



Teknik Pengolahan Big Data: Dari Batch hingga Real-time

Dalam era digital saat ini, organisasi menghadapi tantangan mengelola volume data yang terus meningkat secara eksponensial. Presentasi ini akan mengupas tuntas berbagai teknik pengolahan Big Data, mulai dari pemrosesan batch tradisional hingga analisis real-time yang canggih.

Kita akan menjelajahi framework populer seperti Hadoop MapReduce dan Apache Spark, serta teknik optimasi query dengan Apache Hive dan Impala.



Pengantar Big Data

1 Volume

Big Data ditandai dengan jumlah data yang sangat besar, mencapai petabyte atau bahkan exabyte. Volume data ini jauh melampaui kemampuan sistem database tradisional dan memerlukan infrastruktur terdistribusi untuk pengolahannya yang efisien.

2 Velocity

Kecepatan data dihasilkan dan perlu diproses menjadi tantangan tersendiri. Data streaming dari media sosial, perangkat IoT, dan transaksi online membutuhkan pengolahan real-time untuk menghasilkan insight yang berharga segera.

3 Variety

Data hadir dalam berbagai format - terstruktur (database), semi-terstruktur (XML, JSON), dan tidak terstruktur (teks, gambar, video). Keberagaman ini memerlukan pendekatan pengolahan yang fleksibel dan mampu beradaptasi.

Pemrosesan Batch dengan Hadoop MapReduce



Konsep Dasar

Hadoop MapReduce adalah framework pemrosesan data terdistribusi yang memungkinkan pengolahan dataset berskala besar secara paralel di cluster komputer. Paradigma ini memecah tugas besar menjadi tugas-tugas kecil yang independen.



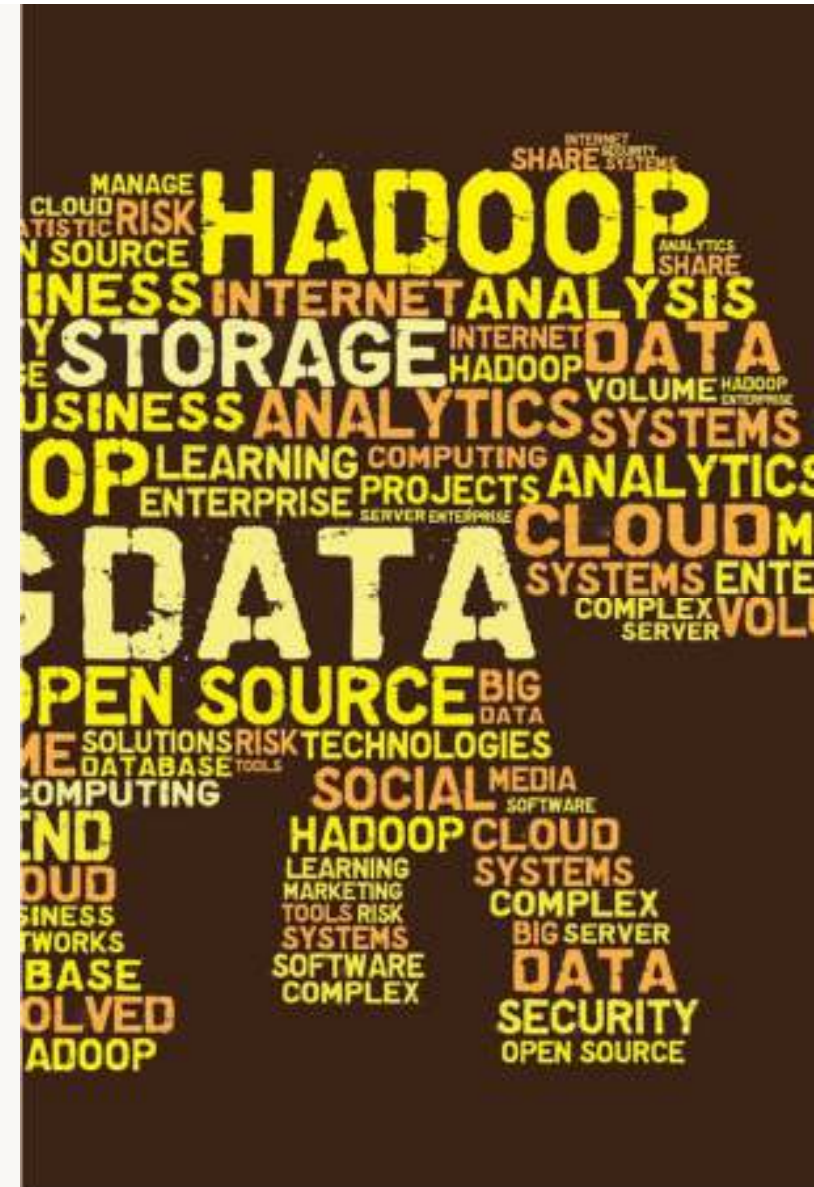
Arsitektur Hadoop

HDFS (Hadoop Distributed File System) menyediakan penyimpanan terdistribusi yang toleran terhadap kegagalan, sementara YARN (Yet Another Resource Negotiator) mengelola sumber daya dan penjadwalan job di cluster.



Kelebihan Batch

Pemrosesan batch sangat efisien untuk analisis retrospektif terhadap dataset historis dalam jumlah sangat besar, dengan fokus pada throughput tinggi daripada latensi rendah.





Cara Kerja Hadoop MapReduce

1

Map

Tahap Map memecah dataset input menjadi pasangan key-value independen. Setiap node dalam cluster memproses subset data secara paralel, menghasilkan output menengah dalam format key-value yang konsisten.

2

Shuffle dan Sort

Sistem secara otomatis mengelompokkan semua nilai dengan key yang sama dari output Map. Data ditransfer antar node (shuffle) dan diurutkan berdasarkan key untuk mempersiapkan tahap Reduce.

3

Reduce

Tahap Reduce mengambil output dari tahap Shuffle dan Sort, kemudian melakukan operasi agregasi atau kalkulasi final untuk menghasilkan output akhir yang disimpan kembali ke HDFS.



Contoh Aplikasi Hadoop MapReduce

Analisis Log Server

MapReduce sangat efektif untuk memproses log server dalam jumlah besar. Tahap Map mengekstrak informasi penting seperti timestamp, kode status, dan URL, sementara Reduce mengagregatnya menjadi statistik bermanfaat seperti jumlah kunjungan per halaman atau tingkat kesalahan.

Perhitungan PageRank

Algoritma PageRank Google awalnya diimplementasikan menggunakan MapReduce. Setiap iterasi algoritma dijalankan sebagai job MapReduce terpisah, dengan Map memproses link antar halaman dan Reduce menghitung skor PageRank baru.

Pemrosesan Data IoT

Data sensor IoT yang terkumpul dalam jumlah masif dapat diproses secara efisien dengan MapReduce untuk menghasilkan insight seperti tren penggunaan energi, deteksi anomali, atau prediksi pemeliharaan peralatan.



Pemrosesan Real-time dengan Apache Spark

1

Kebutuhan Real-time

Meningkatnya kebutuhan organisasi untuk menganalisis data saat dihasilkan mendorong pengembangan sistem pemrosesan real-time. Kasus penggunaan seperti deteksi fraud, rekomendasi dinamis, dan monitoring perangkat memerlukan insight instan.

2

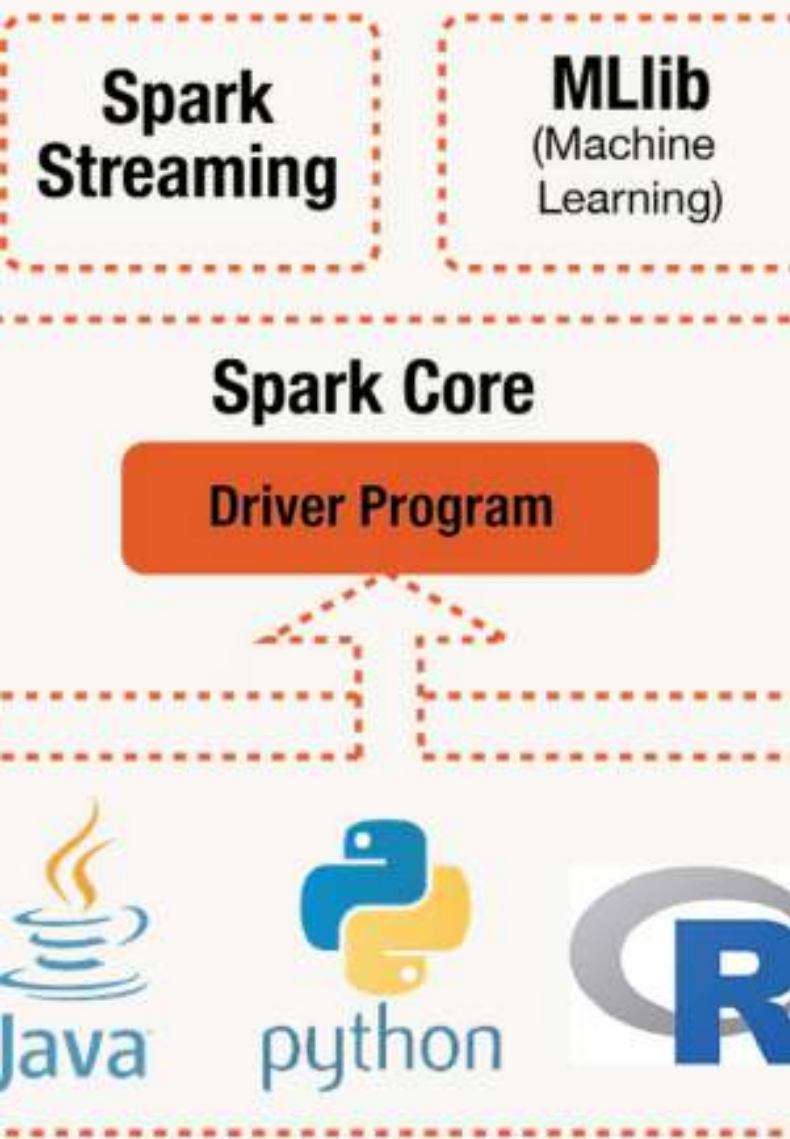
Arsitektur Spark

Apache Spark dibangun di atas konsep Resilient Distributed Datasets (RDD) dan Directed Acyclic Graph (DAG) untuk eksekusi tugas. Arsitektur ini memungkinkan pemrosesan data in-memory yang jauh lebih cepat dibandingkan MapReduce.

3

Spark Streaming

Komponen Spark Streaming memungkinkan pemrosesan data berkelanjutan dengan membagi aliran data menjadi micro-batch yang diproses secara paralel, memberikan ilusi pemrosesan real-time dengan latensi rendah.



Keunggulan Apache Spark

100x

Lebih Cepat

Apache Spark menawarkan kecepatan eksekusi hingga 100 kali lebih cepat dibandingkan Hadoop MapReduce untuk operasi tertentu, terutama untuk algoritma iteratif seperti machine learning dan analisis grafik.

85%

Efisiensi Memori

Kemampuan pemrosesan in-memory Spark mempertahankan data dalam RAM daripada menulis ke disk setelah setiap operasi, secara dramatis mengurangi overhead I/O dan meningkatkan throughput hingga 85% dalam banyak kasus.

4+

Bahasa Pemrograman

Spark menyediakan API untuk Scala, Java, Python, dan R, memungkinkan tim data science menggunakan bahasa yang paling nyaman bagi mereka tanpa harus beralih platform untuk tahap pengolahan data yang berbeda.

Spark Streaming: Konsep DStream

DStream

Discretized Stream (DStream) adalah abstraksi level tinggi yang mewakili aliran data kontinu. Secara internal, DStream diimplementasikan sebagai rangkaian RDD yang masing-masing mewakili data dalam interval waktu tertentu.

Stateful Processing

Spark Streaming mendukung pemrosesan stateful di mana hasil dari batch terdahulu dapat mempengaruhi komputasi batch berikutnya, memungkinkan analisis tren dan deteksi anomali yang lebih kompleks.



Integrasi Sumber Data

Spark Streaming mendukung integrasi dengan berbagai sumber data real-time termasuk Kafka, Flume, Kinesis, dan HDFS, memungkinkan pengolahan streaming data dari hampir semua platform sumber.

Windowing

Operasi windowing memungkinkan komputasi dilakukan pada rentang waktu sliding, seperti menghitung rata-rata data dalam jendela 5 menit terakhir yang diperbarui setiap 1 menit.

RDD (Resilient Distributed Datasets) di Spark

Konsep Dasar RDD

RDD adalah abstraksi data fundamental dalam Spark, mewakili koleksi objek yang terdistribusi dan toleran terhadap kegagalan. RDD dapat dipandang sebagai dataset terdistribusi dengan operasi-operasi yang dapat dijalankan secara paralel pada cluster.

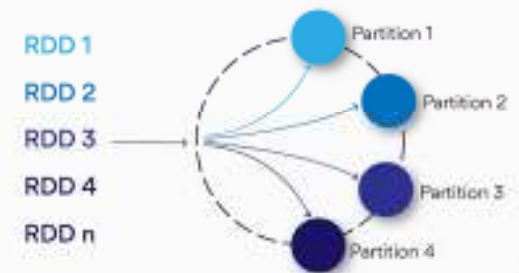
Transformasi dan Aksi

Operasi pada RDD dibagi menjadi transformasi (seperti map, filter, join) yang menghasilkan RDD baru dan lazy evaluation, serta aksi (seperti count, collect, save) yang mengembalikan nilai atau menulis hasil ke penyimpanan.

Persistensi dan Partisi

RDD dapat dipertahankan dalam memori atau disk untuk penggunaan berulang, dengan berbagai level persistensi. Pengaturan partisi yang tepat memungkinkan distribusi data yang seimbang di seluruh cluster untuk pemrosesan optimal.

Partition Process in Resilient Distribution Dataset



DataFrames di Spark

Struktur Data Tabular

DataFrames adalah abstraksi data tingkat tinggi dalam Spark yang mengorganisir data dalam format tabular dengan skema tertentu, mirip dengan tabel database relasional atau dataframe di R/Python. Struktur ini memungkinkan Spark untuk lebih memahami data dan mengoptimalkan eksekusi.

Optimasi Query

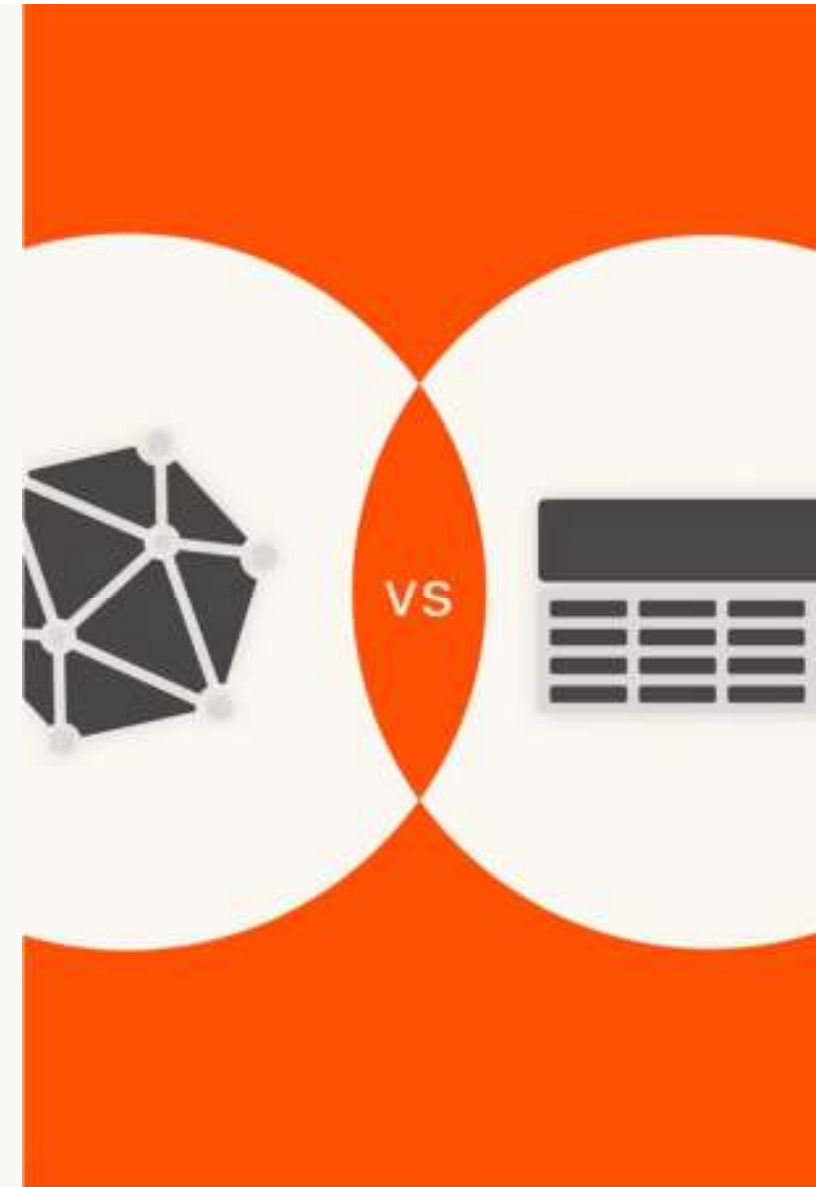
Catalyst Optimizer secara otomatis mengoptimalkan query DataFrames dengan melakukan analisis logis dan fisik untuk menemukan rencana eksekusi terbaik. Ini termasuk pushdown predicate, column pruning, dan berbagai optimasi lain yang meningkatkan performa secara signifikan.

Operasi SQL

Spark SQL memungkinkan penggunaan syntax SQL standar untuk query DataFrames, membuat transisi lebih mudah bagi profesional data yang sudah familiar dengan SQL. Ini juga memungkinkan integrasi dengan tools BI yang menggunakan SQL sebagai bahasa query.

Perbandingan RDD dan DataFrames

Aspek	RDD	DataFrames
Skema Data	Tidak terstruktur	Terstruktur dengan skema
Optimasi	Minimal, dikelola developer	Otomatis via Catalyst Optimizer
Performa	Baik untuk data tidak terstruktur	Superior untuk data terstruktur
API	Fungsional, detail kontrol tinggi	Deklaratif, abstraksi tingkat tinggi
Penggunaan Memori	Lebih tinggi	Lebih efisien (format kolumnar)
Kasus Penggunaan	Transformasi kompleks, unstructured	Analitik, ML, integrasi SQL



Apache Hive: Data Warehousing di Hadoop



Apache Hive adalah data warehouse framework yang dibangun di atas Hadoop, memungkinkan analisis data dan query yang familiar menggunakan bahasa mirip SQL. HiveQL menyediakan abstraksi SQL-like yang diterjemahkan ke job MapReduce, Tez, atau Spark di belakang layar.

Metastore Hive menyimpan semua informasi struktural tentang tabel, partisi, kolom, dan tipe data, memudahkan pengguna untuk mengatur dan mengelola data mereka. Hive sangat cocok untuk analisis batch pada dataset besar dan berintegrasi dengan berbagai tools business intelligence.

Apache Impala: Query SQL Interaktif



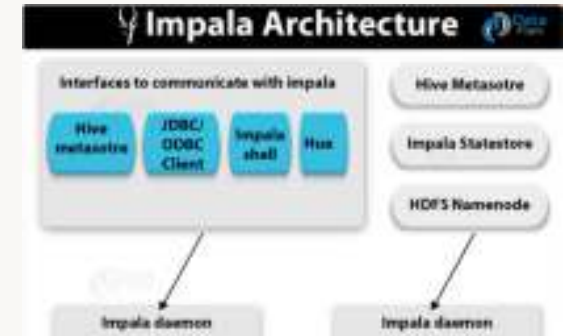
Arsitektur MPP

Impala menggunakan arsitektur Massively Parallel Processing (MPP) yang memungkinkan eksekusi query terdistribusi tanpa menggunakan MapReduce. Dengan daemon yang berjalan pada setiap node dalam cluster, Impala mampu mengakses data HDFS atau HBase secara langsung untuk performa optimal.



Performa vs Hive

Dibandingkan dengan Hive, Impala menawarkan performa query yang jauh lebih cepat, terutama untuk query analitik interaktif. Benchmark menunjukkan Impala seringkali 10-100x lebih cepat untuk query seleksi, agregasi, dan join pada dataset menengah hingga besar.



Optimasi Real-time

Impala menggunakan runtime code generation, predicate pushdown, dan partitioning awareness untuk mengoptimalkan query. Kemampuan runtime adaptif memungkinkan Impala menyesuaikan rencana eksekusi berdasarkan statistik yang dikumpulkan selama eksekusi untuk performa maksimal.

Kesimpulan: Memilih Teknik yang Tepat

1	Batch Processing Ideal untuk analisis historis kompleks
2	Streaming Analytics Untuk insight real-time dan responsif
3	SQL Warehousing Business intelligence dan reporting
4	Optimasi Menyeluruh Kombinasi teknik untuk hasil terbaik

Pemilihan teknik pengolahan Big Data yang tepat harus didasarkan pada karakteristik dan kebutuhan bisnis spesifik. Untuk analisis historis dan ETL berat, pemrosesan batch dengan Hadoop MapReduce masih relevan. Untuk analisis real-time dan pemrosesan streaming, Apache Spark menawarkan performa dan fleksibilitas yang unggul.

Tren masa depan mengarah pada arsitektur hybrid yang menggabungkan pemrosesan batch dan streaming dalam platform terpadu, serta penggunaan lebih banyak solusi serverless dan cloud-native. Perusahaan juga perlu mengantisipasi integrasi teknik AI dan machine learning yang lebih mendalam dengan platform pengolahan data besar mereka.