# Hadoop Ecosystem
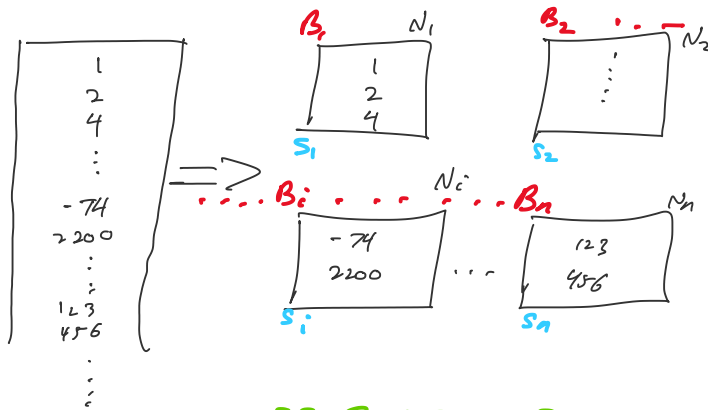Thursday, January 26, 2023   3:22 PM

- Data Processing
  - Data Parallelism
    - Multiple Nodes / Machines (Cluster): Distributed Storage
                                                      Distributed Compute

\* Storage
  - File: Sequence of Bytes (Text) Code & Data
    Typical Filesystem ⇒ Inodes

    DFS ⇒ Blocks



$$S = S_1 + S_2 + \dots S_i + \dots + S_n$$

Linear

\* Fault-Tolerance !                    \* Performance
    ↳ K-safety                              ↳ Storage: Redundancy / IOPS
                                                  Compute: Locality

    Safety                                   Efficiency

\* Compute
  - Parallelism & Concurrency   (HPC ⇒ Schedulers / Task Management)

  - Task / Job: Process ⇒ Scheduling
    - Stream of Bytes: Code & Data (Binary)

  - Pseudo Code ⇒ Adding a list of numbers

    m () {          ⇐ input blocks                    r () {          ⇒ output blocks
        partial-total = 0                                    final-total
        for (each block) {                                   for (each block) {
            read-block()                                         read-block()
          partial-total += record                               final-total += record
        }                                                    }
        return partial-total                                 return final-total
    }

call in parallel
across each node $N_c$ !

*intermediary* blocks

single call for
across all portslate data blocks !

parallel

serial

map

reduce

※ Hadoop cluster

read()  write()

DFS → Blocks → map, map, map → DFS

Shuffle/Sort

read()  write()

DFS ← Blocks → reduce → DFS

※ Cloud : AWS
 - Hadoop (EMR): HDFS (dynamic) / YARN (single tenant)
 - S3: Object Storage

S3 object

| 1 2 | : | : | 123 456 |

Byte Offsets