# HDFS Architecture

Tuesday, January 10, 2023   7:09 AM

- Distributed File Systems (DFS)
  - Overview
    - Storage partitioned across multiple machines
    - Distributed systems
      - Consistency: read() followed after write()
      - Availability: nodes allowed for read/write access
      - Partitioning: node/net split & connectivity

  CAP { (grouping brace around the three points above)

- HDFS
  - DFS
  - Filesystem Abstraction /API over Blob/object Storage
    ↳ S3, Azure... (HTTP/web API)

  - Use Cases
    - Large Files (very!) → GB/TB/PB+
      ↳ not good for too many (small) files
    - Fault Tolerance → Replication of files/blocks across nodes
      ↳ K-safety    k ≥ 2
    - Read-Only / Read-Heavy Workloads (WORM)
      ↳ No edits, overwrites only, appends allowed
    - High Throughput / Not Low-Latency
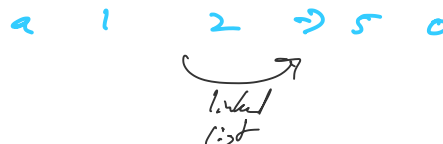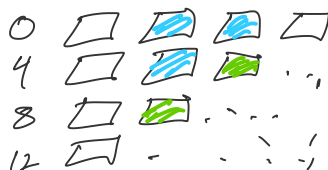      ↳ Batch jobs/tasks → OLAP

  - Blocks
    - Traditional FS          **# Blocks ~ kB**

      Disk                          Disk / Memory (Cache)
      ( Data )                      ( Meta Data )

      0  [ ][ ][ ][ ]              file    start   length
      4  [ ][ ][ ][ ]                2       0       3
      8  [ ][ ][ ][ ]                6       5       2
      12 [ ]

    **# Fragmentation**

      0  [ ][ ][ ][ ]              2    1    2  → 5   0
      4  [ ][ ][ ][ ]
      8  [ ][ ][ ]                           linked
      12 [ ]                                 list
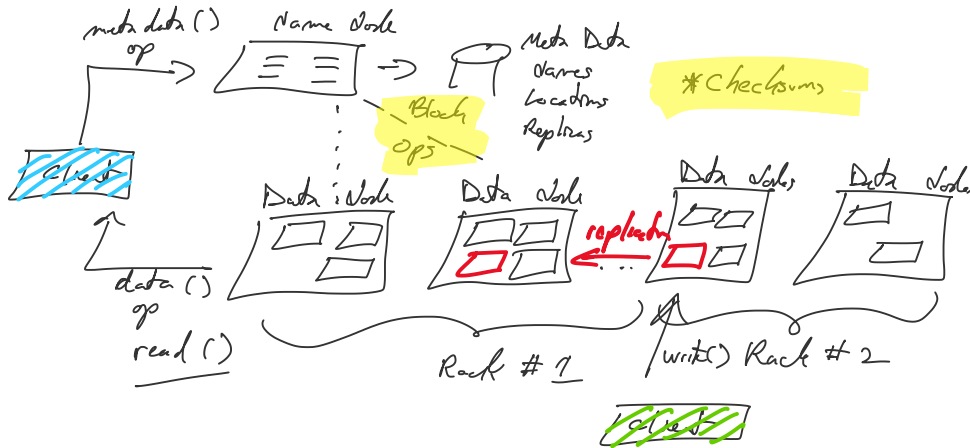
- DFS (HDFS)
  - Blocks       \* Large  ~ 64MB
  - Block - Level Operations
        read / write
        copy / replicate   (k-safety)

\* Replication
  - k = 3 (default) : original + 2copies per block
  - replicas stored different nodes /racks

  - Block size & Replication factor <-> File size



- Nodes
  - Two types of data: Data, Meta Data
      ⤷ 2 node types

  - Name Node (1+) : Meta Data
    - FS
      - Directory Tree
      - Permissions        Cache
      - Block Locations       &
      - checksums ...      Persist

    - Client
      - Block Locations / Data Node Addresses
      - Alternative Locations / Failover

- Data Node (N+) : Data
      - Files consisting of Blocks
      - Blocks are unaware of file-level-memberships
      - Block ops: new blocks, replication, notifications (delete, etc...)
      - Heartbeats to Name Node !

\# write new file:

1. add block to Name Node (Meta Data)
2. write block to Local Data Node
3. background block generation ⇒ replication
   - orchestrated by Name Node w/ Local
     Data Node as source