

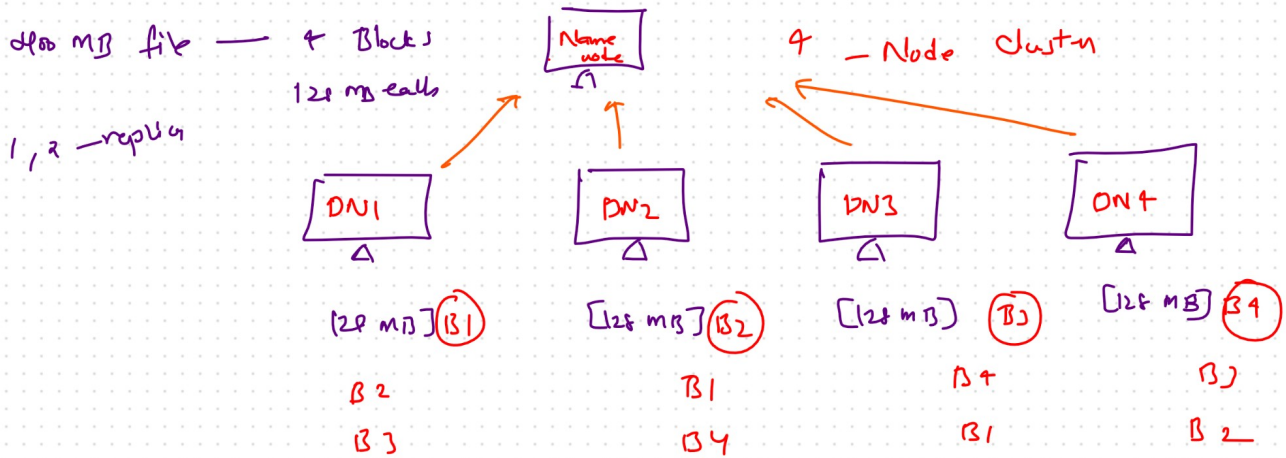
Map Reduce

— Cluster Processing

Multi-threading vs Distributed processing

Multi-threading: parallel processing within a single machine

Map-Reduce: Distributed processing across multiple machines (cluster nodes)



Node Configuration

CPU Resources (4 cores)

RAM

hard

* Core Processing Rules (fundamental principle)

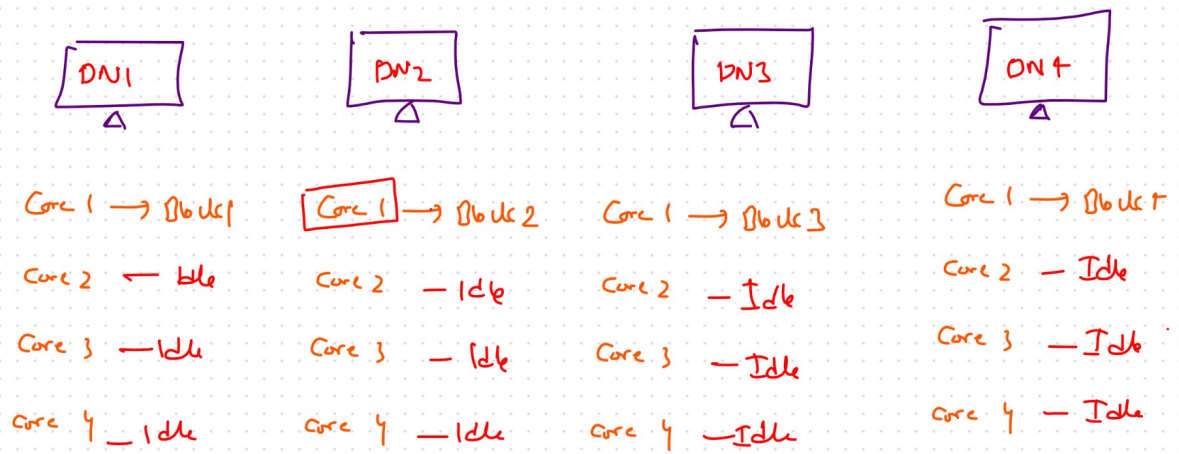
→ One CPU Core can process one mapper job on one block of data

* remember important concept for understanding Map Reduce processing

example 1 : 4-node cluster, each with 4 cores = 16 total cores

file : 400 MB \rightarrow 4 blocks of 128 MB each

Yarn



Ex 2 : Same 16-core cluster

File : large file with (~ 4.1 GB) = 32 blocks

Processing : Can only process 16 blocks at a time

\rightarrow Batch 1 : Process 16 blocks simultaneously (utilizing all 16 cores)
 \downarrow
 \rightarrow Batch 2 : process remaining 16 blocks

Result : Processing happens in batches, takes more time but doesn't fail

Block size Impact on parallelism

Block size \uparrow Parallelism

(1) 128 MB : block size [standard]

1 GB file / 128 MB = 8 blocks = 8 parallel processes

2. Smaller Block size : 64 MB

$$1 \text{ GB File} / 64 \text{ MB} = 16 \text{ blocks} = 16 \text{ parallel processes}$$

✓ Better parallelism

✗ More metadata overhead

3. Very small Block size : 4 MB

16 Cores

$$1 \text{ GB File} / 4 \text{ MB} = \underline{256} \text{ Blocks}$$

✗ Too many batches needed

✗ High metadata overhead

✗ Processing overhead

1 Batch — 16 Blocks

2 Batch — 16 Blocks

⋮

16 Batches — 16 Blocks

256 Blocks

4. Larger Block size : 512 MB

$$1 \text{ GB} / 512 \text{ MB} = 2 \text{ Blocks} = 2 \text{ parallel process}$$

✗ Reduced parallelism

✗ Underutilized cores

Block size selection:

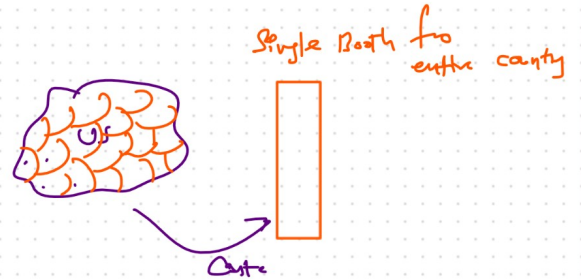
128 MB → well-researched & experimentally proven optimal size

Map Reduce

def: mapReduce is a framework for distributed data processing in Hadoop Ecosystem. helps on the processing aspect of big data operations

Key Concept : Code Locality

ex: Election \rightarrow Polling System



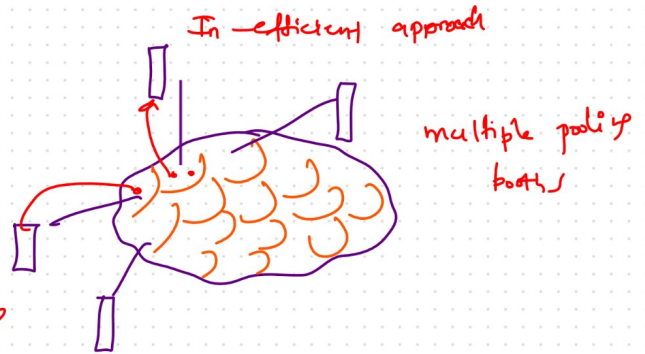
Example 2 Software Distribution

Inefficient Approach

\rightarrow users buy laptops to

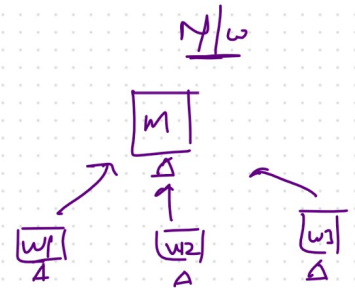
Company headquarters for s/w installation

G.C



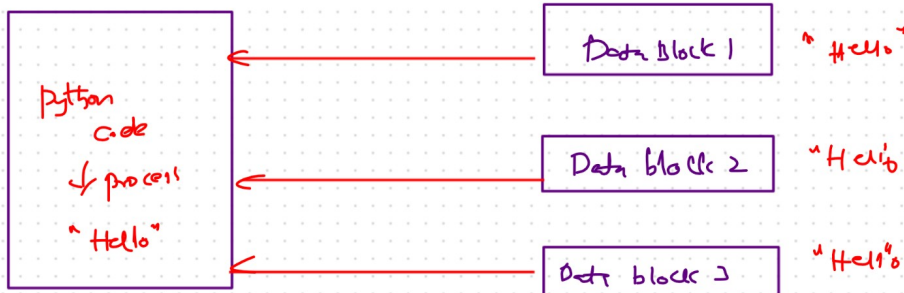
Efficient Approach : Download Software to user's device

Traditional vs Distributed Processing



Master Machine

Worker Nodes

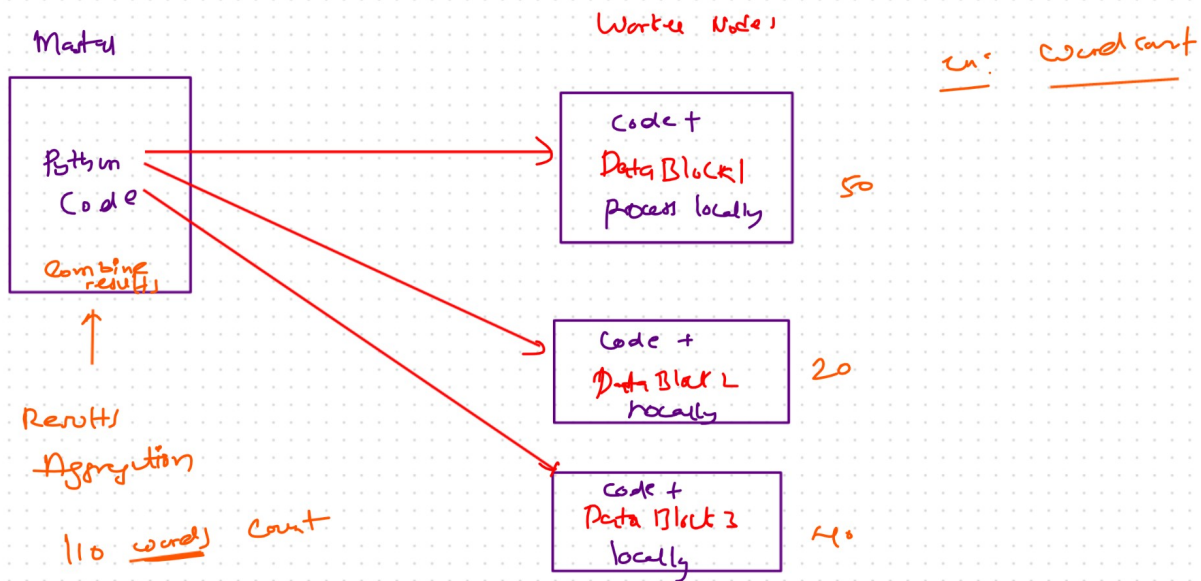


Issues with Traditional Approach :

(1) Data movement : Large Data sets travelling to processing unit (master)

- 3) Single Point Processing: only master machine's resources utilized
- 4) No Data Locality: Data stored locally on workers but processed remotely
- 5) No Local Processing: Worker node CPU/RAM resources wasted
- 6) Scalability Issues: Bottleneck at master node

Distributed Processing Solution



Advantages of Distributed Approach

- 1) Code to Data Nodes: lightweight code (KB/MB) travels instead of heavy data files (GB)
- 2) Data Locality: Data processed where it resides
- 3) Local processing: All worker node's resources utilized
- 4) Parallel processing: multiple machines process simultaneously
- 5) Scalability: Performance improves with more machines