

# Agenda

Learn the architecture of YARN and its key components.

Understand how MapReduce tasks are scheduled and executed using YARN.

Explore the communication between HDFS and YARN.

Work through examples to solidify these concepts.



**Map phase**



**Partition phase**



**Shuffle phase**



**Sort phase**



**Reduce phase**

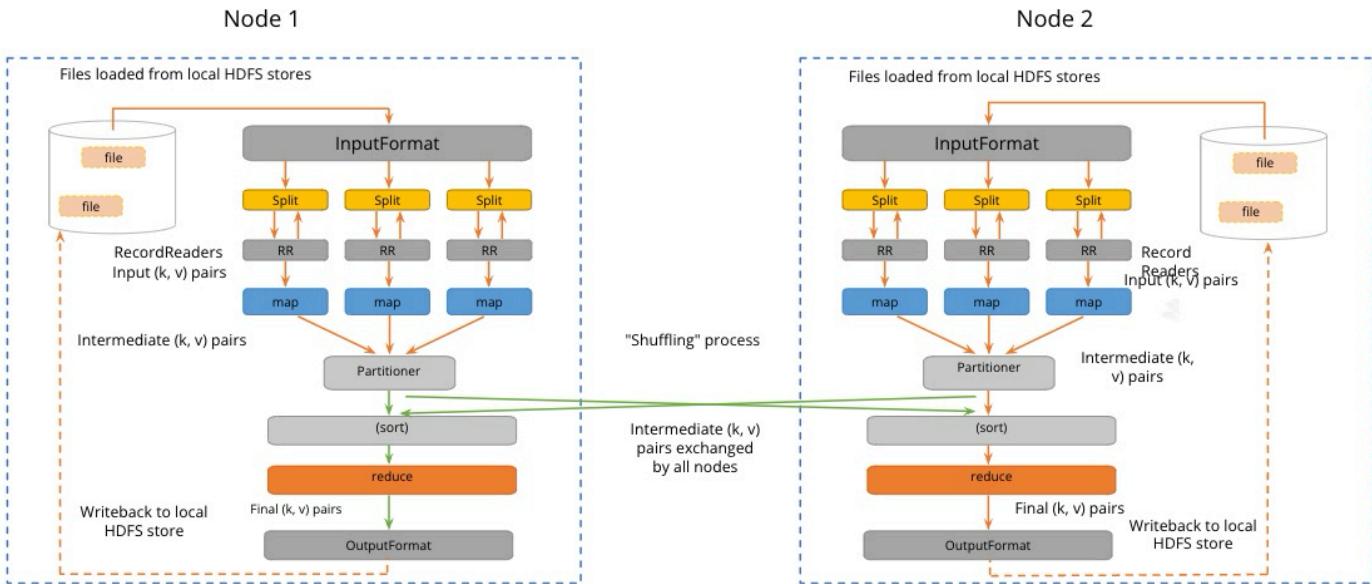
- Reads assigned input split from HDFS
- Parses input into records as key-value pairs
- Applies map function to each record
- Informs master node of its completion

- Each mapper must determine which reducer will receive each of the outputs
- For any key, the destination partition is the same
- Number of partitions = Number of reducers

- Fetches input data from all map tasks for the portion corresponding to the reduce tasks bucket

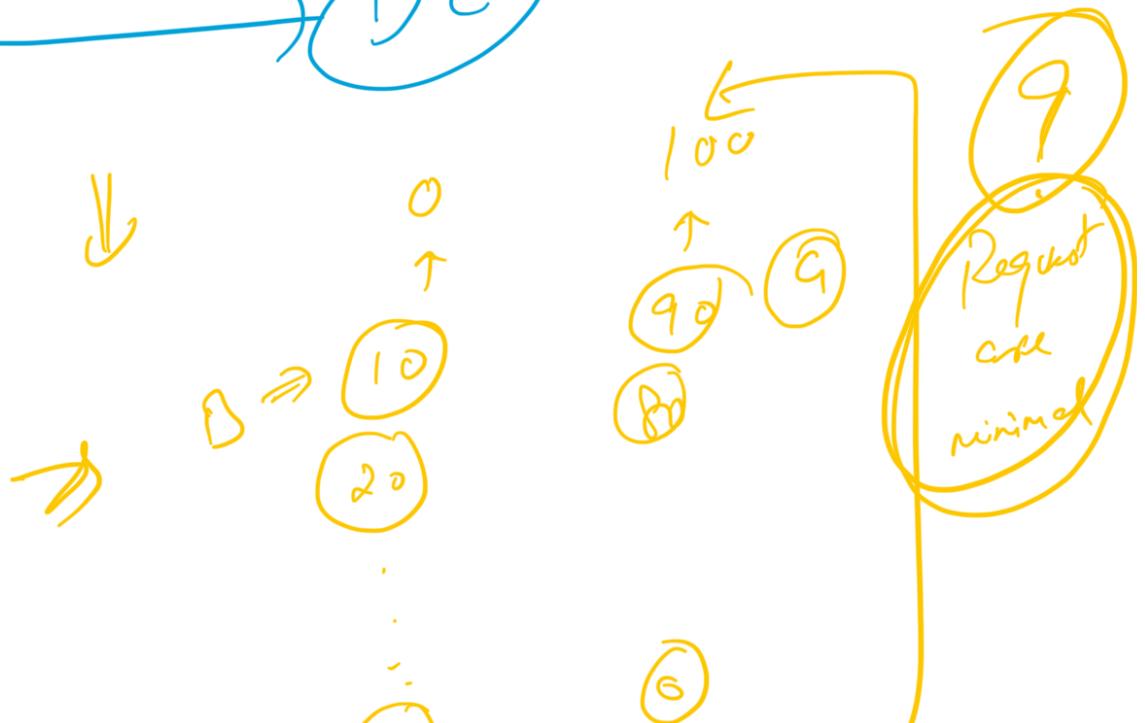
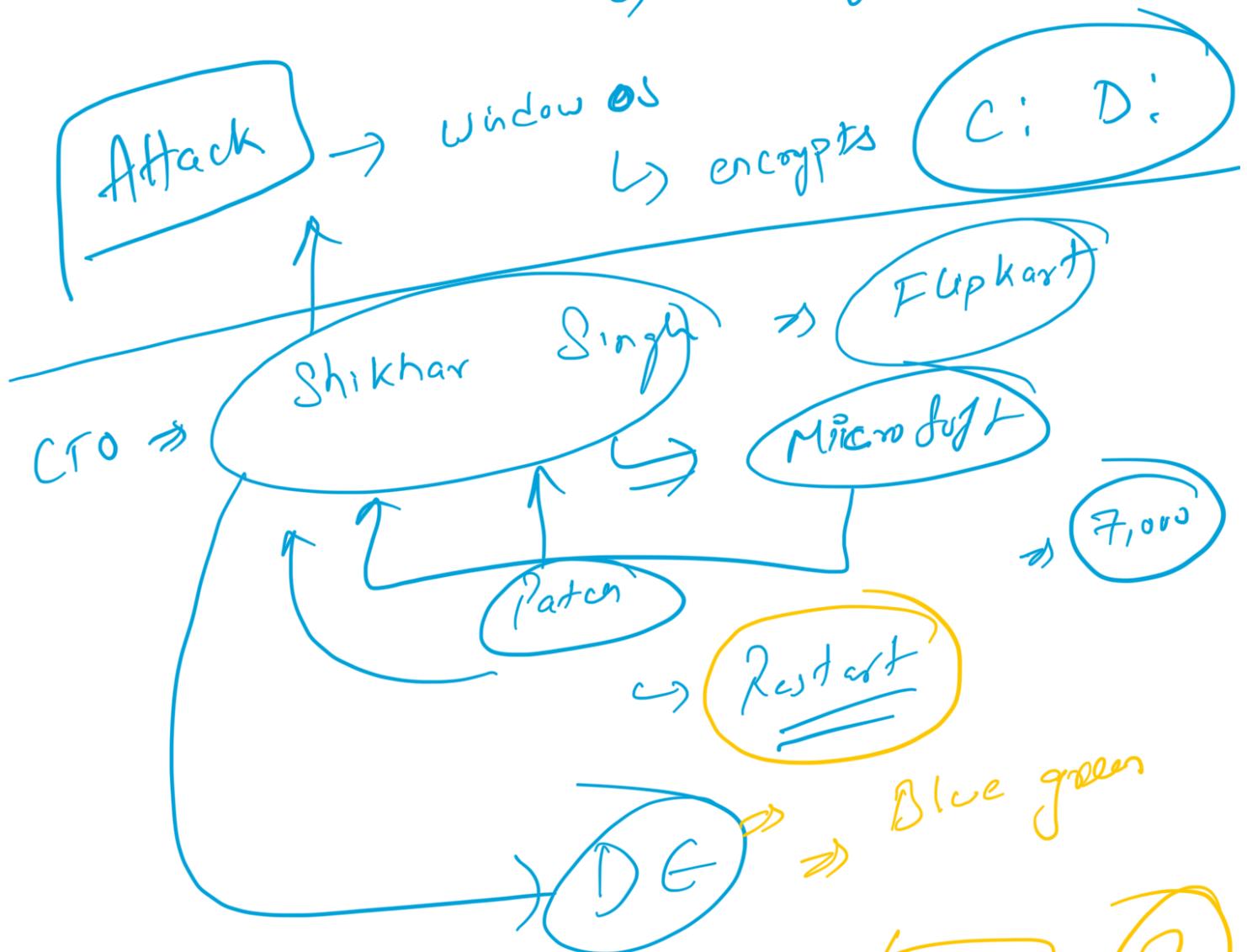
- Merge sorts all map outputs into a single run

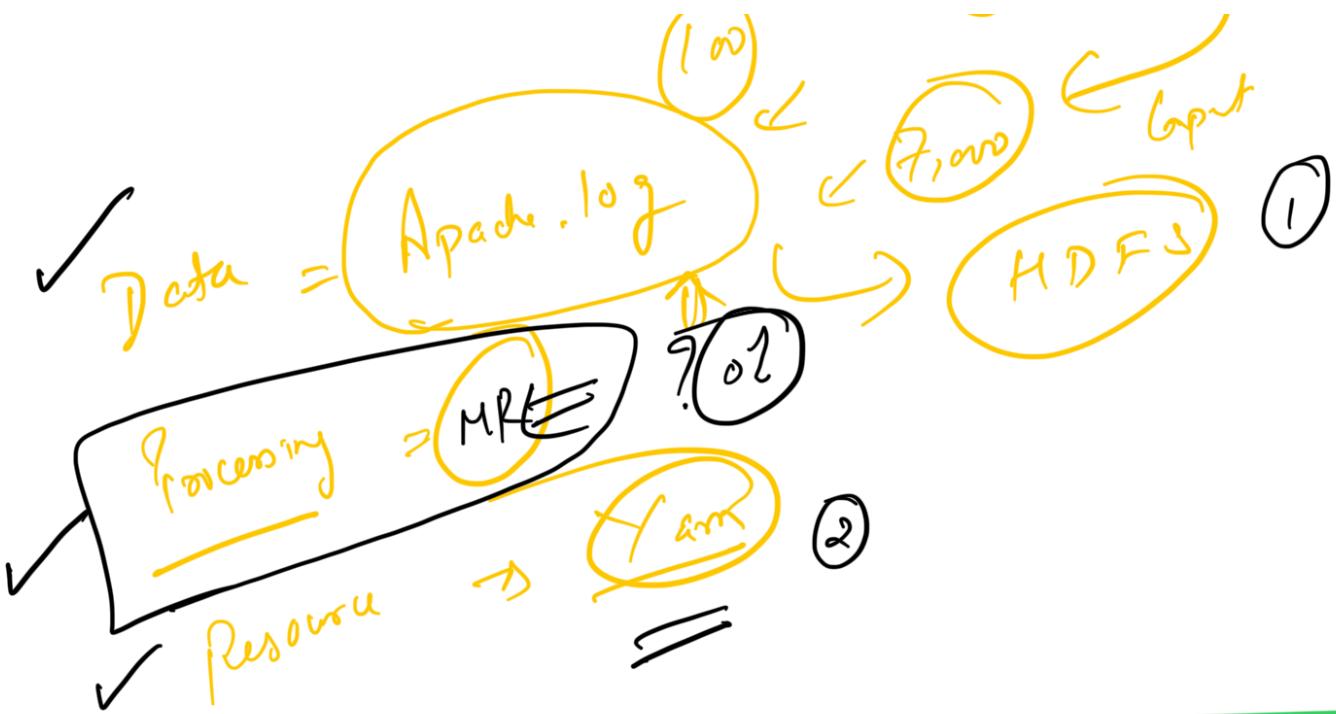
- Applies user-defined reduce function to the merged run



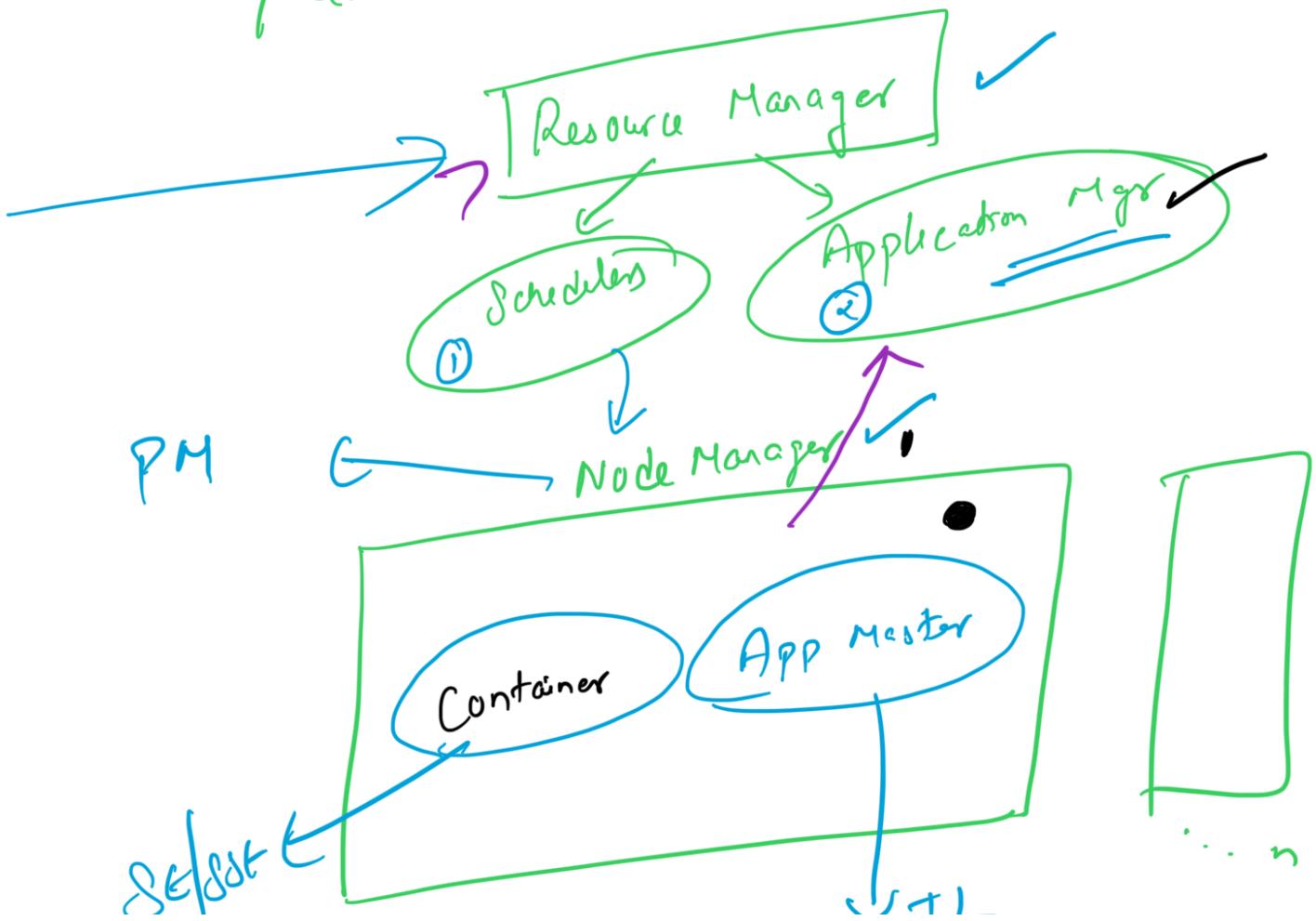
May 2017

↳ WannaCry ransomware



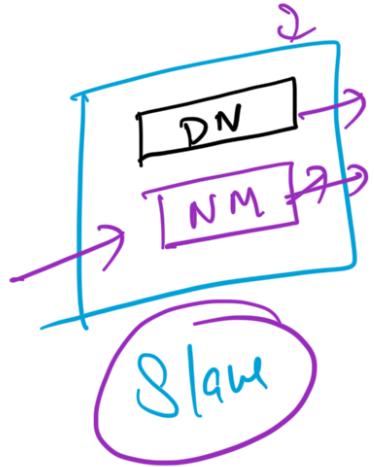


Yarn =



## ① Node Manager (NM)

↳ runs on each node and manages ↓



- ① Container lifecycle
- ② Container dependencies
- ③ Container leases
- ④ Node health | Log import
- ⑤ Reporting for node status to RM.

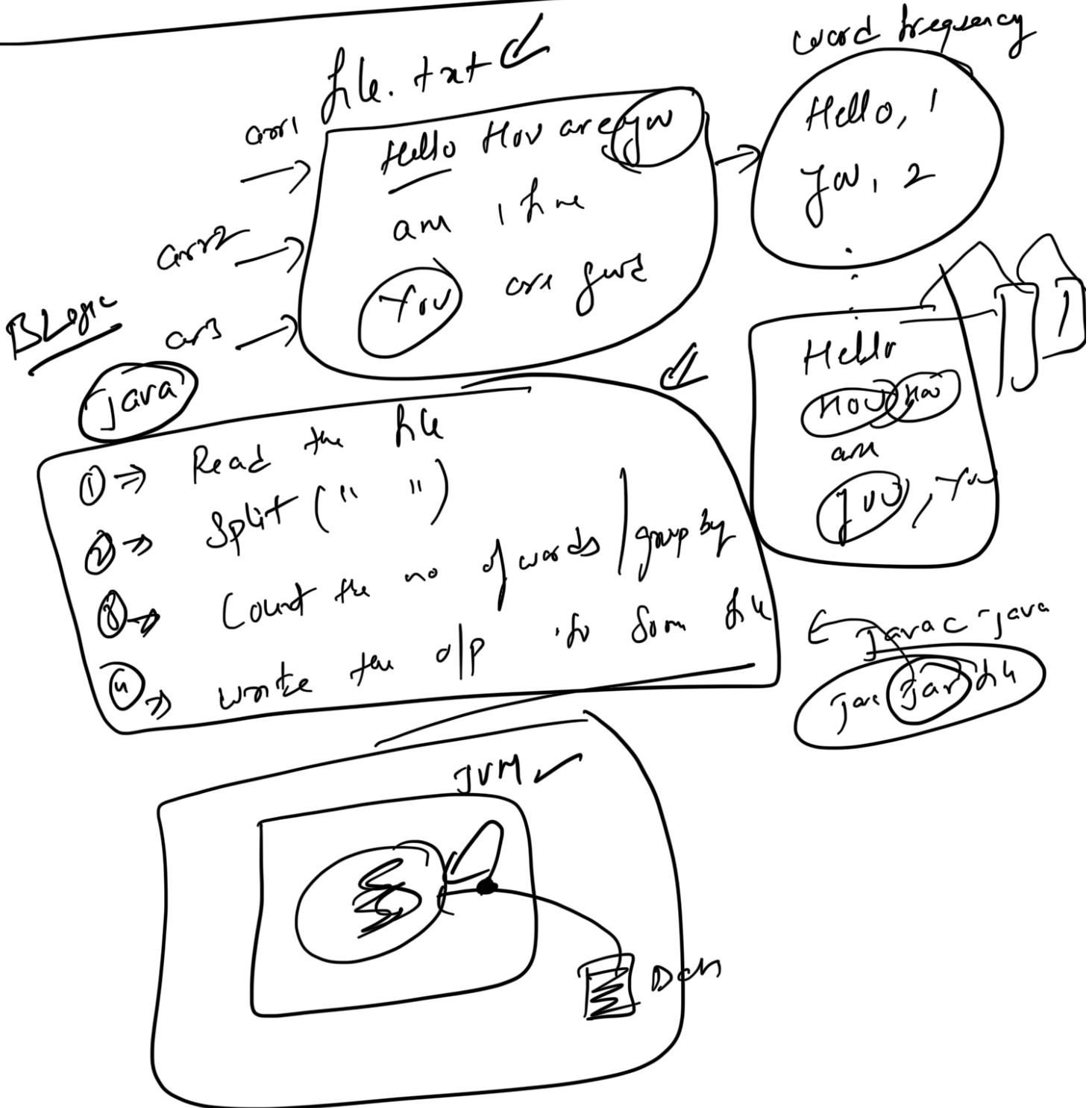
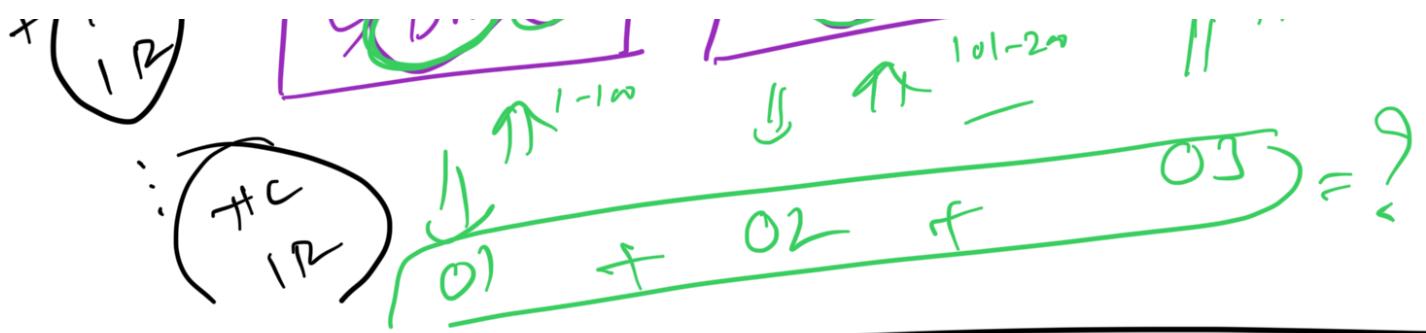
## ② Application Master :-

↳ this is your distributed job.



is a framework specific library, which negotiates resources from RM and works with NM to execute it self with the help of Container.

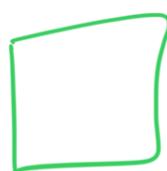




Container  $\rightarrow$  C&R

, , , ,

↳ a collection of set of resources (CPU + RAM) in certain numbers on a specific node.



Job ?

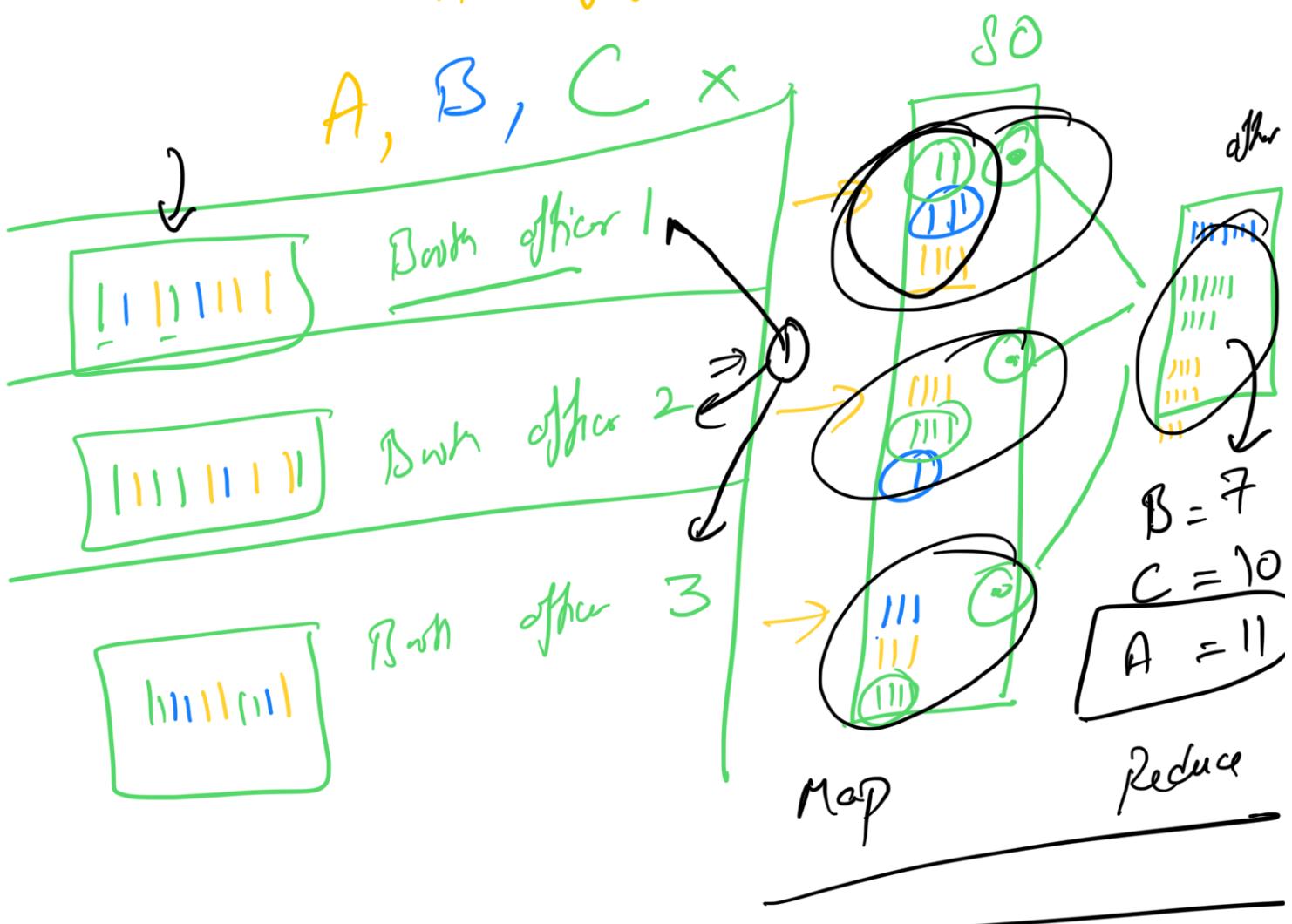
⇒ Distributed Processing

⇒ Map Reduce

↳ is a programming model that simultaneously processes and analyzes huge data sets logically into clusters.

↳ while Map sorts the data, Reduce segregates the data

## Analogy



## Map Reduce

