



Business Intelligence 3 SKS: Big Data & BI

Integrasi BI dengan Hadoop dan Spark

Memahami Business Intelligence dalam Skala Besar

Bab 1: Mengapa Big Data Penting untuk BI?

175

Zettabytes

Proyeksi volume data global pada tahun 2025
menurut IDC

90%

Data Tak Terstruktur

Mayoritas data bisnis yang sulit diproses BI
tradisional

3x

Peningkatan Insight

Kualitas keputusan bisnis dengan analisis Big
Data

Sistem Business Intelligence tradisional menghadapi keterbatasan serius dalam mengolah volume data yang terus berkembang secara eksponensial. Data global diperkirakan mencapai **175 zettabytes pada 2025** menurut IDC, menciptakan tantangan besar bagi sistem BI konvensional.

BI tradisional tidak mampu menangani data dalam skala besar dan pemrosesan real-time yang dibutuhkan bisnis modern. Big Data memberikan solusi dengan memungkinkan analisis mendalam terhadap volume data masif, menghasilkan insight bisnis yang lebih tajam, akurat, dan dapat ditindaklanjuti untuk pengambilan keputusan strategis.

Apa itu Hadoop?

Hadoop adalah framework open-source yang revolusioner untuk penyimpanan dan pemrosesan data besar secara terdistribusi. Dikembangkan oleh Apache Software Foundation, Hadoop memungkinkan organisasi untuk mengelola dan menganalisis petabytes data menggunakan cluster komputer yang terdistribusi.

Komponen Utama Hadoop

- **HDFS (Hadoop Distributed File System)** - sistem penyimpanan terdistribusi yang fault-tolerant
- **MapReduce** - model pemrograman untuk pemrosesan data paralel
- **YARN** - resource manager untuk alokasi sumber daya cluster

Perusahaan teknologi terkemuka seperti **Yahoo dan Facebook** menggunakan Hadoop untuk mengelola petabytes data pengguna mereka setiap hari.





Kelebihan Hadoop untuk BI Skala Besar



Skalabilitas Tinggi

Kemampuan menambahkan ribuan node dengan biaya rendah. Arsitektur scale-out memungkinkan pertumbuhan horizontal tanpa downtime.



Fault Tolerance

Data direplikasi otomatis di multiple node untuk mencegah kehilangan data. Sistem terus berjalan meskipun ada node yang gagal.



Batch Processing

Sangat cocok untuk pemrosesan batch data besar seperti laporan bulanan, analisis historis, dan arsip data jangka panjang.

Hadoop memberikan fondasi yang solid untuk infrastruktur BI skala enterprise dengan kombinasi skalabilitas, keandalan, dan efisiensi biaya yang tidak tertandingi oleh sistem tradisional.



Apa itu Apache Spark?

Apache Spark adalah mesin pemrosesan data **in-memory** yang dirancang untuk kecepatan dan fleksibilitas tinggi. Berbeda dengan Hadoop MapReduce yang menulis data ke disk di setiap tahap, Spark memproses data langsung di RAM, menghasilkan performa yang jauh lebih cepat.



Real-Time Analytics

Mendukung pemrosesan streaming data real-time untuk analisis instan terhadap data yang terus mengalir



MLlib Library

Library machine learning terintegrasi untuk algoritma klasifikasi, regresi, clustering, dan collaborative filtering



GraphX

Framework untuk analisis graf dan komputasi graf paralel, ideal untuk analisis jaringan sosial dan rekomendasi



Iterative Processing

Optimal untuk algoritma iteratif yang memerlukan multiple pass melalui dataset yang sama

Kelebihan Spark dibanding Hadoop MapReduce



Kecepatan Pemrosesan In-Memory

Spark memproses data di RAM, bukan disk, menghasilkan kecepatan **10-100x lebih cepat** dibanding MapReduce untuk workload tertentu. Ideal untuk iterative algorithms dan interactive queries.

Real-Time Streaming Support

Mendukung streaming data real-time dengan Spark Streaming, sangat cocok untuk analisis media sosial, IoT sensor data, dan monitoring sistem yang memerlukan respons cepat.

Integrasi Multi-Bahasa

Mudah diintegrasikan dengan bahasa pemrograman populer seperti Python, Scala, Java, dan R. API yang user-friendly mempercepat development dan deployment aplikasi BI.

Integrasi BI dengan Hadoop dan Spark

Arsitektur Terintegrasi

Hadoop menyediakan **layer penyimpanan data besar** melalui HDFS dan resource management melalui YARN. Spark berjalan di atas infrastruktur Hadoop untuk mempercepat pemrosesan data dengan in-memory computing.

Integrasi ini menciptakan ekosistem yang powerful: Hadoop menangani storage dan orchestration, sementara Spark memberikan kecepatan dan fleksibilitas untuk analytics.



01

Spark on YARN

Mode paling umum - Spark berjalan sebagai aplikasi YARN dengan dynamic resource allocation

02

Spark Standalone

Cluster manager bawaan Spark untuk deployment independen dengan kontrol penuh

03

SIMR (Spark in MapReduce)

Menjalankan Spark job dalam MapReduce job untuk kompatibilitas dengan sistem legacy

Manfaat Integrasi Hadoop & Spark untuk BI

1

Akselerasi Analisis Data

Mempercepat analisis data besar secara real-time dan batch processing. Kombinasi HDFS untuk storage dan Spark untuk processing menghasilkan pipeline analytics yang sangat efisien.

- Query response time berkurang hingga 90%
- Batch job completion 10-100x lebih cepat

2

Efisiensi Biaya Infrastruktur

Mengurangi biaya infrastruktur dengan pemanfaatan cluster terdistribusi menggunakan commodity hardware dibanding sistem proprietary yang mahal.

- TCO (Total Cost of Ownership) lebih rendah
- Scale-out architecture lebih ekonomis

3

Fleksibilitas Platform

Memungkinkan BI tools mengakses data lintas platform dengan fleksibilitas tinggi. Support untuk berbagai data sources dan formats meningkatkan agility bisnis.

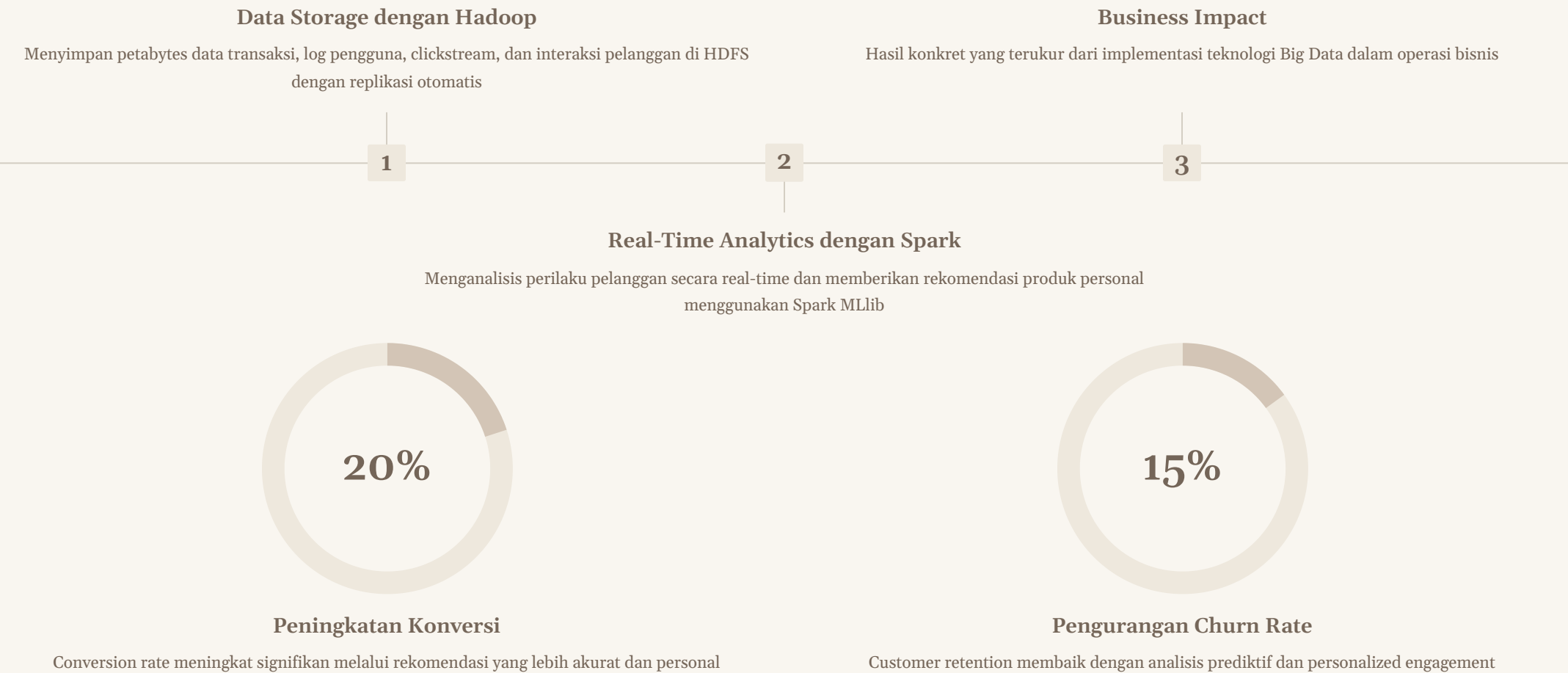
- Integrasi dengan Tableau, Power BI, QlikView
- API standard untuk custom applications



Studi Kasus: Perusahaan E-Commerce

Implementasi Hadoop & Spark untuk Transformasi Digital

Sebuah perusahaan e-commerce terkemuka mengimplementasikan arsitektur Big Data menggunakan Hadoop dan Spark untuk meningkatkan customer experience dan revenue. Implementasi ini menunjukkan kekuatan real-world dari integrasi teknologi Big Data dengan BI.



Tantangan dan Solusi dalam Integrasi BI Big Data

Tantangan Utama

Kompleksitas Pengelolaan

Mengelola cluster distributed yang besar memerlukan expertise khusus dalam konfigurasi, monitoring, dan troubleshooting

Keamanan Data

Memastikan data security, access control, dan compliance di lingkungan distributed yang kompleks

Skill Gap

Kurangnya talent dengan skill khusus untuk mengoptimalkan Spark dan Hadoop secara efektif

Solusi Modern

AWS EMR

Amazon Elastic MapReduce - managed Hadoop framework dengan auto-scaling dan easy deployment

Azure Synapse

Platform analytics terintegrasi dengan built-in Spark dan data warehousing capabilities

Google Dataproc

Fast, easy-to-use, fully managed cloud service untuk running Apache Spark dan Hadoop clusters

Penggunaan **managed services cloud** mengurangi kompleksitas operasional, mempercepat time-to-value, dan memungkinkan tim fokus pada analytics dibanding infrastructure management.

Ringkasan dan Kesimpulan

Big Data + BI = Keputusan Cerdas

Kombinasi Big Data dan Business Intelligence adalah kunci untuk pengambilan keputusan bisnis yang lebih cerdas, cepat, dan berbasis data di era digital

Hadoop & Spark: Fondasi Teknologi

Hadoop dan Spark membentuk fondasi teknologi untuk BI skala besar dengan penyimpanan distributed dan pemrosesan in-memory yang powerful

Integrasi untuk Keunggulan Kompetitif

Integrasi keduanya memberikan kecepatan, skalabilitas, dan fleksibilitas analitik yang dibutuhkan untuk competitive advantage di pasar modern

Poin-Poin Penting

- Volume data global mencapai **175 zettabytes pada 2025** - memerlukan solusi Big Data
- Hadoop menyediakan **storage distributed dan fault tolerance** untuk data besar
- Spark memberikan **kecepatan 10-100x lebih cepat** dengan in-memory processing
- Integrasi menghasilkan **ROI signifikan** - contoh: peningkatan konversi 20%, pengurangan churn 15%
- Cloud managed services seperti **AWS EMR, Azure Synapse, Google Dataproc** mempermudah adopsi

❏ **Langkah Selanjutnya:** Mulai eksplorasi hands-on dengan Apache Spark menggunakan free tier cloud services untuk memahami capabilities dan best practices dalam implementasi BI skala besar.