

SKRIPSI

**STUDI DAN PERBANDINGAN APACHE SPARK SQL DAN
HIVE DALAM KONTEKS ANALISIS BIG DATA**



Lydia Febtriani

NPM: 2014730044

**PROGRAM STUDI TEKNIK INFORMATIKA
FAKULTAS TEKNOLOGI INFORMASI DAN SAINS
UNIVERSITAS KATOLIK PARAHYANGAN
2018**

UNDERGRADUATE THESIS

**STUDY AND COMPARISON OF APACHE SPARK SQL AND
HIVE FOR BIG DATA ANALYSIS**



Lydia Febtriani

NPM: 2014730044

**DEPARTMENT OF INFORMATICS
FACULTY OF INFORMATION TECHNOLOGY AND SCIENCES
PARAHYANGAN CATHOLIC UNIVERSITY
2018**

LEMBAR PENGESAHAN



**STUDI DAN PERBANDINGAN APACHE SPARK SQL DAN
HIVE DALAM KONTEKS ANALISIS BIG DATA**

Lydia Febtriani

NPM: 2014730044

Bandung, 21 Mei 2018

Menyetujui,

Pembimbing Utama

Pembimbing Pendamping



Dr. Veronica Sri Moertini



Gede Karya, M.T., CISA, IPM

Ketua Tim Penguji

Anggota Tim Penguji



Vania Natali, M.T.



Claudio Franciscus, M.T.

Mengetahui,

Ketua Program Studi



Mariskha Tri Adithia, P.D.Eng



PERNYATAAN

Dengan ini saya yang bertandatangan di bawah ini menyatakan bahwa skripsi dengan judul:

STUDI DAN PERBANDINGAN APACHE SPARK SQL DAN HIVE DALAM KONTEKS ANALISIS BIG DATA

adalah benar-benar karya saya sendiri, dan saya tidak melakukan penjiplakan atau pengutipan dengan cara-cara yang tidak sesuai dengan etika keilmuan yang berlaku dalam masyarakat keilmuan.

Atas pernyataan ini, saya siap menanggung segala risiko dan sanksi yang dijatuhkan kepada saya, apabila di kemudian hari ditemukan adanya pelanggaran terhadap etika keilmuan dalam karya saya, atau jika ada tuntutan formal atau non-formal dari pihak lain berkaitan dengan keaslian karya saya ini.

Dinyatakan di Bandung,
Tanggal 21 Mei 2018



Lydia Febtriani
NPM: 2014730044

ABSTRAK

Seiring kemajuan dan pemanfaatan teknologi, data bertambah dengan cepat dalam waktu singkat sehingga ukuran data menjadi sangat besar atau sering disebut *big data*. Agar pengguna bisa mendapatkan informasi dari *big data*, *big data* harus diolah dan dianalisis dengan cepat. Beberapa teknologi yang dapat digunakan untuk mengolah dan menganalisis *big data*, yaitu Apache Hadoop dan Apache Spark. Hadoop adalah sistem terdistribusi yang dapat digunakan untuk mengolah dan menganalisis *big data* dengan cepat. Spark merupakan *framework* komputasi di *cluster*. Spark dioptimalkan untuk berjalan di memori sehingga pemrosesan data lebih cepat dibandingkan komponen komputasi Hadoop, MapReduce.

Hadoop dan Spark memiliki subproyek yang bekerja dengan data terstruktur, yaitu Apache Hive dan Spark SQL. Hive merupakan infrastruktur data *warehouse* berbasis Hadoop, sedangkan Spark SQL adalah *library* Spark untuk bekerja dengan data terstruktur. Hive dan Spark menyediakan pengerjaan kueri dengan atau mirip sintaks SQL. Spark SQL dan Hive yang digunakan untuk tujuan yang sama, yaitu mengolah dan menganalisis *big data* dengan dialek SQL, perlu dipelajari dan diteliti mengenai kesamaan dan perbedaan penggunaan sintaks SQL serta bagaimana performasinya dalam mengerjakan kueri dengan *big data*.

Penelitian ini dilakukan untuk mencari perbedaan Spark dan Hive. Perbedaan tersebut dapat berupa penggunaan sintaks SQL dan bagaimana performa Spark SQL dan Hive dalam mengerjakan kueri. Eksperimen dilakukan untuk melihat performa Spark dan Hive dalam mengerjakan kueri-kueri SELECT. Langkah penelitian yang dilakukan adalah mempelajari konsep dan pemanfaatan Spark SQL dan Hive, mengeksplorasi fitur SQL pada Spark dan Hive, menyiapkan *cluster* Hadoop dan Spark, mengumpulkan set *big data* untuk eksperimen, mengembangkan perangkat lunak berbasis *web* untuk eksperimen, dan melakukan eksperimen perbandingan kinerja Spark SQL dan Hive. Hasil dari penelitian ini berupa perangkat lunak berbasis *web* yang dapat menjalankan perangkat lunak dengan Spark untuk mengerjakan kueri-kueri Spark SQL dan Hive. Perangkat lunak *web* menampilkan daftar kueri untuk eksperimen dan dapat menampilkan hasil kueri serta waktu eksekusi kueri Spark SQL dan Hive. Kueri yang dipilih merupakan kueri-kueri SELECT yang diurutkan dari kueri sederhana hingga kueri yang lebih kompleks. Kueri-kueri yang digunakan sama baik untuk Spark SQL maupun Hive. Eksperimen dilakukan sebanyak tiga kali menggunakan perangkat lunak berbasis *web*. Data yang digunakan terdiri dari dua data, yaitu data Movies dengan ukuran 358 MB dan data Ratings dengan ukuran 2,6 GB. Hasil yang didapat dari eksperimen adalah waktu eksekusi kueri-kueri SELECT menggunakan Spark SQL dan Hive. Berdasarkan tiga kali eksperimen waktu eksekusi kueri-kueri SELECT dengan Spark SQL dan Hive, Spark SQL lebih cepat dalam menjalankan kueri-kueri SELECT dibandingkan Hive.

Kesimpulan pertama yang didapat dari penelitian ini adalah Spark SQL dan Hive digunakan untuk bekerja dengan data terstruktur. Tetapi, Spark SQL memiliki abstraksi RDD yang dirancang untuk penyimpanan data dalam memori sedangkan Hive menggunakan HDFS Hadoop sebagai komponen penyimpanan data. Hive juga mendukung lebih banyak pernyataan SQL dibandingkan Spark SQL. Kesimpulan dari eksperimen-eksperimen waktu eksekusi kueri Spark SQL dan Hive adalah Spark SQL lebih cepat dalam mengerjakan kueri dibandingkan Hive.

Kata-kata kunci: Apache Spark, Spark SQL, Apache Hadoop, Hive, *Big Data*

ABSTRACT

With current advanced technology and its utilization, in short period of time, the size of data has increased which results in what we know as big data. As users should be able to retrieve information from a significant amount of data, the system must be able to process and analyze it faster. Apache Hadoop is a distribution system which can efficiently handle the computation and analysis of big data. Meanwhile, Apache Spark is a computation framework that works on a cluster. One main advantage of Spark, compared to Hadoop's MapReduce, is in their processing speed since the computation runs in memory.

Both Apache Hadoop and Spark's subproject work with structured data, known as Apache Hive and Spark SQL. Hive is a data warehouse infrastructure with Hadoop as its basis, and Spark SQL is the Spark's library for working with structured data. Both Hive and Spark allow queries with SQL syntax or SQL-like syntax. Since both of them have the same purpose of processing and analyzing big data, Spark and Hive can be studied and analyzed further not only in the similarities and differences among usages of SQL syntax but also their performance on processing big data query.

This research aims at comparing Spark and Hive. The comparison can be in the form of SQL syntax usage and their performance on processing queries. The experiment demonstrates the performance of Spark and Hive on processing the SELECT queries. The research studies the concept and application of Spark SQL and Hive, explores the features of SQL in Spark SQL and Hive, prepares Hadoop and Spark clusters, collects big dataset for experiments, develops web-based application, and compares the performance of Spark SQL and Hive. The result of this research is a web application which can run Spark application to execute Spark SQL and Hive queries. The web application shows Spark SQL and Hive query lists for the experiments, the query results, and the execution time. The selected queries are SELECT queries, sorted from the simplest to more complex one. The queries for Spark SQL and Hive are same. The experiments are conducted three times using the web-based application. The data being used for this experiments are the Movies data, which is 358 MB in size, and the Ratings data, which is 2.6 GB in size. The result of this experiment is execution time SELECT queries for Spark SQL and Hive. Based on the experiments of SELECT queries, Spark SQL executes the commands faster than Hive.

The first conclusion of this research is that Spark SQL and Hive can be used to work with structured data. But the Spark SQL has RDD abstraction which is specifically designed for in-memory storage, while Hive uses HDFS Hadoop as the storage component. In contrast, Hive supports more SQL commands than Spark SQL. From the experiments of Spark SQL and Hive's query execution time, the Spark SQL runs faster than Hive.

Keywords: Apache Spark, Spark SQL, Apache Hadoop, Hive, Big Data

*Dipersembahkan untuk orang tua, kedua kakak, dan keluarga
penulis*

KATA PENGANTAR

Puji dan syukur kepada Tuhan Yang Maha Esa atas dukungan-Nya sehingga penulis dapat menyelesaikan skripsi ini tepat pada waktunya. Pada kesempatan ini, penulis ingin menyampaikan rasa terima kasih kepada:

1. Kedua orang tua, kedua kakak, dan keluarga besar penulis yang telah mendukung penulis selama pengerjaan skripsi ini.
2. Ibu Veronica Sri Moertini sebagai dosen pembimbing utama dan Pak Gede Karya sebagai dosen pembimbing serta yang telah membimbing dan membantu penulis dalam mengerjakan skripsi ini.
3. Teman-teman seperjuangan yang saling membantu dan saling curhat, Tania dan Gabriella.
4. Teman-teman yang selalu menyemangati penulis, Vanessa, Tania, Gabriella, Betari.
5. Rekan-rekan mahasiswa bimbingan Ibu Veronica dan rekan-rekan dengan topik skripsi Spark yang telah saling membantu dan saling memberi informasi.
6. Rekan-rekan dengan topik skripsi Hadoop yang saling mengerti karena sama-sama menggunakan *cluster* Hadoop.
7. Rekan-rekan seangkatan Teknik Informatika yang saling memberi semangat.
8. Kakak perempuan dan tante yang selalu mendampingi dan menyemangati penulis.
9. Kakak laki-laki dan sepupu terdekat penulis, Gebbie De Nova, yang telah membantu dalam penulisan judul dan abstrak bahasa inggris.
10. Teman-teman SMA penulis yang memberi semangat ke penulis Jean, Nelly dan teman-teman lain yang tidak bisa disebutkan satu-per-satu.
11. Lagu-lagu Korea yang selalu menemani penulis dalam pengerjaan skripsi ini serta drama dan acara Korea yang menjadi penghibur penulis di saat jenuh dalam pengerjaan skripsi ini.

Akhir kata, penulis menyadari bahwa penelitian dan skripsi ini masih memiliki kekurangan. Penulis memohon maaf bila terdapat kesalahan pada penelitian atau skripsi ini. Penulis juga menerima kritik dan saran sehingga penelitian ini dapat berkembang menjadi lebih baik. Semoga penelitian ini dapat membantu penelitian-penelitian lainnya.

Bandung, Mei 2018

Penulis

DAFTAR ISI

KATA PENGANTAR	xv
DAFTAR ISI	xvii
DAFTAR GAMBAR	xxi
DAFTAR TABEL	xxv
1 PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	1
1.3 Tujuan	2
1.4 Batasan Masalah	2
1.5 Metodologi	2
1.6 Sistematika Pembahasan	3
2 LANDASAN TEORI	5
2.1 <i>Big Data</i>	5
2.2 Scala	6
2.2.1 Variabel	6
2.2.2 <i>String</i>	7
2.2.3 Angka	7
2.2.4 <i>Array</i>	8
2.2.5 Kondisional	8
2.2.6 <i>Loops</i>	9
2.2.7 Fungsi	10
2.3 Structured Query Language (SQL)	11
2.3.1 Deskripsi Structured Query Language (SQL)	11
2.3.2 Basis Data Relasional	11
2.3.3 <i>Primary Key</i> dan <i>Foreign Key</i>	12
2.3.4 Operasi SQL	12
2.3.5 Pernyataan JOIN	18
2.4 Apache Hadoop	20
2.4.1 Deskripsi Hadoop	20
2.4.2 Komponen Inti Hadoop	21
2.5 Apache Hive	28
2.5.1 Deskripsi Apache Hive	28
2.5.2 Operasi-operasi Hive	31
2.6 Apache Spark	40
2.6.1 Deskripsi Apache Spark	40
2.6.2 Spark Core	42
2.6.3 Apache Spark SQL	45

3	STUDI EKSPLORASI	53
3.1	Instalasi dan Konfigurasi	53
3.1.1	Instalasi dan Konfigurasi Hive	53
3.1.2	Instalasi dan Konfigurasi Spark	53
3.1.3	Instalasi dan Konfigurasi IntelliJ IDEA	54
3.2	Eksperimen	55
3.2.1	Eksperimen Spark Sederhana	55
3.2.2	Eksperimen Spark SQL Menggunakan Akka	59
3.2.3	Eksperimen Spark SQL Menggunakan HttpServer	60
3.2.4	Eksperimen Hive	61
4	ANALISIS DAN PERANCANGAN	73
4.1	Analisis Perangkat Lunak	73
4.2	Analisis Data Studi Kasus	73
4.3	Analisis Praolah Data	74
4.3.1	Praolah Set Data <i>Movies</i>	74
4.3.2	Praolah <i>Big Data</i>	75
4.4	Analisis Kueri	75
4.5	Perancangan Perangkat Lunak dengan Spark	77
4.5.1	Diagram Kelas	78
4.5.2	Deskripsi Kelas dan <i>Method</i>	78
4.6	Perancangan Perangkat Lunak <i>Web</i>	80
4.6.1	Diagram Use Case	80
4.6.2	Diagram Kelas	82
4.6.3	Perancangan Tampilan Antarmuka	85
5	IMPLEMENTASI DAN EKSPERIMEN	89
5.1	Deskripsi Lingkungan Perangkat yang Digunakan	89
5.1.1	Lingkungan Implementasi Perangkat Keras	89
5.1.2	Lingkungan Implementasi Perangkat Lunak	89
5.2	Implementasi Spark dan Spark SQL	90
5.2.1	Implementasi SparkContext	90
5.2.2	Implementasi RDD	90
5.2.3	Implementasi Spark SQL	92
5.3	Implementasi Hive	93
5.3.1	Implementasi Membuat Tabel dan <i>Load Data</i>	93
5.3.2	Implementasi Koneksi JDBC HiveServer2	94
5.3.3	Implementasi Kueri Hive dengan Koneksi JDBC	95
5.4	Implementasi Modul Membaca Kueri dan Menulis Hasil Kueri ke HDFS	96
5.5	Pengujian Kueri	98
5.6	Pengujian Hasil Kueri	101
5.6.1	Pengujian Hasil Kueri Set Data <i>Movies</i>	101
5.6.2	Pengujian Hasil Kueri Set Data <i>Ratings</i>	115
5.7	Eksperimen Waktu Eksekusi Kueri	129
5.7.1	Eksperimen Waktu Eksekusi dengan Set Data <i>Movies</i>	129
5.7.2	Eksperimen Waktu Eksekusi dengan Set Data <i>Ratings</i>	131
5.7.3	Eksperimen Waktu Eksekusi Kueri JOIN	133
5.7.4	Kesimpulan Hasil Eksperimen Waktu Eksekusi Spark SQL dan Hive	134
6	KESIMPULAN DAN SARAN	135
6.1	Kesimpulan	135
6.2	Saran	136

DAFTAR REFERENSI	137
A KONFIGURASI HIVE	139
B KODE PROGRAM UNTUK EKSPERIMEN	141
C KODE PROGRAM UNTUK DATASET CLEANER	157
D KODE PROGRAM UNTUK PERANGKAT LUNAK DENGAN SPARK	159
E KODE PROGRAM UNTUK PERANGKAT LUNAK WEB	165

DAFTAR GAMBAR

2.1	Arsitektur Hadoop [1]	22
2.2	Arsitektur HDFS dan Komunikasi NameNode dengan DataNode [1]	24
2.3	Proses Membaca dari HDFS [2]	24
2.4	Proses Menulis ke HDFS [2]	25
2.5	Arsitektur MapReduce dan Komunikasi JobTracker dengan TaskTracker [3]	27
2.6	Algoritma Word Count dengan MapReduce [4]	28
2.7	Arsitektur Hive [4]	29
2.8	Arsitektur Spark [5]	41
3.1	Eksperimen Kedua Spark Menggunakan Spark-Shell (1)	56
3.2	Eksperimen Kedua Spark Menggunakan Spark-Shell (2)	56
3.3	SparkSQL-IntelliJ-1: println(df.count)	57
3.4	SparkSQL-IntelliJ-1: df.show()	57
3.5	SparkSQL-IntelliJ-1: SELECT * FROM tabel WHERE Region='ar' AND Position='1'	58
3.6	SparkSQL-IntelliJ-1: SELECT Region FROM tabel	58
3.7	SparkSQL-IntelliJ-1: SELECT COUNT(DISTINCT(Region)) FROM tabel	59
3.8	SparkSQL-IntelliJ-1: spark.catalog.listTables.show	59
3.9	Tampilan Web Eksperimen Spark SQL dengan Akka	59
3.10	Tampilan Utama Web Spark SQL Menggunakan HttpServer	61
3.11	Hive-Interaktif: SHOW DATABASES	62
3.12	Hive-Interaktif: DESCRIBE DATABASE EXTENDED mydb	62
3.13	Hive-Interaktif: DESCRIBE mytable	63
3.14	Hive-Interaktif: SELECT * FROM mytable	64
3.15	Hive-Interaktif: SELECT DISTINCT huruf FROM mytable	64
3.16	Hive-Interaktif: SELECT * FROM mytable WHERE nomor<3	65
3.17	Hive-Interaktif: Subkueri dengan Format WITH	65
3.18	Hive-Interaktif: Subkueri di FROM	66
3.19	Hive-Interaktif: Subkueri di WHERE	67
3.20	Hive-Noninteraktif: SELECT * FROM movies	68
3.21	Hive-Noninteraktif: SELECT * FROM movies WHERE year = 1997	68
3.22	Hive-Noninteraktif: Subkueri di FROM	69
3.23	Hive-LOAD: SELECT * FROM movies2 LIMIT 5	70
3.24	Hive-LOAD: SELECT id FROM movies2 WHERE title='Toy Story (1995)'	70
3.25	Hive-LOAD: SELECT title FROM movies2 WHERE id=1	71
3.26	Hive-LOAD: SELECT genres FROM movies2 WHERE id=1	71
3.27	Hive-LOAD: SELECT COUNT(id) FROM movies2	72
4.1	Diagram Kelas Perangkat Lunak dengan Spark	78
4.2	Diagram Use Case Perangkat Lunak Web	80
4.3	Diagram Kelas Perangkat Lunak Web	83
4.4	Rancangan Tampilan Halaman Utama	86
4.5	Rancangan Tampilan Halaman Daftar Kueri	86
4.6	Rancangan Tampilan Halaman Hasil Kueri Spark SQL	87

4.7	Rancangan Tampilan Halaman Hasil Kueri Hive	87
5.1	Tampilan Awal Perangkat Lunak <i>Web</i>	99
5.2	Daftar Kueri untuk Set Data Movies	99
5.3	Tampilan Daftar Kueri untuk Set Data Movies	100
5.4	Daftar Kueri untuk Set Data Ratings	100
5.5	Tampilan Daftar Kueri untuk Set Data Ratings	101
5.6	Hasil Kueri 1 Spark SQL dengan Set Data Movies	102
5.7	Hasil Kueri 1 Hive dengan Set Data Movies	102
5.8	Hasil Kueri 2 Spark SQL dengan Set Data Movies	103
5.9	Hasil Kueri 2 Hive dengan Set Data Movies	103
5.10	Hasil Kueri 3 Spark SQL dengan Set Data Movies	104
5.11	Hasil Kueri 3 Hive dengan Set Data Movies	104
5.12	Hasil Kueri 4 Spark SQL dengan Set Data Movies	105
5.13	Hasil Kueri 4 Hive dengan Set Data Movies	105
5.14	Hasil Kueri 5 Spark SQL dengan set data Movies	106
5.15	Hasil Kueri 5 Hive dengan Set Data Movies	106
5.16	Hasil Kueri 6 Spark SQL dengan Set Data Movies	107
5.17	Hasil Kueri 6 Hive dengan Set Data Movies	107
5.18	Hasil Kueri 7 Spark SQL dengan Set Data Movies	108
5.19	Hasil Kueri 7 Hive dengan Set Data Movies	108
5.20	Hasil Kueri 8 Spark SQL dengan Set Data Movies	109
5.21	Hasil Kueri 8 Hive dengan Set Data Movies	109
5.22	Hasil Kueri 9 Spark SQL dengan Set Data Movies	110
5.23	Hasil Kueri 9 Hive dengan Set Data Movies	110
5.24	Hasil Kueri 10 Spark SQL dengan Set Data Movies	111
5.25	Hasil Kueri 10 Hive dengan Set Data Movies	111
5.26	Hasil Kueri 11 Spark SQL dengan Set Data Movies	112
5.27	Hasil Kueri 11 Hive dengan Set Data Movies	112
5.28	Hasil Kueri 12 Spark SQL dengan Set Data Movies	113
5.29	Hasil Kueri 12 Hive dengan Set Data Movies	113
5.30	Hasil Kueri 13 Spark SQL dengan Set Data Movies	114
5.31	Hasil Kueri 13 Hive dengan Set Data Movies	114
5.32	Hasil Kueri 14 Spark SQL dengan Set Data Movies	115
5.33	Hasil Kueri 14 Hive dengan Set Data Movies	115
5.34	Hasil Kueri 1 Spark SQL dengan Set Data Ratings	116
5.35	Hasil Kueri 1 Hive dengan Set Data Ratings	116
5.36	Hasil Kueri 2 Spark SQL dengan Set Data Ratings	117
5.37	Hasil Kueri 2 Hive dengan Set Data Ratings	117
5.38	Hasil Kueri 3 Spark SQL dengan Set Data Ratings	118
5.39	Hasil Kueri 3 Hive dengan Set Data Ratings	118
5.40	Hasil Kueri 4 Spark SQL dengan Set Data Ratings	119
5.41	Hasil Kueri 4 Hive dengan Set Data Ratings	119
5.42	Hasil Kueri 5 Spark SQL dengan Set Data Ratings	120
5.43	Hasil Kueri 5 Hive dengan Set Data Ratings	120
5.44	Hasil Kueri 6 Spark SQL dengan Set Data Ratings	121
5.45	Hasil Kueri 6 Hive dengan Set Data Ratings	121
5.46	Hasil Kueri 7 Spark SQL dengan Set Data Ratings	122
5.47	Hasil Kueri 7 Hive dengan Set Data Ratings	122
5.48	Hasil Kueri 8 Spark SQL dengan Set Data Ratings	123
5.49	Hasil Kueri 8 Hive dengan Set Data Ratings	123
5.50	Hasil Kueri 9 Spark SQL dengan Set Data Ratings	124

5.51 Hasil Kueri 9 Hive dengan Set Data Ratings	124
5.52 Hasil Kueri 10 Spark SQL dengan Set Data Ratings	125
5.53 Hasil Kueri 10 Hive dengan Set Data Ratings	125
5.54 Hasil Kueri 11 Spark SQL dengan Set Data Ratings	126
5.55 Hasil Kueri 11 Hive dengan Set Data Ratings	126
5.56 Hasil Kueri 12 Spark SQL dengan Set Data Ratings	127
5.57 Hasil Kueri 12 Hive dengan Set Data Ratings	127
5.58 Hasil Kueri 13 Spark SQL dengan Set Data Ratings	128
5.59 Hasil Kueri 13 Hive dengan Set Data Ratings	128
5.60 Hasil Kueri 14 Spark SQL dengan Set Data Ratings	129
5.61 Hasil Kueri 14 Hive dengan Set Data Ratings	129
5.62 Grafik Perbandingan Waktu Eksekusi Kueri dengan Set Data Movies	131
5.63 Grafik Perbandingan Waktu Eksekusi Kueri dengan Set Data Ratings	133
5.64 Grafik Perbandingan Waktu Eksekusi Kueri JOIN	134

DAFTAR TABEL

2.1	Tipe Variabel Scala [6]	6
2.2	Tipe Variabel Numerik Scala [6]	7
2.3	Operator Kondisional Scala [6]	9
2.4	Tabel Customers [7]	12
2.5	Tabel Orders [7]	12
2.6	Tabel Customers [7]	14
2.7	Output Pernyataan SELECT * FROM Customer [7]	14
2.8	Tabel Clients [7]	15
2.9	Tabel Clients Setelah Pernyataan INSERT [7]	16
2.10	Visualisasi Baris yang Akan Dimasukkan ke Tabel Clients [7]	16
2.11	Table NewClients [7]	16
2.12	Table Clients Setelah Pernyataan INSERT dengan SELECT [7]	17
2.13	Visualisasi Baris Pertama Tabel Clients [7]	18
2.14	Table Orders [7]	18
2.15	Tabel Setelah Pernyataan INNER JOIN [7]	19
2.16	Tabel Refunds [7]	19
2.17	Tabel Hasil Kueri LEFT JOIN [7]	20
3.1	Tabel Data yang Dimasukkan ke Tabel <i>movies</i>	67
5.1	Tabel Perbandingan Waktu Eksekusi Rata-rata Spark SQL dan Hive untuk Set Data Movies	130
5.2	Tabel Perbandingan Waktu Eksekusi Rata-rata Spark SQL dan Hive untuk Set Data Ratings	132
5.3	Tabel Perbandingan Waktu Eksekusi Rata-rata Spark SQL dan Hive untuk Kueri JOIN	133

BAB 1

PENDAHULUAN

Pada bab ini dibahas tentang latar belakang penelitian dilakukan, rumusan masalah, dan tujuan dari penelitian ini. Selain itu batasan-batasan masalah, metodologi, dan sistematika pembahasan penelitian ini juga dijelaskan pada bab ini.

1.1 Latar Belakang

Seiring kemajuan dan pemanfaatan teknologi, data bertambah dengan cepat dalam waktu singkat sehingga ukuran data menjadi sangat besar. Ukuran data tersebut mencapai puluhan, ratusan *gigabyte*, hingga *petabyte*. Kumpulan data dengan ukuran besar tersebut disebut dengan *big data*. *Big data* perlu diolah dan dianalisis untuk memudahkan pengguna dalam memahami dan menyelesaikan masalahnya.

Beberapa teknologi, seperti Apache Hadoop dan Apache Spark, dibangun untuk mengolah dan menganalisis *big data*. Apache Hadoop merupakan *framework open source* untuk menulis dan menjalankan aplikasi terdistribusi yang memproses *big data*. Hadoop memiliki komponen penyimpanan, yaitu Hadoop Distributed File System (HDFS), dan komponen untuk komputasi, yaitu MapReduce. Hadoop mempartisi data dan melakukan komputasi *big data* secara paralel sehingga analisis dapat diselesaikan dengan cepat. Apache Hive merupakan infrastruktur data *warehouse* berbasis Apache Hadoop yang menyediakan dialek yang mirip dengan SQL. Dialek yang disediakan oleh Hive disebut HiveQL atau HQL.

Apache Spark merupakan *framework* komputasi di *cluster* untuk memproses dan menganalisis *big data*. Spark dioptimalkan untuk berjalan di memori sehingga membantu pemrosesan data lebih cepat dari pada MapReduce Hadoop. Spark memiliki komponen utama untuk menyimpan data dalam memori, yaitu RDD. Spark juga memiliki *library* untuk bekerja dengan data terstruktur, seperti data CSV, JSON, dan lainnya, yaitu Spark SQL. Spark SQL mendukung penggunaan bahasa SQL dan HiveQL.

Spark SQL dan Hive digunakan untuk tujuan yang sama, yaitu mengolah dan menganalisis *big data* dengan dialek SQL. Berdasarkan hal tersebut, Spark SQL dan Hive perlu dipelajari lebih lanjut. Penelitian ini diharapkan dapat membandingkan kinerja Spark SQL dan Hive dalam mengolah dan menganalisis *big data*.

1.2 Rumusan Masalah

Berikut rumusan masalah dari penelitian ini:

1. Bagaimana konsep dan pemanfaatan Spark SQL?
2. Bagaimana konsep dan pemanfaatan Hive?
3. Bagaimana perbandingan Spark SQL dan Hive dalam sintaks penulisan SQL dan kinerjanya untuk menganalisis *big data*?
4. Bagaimana rancangan dan implementasi perangkat lunak untuk Spark SQL dan Hive?

1.3 Tujuan

Berikut tujuan dilakukannya penelitian ini:

1. Mempelajari konsep dan pemanfaatan Spark SQL.
2. Mempelajari konsep dan pemanfaatan Hive.
3. Melakukan studi untuk membandingkan sintaks penulisan SQL pada Spark SQL dan Hive.
4. Melakukan eksperimen untuk membandingkan kinerja dalam menganalisis *big data* pada Spark SQL dan Hive.
5. Merancang dan mengimplementasikan perangkat lunak untuk menganalisis *big data* dengan Spark SQL dan Hive.

1.4 Batasan Masalah

Batasan masalah dari penelitian ini antara lain:

1. Data studi kasus yang digunakan untuk analisis merupakan data yang dapat diambil secara bebas di internet.
2. Data studi kasus terdiri dari 2 data berformat CSV, yaitu data *movies* dan data *ratings*.

1.5 Metodologi

Metodologi yang digunakan untuk penelitian ini antara lain:

1. Mempelajari konsep Spark, Hive, serta perintah SQL pada Spark SQL dan Hive.
2. Mempelajari bahasa pemrograman Scala.
3. Menginstalasi dan mengkonfigurasi Spark dan Hive *standalone* dan pada *cluster* Hadoop.
4. Mencoba perintah-perintah Spark dan Spark SQL.
5. Menulis program-program Spark SQL untuk mempelajari kemampuan Spark SQL.
6. Menulis skrip HiveQL untuk mempelajari kemampuan Hive.
7. Mencari dan mengumpulkan beberapa data studi kasus (berformat teks, CSV) dan menyimpannya ke HDFS.
8. Merancang antarmuka perangkat lunak.
9. Mengimplementasi perangkat lunak.
10. Menguji performa perangkat lunak dalam menganalisis *big data*.
11. Menulis dokumen skripsi.

1.6 Sistematika Pembahasan

- Bab 1 Pendahuluan
Bab 1 membahas tentang latar belakang penelitian dilakukan, rumusan masalah, dan tujuan dari penelitian ini. Selain itu, batasan-batasan masalah dari penelitian yang dilakukan, metodologi penelitian, dan sistematika pembahasan dalam penelitian ini juga dijelaskan pada bab ini.
- Bab 2 Landasan Teori
Bab 2 membahas tentang hasil studi literatur yang dilakukan. Bab ini terdiri dari penjelasan singkat mengenai *big data*, perintah-perintah, fungsi-fungsi dari bahasa Scala, deskripsi singkat SQL, basis data relasional, *primary* dan *foreign key*, serta operasi-operasi SQL. Selain itu, dijelaskan juga mengenai deskripsi dan komponen-komponen inti Hadoop, deskripsi dan operasi-operasi Hive, deskripsi dan komponen utama Spark, serta Spark SQL dan operasi-operasinya.
- Bab 3 Studi Eksplorasi
Bab 3 membahas tentang cara instalasi dan konfigurasi Hive dan Spark, serta eksperimen-eksperimen Spark SQL dan Hive yang dilakukan. Eksperimen-eksperimen yang dilakukan berupa eksperimen sederhana melalui *console* maupun IDE.
- Bab 4 Analisis dan Perancangan
Bab 4 membahas tentang analisis yang dilakukan mengenai perangkat lunak, data studi kasus, praolah data, dan kueri. Selain itu, bab ini juga berisi mengenai perancangan perangkat lunak dengan Spark beserta diagram kelas dan deskripsinya, serta perancangan perangkat lunak *web* beserta diagram *use case* dan diagram kelasnya.
- Bab 5 Implementasi dan Eksperimen
Bab 5 membahas tentang deskripsi lingkungan perangkat yang digunakan, implementasi Spark dan Spark SQL, implementasi Hive, implementasi modul membaca kueri dan menulis hasil kueri ke HDFS. Selain itu, bab ini juga berisi pengujian kueri dan hasil kueri yang dilakukan. Eksperimen waktu eksekusi kueri dan kesimpulannya juga dibahas di bab ini.
- Bab 6 Kesimpulan dan Saran
Bab 6 membahas tentang kesimpulan dari penelitian yang dilakukan. Bab ini juga berisi saran penulis untuk pengembangan perangkat lunak dengan Spark SQL selanjutnya.