yang ditaruh di Hadoop. Hasilnya, tidak butuh lagi mentransfer kumpulan data ke sistem spesial ataupun format kepemilikan buat melaksanakan analisis. Pencarian interaktif serta pemrosesan informasi skala besar (lewat MapReduce) bisa dicoba secara bertepatan pada sistem yang sama dengan memakai data serta metadata yang sama. Apache Spark adalah framework yang bisa digunakan buat mengambil informasi dari bermacam sumber, memprosesnya, dan setelah itu menaruh informasi yang sudah diproses di penyimpanan informasi buat dianalisis. Data engineer bisa membuat aplikasi alur pemrosesan Big Data memakai fitur-fitur Apache Spark. Apache Spark terintegrasi dengan beberapa bahasa pemrograman agar dapat mengubah himpunan data yang terdistribusi[5].

3. METODELOGI PENELITIAN

1.1 Jenis Penelitian

Penelitian pustaka (Library Research), yaitu penelitian yang dilakukan dengan menggunakan beberapa artikel sebagai referensi untuk penulisan.

3.2 Alat dan Bahan Penelitian

3.2.1 Perangkat Lunak (Software)

1. Virtual Box
2. Apache Impala
3. Apache Spark
4. Apache Hive
5. Cloudera
6. Hue Manager

3.2.2 Bahan Penelitian

1. Covid_19_indonesia_time_series_all.csv

3.2.3 Sistem Operasi Yang Digunakan

1. CentOS 6.7
2. Ubuntu 18.04

3.2.4 Perangkat Keras (Hardware)

1. Processor i5
2. Memory DDR4 8 GB
3. LocalDrive 20 GB

3.3 Jenis dan Pendekatan

Membandingkan performa apache Impala dan apache Spark dalam pengeksekusian kueri berupa data set “`covid_19_indonesia_time_series_all`” yang berisikan data sebanyak 32823 kolom. Penelitian ini akan menunjukkan perbedaan secara signifikan pada performa antara kedua software tersebut dalam mengeksekusi kueri.

3.4 Tahap Penelitian

Beberapa tahapan yang dilaksanakan selama penelitian, sebagai berikut:

- Langkah mengeksekusi query menggunakan Apache Impala dengan software cloudera manager

```
[cloudera@quickstart ~]$ hostname
quickstart.cloudera
[cloudera@quickstart ~]$ hdfs dfs -ls /
ls: NameNode still not started
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 6 items
drwxrwxrwx - hdfs supergroup      0 2017-10-24 00:15 /benchmarks
drwxr-xr-x - hbase supergroup      0 2023-10-10 10:06 /hbase
drwxr-xr-x - solr solr              0 2017-10-24 00:18 /solr
drwxrwxrwt - hdfs supergroup      0 2023-10-08 15:39 /tmp
drwxr-xr-x - hdfs supergroup      0 2017-10-24 00:17 /user
drwxr-xr-x - hdfs supergroup      0 2017-10-24 00:17 /var
[cloudera@quickstart ~]$ sudo /home/cloudera/cloudera-manager --express --force
[QuickStart] Shutting down CDH services via init scripts...
kafka-server: unrecognized service
[QuickStart] Disabling CDH services on boot...
error reading information on service kafka-server: No such file or directory
[QuickStart] Starting Cloudera Manager server...
[QuickStart] Waiting for Cloudera Manager API...
[QuickStart] Starting Cloudera Manager agent...
[QuickStart] Configuring deployment...
Submitted jobs: 502
[QuickStart] Deploying client configuration...
Submitted jobs: 503
[QuickStart] Starting Cloudera Management Service...
Submitted jobs: 511
[QuickStart] Enabling Cloudera Manager daemons on boot...

Success! You can now log into Cloudera Manager from the QuickStart VM's browser:

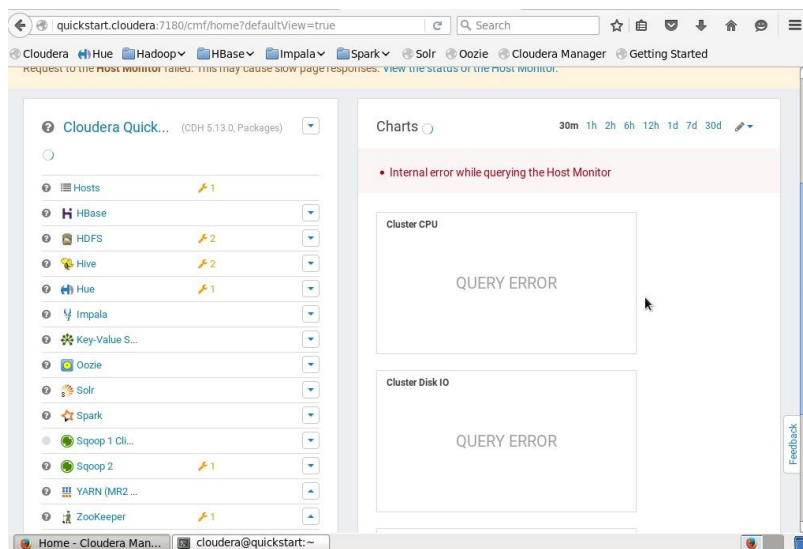
http://quickstart.cloudera:7180

Username: cloudera
Password: cloudera

[cloudera@quickstart ~]$
```

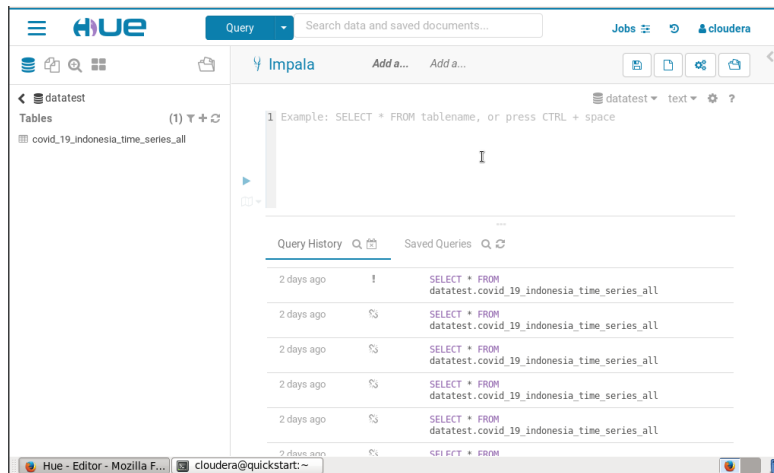
Gambar 1. Cloudera Manager

Setelah menyalakan Virtual Machine, kita jalankan command `hdfs dfs /` untuk memberi permission terhadap sistem. Dilanjutkan dengan command `sudo /home/cloudera/cloudera-manager --express --force` untuk mengaktifkan client cloudera manager. Saat client sudah aktif, akan muncul link serta username dan password yang akan kita gunakan untuk login ke cloudera manager.



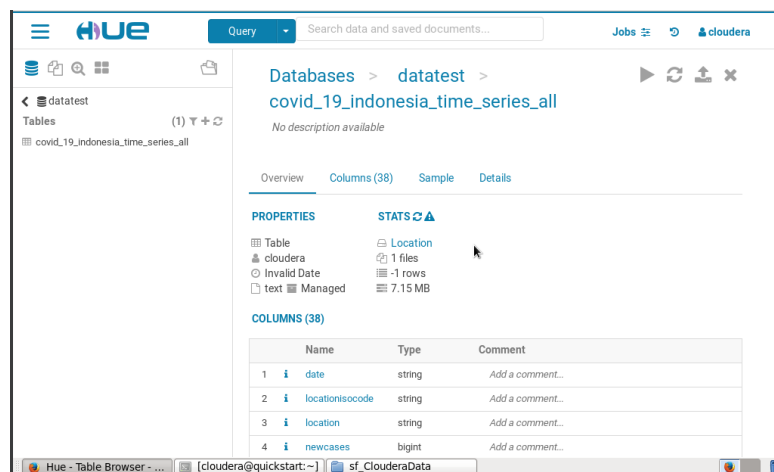
Gambar 2. Impala

Setelah login, kita harus menjalankan program HBase, HDFS, Hive, Impala, YARN, Zookeeper dan Hue yang diperlukan untuk mengeksekusi query dengan Impala. Pengeksekusian query akan dilakukan di halaman Hue yang perlu kita buka terlebih dahulu.



Gambar 3. Membuka Hue

Saat kita membuka Hue, akan muncul tampilan seperti diatas. Kolom ditengah berfungsi sebagai tempat menaruh command untuk pengeksekusian query. Kita perlu memasukkan dataset dengan cara menambahkannya kedalam table, tekan simbol “+” pada bagian kiri.



Gambar 4. Dataset penelitian

Setelah kita menambahkan dataset, akan muncul preview seperti diatas. Kita perlu memberi izin terlebih dahulu kepada sistem untuk bisa mengakses data dalam dataset tersebut dengan menggunakan command \$hdfs dfs https://:cloudera/user/cloudera.


```
scala> import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.SparkSession

scala> val session = SparkSession.builder().appName("app_name").master("local").getOrCreate()
23/10/11 08:02:19 WARN SparkSession: Using an existing Spark session; only runtime SQL configurations will take effect.
session: org.apache.spark.sql.SparkSession = org.apache.spark.sql.SparkSession@3c01915e

scala> val training = session.read.format("csv").load("/home/vboxuser/Downloads/covid_19_indonesia_time_series_all.csv")
training: org.apache.spark.sql.DataFrame = [_c0: string, _c1: string ... 36 more fields]

scala> training.show()
23/10/11 08:03:14 WARN SparkStringUtils: Truncated the string representation of a plan since it was too large. This behavior can be adjusted by setting 'spark.sql.debug.maxToStringFields'.
-----
|_c0|_c1|_c2|_c3|_c4|_c5|_c6|_c7|_c8|_c9|_c10|_c11|_c12|_c13|_c14|_c15|_c16|_c17|_c18|_c19|_c20|_c21|_c22|_c23|_c24|_c25|_c26|_c27|_c28|_c29|_c30|_c31|_c32|_c33|
-----
|Date|Location ISO Code|Location|New Cases|New Deaths|New Recovered|New Active Cases|Total Cases|Total Deaths|Total Recovered|Total Active Cases|Location Level|City or Regency|Province|Country|Continent|Island|Time Zone|Special Status|Total Regencies|Total Cities|Total Districts|Total Urban Villages|Total Rural Villages|Area (Km2)|Population|Population Density|Longitude|Latitude|New Cases per M...|Total Cases per M...|New Deaths per M...|Total Deaths per ...|Case Fatality Rate|Case Recovered Rate|Growth Factor of ...|Growth Factor of ...|
|3/1/2020|ID-JK|DKI Jakarta|Asia|Jawa|UTC+07:00|Daerah Khusus Ibu...|1|39|20|44|75|-56|Province|NULL|NULL|0.18|18086145|16334.31|106.8361183|-6.204698991|NULL|0.18|664|10846145|51.28%|192.31%|NULL|NULL|0.18|2|41|20|44|75|-54|Province|NULL|NULL|0.18|16334.31|106.8361183|-6.204698991|0.18|3.78|0.00|0.00|
-----
```

Gambar 7. Eksekusi Query

Untuk mulai mengeksekusi query, kita harus memulai spark session terlebih dahulu dengan command `import org.apache.spark.sql.SparkSession`. Lalu kita build sesi spark dengan command `val session = SparkSession.builder().appName("app_name").master("local").getOrCreate()`. Setelah itu spark session akan aktif. Lalu kita masukan dataset yang berupa csv dengan command `val training = session.read.format("csv").load("/home/vboxuser/Downloads/covid_19_indonesia_time_series_all.csv")` dimana "training" adalah nama database kita. Lalu kita panggil database sesuai namanya dengan command `<database>.show()`. Setelah itu spark akan mulai mengeksekusi data yang terdapat pada file csv yang telah kita tuliskan pathnya.

```
root@Ubuntu: /home/vboxuser
File Edit View Search Terminal Help
416| 98| 7230| 8488| 74953| 1916907| 265185520| 138.34| 113.921327| -
0.789275| 0.00%| 0.00| 1.00| 0.01| 1.00| 0.00| 0.00| 0.00| 0.00| 0.00| 0.00|
|3/5/2020|-58|ID-JB|Jawa Barat|1|0|0|0|5312|35378|45161325|1276.55|107.6037083|-6.9
10|91|627|645|0.07|0.00|0.00|0.02|0.00|0.00|0.00|0.00|
20432083|2000.00%|1.00|0.07|1.00|0.00|0.00|0.02|0.00|0.00|0.00|0.00|
|3/5/2020|0|ID-RI|Riau|0|0|0|0|0|0|0|0|1|0|1|
10|2|169|268|1591|87024|6074100|69.80|101.8051092|0.5
11647851|100.00%|1.00|0.16|1.00|0.00|0.00|0.00|0.00|0.00|0.00|0.00|
|3/6/2020|-115|ID-BT|Banten|1|0|1|1|1238|9663|10722374|1109.64|106.1090043|-6.4
4|4|155|313|0.09|0.00|0.00|0.47|0.00|0.00|0.00|0.00|
56736388|11100.00%|NULL|NULL|NULL|NULL|NULL|NULL|0.47|0.00|0.00|0.00|
|3/6/2020|-51|ID-JK|DKI Jakarta|0|0|0|0|0|0|45|21|75|
1|5|44|267|0.00|0.00|0.00|1.94|664|10846145|16334.31|106.8361183|-6.2
04698991|166.67%|1.00|4.15|0.00|0.00|0.00|1.94|664|10846145|16334.31|106.8361183|-6.2
|3/6/2020|4|IDN|Indonesia|2|0|0|0|2|4|0|0|0|0|
416| 98| 7230| 8488| 74953| 1916907| 265185520| 138.34| 113.921327| -
0.789275| 0.00%| 0.01| 1.00| 0.01| 1.00| 0.00| 0.00| 0.00| 0.00| 0.00| 0.00|
-----
only showing top 20 rows
Time taken: 460 ms
```

Gambar 8. Preview isi data

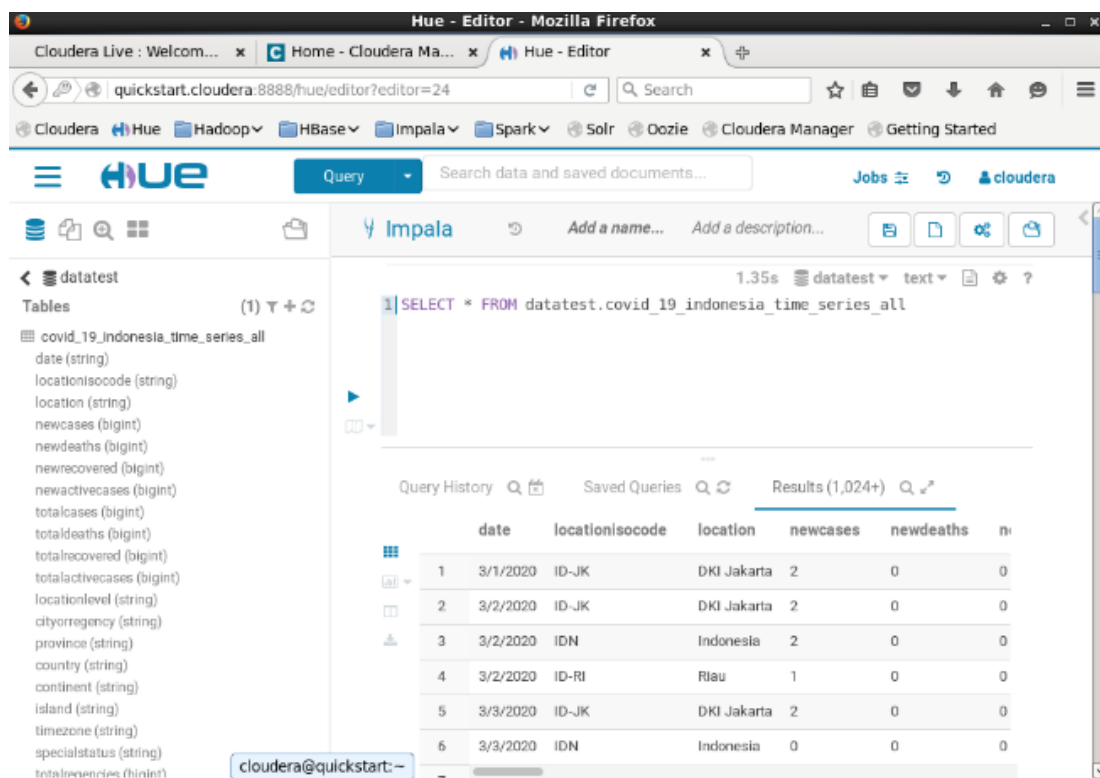
Setelah selesai, spark akan menunjukkan preview isi data dan waktu yang diperlukan untuk menyelesaikan query

4. HASIL PENELITIAN

Dalam percobaan ini, kami menggunakan file CSV yang berisi data "covid_19_indonesia_time_series_all" dengan 32,823 kolom. Penelitian ini bertujuan untuk

membandingkan kinerja eksekusi query antara Apache Impala dan Apache Spark pada kluster big data. Pemilihan Apache Impala didasarkan pada keefektifannya sebagai framework open-source untuk menyimpan dataset dalam jumlah besar dan kemampuannya memproses data dari ukuran gigabyte hingga petabyte secara efisien. Sebagai pesaing utama, Apache Spark dipilih karena kemampuannya mengelola dan memproses data besar dengan memanfaatkan prosesor multi-core serta melakukan pemrosesan data secara distribusi. Hasil percobaan menunjukkan bahwa Apache Impala unggul dalam eksekusi query pada dataset dengan skala besar, terutama dalam hal efisiensi dan kecepatan. Meskipun demikian, Apache Spark tetap memberikan kinerja yang baik, terutama saat menangani query pada dataset yang lebih kompleks dan memanfaatkan prosesor multi-core. Temuan ini memberikan wawasan penting bagi pengembang dan analis data dalam memilih platform yang sesuai dengan kebutuhan spesifik mereka, terutama dalam konteks pengolahan dan analisis data pada skala besar.

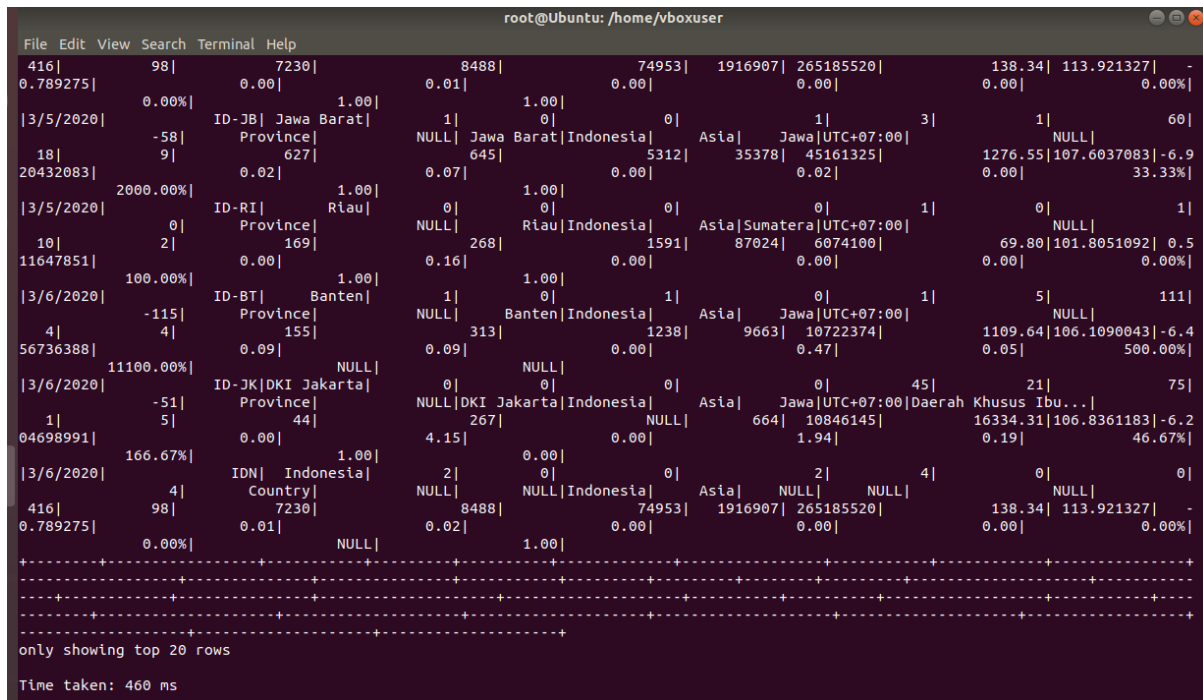
1.Impala



Gambar 9. Waktu eksekusi

Waktu yang dibutuhkan untuk menampilkan file csv pada impala dicloudera adalah 1.35s.

2. Spark



Gambar 10. Proses pada apache spark

Di Spark, waktu yang dibutuhkan untuk menghasilkan tabel seperti di atas adalah sekitar 460 milidetik atau sekitar 0,46 detik untuk mengeksekusinya.

5. KESIMPULAN

Berdasarkan hasil penelitian ini, dapat disimpulkan bahwa Apache Spark unggul dibandingkan Apache Impala dalam kecepatan eksekusi query pada dataset besar. Apache Spark menunjukkan kinerja eksekusi yang lebih cepat, terutama saat memanfaatkan pemrosesan paralel dan distribusi yang optimal. Meskipun Apache Impala efisien dalam menangani dataset besar, Apache Spark memberikan performa eksekusi yang lebih baik dan efisien, terutama dalam pemrosesan data pada skala yang lebih besar. Meskipun Apache Impala memiliki keunggulan dalam mendukung pemrosesan paralel langsung tanpa melalui MapReduce, Apache Spark tetap kompetitif, terutama dengan memanfaatkan prosesor multi-core untuk mempercepat pemrosesan query pada dataset besar. Oleh karena itu, dapat disimpulkan bahwa Apache Spark menawarkan kinerja eksekusi yang lebih cepat dibandingkan dengan Apache Impala, dengan efisiensi dan kecepatan sebagai faktor utama dalam memilih framework untuk menangani query pada data besar.

DAFTAR PUSTAKA

- [1]. Armbrust, M., Xin, R. S., Lian, C., Huai, Y., Liu, D., Bradley, J. K., ... & Zaharia, M. (2015, May). Spark sql: Relational data processing in spark. In Proceedings of the 2015 ACM SIGMOD international conference on management of data (pp. 1383-1394).

- [2]. Bittorf, M. K. A. B. V., Bobrovytsky, T., Erickson, C. C. A. C. J., Hecht, M. G. D., Kuff, M. J. I. J. L., Leblang, D. K. A., ... & Yoder, M. M. (2015, January). Impala: A modern, open-source SQL engine for Hadoop. In Proceedings of the 7th biennial conference on innovative data systems research (pp. 1-10).
- [3]. Wibawa, C., Wirawan, S., Mustikasari, M., & Anggraeni, D. T. (2022). KOMPARASI KECEPATAN HADOOP MAPREDUCE DAN APACHE SPARK DALAM MENGOLAH DATA TEKS. *Jurnal Ilmiah MATRIK*, 24(1), 10-20.
- [4]. Aminudin, A., & Cahyono, E. B. (2019). Pengukuran Performa Apache Spark dengan Library H2O Menggunakan Benchmark Hibench Berbasis Cloud Computing. *Jurnal Teknologi Informasi dan Ilmu Komputer (JTIIK)*, 6(5), 519-526.
- [5]. Cloudera. (n.d.). Impala Documentation. Retrieved from <https://docs.cloudera.com/documentation/enterprise/6/6.3/topics/impala.html>