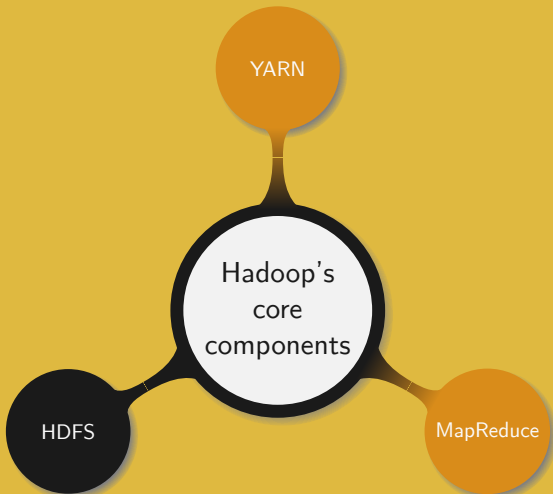


Methods and tools for big data

2. Hadoop's core components

Manuel – Summer 2022



Regular filesystems on a computer:

- Partition
- Hard drive
- LVM

Distributed filesystems:

- Spans several computers
- Has to deal with potential network issues

Idea behind HDFS summarised on slide 1.17

Blocks in HDFS:

- Default size of 128 MB
- Files smaller than a block size do not occupy the whole block
- A file can be larger than a whole disk
- Data and metadata handled separately
- Easy to implement fault tolerance and availability

Two types of node in a cluster:

- Namenode:
 - Maintains FS tree and metadata for all files and directories
 - Locally stores information in namespace image and edit log
 - Knows on which datanodes the blocks of a file are located
- Datanode:
 - Stores and retrieve blocks
 - Regularly reports the list of stored blocks to namenode
 - Can store certain blocks in cache

A namenode has no persistent copy of where blocks are:

- Each datanode announces the blocks it has
- All the information is kept in memory by the namenode
- When a write occurs an entry is added to the edit log

What to do if the namenode fails?

A namenode stores all the blocks of all the files in its memory:

- Assume 1 GB of memory for 1 million blocks
- 200 nodes cluster, 24 TB each: ~ 12 GB of memory
- Cluster at Yahoo!: 25 PB $\rightarrow \sim 64$ GB of memory
- Cluster at Facebook: 60 PB $\rightarrow \sim 156$ GB of memory

How about having more namenodes?

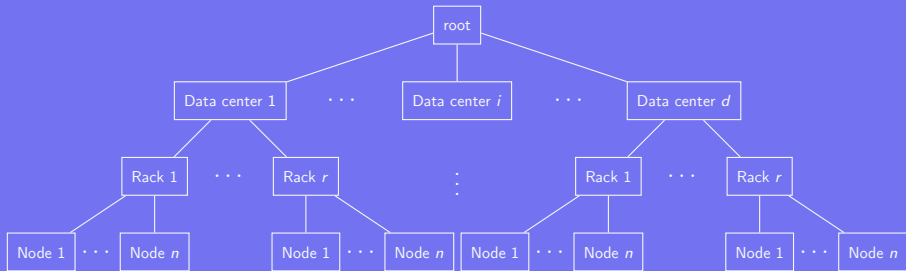
Allowing more namenodes:

- Split the filesystem over several independent namenodes
- Each namenode has a namespace
- Each namespace has its own pool of blocks
- A namespace with a block pool is called namespace volume
- A datanode is not attached to a specific namespace volume

Two namenodes in an active-passive mode:

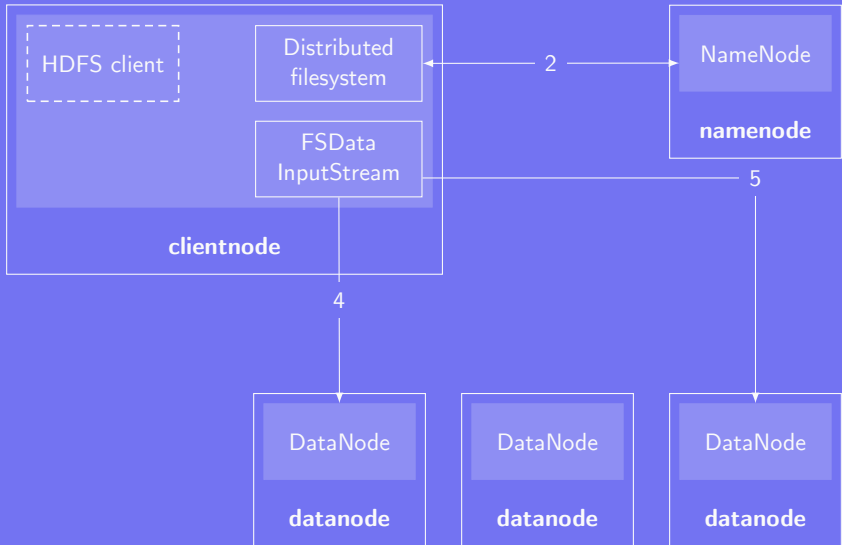
- Passive node takes over in case of failure of the active one
- The two namenodes share the same edit log
- Only the active namenode can write to the edit log
- Passive namenode reads entries when written in edit log
- Datanodes send block reports to both namenodes
- Passive namenode also works as secondary namenode
- Clients must be configured to handle namenode failures

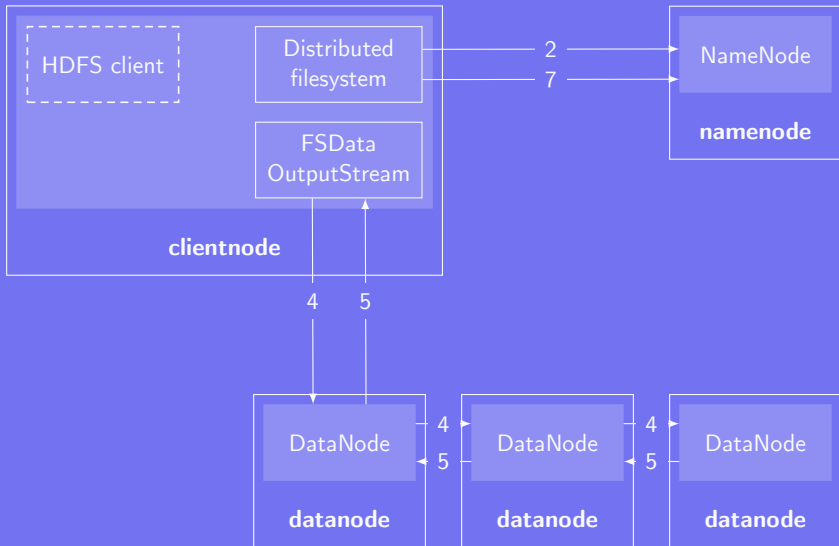
The goal is to evenly spread blocks across the whole cluster:



Replicas' location:

- First: same node as the client
- Second: random, different rack from the first
- Third: same rack as the second but different node
- Others: random nodes in the cluster





When write on a data node fails:

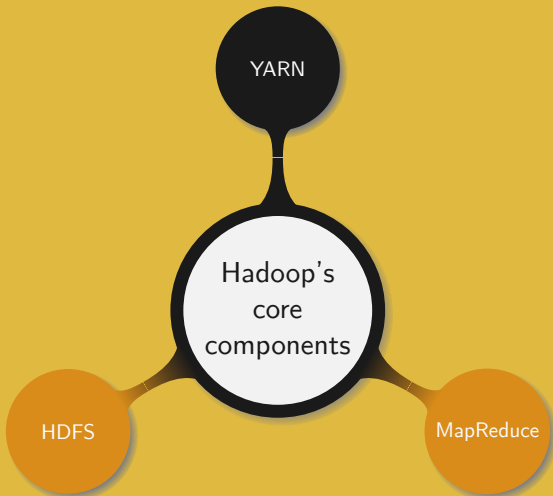
- ① Close the pipeline
- ② Add packets in front of the acknowledgment queue
- ③ Inform the name node of the failing data node
- ④ Remove the faulty data node from the pipeline
- ⑤ Construct a new pipeline using only the healthy data nodes
- ⑥ Complete the writing of the block across the pipeline
- ⑦ Arrange for the replication of under-replicated blocks

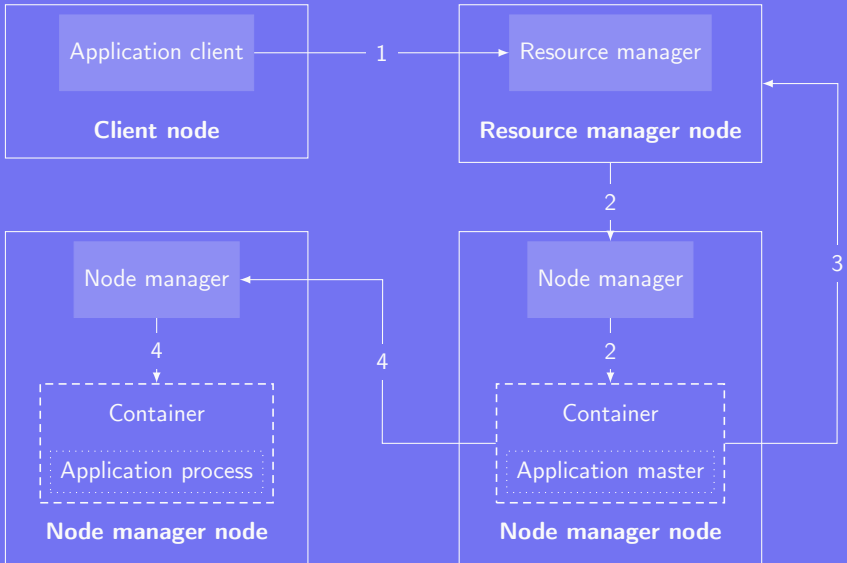
Alternative filesystems built on top of Hadoop abstract filesystem:

- HAR: pack HDFS files into an archive
- View: used to create mount points for federated FS
- FTP: FS backed by an FTP server
- S3: FS based on Amazon S3
- Azure: FS developed by Microsoft Azure

Interfaces to access HDFS data:

- libhdfs: C library with API similar to the Java one
- NFS: use Hadoop NFS gateway
- FUSE: based on libhdfs





A request has two components:

- Amount of resources for each container
- Preferred location of the containers

Common strategy for a request:

- Announce all the resources requests at the beginning
- Dynamically request resources based on the needs

YARN can be used in three ways by applications:

- One application per user job: simplest model
- One application per user session:
 - Containers can be reused between jobs
 - Possibility to cache data between jobs
- Long-running application shared among users:
 - Application master is always “on”
 - Application master acts as a coordinator for other applications

Three schedulers available in YARN:

- FIFO: request served one by one in a queue
- Capacity:
 - Define queues based on the “size” of the jobs to complete
 - All the jobs start early
 - Resources are wasted when unused
- Fair:
 - Resources are dynamically balanced over all the jobs
 - All the resources are fully used
 - Delay due the resource reallocation

Capacity scheduler setup:

- Each queue is handled as a FIFO
- Queue elasticity
 - A single job cannot exceed the capacity of the queue
 - The capacity can be exceeded when several jobs wait and resources are available
- Containers are not preempted
- Can control:
 - The number of resources per user/application
 - The number of applications that can be run at a time
 - Access Control Lists (ACL) on the queues

The challenge is to find a reasonable trade-off for the capacity

Fair scheduler queues:

- One or more queues allowed:
 - Single queue: resources are fairly shared among all applications
 - Several queues, each having:
 - Its own scheduling policy
 - A max/min resources and number of applications
- Queues can be precisely configured using an allocation file
- By default a queue is dynamically created for each user

Preemption in the fair scheduler:

- Efficiency is reduced as killed containers must be restarted
- Two timeout settings t_1 and t_2 to trigger preemption:
 - Minimum share: if a queue waits longer than t_1
 - Fair share: if a queue remains below half of its fair share for longer than t_2
- Timeouts can be set per queue or globally

Locality problem when scheduling:

- An application requests a specific node
- The node is busy
- Should the application wait for the node or loosen its request?

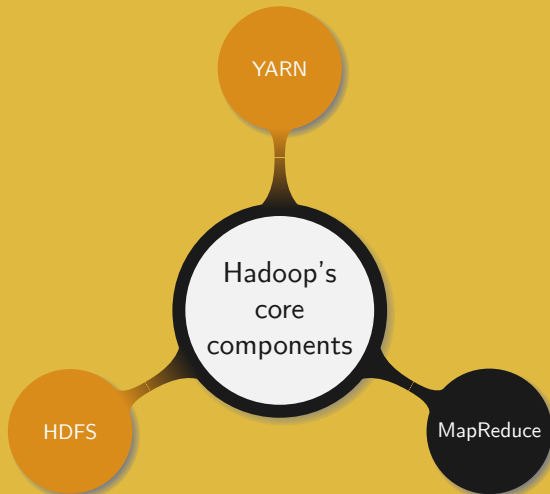
YARN schedulers' approach:

- Every second each node manager sends a heartbeat reporting the running containers and available resources
- Capacity scheduler: wait for a predefined number of heartbeats before loosening the requirement
- Fair scheduler: wait for a predefined portion of nodes in the cluster to offer opportunities before loosening the requirement

How to fairly share resources between applications when they do not use the same type of resources?

Basic idea for two applications:

- Consider the proportion of resources requested for a container by each application
- Call the largest proportion the dominant resource and use it as measure of cluster usage
- Proportionally offer less containers to the more demanding application



Parties involved in a MapReduce job:

- A client which initiates the job
- YARN resource manager
- YARN node manager
- MapReduce application master
- HDFS

Starting a MapReduce job:

- 1 Request a new application ID to the resource manager
- 2 Check the job parameters
- 3 Split the job into subtasks
- 4 Copy the splits and other necessary information to run the job onto the shared FS
- 5 Effectively submit the job on the resource manager

Running a MapReduce job:

- ① YARN scheduler allocates a container
- ② Application master launched by the resource manager
- ③ Setup the tasks
- ④ Retrieve the splits from the shared FS
- ⑤ Create a Map task for each split and specify the number of tasks for the Reduce part
- ⑥ Resources for
 - a Small tasks: run on the same node
 - b Large tasks: contact the resource manager for more containers
 - i Request resources for the maps (high priority)
 - ii Request resources for the reducers when enough maps have completed
- ⑦ Locate the data on the distributed FS and start the task

Potential points of failure:

- Task failure
- Application master failure
- Node manager failure: YARN level (no heartbeat received)
- Resource manager failure: YARN level (high availability mode)

Protection and progress monitoring:

- Tasks are run in a separate JVM: avoid crashing namenode
- Each task has a status and some counters
- Each task reports its progress to the application master
- Client application polls the application master every second to retrieve the latest status

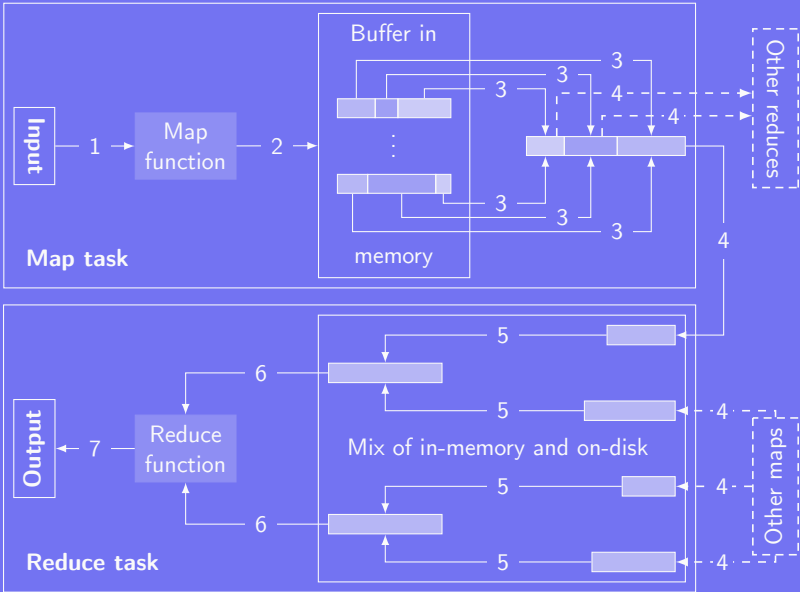
Task failure occurrences and solutions:

- The map or reduce task fails:
 - When receiving a failure notice the application master marks the task as failed
 - The container is freed and resources released
- The JVM crashes: the node manager notices the application manager of the failure
- A task hangs:
 - Tasks marked as failed if no report is received
 - The JVM is killed by the application master

Failed tasks are rescheduled on a different nodemanager

Failure detection and recovery:

- The application master sends periodic heartbeats to the resource manager
- On failure a new instance is run on a new container
- The tasks progress is known so it is possible to resume without re-running the completed tasks
- Each task caches the application master's address
- Use a timeout after which the resource manager is contacted for the new application's master address

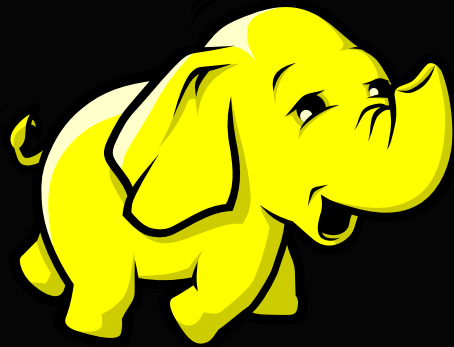


Optimized configuration setup:

- Provide shuffle with as much memory as possible
- Keep enough memory for map and reduce functions
- Optimize the code with respect to memory consumption
- Minimize the number of spills for the map part
- As much as possible keep intermediate reduce data in memory

Speculative execution:

- A task is detected as much slower than average
- Re-run it on a different node
- Kill all the other duplicates as soon as one completes



Thank you!