

- Apache Hadoop : រំពោលការងាររបស់វិចាល់ដែលអាចបង្កើតឡើង Big Data នៃវឌ្ឍនភាព distributed
 - It Hadoop បានបង្កើតឡើងនូវការងារនៃវិចាល់ទីផ្សារ music streaming, video, social media sentiment analysis និង big data warehouse
 - Stakeholder analysis នូវការងារ Real time
 - នូវតម្លៃដែលអាចបង្កើតឡើង costs ទូទៅឡើង និងការបង្កើតឡើងនូវការងារដែលមានការងារ "cold" date. និងនូវការងារនៃវិចាល់នៃ Hadoop

La Hadoop Distributed File System (HDFS): Բայց յօւնչող բարեկարգ

ห้องเรียน Hardware ประกอบด้วยคอมพิวเตอร์ Network ที่ต่ออยู่กับเครือข่าย LAN และมีไฟ LED

→ find Position file zu mainnode. Darin access 19: nimmt mehr parallel

→ trace file block von node mit den benötigten Daten für mehrere Minuten

→ HTTPS 2. man mit https einrichten, https server aktivieren & https zertifikat installieren

ใน Server หนึ่ง ๆ ก็จะมี Cluster มีชื่อเดียวกัน

ទាញយកព័ត៌មានទាំងអស់នៅក្នុងបច្ចេកទេសទិន្នន័យនូវ Big Data

↳ ก็ติ : • Recover ໄດ້ເງື່ອນໄຫວ້າມີການ Hard ware ຢັງ HDFS ເປົ້າມີໂຄງການ

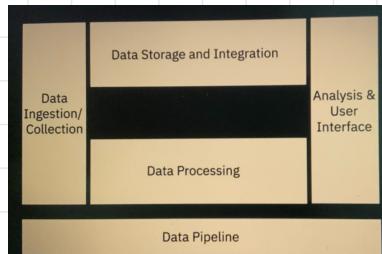
Digitized by srujanika@gmail.com

- ↳ fast: • Recover lost data quickly because Hardware has HDFS redundancy feature
that stores data in multiple copies in different nodes
 - handles Big Data very well in node level in Cluster world
 - handles loss: minimum Hardware failure in multi-level Operation system fail safe

- **Apache Hive:** Hive is a data warehouse software built on top of Hadoop, designed to handle large-scale data processing. It provides a SQL-like interface for querying data stored in various formats (e.g., CSV, JSON, Parquet) across distributed storage systems like HDFS.
 - ↳ **Hive Data Model:** The data is organized into tables, which are composed of partitions. Each partition contains data for a specific subset of the table's columns. This allows for efficient parallel processing of large datasets.
 - ↳ **SQL-like Query Language:** Hive supports a query language called HiveQL, which is based on SQL. This makes it easier for analysts and developers to work with data stored in various formats.
 - ↳ **MapReduce Integration:** Hive is built on top of the MapReduce framework, which provides a distributed processing engine for large datasets. This integration allows for efficient data processing and analysis.
 - ↳ **Apache Hadoop:** Hive is built on top of the Hadoop framework, which provides a distributed file system and a distributed processing engine. This integration allows for efficient data processing and analysis.

Third Module

Layer of a Data Platform Architecture: • Data



- Diagram of a Data Platform Architecture:**

```

graph TD
    A[Data Ingestion/Collection] --- B[Data Storage and Integration]
    B --- C[Data Processing]
    C --- D[Analysis & User Interface]
    D --- E[Data Pipeline]

```

 - **Data Ingestion/Collection layer** → main task: identify data source, transfer data instance into data platform in streaming, batch or both modes
 - **Streaming data collection to metadata repository**
 - **How in data platform work to store data: Integrate**
 - **Data Storage and Integration layer:**
 - ↳ from data to logical processing (e.g. transformation) then transformation to merge data in logical or physical with data to intermediate streaming or batch mode
 - ↳ **Storage layer** also includes, summarize, aggregate information in option
 - **Data Processing** → Data Validation, Transformations, etc. to logical mapping from front layer
 - ↳ layer between raw data in batch to streaming mode in storage
 - ↳ transform data to DB or DS API
 - ↳ Transformation involves in this layer
 - **entity structure** transforming data in schema entities like minification, max complex representation data, Table alias, Joins or Unions
 - **no Normalization** in focusing clean database (in fact data in fact is normalized) for analytical purpose
 - **De-normalization** creating new table(s) in table wise in fact every entity want to report is Analyze table is derived from this
 - **Data Cleaning** in term clean data into fit for analysis into downstream application
 - ↳ Storage or processing on this layer include storage layer below
 - **storage:** Relational data management, HDFS now part of data processing via Spark
 - **No layer** processing engine store their own transformation load into storage to Database
 - **Data Pipeline** → Data Ingestion, Data Storage or Data Processing or Data pipeline to move to Data pipeline
 - ↳ responsibilities to maintain relationships between

Designing of Data Stores:

- ↳ Data store into Data Repositories do data can collect, organize into isolation into Analyze, Planning
- ↳ Data Repository or Database, Data Warehouse, Data Mart, Big data store or Data lake
- ↳ Consideration for designing a data store:
 - **Type of data:** Type of data, Volume data, Implementation, Storage consideration, security, Privacy, Governance data
 - **Type of Data:** Relational Database (RDBMSes) or Non-Relational (NoSQL)
 - RDBMSes implement data in structure of tables into schema into
 - NoSQL implement data into schema or free-form implementation into
 - key value
 - Document based: Implement complex search query into document storage
 - Column based
 - Graph based: Implement data in terms of analytics queries using
- **Volume of Data:** Data lake, Big Data Store
 - **Data lake:** Integrate into single raw data from various source into platform ADLS into NOSQL
 - **Big Data Store:** into data storage, processing, handling into various distributed processing
 - Split large file into computer into nodes, into parallel processing
 - Analyze into multiple node into data mining
- **Intend to use:** 
 - Transactional System: Implement data in term of update, update no update
 - Analytical System: do complex queries into analyzing, processing data in transaction system into decision making
 - Schema design, indexing, or partitioning strategy into optimization into system performance data quality
 - Intend to use data • do with optimization like scalability, trans. capacity, load balancing, fault tolerance, availability
 - Normalization: Implement into data base into data redundancy, data inconsistency, data inconsistency
- **Storage Considerations:** from the perspective of storage
 - **Performance:** throughput (throughput), latency
 - **Throughput:** Implement吞吐量 into storage
 - **Latency:** Implement 延迟 into storage
 - **Availability:** Storage availability into data protection, backup, disaster recovery
 - **Integrity:** Data integrity, system corruption, fragmentation, inconsistency
 - **Recoverability:** storage availability, maximum recover data into protection into data
- **Privacy, Security, and Governance:** Access control, Multi-zone Encryption, Data Management, Monitoring System
 - **Regulations:** GDPR, CCPA, HIPAA

Security: → CIA Triad မြန်မာစုနစ်အတွက် ၂ ခုခု

- **C:** Control မှုပ်နည်း၊ **I:** Integrity မှုပ်နည်း၊ **A:** Availability မှုပ်နည်း

↳ ၁၀. CIA Triad မြန်မာစုနစ်: Infrastructure Security, Network Security, Application Security, Data Security

• Physical Infrastructure Security:

- Access to the perimeter (perimeter), network segmentation, authentication
- Authentication
- အိမ်အော်မှုပ်နည်း 24 hr.
- အိမ်အော်မှုပ်နည်း၊ အိမ်အော်မှုပ်နည်း၊ အိမ်အော်မှုပ်နည်း
- ရန်အိမ်အော်မှုပ်နည်း 10.000
- အိမ်အော်မှုပ်နည်းအတွက် အိမ်အော်မှုပ်နည်း၊ အိမ်အော်မှုပ်နည်း

• Network Security: အိမ်အော်မှုပ်နည်း၊ အိမ်အော်မှုပ်နည်း

- Firewall မှုပ်နည်း၊ Network flow monitoring
- Network Access Control မှုပ်နည်း၊ အိမ်အော်မှုပ်နည်း
- Network segmentation မှုပ်နည်း၊ Network into data silos
- Security protocols မှုပ်နည်း၊ SSL/TLS, IPsec, Kerberos
- Intrusion Detection မှုပ်နည်း၊ Intrusion Prevention မှုပ်နည်း

• Application Security: အုပ်စုဆောင်ရွက်မှုပ်နည်း၊ Application Monitoring

- Threat model မှုပ်နည်း၊ အုပ်စုဆောင်ရွက်မှုပ်နည်း
- Secure design မှုပ်နည်း
- Secure coding မှုပ်နည်း၊ လုပ်ငန်းလုပ်ငန်း
- Security testing မှုပ်နည်း၊ ဖော်ဆောင်ရွက်မှုပ်နည်း

• Data Security

- Data at rest storage မှုပ်နည်း၊ အိမ်အော်မှုပ်နည်း
- Data Authorization မှုပ်နည်း၊ ဂျပ်ဆောင်ရွက်မှုပ်နည်း၊ မှုပ်နည်း၊ မှုပ်နည်း၊ Token
- Data in storage မှုပ်နည်း၊ Data warehouse မှုပ်နည်း၊ Data湖 မှုပ်နည်း
- Data in transit မှုပ်နည်း၊ https, SSL မှုပ်နည်း