# Hadoop Core

Tuesday, January 10, 2023   7:09 AM

- Requirements
  - Data Capture
    - Handle ingestion of both structured/unstructured data from multiple data sources (quickly/efficiently)
    - Handle both batch & online use cases
      - warehouse reporting
      - realtime/streaming

  - Data Processing
    - ETL: Extract/Transform/Load — Batch
    - OLAP: Aggregations, Queries, Reports (Analytics)
      - Ad-Hoc (Interactive)

  - Data Exchange
    - Integration: Interoperability with other systems
    - Sharing: Data/Results across systems & students

  - Easy to Operate
    - Administration: Provisioning, Scaling, Diagnosing, Monitoring, Managing, Maintaining
    - Operations: Availability, Reliability, Dirt Testing ....

* Conceptually

    Temporal

  scale { • Compute: Code (Task) is "too large" ~ run distributed across a system
         • Storage: Data (Memory) is "too large" - store distributed across a system

    Spatial

- Categories
  - MPP (Massively Parallel Processing)
    - Hadoop ✓
  * Row stores (OLTP)
    - ELK: Elasticsearch, Logstash, Kibana (Lucene/Solr) ⇒ Document Search/IR
    - Postgres + PostGIS & Friends ⇒ Geospatial
    - Clickhouse & Friends → Time Series
  - NoSQL (Non-Relational)
    - CouchDB, MongoDB ⇒ Document DB
    - Cassandra, HBase ⇒ Wide-Column DB
    - Memcached, Redis ⇒ Key-Value DB
    - BigTable, Dynamo DB, Cosmos DB → Multi Paradigm DB
    - Neo4j, Gremlin ⇒ Graph DB
  - Column Stores (OLAP)
    - C Store ⇒ Vertica
    - Redshift

- Monet DB
- kdb +

\* HPC Systems
OpenMPI & Friends → Supercomputing / Scientific Computing
· PFS, OpenMP, OpenACC ....

- Hadoop
  · Concepts (Files/Blocks)
    - Storage: Data loaded into system → Distribute across nodes
    - Compute: Code (Jobs/Tasks) run on system → Distribute across nodes

Where the data is located!

\* Locality

  · Techniques                    CPU → Net → HDD
    - COTS Hardware : Linux / x86-64 / Ethernet / Disk     Open Source  (Dev)
                                      (PCIE)
    - API : Understandable / Usable by most programmers    Open Access  (User)

  \* Hadoop Core & Hadoop Ecosystem are two different things!

  · Characteristics
    - Scalable
    - Reliability
    - Flexibility
    - Economical

  · History
    - Google
      · Crawl / Index documents
      · PageRank Calculations → Linear Algebra / Matrix-Vectors
    - Proprietary
      · GFS
      · MR

  \* Open-Source Implementation
    - Apache
    - Doug Cutting

- Uptake
  - IT/Cost
    - Storage Costs
    - ETL Efficiency ==(TTM)==
    - EDW Optimization
    - Archiving
  - Business/Revenue
    - Data Science
    - EDA ==  "Business Analytics"==
    - Predictive Analytics / Data Mining

- Hadoop Core
  - Architecture
    - Common : Libraries, Utilities, etc...
    - **Storage** — HDFS : FOSS implementation of GFS ⇒ <u>DFS</u>
      - ↳ runs across all Hadoop cluster nodes!
    - **Compute** — YARN : Cluster resource management & monitoring
      - ↳ utilization / availability ⇒ intelligent (efficient) job (task) <u>scheduling</u>
    - Zookeeper : Centralized service holder cluster state data
      - ↳ DHT serves as KV-store