

Hadoop Multi-Cluster Lab

Mahnoor Adeel - 26913 | Zuha Aqib - 26106

Purpose

This lab sets up two independent Hadoop clusters for learning HDFS administration, YARN operations, and inter-cluster data transfer using *DistCp*

Note: DistCp (Distributed Copy) is a Hadoop tool used for large-scale data copying between HDFS clusters (or between HDFS and compatible file systems such as S3, Azure Blob, etc.). It's designed to efficiently copy massive datasets in a parallel, fault-tolerant manner using MapReduce.

Prerequisites

- Ubuntu host with Docker and Docker Compose installed
- At least 16 GB RAM available
- Ports 9870, 8088, 9970, 9088 should be free on your host

Setup Instructions – First we will create Cluster 1 and then Cluster 2

Step 1: Create Lab Directory

- `mkdir -p ~/hadoop-multicloud-lab`
- `cd ~/hadoop-multicloud-lab`

Step 2: Create docker-compose-cluster1.yml

Save this file in `~/hadoop-multicloud-lab/docker-compose-cluster1.yml`:

- `version: "3.8"`
-
- `services:`
- `namenode-c1:`
- `image: bde2020/hadoop-namenode:2.0.0-hadoop3.2.1-java8`
- `container_name: namenode-c1`
- `hostname: namenode-c1`
- `ports:`
- `- "9870:9870"`
- `- "8020:8020"`
- `volumes:`

```
●      - ./cluster1-data/namenode:/hadoop/dfs/name
●      environment:
●          - CLUSTER_NAME=hadoop-cluster-1
●          - CORE_CONF_fs_defaultFS=hdfs://namenode-c1:8020
●      networks:
●          - cluster1-net
●          - shared-net
●      restart: unless-stopped
●
●      datanode1-c1:
●          image: bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8
●          container_name: datanode1-c1
●          hostname: datanode1-c1
●          depends_on:
●              - namenode-c1
●          ports:
●              - "9864:9864"
●          volumes:
●              - ./cluster1-data/datanode1:/hadoop/dfs/data
●          environment:
●              - CORE_CONF_fs_defaultFS=hdfs://namenode-c1:8020
●      networks:
●          - cluster1-net
●          - shared-net
●      restart: unless-stopped
●
●      datanode2-c1:
●          image: bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8
●          container_name: datanode2-c1
●          hostname: datanode2-c1
●          depends_on:
●              - namenode-c1
●          ports:
●              - "9865:9864"
●          volumes:
●              - ./cluster1-data/datanode2:/hadoop/dfs/data
●          environment:
●              - CORE_CONF_fs_defaultFS=hdfs://namenode-c1:8020
●      networks:
●          - cluster1-net
●          - shared-net
●      restart: unless-stopped
●
●      datanode3-c1:
```

```
●  image: bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8
●  container_name: datanode3-c1
●  hostname: datanode3-c1
●  depends_on:
●    - namenode-c1
●  ports:
●    - "9866:9864"
●  volumes:
●    - ./cluster1-data/datanode3:/hadoop/dfs/data
●  environment:
●    - CORE_CONF_fs_defaultFS=hdfs://namenode-c1:8020
●  networks:
●    - cluster1-net
●    - shared-net
●  restart: unless-stopped
●
●  resourcemanager-c1:
●  image: bde2020/hadoop-resourcemanager:2.0.0-hadoop3.2.1-java8
●  container_name: resourcemanager-c1
●  hostname: resourcemanager-c1
●  ports:
●    - "8088:8088"
●  depends_on:
●    - namenode-c1
●  environment:
●    - CORE_CONF_fs_defaultFS=hdfs://namenode-c1:8020
●    - YARN_CONF_yarn_resourcemanager_hostname=resourcemanager-c1
●  networks:
●    - cluster1-net
●    - shared-net
●  restart: unless-stopped
●
●  nodemanager1-c1:
●  image: bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8
●  container_name: nodemanager1-c1
●  hostname: nodemanager1-c1
●  depends_on:
●    - resourcemanager-c1
●  ports:
●    - "8042:8042"
●  environment:
●    - CORE_CONF_fs_defaultFS=hdfs://namenode-c1:8020
●    - YARN_CONF_yarn_resourcemanager_hostname=resourcemanager-c1
●    - YARN_CONF_yarn_nodemanager_aux__services=mapreduce_shuffle
```

```

●   - YARN_CONF_yarn_nodemanager_resource_memory_mb=2048
●   - YARN_CONF_yarn_nodemanager_resource_cpu_vcores=2
●   networks:
●     - cluster1-net
●     - shared-net
●   restart: unless-stopped
●
●   nodemanager2-c1:
●     image: bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8
●     container_name: nodemanager2-c1
●     hostname: nodemanager2-c1
●     depends_on:
●       - resourcemanager-c1
●     ports:
●       - "8043:8042"
●     environment:
●       - CORE_CONF_fs_defaultFS=hdfs://namenode-c1:8020
●       - YARN_CONF_yarn_resourcemanager_hostname=resourcemanager-c1
●       - YARN_CONF_yarn_nodemanager_aux_services=mapreduce_shuffle
●       - YARN_CONF_yarn_nodemanager_resource_memory_mb=2048
●       - YARN_CONF_yarn_nodemanager_resource_cpu_vcores=2
●     networks:
●       - cluster1-net
●       - shared-net
●   restart: unless-stopped
●
●   networks:
●     cluster1-net:
●     driver: bridge
●     shared-net:
●     external: true

```

Task: Summarize the contents of this file in a few sentences

The compose file defines a single-cluster Hadoop setup. It launches namenode, datanodes, resourcemanager, and nodemanagers. All the containers belong to Hadoop Cluster 1. They are connected through two Docker networks — cluster1-net, for in-cluster communication, and shared-net, for communication with other clusters.

- The Namenode runs the image `bde2020/hadoop-namenode:2.0.0-hadoop3.2.1-java8` and exposes ports 9870 for the HDFS web UI and 8020 for HDFS RPC communication.
- The Datanodes the image `bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8` and expose ports 9864, 9865, and 9866 respectively on the host.
- The ResourceManager, running the image `bde2020/hadoop-resourcemanager:2.0.0-hadoop3.2.1-java8`, is accessible on port 8088, providing the YARN web interface.

- *The NodeManagers use the image bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8 and expose ports 8042 and 8043 respectively for their monitoring UIs.*

Step 3: Create docker-compose-cluster2.yml

Save this file in ~/hadoop-multicloud-lab/docker-compose-cluster2.yml:

```

● version: "3.8"
  ○
● services:
● namenode-c2:
●   image: bde2020/hadoop-namenode:2.0.0-hadoop3.2.1-java8
●   container_name: namenode-c2
●   hostname: namenode-c2
●   ports:
●     - "9970:9870"
●     - "8021:8020"
●   volumes:
●     - ./cluster2-data/namenode:/hadoop/dfs/name
●   environment:
●     - CLUSTER_NAME=hadoop-cluster-2
●     - CORE_CONF_fs_defaultFS=hdfs://namenode-c2:8020
●   networks:
●     - cluster2-net
●     - shared-net
●   restart: unless-stopped
●
● datanode1-c2:
●   image: bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8
●   container_name: datanode1-c2
●   hostname: datanode1-c2
●   depends_on:
●     - namenode-c2
●   ports:
●     - "9964:9864"
●   volumes:
●     - ./cluster2-data/datanode1:/hadoop/dfs/data
●   environment:
●     - CORE_CONF_fs_defaultFS=hdfs://namenode-c2:8020
●   networks:
●     - cluster2-net
●     - shared-net
●   restart: unless-stopped

```

-
- datanode2-c2:
 - image: bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8
 - container_name: datanode2-c2
 - hostname: datanode2-c2
 - depends_on:
 - namenode-c2
 - ports:
 - "9965:9864"
 - volumes:
 - ./cluster2-data/datanode2:/hadoop/dfs/data
 - environment:
 - CORE_CONF_fs_defaultFS=hdfs://namenode-c2:8020
 - networks:
 - cluster2-net
 - shared-net
 - restart: unless-stopped
 -
- datanode3-c2:
 - image: bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8
 - container_name: datanode3-c2
 - hostname: datanode3-c2
 - depends_on:
 - namenode-c2
 - ports:
 - "9966:9864"
 - volumes:
 - ./cluster2-data/datanode3:/hadoop/dfs/data
 - environment:
 - CORE_CONF_fs_defaultFS=hdfs://namenode-c2:8020
 - networks:
 - cluster2-net
 - shared-net
 - restart: unless-stopped
 -
- resourcemanager-c2:
 - image: bde2020/hadoop-resourcemanager:2.0.0-hadoop3.2.1-java8
 - container_name: resourcemanager-c2
 - hostname: resourcemanager-c2
 - ports:
 - "9088:8088"
 - depends_on:
 - namenode-c2
 - environment:

- - CORE_CONF_fs_defaultFS=hdfs://namenode-c2:8020
- - YARN_CONF_yarn_resourcemanager_hostname=resourcemanager-c2
- networks:
 - - cluster2-net
 - - shared-net
- restart: unless-stopped
-
- nodemanager1-c2:
 - image: bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8
 - container_name: nodemanager1-c2
 - hostname: nodemanager1-c2
 - depends_on:
 - - resourcemanager-c2
 - ports:
 - - "8142:8042"
 - environment:
 - - CORE_CONF_fs_defaultFS=hdfs://namenode-c2:8020
 - - YARN_CONF_yarn_resourcemanager_hostname=resourcemanager-c2
 - - YARN_CONF_yarn_nodemanager_aux__services=mapreduce_shuffle
 - - YARN_CONF_yarn_nodemanager_resource_memory__mb=2048
 - - YARN_CONF_yarn_nodemanager_resource_cpu__vcores=2
 - networks:
 - - cluster2-net
 - - shared-net
 - restart: unless-stopped
-
- nodemanager2-c2:
 - image: bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8
 - container_name: nodemanager2-c2
 - hostname: nodemanager2-c2
 - depends_on:
 - - resourcemanager-c2
 - ports:
 - - "8143:8042"
 - environment:
 - - CORE_CONF_fs_defaultFS=hdfs://namenode-c2:8020
 - - YARN_CONF_yarn_resourcemanager_hostname=resourcemanager-c2
 - - YARN_CONF_yarn_nodemanager_aux__services=mapreduce_shuffle
 - - YARN_CONF_yarn_nodemanager_resource_memory__mb=2048
 - - YARN_CONF_yarn_nodemanager_resource_cpu__vcores=2
 - networks:
 - - cluster2-net
 - - shared-net
 - restart: unless-stopped

- networks:
- cluster2-net:
- driver: bridge
- shared-net:
- external: true

Task: Summarize the contents of this file in a few sentences

All the containers in this setup belong to Hadoop Cluster 2. They are connected through two networks, cluster2-net, for internal cluster communication, and shared-net for communication with other clusters.

Cluster 2 is almost identical to Cluster 1 in setup and configuration, except that it uses different port mappings to avoid conflicts when both clusters run simultaneously.

Step 4: Create Shared Network

- docker network create shared-net

Step 5: Start Cluster 1

- cd ~/hadoop-multicloud-lab
- docker compose -f docker-compose-cluster1.yml up -d
- sleep 15
- docker ps --filter "name=c1"

Task: Why the sleep 15 command?

This gives Cluster 1 containers enough time to start up and initialize properly. Otherwise, Docker might still be initializing containers, and docker ps command could show incomplete or missing results.

Task: Take a screenshot showing all Cluster 1 containers running

```
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ docker compose -f docker-compose-cluster1.yml up -d
sleep 15
docker ps --filter "name=c1"
WARN[0000] /home/madeel/mahnoor-hadoop-multicloud-lab/docker-compose-cluster1.yml: the attribute 'version' is obsolete, it will be ignored, please remove it to avoid potential confusion
[+] Running 8/8
✓ Network mahnoor-hadoop-multicloud-lab_cluster1-net Created          0.1s
✓ Container namenode-c1 Started           0.6s
✓ Container datanode1-c1 Started          1.0s
✓ Container datanode3-c1 Started          0.9s
✓ Container resourcemanager-c1 Started      1.1s
✓ Container datanode2-c1 Started          1.1s
✓ Container nodemanager1-c1 Started        1.5s
✓ Container nodemanager2-c1 Started        1.9s
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ docker ps --filter "name=c1"
CONTAINER ID IMAGE NAMES COMMAND CREATED STATUS PORTS
195e0680b9ff bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." About a minute ago Up About a minute (healthy) 0.0.0.0:8043->8042/tcp, [::]:8043->8042/
p nodemanager2-c1
3336163ff5c2 bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." About a minute ago Up About a minute (healthy) 0.0.0.0:8042->8042/tcp, [::]:8042->8042/
p nodemanager1-c1
28188a288cb2 bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." About a minute ago Up About a minute (healthy) 0.0.0.0:9866->9864/tcp, [::]:9866->9864/
p datanode3-c1
80096806382b bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." About a minute ago Up About a minute (healthy) 0.0.0.0:9864->9864/tcp, [::]:9864->9864/
p datanode1-c1
9dbb5bb36ac0 bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." About a minute ago Up About a minute (healthy) 0.0.0.0:9865->9864/tcp, [::]:9865->9864/
p datanode2-c1
1f0304a868c8 bde2020/hadoop-resourcemanager:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." About a minute ago Up About a minute (healthy) 0.0.0.0:8088->8088/tcp, [::]:8088->8088/
p resourcemanager-c1
06a6b8c75d8f bde2020/hadoop-namenode:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." About a minute ago Up About a minute (healthy) 0.0.0.0:8020->8020/tcp, [::]:8020->8020/
p ,0.0.0.0:9870->9870/tcp, [::]:9870->9870/tcp namenode-c1
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$
```

Step 6: Start Cluster 2

- docker compose -f docker-compose-cluster2.yml up -d
- sleep 15
- docker ps --filter "name=c2"

Task: Take a screenshot showing all Cluster 2 containers running

```
service in your compose file, you can run this command with the --remove-orphan flag to clean it up.
[+] Running 8/8
✓ Network mahnoor-hadoop-multicuster-lab_cluster2-net  Created
✓ Container namenode-c2  Started
✓ Container datanode1-c2  Started
✓ Container resourcemanager-c2  Started
✓ Container datanode2-c2  Started
✓ Container datanode3-c2  Started
✓ Container nodemanager1-c2  Started
✓ Container nodemanager2-c2  Started
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ docker ps --filter "name=c2"
CONTAINER ID IMAGE NAMES COMMAND CREATED STATUS PORTS
5fd6383c7e30 bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." 34 seconds ago Up 32 seconds (healthy) 0.0.0.0:8142->8042/tcp, [::]:8142->8042/
p nodemanager1-c2
5e2e1ca3bf3b bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." 34 seconds ago Up 32 seconds (healthy) 0.0.0.0:8143->8042/tcp, [::]:8143->8042/
p nodemanager2-c2
489b5720bfa3 bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." 34 seconds ago Up 33 seconds (healthy) 0.0.0.0:9964->9864/tcp, [::]:9964->9864/tcp
p datanode1-c2
c6dc1308cb39 bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." 34 seconds ago Up 33 seconds (healthy) 0.0.0.0:9966->9864/tcp, [::]:9966->9864/tcp
p datanode3-c2
96b948edd9ff bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." 34 seconds ago Up 32 seconds (healthy) 0.0.0.0:9965->9864/tcp, [::]:9965->9864/tcp
p datanode2-c2
0ab87b503c0c bde2020/hadoop-resourcemanager:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." 34 seconds ago Up 33 seconds (healthy) 0.0.0.0:9088->8088/tcp, [::]:9088->8088/tcp
p resourcemanager-c2
3c9ae45c5934 bde2020/hadoop-namenode:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..." 34 seconds ago Up 33 seconds (healthy) 0.0.0.0:8021->8020/tcp, [::]:8021->8020/tcp, 0.0
,0.0.0.0:9870->9870/tcp, [::]:9870->9870/tcp namenode-c2
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$
```

Step 7: Verify Network Connectivity

- docker exec namenode-c1 ping -c 3 namenode-c2
- docker exec namenode-c2 ping -c 3 namenode-c1

Task: Screenshot showing successful ping between clusters – also explain the outputs.

```
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ docker exec namenode-c1 ping -c 3 namenode-c2
PING namenode-c2 (172.19.0.9) 56(84) bytes of data.
64 bytes from namenode-c2.shared-net (172.19.0.9): icmp_seq=1 ttl=64 time=0.167 ms
64 bytes from namenode-c2.shared-net (172.19.0.9): icmp_seq=2 ttl=64 time=0.059 ms
64 bytes from namenode-c2.shared-net (172.19.0.9): icmp_seq=3 ttl=64 time=0.079 ms

--- namenode-c2 ping statistics ---
3 packets transmitted, 3 received, 0% packet loss, time 2025ms
rtt min/avg/max/mdev = 0.059/0.101/0.167/0.048 ms
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec namenode-c2 ping -c 3 namenode-c1
PING namenode-c1 (172.19.0.2) 56(84) bytes of data.
64 bytes from namenode-c1.shared-net (172.19.0.2): icmp_seq=1 ttl=64 time=0.056 ms
64 bytes from namenode-c1.shared-net (172.19.0.2): icmp_seq=2 ttl=64 time=0.072 ms
64 bytes from namenode-c1.shared-net (172.19.0.2): icmp_seq=3 ttl=64 time=0.060 ms

--- namenode-c1 ping statistics ---
3 packets transmitted, 3 received, 0% packet loss, time 2071ms
rtt min/avg/max/mdev = 0.056/0.062/0.072/0.011 ms
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$
```

Web UI Access

Cluster 1 UIs:

- **NameNode:** <http://localhost:9870>
- **ResourceManager:** <http://localhost:8088>
- **DataNode 1:** <http://localhost:9864>

Cluster 2 UIs:

- **NameNode:** <http://localhost:9970>
- **ResourceManager:** <http://localhost:9088>
- **DataNode 1:** <http://localhost:9964>

Task: Open both NameNode UIs side-by-side showing the cluster overview. Explain the outputs.

Overview 'namenode-c1:8020' (active)

Started:	Wed Oct 22 22:02:24 +0500 2025
Version:	3.2.1, rb3cb6b467e22ea32b380ff4b7601d07e0bf3842
Compiled:	Tue Sep 10 20:56:00 +0500 2019 by rohitsharmaks from branch-3.2.1
Cluster ID:	CID-0110259-865a-4105-a1e1-8092923a5789
Block Pool ID:	BP-1865643019-172.19.0.2-1761152541275

Summary

Security is off.
Safemode is off.
1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem objects.
Heap Memory used 126.45 MB of 252.5 MB Heap Memory. Max Heap Memory is 1.72 GB.
Non Heap Memory used 47.94 MB of 49.44 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.
Configured Capacity: 136.6 GB
Configured Remote Capacity: 0 B
DFS Used: 84 KB (0%)
Non DFS Used: 105.01 GB
DFS Remaining: 24.56 GB (17.98%)
Block Pool Used: 84 KB (0%)
DataNodes usages% (Min/Median/Max/stdDev): 0.00% / 0.00% / 0.00% / 0.00%
Live Nodes : 3 (Decommissioned: 0, In Maintenance: 0)
Dead Nodes : 0 (Decommissioned: 0, In Maintenance: 0)
Decommissioning Nodes : 0
Entering Maintenance Nodes : 0
Total Datanode Volume Failures : 0 (0 B)
Number of Under-Replicated Blocks : 0
Number of Blocks Pending Deletion (including replicas) : 0
Block Deletion Start Time : Wed Oct 22 22:02:24 +0500 2025
Last Checkpoint Time : Wed Oct 22 22:02:21 +0500 2025
Enabled Erasure Coding Policies : RS-6-3-1024k

Overview 'namenode-c2:8020' (active)

Started:	Wed Oct 22 22:04:40 +0500 2025
Version:	3.2.1, rb3cb6b467e22ea32b380ff4b7601d07e0bf3842
Compiled:	Tue Sep 10 20:56:00 +0500 2019 by rohitsharmaks from branch-3.2.1
Cluster ID:	CID-922ef1232-481c-4182-bf6d-9e6c6df8beda
Block Pool ID:	BP-1899536788-172.19.0.9-1761152677037

Summary

Security is off.
Safemode is off.
1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem objects.
Heap Memory used 95.19 MB of 351.5 MB Heap Memory. Max Heap Memory is 1.72 GB.
Non Heap Memory used 47.34 MB of 48.38 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.
Configured Capacity: 136.6 GB
Configured Remote Capacity: 0 B
DFS Used: 84 KB (0%)
Non DFS Used: 105.06 GB
DFS Remaining: 24.51 GB (17.94%)
Block Pool Used: 84 KB (0%)
DataNodes usages% (Min/Median/Max/stdDev): 0.00% / 0.00% / 0.00% / 0.00%
Live Nodes : 3 (Decommissioned: 0, In Maintenance: 0)
Dead Nodes : 0 (Decommissioned: 0, In Maintenance: 0)
Decommissioning Nodes : 0
Entering Maintenance Nodes : 0
Total Datanode Volume Failures : 0 (0 B)
Number of Under-Replicated Blocks : 0
Number of Blocks Pending Deletion (including replicas) : 0
Block Deletion Start Time : Wed Oct 22 22:04:40 +0500 2025
Last Checkpoint Time : Wed Oct 22 22:04:37 +0500 2025
Enabled Erasure Coding Policies : RS-6-3-1024k

Lab Exercises

Exercise 1: Create and Upload Your Own Dataset to Cluster 1

Task: Create a custom dataset file on your local machine, then upload it to HDFS on Cluster 1. Show proof of transfer and give some description about this dataset (it must be a big dataset or from amongst the data sources I uploaded).

What to do:

1. Create a text file with at least 20 lines of data (use any text editor or echo commands)
2. Upload this file to HDFS under /user/yourusername/dataset/
3. Verify the upload was successful by listing the directory
4. Display the contents of your uploaded file

Task: Screenshot of the terminal output of your file listing and content display commands (**Hint:** You'll need to use `hdfs dfs -mkdir`, `hdfs dfs -put`, `hdfs dfs -ls`, and `hdfs dfs -cat`)

```
root@namenode-c1:/# hdfs dfs -put /tmp/mydataset.txt /user/madeel/dataset/
2025-10-22 17:45:43,494 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
root@namenode-c1:/# hdfs dfs -ls /user/madeel/dataset
Found 1 items
-rw-r--r-- 3 root supergroup          618 2025-10-22 17:45 /user/madeel/dataset/mydataset.txt
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs dfs -cat /user/madeel/dataset/mydataset.txt
2025-10-22 17:46:59,305 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
1, Alice, 25, Data Scientist
2, Bob, 30, Machine Learning Engineer
3, Charlie, 28, Data Analyst
4, David, 35, Software Developer
5, Emma, 22, Research Intern
6, Frank, 40, System Architect
7, Grace, 26, Database Administrator
8, Henry, 33, Hadoop Engineer
9, Ivy, 27, Business Analyst
10, Jack, 29, Statistician
11, Karen, 31, Cloud Engineer
12, Leo, 24, Data Engineer
13, Mia, 32, Research Scientist
14, Noah, 38, Product Manager
15, Olivia, 21, Student Researcher
16, Peter, 34, AI Specialist
17, Quinn, 25, Data Associate
18, Rose, 28, ML Researcher
19, Steve, 29, DevOps Engineer
20, Tina, 23, Junior Data Analyst
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$
```

Exercise 2: Check File Replication and Block Distribution

Task: Examine how your uploaded file is replicated across DataNodes.

What to do:

1. Use HDFS commands to get detailed information about your file (replication factor, block size)
2. Check which DataNodes are storing copies of your file
3. Take note of the block IDs and their locations

Screenshot of the output displaying file status and block locations

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -ls /user/madeel/dataset/mydataset.txt
-rw-r--r-- 3 root supergroup 618 2025-10-22 17:45 /user/madeel/dataset/mydataset.txt
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$
```

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs fsck /user/madeel/dataset/mydataset.txt -files -blocks -locations
Connecting to namenode via http://namenode-c1:9870/fsck?gi=root&files=1&blocks=1&locations=1&path=%2Fuser%2Fmadeel%2Fdataset%2Fmydataset.txt
FSCK started by root (auth:SIMPLE) from /172.19.0.2 for path /user/madeel/dataset/mydataset.txt at Wed Oct 22 17:51:41 UTC 2025
/usr/madeel/dataset/mydataset.txt 618 bytes, replicated: replication=3, 1 block(s): OK
0. BP-1865643019-172.19.0.2-1761152541275:blk_1073741825_1001 len=618 Live_repl=3 [DatanodeInfoWithStorage[172.20.0.3:9866,DS-44f5d53e-18f2-401e-8518-420da57cc385,1
rage[172.20.0.4:9866,DS-ab2ca4fa-4d34-480f-b093-f0615e4b82cf,DISK], DatanodeInfoWithStorage[172.20.0.5:9866,DS-404605dc-1773-42fd-8fc7-71ac0c360b59,DISK]]]

Status: HEALTHY
Number of data-nodes: 3
Number of racks: 1
Total dirs: 0
Total symlinks: 0

Replicated Blocks:
Total size: 618 B
Total files: 1
Total blocks (validated): 1 (avg. block size 618 B)
Minimally replicated blocks: 1 (100.0 %)
Over-replicated blocks: 0 (0.0 %)
Under-replicated blocks: 0 (0.0 %)
Mis-replicated blocks: 0 (0.0 %)
Default replication factor: 3
Average block replication: 3.0
Missing blocks: 0
Corrupt blocks: 0
Missing replicas: 0 (0.0 %)
```

Then verify in UI: Open <http://localhost:9870>, browse to your file, and compare what you see in the UI with the command output (*Hint: Use hdfs fsck*

/user/yourname/dataset/yourfile -files -blocks -locations)

The screenshot shows the HDFS web UI with the URL `/user/madeel/dataset` in the address bar. The page title is "Browse Directory". A search bar at the top right contains the placeholder "Search:". Below the search bar is a table header with columns: "Show", "Permission", "Owner", "Group", "Size", "Last Modified", "Replication", "Block Size", and "Name". A single entry is listed: "mydataset.txt" with a size of 618 B, last modified on Oct 22 22:45, replication factor of 3, and a block size of 128 MB. The "Name" column shows a small icon of a document. At the bottom left, it says "Showing 1 to 1 of 1 entries". On the right, there are "Previous" and "Next" buttons, with "1" highlighted. The footer of the main page says "Hadoop, 2019."

This screenshot shows the same HDFS web UI interface, but a modal dialog box is open over the "Browse Director" page. The dialog is titled "File information - mydataset.txt". It contains tabs for "Download", "Head the file (first 32K)", and "Tail the file (last 32K)". The "Download" tab is selected. Inside the dialog, under "Block information", it says "Block 0" (with a dropdown arrow). Below that, it lists: "Block ID: 1073741825", "Block Pool ID: BP-1865643019-172.19.0.2-1761152541275", "Generation Stamp: 1001", and "Size: 618". Under "Availability:", there is a bulleted list: "• datanode3-c1", "• datanode2-c1", and "• datanode1-c1". At the bottom right of the dialog is a "Close" button.

Both the command-line output and the HDFS web UI show consistent information about the file, including its size, replication factor, and the DataNodes storing each block. The web UI shows in a more easy-to-read format.

Exercise 3: Test Fault Tolerance - DataNode Failure

Task: Stop one DataNode and observe how the cluster maintains data availability.

What to do:

1. Check current cluster status (number of live nodes)
2. Stop datanode2-c1
3. Wait 30 seconds, then check cluster status again
4. Try to read your file - it should still be accessible!
5. Check which DataNodes are now serving the file
6. Restart the DataNode and verify it rejoins the cluster

Task: Screenshot cluster status before DataNode failure (all 3 nodes live)

```
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ docker exec -it namenode-c1 hdfs dfsadmin -report
Configured Capacity: 146675011584 (136.60 GB)
Present Capacity: 26282741760 (24.48 GB)
DFS Remaining: 26281783296 (24.48 GB)
DFS Used: 958464 (936 KB)
DFS Used%: 0.00%
Replicated Blocks:
    Under replicated blocks: 0
    Blocks with corrupt replicas: 0
    Missing blocks: 0
    Missing blocks (with replication factor 1): 0
    Low redundancy blocks with highest priority to recover: 0
    Pending deletion blocks: 0
Erasure Coded Block Groups:
    Low redundancy block groups: 0
    Block groups with corrupt internal blocks: 0
    Missing block groups: 0
    Low redundancy blocks with highest priority to recover: 0
    Pending deletion blocks: 0
-----
Live datanodes (3):
Name: 172.20.0.3:9866 (datanode3-c1.mahnoor-hadoop-multicuster-lab_cluster1-net)
Hostname: datanode3-c1
Decommission Status : Normal
Configured Capacity: 48891670528 (45.53 GB)
DFS Used: 319488 (312 KB)
Non DFS Used: 37613965312 (35.03 GB)
DFS Remaining: 8760594432 (8.16 GB)
```

Task: Screenshot cluster status after failure (2 nodes live) and successful file read despite the failure

```
Live datanodes (2):
```

```
Name: 172.20.0.3:9866 (datanode3-c1.mahnoor-hadoop-multicloud-lab_cluster1-net)
Hostname: datanode3-c1
Decommission Status : Normal
Configured Capacity: 48891670528 (45.53 GB)
DFS Used: 319488 (312 KB)
Non DFS Used: 37613969408 (35.03 GB)
DFS Remaining: 8760590336 (8.16 GB)
DFS Used%: 0.00%
DFS Remaining%: 17.92%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Sun Nov 02 07:32:25 UTC 2025
Last Block Report: Sun Nov 02 03:19:18 UTC 2025
Num of Blocks: 4
```

```
Name: 172.20.0.4:9866 (datanode1-c1.mahnoor-hadoop-multicloud-lab_cluster1-net)
Hostname: datanode1-c1
Decommission Status : Normal
Configured Capacity: 48891670528 (45.53 GB)
DFS Used: 319488 (312 KB)
Non DFS Used: 37613969408 (35.03 GB)
DFS Remaining: 8760590336 (8.16 GB)
DFS Used%: 0.00%
DFS Remaining%: 17.92%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Sun Nov 02 07:32:23 UTC 2025
Last Block Report: Sun Nov 02 05:13:14 UTC 2025
Num of Blocks: 4
```

```
Dead datanodes (1):
```

```
Name: 172.20.0.5:9866 (172.20.0.5)
Hostname: datanode2-c1
Decommission Status : Normal
Configured Capacity: 48891670528 (45.53 GB)
DFS Used: 319488 (312 KB)
Non DFS Used: 37613965312 (35.03 GB)
DFS Remaining: 8760594432 (8.16 GB)
DFS Used%: 0.00%
DFS Remaining%: 17.92%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Sun Nov 02 07:17:52 UTC 2025
Last Block Report: Sun Nov 02 04:55:04 UTC 2025
Num of Blocks: 4
```

File Still accessible:

```
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ docker exec -it namenode-c1 hdfs dfs -cat /user/madeel/dataset/mydataset.txt
2025-11-02 07:34:50,260 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted =
1, Alice, 25, Data Scientist
2, Bob, 30, Machine Learning Engineer
3, Charlie, 28, Data Analyst
4, David, 35, Software Developer
5, Emma, 22, Research Intern
6, Frank, 40, System Architect
7, Grace, 26, Database Administrator
8, Henry, 33, Hadoop Engineer
9, Ivy, 27, Business Analyst
10, Jack, 29, Statistician
11, Karen, 31, Cloud Engineer
12, Leo, 24, Data Engineer
13, Mia, 32, Research Scientist
14, Noah, 38, Product Manager
15, Olivia, 21, Student Researcher
16, Peter, 34, AI Specialist
17, Quinn, 25, Data Associate
18, Rose, 28, ML Researcher
19, Steve, 29, DevOps Engineer
20, Tina, 23, Junior Data Analyst
```

```

madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs fsck /user/madeel/dataset/mydataset.txt -files -blocks -locations
Connecting to namenode via http://namenode-c1:9870/fsck?ugi=root&files=1&blocks=1&locations=1&path=%2Fuser%2Fmadeel%2Fdataset%2Fmydataset.txt
FSCK started by root (auth:SIMPLE) from /172.19.0.2 for path /user/madeel/dataset/mydataset.txt at Sun Nov 02 08:00:45 UTC 2025
/user/madeel/dataset/mydataset.txt 618 bytes, replicated: replication=3, 1 block(s): Under replicated BP-1865643019-172.19.0.2-1761152541275:blk_1073741825_1001. Tar
and 2 live replica(s), 0 decommissioned replica(s), 0 decommissioning replica(s).
0. BP-1865643019-172.19.0.2-1761152541275:blk_1073741825_1001 len=618 Live_repl=2 [DatanodeInfoWithStorage[172.20.0.3:9866,DS-44f5d53e-18f2-401e-8518-420da57cc385,DI
rage[172.20.0.4:9866,DS-ab2ca4fa-4d34-480f-b093-f0615e4b82cf,DISK]]]

Status: HEALTHY
Number of data-nodes: 2
Number of racks: 1
Total dirs: 0
Total symlinks: 0

Replicated Blocks:
Total size: 618 B
Total files: 1
Total blocks (validated): 1 (avg. block size 618 B)
Minimally replicated blocks: 1 (100.0 %)
Over-replicated blocks: 0 (0.0 %)
Under-replicated blocks: 1 (100.0 %)
Mis-replicated blocks: 0 (0.0 %)
Default replication factor: 3
Average block replication: 2.0
Missing blocks: 0
Corrupt blocks: 0
Missing replicas: 1 (33.333332 %)

Erasure Coded Block Groups:
Total size: 0 B
Total files: 0
Total block groups (validated): 0
Minimally erasure-coded block groups: 0

```

Screenshot

4: Screenshot of the cluster status (*Hint: Use hdfs dfsadmin -report and docker stop/start datanode2-c1*)

```

madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs dfsadmin -report
Configured Capacity: 146675011584 (136.60 GB)
Present Capacity: 26282729472 (24.48 GB)
DFS Remaining: 26281771008 (24.48 GB)
DFS Used: 958464 (936 KB)
DFS Used%: 0.00%
Replicated Blocks:
    Under replicated blocks: 0
    Blocks with corrupt replicas: 0
    Missing blocks: 0
    Missing blocks (with replication factor 1): 0
    Low redundancy blocks with highest priority to recover: 0
    Pending deletion blocks: 0
Erasure Coded Block Groups:
    Low redundancy block groups: 0
    Block groups with corrupt internal blocks: 0
    Missing block groups: 0
    Low redundancy blocks with highest priority to recover: 0
    Pending deletion blocks: 0
-----
Live datanodes (3):

```

Exercise 4: Create Multiple Files with Different Sizes

Task: Create and upload files of different sizes to see how HDFS handles them.

What to do:

1. Create three files:
 - o Small file: ~100 lines (~5KB)
 - o Medium file: ~5000 lines (~250KB)
 - o Large file: ~20000 lines (~1MB)
2. Upload all three to HDFS under /user/yourname/sizetest/
3. Check the block distribution for each file
4. Observe: How many blocks does each file have?

Task: Screenshot the fsck output for all three files, highlighting the number of blocks each file

uses (Hint: Use a for loop to generate files: for i in {1..5000}; do echo "Line \$i content"; done > medium.txt)

Then we generated SMALL MEDIUM LARGE text files

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ nn1 'cd /tmp \
&& for i in {1..100}; do echo "Line $i content"; done > small.txt \
&& for i in {1..5000}; do echo "Line $i content"; done > medium.txt \
&& for i in {1..20000}; do echo "Line $i content"; done > large.txt \
&& ls -lh small.txt medium.txt large.txt \
&& echo "Lines:" \
&& wc -l small.txt medium.txt large.txt'
-rw-r--r-- 1 root root 361K Nov 1 22:32 large.txt
-rw-r--r-- 1 root root 87K Nov 1 22:32 medium.txt
-rw-r--r-- 1 root root 1.6K Nov 1 22:32 small.txt
Lines:
 100 small.txt
 5000 medium.txt
 20000 large.txt
 25100 total
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$
```

Then we uploaded on HDFS

```
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ nnm '
# put Hadoop on PATH
for C in /opt/hadoop* /usr/local/hadoop; do
    if [ -x "$C/bin/hdfs" ]; then export HADOOP_HOME="$C"; break; fi
done
export PATH="$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$PATH"

# make your target dir and upload
hdfs dfs -mkdir -p /user/madeel/sizetest
hdfs dfs -put -f /tmp/small.txt /user/madeel/sizetest/
hdfs dfs -put -f /tmp/medium.txt /user/madeel/sizetest/
hdfs dfs -put -f /tmp/large.txt /user/madeel/sizetest/

# verify upload
hdfs dfs -ls -h /user/madeel/sizetest
'
2025-11-01 22:44:53,272 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2025-11-01 22:44:58,396 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2025-11-01 22:45:00,442 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
Found 3 items
-rw-r--r--  3 madeel supergroup   204.0 K 2025-11-01 22:45 /user/madeel/sizetest
/large.txt
-rw-r--r--  3 madeel supergroup    47.7 K 2025-11-01 22:44 /user/madeel/sizetest
/medium.txt
-rw-r--r--  3 madeel supergroup     792 2025-11-01 22:44 /user/madeel/sizetest
/small.txt
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$
```

Then we observed block distribution for SMALL

```
Status: HEALTHY
Number of data-nodes: 3
Number of racks: 1
Total dirs: 0
Total symlinks: 0

Replicated Blocks:
Total size: 792 B
Total files: 1
Total blocks (validated): 1 (avg. block size 792 B)
Minimally replicated blocks: 1 (100.0 %)
Over-replicated blocks: 0 (0.0 %)
Under-replicated blocks: 0 (0.0 %)
Mis-replicated blocks: 0 (0.0 %)
Default replication factor: 3
Average block replication: 3.0
Missing blocks: 0
Corrupt blocks: 0
Missing replicas: 0 (0.0 %)

Erasure Coded Block Groups:
Total size: 0 B
Total files: 0
Total block groups (validated): 0
Minimally erasure-coded block groups: 0
Over-erasure-coded block groups: 0
Under-erasure-coded block groups: 0
Unsatisfactory placement block groups: 0
Average block group size: 0.0
Missing block groups: 0
Corrupt block groups: 0
Missing internal blocks: 0
```

THEN MEDIUM

```
STATUS: HEALTHY
Number of data-nodes: 3
Number of racks: 1
Total dirs: 0
Total symlinks: 0

Replicated Blocks:
Total size: 48893 B
Total files: 1
Total blocks (validated): 1 (avg. block size 48893 B)
Minimally replicated blocks: 1 (100.0 %)
Over-replicated blocks: 0 (0.0 %)
Under-replicated blocks: 0 (0.0 %)
Mis-replicated blocks: 0 (0.0 %)
Default replication factor: 3
Average block replication: 3.0
Missing blocks: 0
Corrupt blocks: 0
Missing replicas: 0 (0.0 %)

Erasure Coded Block Groups:
Total size: 0 B
Total files: 0
Total block groups (validated): 0
Minimally erasure-coded block groups: 0
Over-erasure-coded block groups: 0
Under-erasure-coded block groups: 0
Unsatisfactory placement block groups: 0
Average block group size: 0.0
Missing block groups: 0
Corrupt block groups: 0
Missing internal blocks: 0
```

THEN LARGE

```
Number of data-nodes: 3
Number of racks: 1
Total dirs: 0
Total symlinks: 0

Replicated Blocks:
Total size: 208894 B
Total files: 1
Total blocks (validated): 1 (avg. block size 208894 B)
Minimally replicated blocks: 1 (100.0 %)
Over-replicated blocks: 0 (0.0 %)
Under-replicated blocks: 0 (0.0 %)
Mis-replicated blocks: 0 (0.0 %)
Default replication factor: 3
Average block replication: 3.0
Missing blocks: 0
Corrupt blocks: 0
Missing replicas: 0 (0.0 %)

Erasure Coded Block Groups:
Total size: 0 B
Total files: 0
Total block groups (validated): 0
Minimally erasure-coded block groups: 0
Over-erasure-coded block groups: 0
Under-erasure-coded block groups: 0
Unsatisfactory placement block groups: 0
Average block group size: 0.0
Missing block groups: 0
Corrupt block groups: 0
Missing internal blocks: 0
FSCK ended at Sat Nov 01 22:51:21 UTC 2025 in 2 milliseconds
```

OBSERVATION: EACH FILE IS 1 BLOCK BECAUSE ALL FILES ARE <128MB

Exercise 5: Inter-Cluster Data Replication with DistCp

Task: Copy your dataset from Cluster 1 to Cluster 2 using DistCp.

What to do:

1. Use DistCp to copy /user/yourname/ from Cluster 1 to Cluster 2
2. Verify the data exists in Cluster 2 by listing the directory
3. Compare file sizes and checksums between clusters
4. Read a file from Cluster 2 to confirm data integrity

Task: Screenshot of the DistCp command execution and its output

Task: Screenshot of file listings from both clusters side-by-side proving successful replication

(**Hint:** DistCp syntax: hadoop distcp hdfs://namenode-c1:8020/source hdfs://namenode-c2:8020/destination)

Copied

```
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ docker exec -e HADOOP_USER_NAME=madeel -it namenode-c1 bash -lc
for C in /opt/hadoop* /usr/local/hadoop; do
    if [ -x "$C/bin/hadoop" ]; then export HADOOP_HOME="$C"; break; fi
done
export PATH="$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$PATH"

hadoop distcp -update -pugp \
  hdfs://namenode-c1:8020/user/madeel \
  hdfs://namenode-c2:8020/user/madeel
```

```
2025-11-01 23:04:18,324 INFO mapred.LocalJobRunner:
2025-11-01 23:04:18,337 INFO mapred.Task: Task:attempt_local1568786143_0001_m_000000_0 is done. And is
in the process of committing
2025-11-01 23:04:18,339 INFO mapred.LocalJobRunner:
2025-11-01 23:04:18,339 INFO mapred.Task: Task attempt_local1568786143_0001_m_000000_0 is allowed to c
ommit now
2025-11-01 23:04:18,341 INFO output.FileOutputCommitter: Saved output of task 'attempt_local1568786143
_0001_m_000000_0' to file:/tmp/hadoop/mapred/staging/madeel974487544/.staging/_distcp-40894480/_logs
2025-11-01 23:04:18,342 INFO mapred.LocalJobRunner: Copying hdfs://namenode-c1:8020/user/madeel/sizetest
/st/medium.txt to hdfs://namenode-c2:8020/user/madeel/sizetest/medium.txt
2025-11-01 23:04:18,342 INFO mapred.Task: Task 'attempt_local1568786143_0001_m_000000_0' done.
2025-11-01 23:04:18,349 INFO mapred.Task: Final Counters for attempt local1568786143 0001 m 000000 0:
```

Now listing directory

```

madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ # Cluster 1
docker exec -e HADOOP_USER_NAME=madeel -it namenode-c1 bash -lc '
  export PATH=/opt/hadoop-3.2.1/bin:$PATH
  hdfs dfs -ls -h /user/madeel/sizetest
'

# Cluster 2
docker exec -e HADOOP_USER_NAME=madeel -it namenode-c2 bash -lc '
  export PATH=/opt/hadoop-3.2.1/bin:$PATH
  hdfs dfs -ls -h /user/madeel/sizetest
'

Found 3 items
-rw-r--r--  3 madeel supergroup   204.0 K 2025-11-01 22:45 /user/madeel/sizetest/large.txt
-rw-r--r--  3 madeel supergroup   47.7 K 2025-11-01 22:44 /user/madeel/sizetest/medium.txt
-rw-r--r--  3 madeel supergroup   792 2025-11-01 22:44 /user/madeel/sizetest/small.txt
Found 3 items
-rw-r--r--  3 madeel supergroup   204.0 K 2025-11-01 23:03 /user/madeel/sizetest/large.txt
-rw-r--r--  3 madeel supergroup   47.7 K 2025-11-01 23:03 /user/madeel/sizetest/medium.txt
-rw-r--r--  3 madeel supergroup   792 2025-11-01 23:03 /user/madeel/sizetest/small.txt

```

Now compare filesize plus checksum

```

madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ # Cluster 1
docker exec -e HADOOP_USER_NAME=madeel -it namenode-c1 bash -lc '
  export PATH=/opt/hadoop-3.2.1/bin:$PATH
  for f in small.txt medium.txt large.txt; do
    echo "C1 $f"; hdfs dfs -checksum /user/madeel/sizetest/$f; done
'

# Cluster 2
docker exec -e HADOOP_USER_NAME=madeel -it namenode-c2 bash -lc '
  export PATH=/opt/hadoop-3.2.1/bin:$PATH
  for f in small.txt medium.txt large.txt; do
    echo "C2 $f"; hdfs dfs -checksum /user/madeel/sizetest/$f; done
'

C1 small.txt
/user/madeel/sizetest/small.txt MD5-of-0MD5-of-512CRC32C      000002000000000000000000000000a7ef0e7caee0b4
ae3c7e6b1a0b1211b7
C1 medium.txt
/user/madeel/sizetest/medium.txt      MD5-of-0MD5-of-512CRC32C      000002000000000000000000000000934eaf
d82aadcca53c620902631fd839
C1 large.txt
/user/madeel/sizetest/large.txt MD5-of-0MD5-of-512CRC32C      00000200000000000000000000000075594de2149d8d
14f669f1d127fb3894
C2 small.txt
/user/madeel/sizetest/small.txt MD5-of-0MD5-of-512CRC32C      000002000000000000000000000000a7ef0e7caee0b4
ae3c7e6b1a0b1211b7
C2 medium.txt
/user/madeel/sizetest/medium.txt      MD5-of-0MD5-of-512CRC32C      000002000000000000000000000000934eaf
d82aadcca53c620902631fd839
C2 large.txt
/user/madeel/sizetest/large.txt MD5-of-0MD5-of-512CRC32C      00000200000000000000000000000075594de2149d8d
14f669f1d127fb3894
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ █

```

Then listed contents

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -e HADOOP_USER_NAME=madeel -it namenod
e-c2 bash -lc '
  export PATH=/opt/hadoop-3.2.1/bin:$PATH
  hdfs dfs -head /user/madeel/sizetest/small.txt
'

2025-11-01 23:11:52,836 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrus
t = false, remoteHostTrusted = false
Line 1
Line 2
Line 3
Line 4
Line 5
Line 6
Line 7
Line 8
Line 9
Line 10
Line 11
Line 12
Line 13
Line 14
Line 15
Line 16
Line 17
Line 18
Line 19
Line 20
Line 21
Line 22
... 22
```

Exercise 6: Create Directory Structure and Organize Files

Task: Create a proper directory hierarchy and organize files in Cluster 2.

What to do:

1. In Cluster 2, create a directory structure: /production/data/year2024/month09/
2. Move your replicated files into this organized structure
3. Create a separate /backup/ directory
4. Copy (not move) one file to the backup location

Task: Screenshot of your complete directory tree structure using recursive listing

```
madeel@ddacourse:~/mannoor-hadoop-multicloud-lab$ docker exec -it namenode-c2 bash
root@namenode-c2:/# hdfs dfs -mkdir -p /production/data/year2024/month09/
root@namenode-c2:/#
root@namenode-c2:/# hdfs dfs -ls -R /production/
drwxr-xr-x  3 root supergroup          0 2025-11-02 08:15 /production/data
drwxr-xr-x  3 root supergroup          0 2025-11-02 08:15 /production/data/year2024
drwxr-xr-x  3 root supergroup          0 2025-11-02 08:15 /production/data/year2024/month09
root@namenode-c2:/# hdfs dfs -mv /user/nadeel/* /production/data/year2024/month09/
root@namenode-c2:/# hdfs dfs -ls -R /production/data/year2024/month09/
drwxr-xr-x  3 madeel supergroup        0 2025-11-01 23:03 /production/data/year2024/month09/dataset
-rw-r--r--  3 madeel supergroup        618 2025-11-01 23:03 /production/data/year2024/month09/dataset/mydataset.txt
drwxr-xr-x  3 madeel supergroup        0 2025-11-01 23:03 /production/data/year2024/month09/sizetest
-rw-r--r--  3 madeel supergroup        208894 2025-11-01 23:03 /production/data/year2024/month09/sizetest/large.txt
-rw-r--r--  3 madeel supergroup        48893 2025-11-01 23:03 /production/data/year2024/month09/sizetest/medium.txt
-rw-r--r--  3 madeel supergroup        792 2025-11-01 23:03 /production/data/year2024/month09/sizetest/small.txt
root@namenode-c2:/# hdfs dfs -mkdir -p /backup/
root@namenode-c2:/# hdfs dfs -cp /production/data/year2024/month09/sizetest/small.txt /backup/
2025-11-02 08:15:56,571 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2025-11-02 08:15:56,909 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
root@namenode-c2:/# hdfs dfs -ls /backup/
Found 1 items
-rw-r--r--  3 root supergroup        792 2025-11-02 08:15 /backup/small.txt
root@namenode-c2:/#
root@namenode-c2:/# hdfs dfs -ls -R /
drwxr-xr-x  3 root supergroup        0 2025-11-02 08:15 /backup
-rw-r--r--  3 root supergroup        792 2025-11-02 08:15 /backup/small.txt
drwxr-xr-x  3 root supergroup        0 2025-11-02 08:15 /production
drwxr-xr-x  3 root supergroup        0 2025-11-02 08:15 /production/data
drwxr-xr-x  3 root supergroup        0 2025-11-02 08:15 /production/data/year2024
drwxr-xr-x  3 root supergroup        0 2025-11-02 08:15 /production/data/year2024/month09
drwxr-xr-x  3 madeel supergroup      0 2025-11-01 23:03 /production/data/year2024/month09/dataset
-rw-r--r--  3 madeel supergroup      618 2025-11-01 23:03 /production/data/year2024/month09/dataset/mydataset.txt
drwxr-xr-x  3 madeel supergroup      0 2025-11-01 23:03 /production/data/year2024/month09/sizetest
-rw-r--r--  3 madeel supergroup      208894 2025-11-01 23:03 /production/data/year2024/month09/sizetest/large.txt
-rw-r--r--  3 madeel supergroup      48893 2025-11-01 23:03 /production/data/year2024/month09/sizetest/medium.txt
-rw-r--r--  3 madeel supergroup      792 2025-11-01 23:03 /production/data/year2024/month09/sizetest/small.txt
drwxr-xr-x  3 root supergroup        0 2025-11-01 23:03 /user
drwxr-xr-x  3 madeel supergroup      0 2025-11-02 08:15 /user/madeel
```

(Hint: Use `hdfs dfs -mkdir -p` for nested directories and `hdfs dfs -mv` and `hdfs dfs -cp`)

Exercise 7: Test NodeManager Failure and YARN Recovery

Task: Simulate a NodeManager failure and observe YARN's response.

What to do:

1. Check all NodeManagers are running in Cluster 1
2. Stop nodemanager1-c1 container
3. Check YARN node status - how many are active now?
4. Try submitting a simple test job (calculate π using the example jar)
5. Observe: Does the job run with only one NodeManager?
6. Restart the NodeManager

Task: Screenshot of YARN node list before and after NodeManager failure

```
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ docker exec -it resourcemanager-c1 yarn node -list
2025-11-02 08:20:16,197 INFO client.RMProxy: Connecting to ResourceManager at resourcemanager-c1/172.19.0.4:8032
Total Nodes:2
  Node-Id          Node-State Node-Http-Address      Number-of-Running-Containers
nodemanager1-c1:34671      RUNNING  nodemanager1-c1:8042           0
nodemanager2-c1:33475      RUNNING  nodemanager2-c1:8042           0
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ █
```

Task: Screenshot of a successfully completed job running with reduced resources

```
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ sleep 600 && docker exec -it resourcemanager-c1 yarn node -list
2025-11-02 08:43:30,885 INFO client.RMProxy: Connecting to ResourceManager at resourcemanager-c1/172.19.0.4:8032
Total Nodes:1
  Node-Id          Node-State Node-Http-Address      Number-of-Running-Containers
nodemanager2-c1:33475      RUNNING  nodemanager2-c1:8042           0
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ █
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ docker exec -it resourcemanager-c1 yarn jar /opt/hadoop-3.2.1/share/hadoop/mapreduce/hadoop-mapreduce-examples-3
Number of Maps = 2
Samples per Map = 100
2025-11-02 08:45:11,470 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
Wrote input for Map #0
2025-11-02 08:45:11,658 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
Wrote input for Map #1
Starting Job
2025-11-02 08:45:11,760 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-11-02 08:45:11,800 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab$ █
```

Bytes written: 1

Job Finished in 1.434 seconds

Estimated value of Pi is 3.120000000000000000000000000000000

madeel@bdacourse:~/mahnoor-hadoop-multicuster-lab\$ █

When `nodemanager1-c1` was stopped, YARN automatically marked the node as lost but continued to run jobs using the `nodemanager2-c1`. The π calculation job completed successfully despite reduced resources, demonstrating YARN's fault tolerance.

Restarted nodemanager1-c1:

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker start nodemanager1-c1
nodemanager1-c1
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ sleep 600 && docker exec -it resourcemanager-c1 yarn node -list
2025-11-02 08:58:17,956 INFO client.RMProxy: Connecting to ResourceManager at resourcemanager-c1/172.19.0.4:8032
Total Nodes:2
  Node-Id          Node-State Node-Http-Address      Number-of-Running-Containers
nodemanager2-c1:33475      RUNNING  nodemanager2-c1:8042          0
nodemanager1-c1:33003      RUNNING  nodemanager1-c1:8042          0
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$
```

(Hint: Use `yarn node -list` and the example: `yarn jar /opt/hadoop-3.2.1/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.2.1.jar pi 2 100`)

Exercise 8: Change File Replication Factor

Task: Modify the replication factor of specific files and observe the effect.

What to do:

1. Check current replication of one of your files (should be 3)
 2. Change replication factor to 2
 3. Wait and verify the change took effect
 4. Change it to 1, then back to 3
 5. Observe in UI: Watch blocks being added/removed from DataNodes

Task: Screenshot of setrep command and before/after status of the file

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -stat %r /user/madeel/dataset/mydataset.txt  
3  
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -setrep -w 2 /user/madeel/dataset/mydataset.txt  
Replication 2 set: /user/madeel/dataset/mydataset.txt  
Waiting for /user/madeel/dataset/mydataset.txt ...  
WARNING: the waiting time may be long for DECREASING the number of replications.  
. done  
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -stat %r /user/madeel/dataset/mydataset.txt  
2  
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs dfs -setrep -w 1 /user/madeel/dataset/mydataset.txt
Replication 1 set: /user/madeel/dataset/mydataset.txt
Waiting for /user/madeel/dataset/mydataset.txt ...
WARNING: The waiting time may be long for DECREASING the number of replications.
. done
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs dfs -stat %r /user/madeel/dataset/mydataset.txt
1
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs dfs -setrep -w 3 /user/madeel/dataset/mydataset.txt
Replication 3 set: /user/madeel/dataset/mydataset.txt
Waiting for /user/madeel/dataset/mydataset.txt... done
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs dfs -stat %r /user/madeel/dataset/mydataset.txt
3
```

Task: Screenshot of NameNode UI displaying the file's current replication status

Browse Directory

File List						
Actions		Permission	Owner	Group	Size	Last Modified
<input type="checkbox"/>		-rw-r--r--	madeel	supergroup	618 B	Oct 22 22:45
<input type="checkbox"/>						

Browse Directory

/user/madeel/dataset								Go!									
Show	25	▼	entries							Search:	<input type="text"/>						
<input type="checkbox"/>		Permission		Owner		Group		Size		Last Modified		Replication		Block Size		Name	
		-rw-r--r--		madeel		supergroup		618 B		Oct 22 22:45		3		128 MB		mydataset.txt	

Showing 1 to 1 of 1 entries

Previous 1 Next

Hint: Use `hdfs dfs -setrep -w 2 /path/to/file` (the `-w` flag waits for replication to complete)

Exercise 9: Create and Test HDFS Snapshots

Task: Create snapshots to preserve data state and test restoration.

What to do:

1. Enable snapshots on your data directory
2. Create a snapshot named "backup_v1"
3. Modify or delete one of your files
4. List available snapshots
5. Restore the deleted/modified file from the snapshot
6. Verify the restoration was successful

Task: Screenshot of creation and the snapshot directory listing

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -createSnapshot /user/madeel/dataset backup_v1
Created snapshot /user/madeel/dataset/.snapshot/backup_v1
```

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -ls /user/madeel/dataset/.snapshot
Found 1 items
drwxr-xr-x  - madeel supergroup          0 2025-11-02 09:09 /user/madeel/dataset/.snapshot/backup_v1
```

File Deleted:

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -rm /user/madeel/dataset/mydataset.txt
Deleted /user/madeel/dataset/mydataset.txt
```

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -ls /user/madeel/dataset
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$
```

Task: Screenshot of successful file restoration from snapshot

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -cp /user/madeel/dataset/.snapshot/backup_v1/mydataset.txt /user/madeel/dataset
2025-11-02 09:13:36,510 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2025-11-02 09:13:36,605 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -ls /user/madeel/dataset
Found 1 items
-rw-r--r--  3 root supergroup      618 2025-11-02 09:13 /user/madeel/dataset/mydataset.txt
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -cat /user/madeel/dataset/mydataset.txt | head
2025-11-02 09:14:00,357 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
1, Alice, 25, Data Scientist
2, Bob, 30, Machine Learning Engineer
3, Charlie, 28, Data Analyst
4, David, 35, Software Developer
5, Emma, 22, Research Intern
6, Frank, 40, System Architect
7, Grace, 26, Database Administrator
8, Henry, 33, Hadoop Engineer
9, Ivy, 27, Business Analyst
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$
```

Hint: Use `hdfs dfsadmin -allowSnapshot`, `hdfs dfs -createSnapshot`, and `restore` by copying from `.snapshot` directory

Exercise 10: Run HDFS Balancer

Task: Check data distribution and run the balancer to even it out.

What to do:

1. Check current storage distribution across all DataNodes
2. Note which DataNode has the most/least data
3. Run the HDFS balancer with 10% threshold
4. Monitor the balancer progress
5. After completion, check distribution again

Task: Screenshot of DataNode storage distribution before balancing

(Already balanced)

The data distribution across DataNodes was already quite balanced, with storage usage around 312 KB, 312 KB, and 320 KB respectively.

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs dfsadmin -report | grep "Name\|DFS Used"
DFS Used: 967922 (945.24 KB)
DFS Used%: 0.00%
Name: 172.20.0.3:9866 (datanode3-c1.mahnoor-hadoop-multicluster-lab_cluster1-net)
DFS Used: 320121 (312.62 KB)
Non DFS Used: 37625699719 (35.04 GB)
DFS Used%: 0.00%
Name: 172.20.0.4:9866 (datanode1-c1.mahnoor-hadoop-multicluster-lab_cluster1-net)
DFS Used: 320121 (312.62 KB)
Non DFS Used: 37625699719 (35.04 GB)
DFS Used%: 0.00%
Name: 172.20.0.5:9866 (datanode2-c1.mahnoor-hadoop-multicluster-lab_cluster1-net)
DFS Used: 327680 (320 KB)
Non DFS Used: 37625692160 (35.04 GB)
DFS Used%: 0.00%
```

Task: Screenshot of balancer execution output and final distribution

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs balancer -threshold 5
2025-11-02 09:30:29,942 INFO balancer.Balancer: Using a threshold of 5.0
2025-11-02 09:30:29,943 INFO balancer.Balancer: namenodes = [hdfs://namenode-c1:8020]
2025-11-02 09:30:29,945 INFO balancer.Balancer: parameters = Balancer.BalancerParameters [BalancingPolicy.Node, threshold = 5.0, max idle iteration = 5, #excluded nodes = 0, #source nodes = 0, #blockpools = 0, run during upgrade = false]
2025-11-02 09:30:29,945 INFO balancer.Balancer: included nodes = []
2025-11-02 09:30:29,945 INFO balancer.Balancer: excluded nodes = []
2025-11-02 09:30:29,946 INFO balancer.Balancer: source nodes = []
Time Stamp           Iterations Bytes Already Moved Bytes Left To Move Bytes Being Moved
2025-11-02 09:30:30,696 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2025-11-02 09:30:30,788 INFO balancer.Balancer: dfs.balancer.movedInWidth = 5400000 (default=5400000)
2025-11-02 09:30:30,788 INFO balancer.Balancer: dfs.balancer.moverThreads = 1000 (default=1000)
2025-11-02 09:30:30,788 INFO balancer.Balancer: dfs.balancer.dispatcherThreads = 200 (default=200)
2025-11-02 09:30:30,788 INFO balancer.Balancer: dfs.datanode.balance.max.concurrent.moves = 50 (default=50)
2025-11-02 09:30:30,788 INFO balancer.Balancer: dfs.balancer.getLocks.size = 2147483648 (default=2147483648)
2025-11-02 09:30:30,788 INFO balancer.Balancer: dfs.balancer.getBlock.min-block.size = 10485760 (default=10485760)
2025-11-02 09:30:30,795 INFO balancer.Balancer: dfs.balancer.max.size-to-move = 10737418240 (default=10737418240)
2025-11-02 09:30:30,795 INFO balancer.Balancer: dfs.blocksize = 134217728 (default=134217728)
2025-11-02 09:30:30,808 INFO net.NetworkTopology: Adding a new node: /default-rack/172.20.0.4:9866
2025-11-02 09:30:30,808 INFO net.NetworkTopology: Adding a new node: /default-rack/172.20.0.3:9866
2025-11-02 09:30:30,808 INFO net.NetworkTopology: Adding a new node: /default-rack/172.20.0.5:9866
2025-11-02 09:30:30,810 INFO balancer.Balancer: 0 over-utilized: []
2025-11-02 09:30:30,810 INFO balancer.Balancer: 0 underutilized: []
The cluster is balanced. Exiting...
Nov 2, 2025 9:30:30 AM          0          0 B          0 B          0 B
Nov 2, 2025 9:30:30 AM Balancing took 1.026 seconds
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs dfsadmin -report | grep "Name\|DFS Used"
DFS Used: 983084 (960.04 KB)
DFS Used%: 0.00%
Name: 172.20.0.3:9866 (datanode3-c1.mahnoor-hadoop-multicluster-lab_cluster1-net)
DFS Used: 327702 (320.02 KB)
Non DFS Used: 37636775914 (35.05 GB)
DFS Used%: 0.00%
Name: 172.20.0.4:9866 (datanode1-c1.mahnoor-hadoop-multicluster-lab_cluster1-net)
DFS Used: 327680 (320 KB)
Non DFS Used: 37636767744 (35.05 GB)
DFS Used%: 0.00%
Name: 172.20.0.5:9866 (datanode2-c1.mahnoor-hadoop-multicluster-lab_cluster1-net)
DFS Used: 327702 (320.02 KB)
Non DFS Used: 37636775914 (35.05 GB)
DFS Used%: 0.00%
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$
```

Hint: Use `hdfs dfsadmin -report | grep "DFS Used"` and `hdfs balancer -threshold 10`

Exercise 11: Monitor and Compare Both Clusters

Task: Generate a comprehensive status report for both clusters.

What to do:

1. For both Cluster 1 and Cluster 2, collect:
 - o Total number of files and directories
 - o Total storage used
 - o Number of live DataNodes
 - o Number of active NodeManagers
 - o Available vs. used YARN memory
2. Create a comparison document showing these metrics side-by-side

Task: Screenshot of command outputs collecting metrics from both clusters

Task: Screenshot of your completed comparison report (can be a text file, spreadsheet, or document)

Hint: Use `hdfs dfs -count -q /`, `hdfs dfsadmin -report`, and `yarn node -lis`

RUN EACH CLUSTER

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c1 hdfs dfs -count -q /
9223372036854775807 9223372036854775795          none           inf         7          5        259815  /
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ S
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c2 hdfs dfs -count -q /
9223372036854775807 9223372036854775792          none           inf        10          5        259989  /
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$
```

NOW SHOW OUTPUT RESULT

CLUSTER1

```
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ docker exec -it namenode-c1 hdfs dfsadmin -report
Configured Capacity: 146675011584 (136.60 GB)
Present Capacity: 26214297600 (24.41 GB)
DFS Remaining: 26213314560 (24.41 GB)
DFS Used: 983040 (960 KB)
DFS Used%: 0.00%
Replicated Blocks:
    Under replicated blocks: 0
    Blocks with corrupt replicas: 0
    Missing blocks: 0
    Missing blocks (with replication factor 1): 0
    Low redundancy blocks with highest priority to recover: 0
    Pending deletion blocks: 0
Erasure Coded Block Groups:
    Low redundancy block groups: 0
    Block groups with corrupt internal blocks: 0
    Missing block groups: 0
    Low redundancy blocks with highest priority to recover: 0
    Pending deletion blocks: 0
-----
Live datanodes (3):
Name: 172.20.0.3:9866 (datanode3-c1.mahnoor-hadoop-multicloud-lab_cluster1-net)
Hostname: datanode3-c1
Decommission Status : Normal
Configured Capacity: 48891670528 (45.53 GB)
DFS Used: 327680 (320 KB)
Non DFS Used: 37636780032 (35.05 GB)
DFS Remaining: 8737771520 (8.14 GB)
DFS Used%: 0.00%
DFS Remaining%: 17.87%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
```

```
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Sun Nov 02 09:51:45 UTC 2025
Last Block Report: Sun Nov 02 05:13:14 UTC 2025
Num of Blocks: 5

Name: 172.20.0.5:9866 (datanode2-c1.mahnoor-hadoop-multicloud-lab_cluster1-net)
Hostname: datanode2-c1
Decommission Status : Normal
Configured Capacity: 48891670528 (45.53 GB)
DFS Used: 327680 (320 KB)
Non DFS Used: 37636780032 (35.05 GB)
DFS Remaining: 8737771520 (8.14 GB)
DFS Used%: 0.00%
DFS Remaining%: 17.87%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Sun Nov 02 09:51:45 UTC 2025
Last Block Report: Sun Nov 02 09:49:09 UTC 2025
Num of Blocks: 5
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ █
```

CLUSTER 2

```
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ docker exec -it namenode-c2 hdfs dfsadmin -report
Configured Capacity: 146675011584 (136.60 GB)
Present Capacity: 26214248448 (24.41 GB)
DFS Remaining: 26213265408 (24.41 GB)
DFS Used: 983040 (960 KB)
DFS Used%: 0.00%
Replicated Blocks:
    Under replicated blocks: 0
    Blocks with corrupt replicas: 0
    Missing blocks: 0
    Missing blocks (with replication factor 1): 0
    Low redundancy blocks with highest priority to recover: 0
    Pending deletion blocks: 0
Erasure Coded Block Groups:
    Low redundancy block groups: 0
    Block groups with corrupt internal blocks: 0
    Missing block groups: 0
    Low redundancy blocks with highest priority to recover: 0
    Pending deletion blocks: 0
-----
Live datanodes (3):
Name: 172.21.0.4:9866 (datanode1-c2.mahnoor-hadoop-multicloud-lab_cluster2-net)
Hostname: datanode1-c2
Decommission Status : Normal
Configured Capacity: 48891670528 (45.53 GB)
DFS Used: 327680 (320 KB)
Non DFS Used: 37636796416 (35.05 GB)
DFS Remaining: 8737755136 (8.14 GB)
DFS Used%: 0.00%
DFS Remaining%: 17.87%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
```

```
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Sun Nov 02 09:53:24 UTC 2025
Last Block Report: Sun Nov 02 05:06:24 UTC 2025
Num of Blocks: 5

Name: 172.21.0.6:9866 (datanode2-c2.mahnoor-hadoop-multicloud-lab_cluster2-net)
Hostname: datanode2-c2
Decommission Status : Normal
Configured Capacity: 48891670528 (45.53 GB)
DFS Used: 327680 (320 KB)
Non DFS Used: 37636796416 (35.05 GB)
DFS Remaining: 8737755136 (8.14 GB)
DFS Used%: 0.00%
DFS Remaining%: 17.87%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 1
Last contact: Sun Nov 02 09:53:25 UTC 2025
Last Block Report: Sun Nov 02 08:12:43 UTC 2025
Num of Blocks: 5
```

```
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ █
```

FOR YARN MANAGERS

```

madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it resourcemanager-c1 bash -lc 'curl -s http://localhost:8088/ws/v1/cluster/metrics'
{"clusterMetrics":{"appsSubmitted":0,"appsCompleted":0,"appsPending":0,"appsRunning":0,"appsFailed":0,"appsKilled":0,"reservedMB":0,"availableMB":4096,"allocatedMB":0,"reservedVirtualCores":0,"availableVirtualCores":4,"allocatedVirtualCores":0,"containersAllocated":0,"containersReserved":0,"containersPending":0,"totalMB":4096,"totalVirtualCores":4,"totalNodes":3,"lostNodes":1,"unhealthyNodes":0,"decommissioningNodes":0,"decommissionedNodes":0,"rebootedNodes":0,"activeNodes":2,"shutdownNodes":0,"totalUsedResourcesAcrossPartition":[{"memory":0,"vCores":0,"resourceInformations":[{"resourceInformation":[{"maximumAllocation":9223372036854775807,"minimumAllocation":0,"name":"memory-mb","resourceType":"COUNTABLE","units":"Mi","value":0}, {"maximumAllocation":9223372036854775807,"minimumAllocation":0,"name":"vcores","resourceType":"COUNTABLE","units":"","value":0}]}]}, "totalClusterResourcesAcrossPartition":{"memory":4096,"vCores":4,"resourceInformations":[{"resourceInformation":[{"maximumAllocation":9223372036854775807,"minimumAllocation":0,"name":"memory-mb","resourceType":"COUNTABLE","units":"Mi","value":4096}, {"maximumAllocation":9223372036854775807,"minimumAllocation":0,"name":"vcores","resourceType":"COUNTABLE","units":"","value":4}]}]}]}madeel@bdacourse:~/mahnoor-docker exec -it resourcemanager-c2 bash -lc 'curl -s http://localhost:8088/ws/v1/cluster/metrics'
{"clusterMetrics":{"appsSubmitted":0,"appsCompleted":0,"appsPending":0,"appsRunning":0,"appsFailed":0,"appsKilled":0,"reservedMB":0,"availableMB":4096,"allocatedMB":0,"reservedVirtualCores":0,"availableVirtualCores":4,"allocatedVirtualCores":0,"containersAllocated":0,"containersReserved":0,"containersPending":0,"totalMB":4096,"totalVirtualCores":4,"totalNodes":2,"lostNodes":0,"unhealthyNodes":0,"decommissioningNodes":0,"decommissionedNodes":0,"rebootedNodes":0,"activeNodes":2,"shutdownNodes":0,"totalUsedResourcesAcrossPartition":[{"memory":0,"vCores":0,"resourceInformations":[{"resourceInformation":[{"maximumAllocation":9223372036854775807,"minimumAllocation":0,"name":"memory-mb","resourceType":"COUNTABLE","units":"Mi","value":0}, {"maximumAllocation":9223372036854775807,"minimumAllocation":0,"name":"vcores","resourceType":"COUNTABLE","units":"","value":0}]}]}, "totalClusterResourcesAcrossPartition":{"memory":4096,"vCores":4,"resourceInformations":[{"resourceInformation":[{"maximumAllocation":9223372036854775807,"minimumAllocation":0,"name":"memory-mb","resourceType":"COUNTABLE","units":"Mi","value":4096}, {"maximumAllocation":9223372036854775807,"minimumAllocation":0,"name":"vcores","resourceType":"COUNTABLE","units":"","value":4}]}]}]}madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$
```

CLUSTER1

```

madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it resourcemanager-c1 yarn node -list
2025-11-02 09:55:44,118 INFO client.RMProxy: Connecting to ResourceManager at resourcemanager-c1/172.19.0.4:8032
Total Nodes:2
      Node-Id          Node-State Node-Http-Address      Number-of-Running-Containers
nodeManager2-c1:33475        RUNNING nodemanager2-c1:8042                      0
nodeManager1-c1:33003        RUNNING nodemanager1-c1:8042                      0
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$
```

CLUSTER 2

```

madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it resourcemanager-c2 yarn node -list
2025-11-02 09:56:20,429 INFO client.RMProxy: Connecting to ResourceManager at resourcemanager-c2/172.21.0.3:8032
Total Nodes:2
      Node-Id          Node-State Node-Http-Address      Number-of-Running-Containers
nodeManager2-c2:41767        RUNNING nodemanager2-c2:8042                      0
nodeManager1-c2:39657        RUNNING nodemanager1-c2:8042                      0
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$
```

COMPARISION TABLE

Metric	Cluster 1 (C1)	Cluster 2 (C2)
HDFS Configured Capacity (Total)	136.60 GB	136.60 GB

Present Capacity (post-overheads)	24.41 GB (26,214,297,600 B)	24.41 GB (26,214,248,448 B)
DFS Remaining	24.41 GB (26,213,314,560 B)	24.41 GB (26,213,265,408 B)
DFS Used	960 KB (983,040 B)	960 KB (983,040 B)
DFS Used %	0.00%	0.00%
Live DataNodes	3	3
Active YARN NodeManagers	2 (RUNNING)	2 (RUNNING)
Running Containers (YARN)	0	0
YARN allocatedMB	0	0
YARN availableMB	4096	4096
YARN totalMB	4096	4096

Cluster	DataNode	Configured Capacity	DFS Used	Non-DFS Used	DFS Remaining	Num of Blocks	Last contact (UTC)
C1	datanode 3-c1	45.53 GB	320 KB	35.05 GB	8.14 GB	5	Sun Nov 02 09:51:43 2025
C1	datanode 1-c1	45.53 GB	320 KB	35.05 GB	8.14 GB	5	Sun Nov 02 09:51:45 2025
C1	datanode 2-c1	45.53 GB	320 KB	35.05 GB	8.14 GB	5	Sun Nov 02 09:51:45 2025
C2	datanode 1-c2	45.53 GB	320 KB	35.05 GB	8.14 GB	5	Sun Nov 02 09:53:25 2025

C2	datanode 3-c2	45.53 GB	320 KB	35.05 GB	8.14 GB	5	Sun Nov 02 09:53:24 2025
C2	datanode 2-c2	45.53 GB	320 KB	35.05 GB	8.14 GB	5	Sun Nov 02 09:53:25 2025

Exercise 12: Cross-Cluster Data Movement Challenge

Task: Create a workflow that synchronizes specific data between clusters.

What to do:

1. Create a new file in Cluster 1 under /user/yourname/updates/
2. Use DistCp to sync only new/changed files to Cluster 2
3. Verify only the new file was copied
4. Modify a file in Cluster 1
5. Re-run sync and verify the update propagated

Task: Screenshot the DistCp sync operation and verification of selective copying

Hint: Use DistCp with -update flag for incremental sync

GO INSIDE NAMENODE EACH CLUSTER IN SEPERATE TERMINALS



The screenshot shows two separate terminal windows. Both windows have a title bar 'madeel@bdacourse: ~/mahnoor-hadoop-multicuster-lab'. The top window's menu bar includes 'File', 'Edit', 'View', 'Search', 'Terminal', and 'Help'. The bottom window's menu bar includes 'File', 'Edit', 'View', 'Search', 'Terminal', and 'Help'. In both windows, the command 'docker exec -it namenode-c1 bash' is being typed into the terminal. The cursor is visible at the end of the command in the bottom window.

Create a folder

```
root@namenode-c1:/# hdfs dfs -mkdir -p /user/madeel/updates/
```

And create a file and upload to hdfs

```
root@namenode-c1:/# echo "This is a new update file for sync test (v1)." > new_update.txt
root@namenode-c1:/# hdfs dfs -put -f new_update.txt /user/madeel/updates/
2025-11-02 10:36:04,946 INFO sasl.SaslDataTransferClient: SASL encryption trust
check: localHostTrusted = false, remoteHostTrusted = false
root@namenode-c1:/#
```

And verify it is in hdfs

```
root@namenode-c1:/# hdfs dfs -ls -R $SRC_URI/user/madeel/updates/
-rw-r--r-- 3 root supergroup 46 2025-11-02 10:36 hdfs://namenode-c1:80
20/user/madeel/updates/new_update.txt
root@namenode-c1:/#
root@namenode-c1:/# hdfs dfs -stat "%n | %b bytes | %y" $SRC_URI/user/madeel/updates/new_update.txt
new_update.txt | 46 bytes | 2025-11-02 10:36:05
root@namenode-c1:/#
root@namenode-c1:/#
```

Use distcp

```

root@namenode-c1:/# hadoop distcp -pb \
>   $SSRC_URI/user/madeel/updates \
>   $SDST_URI/user/madeel/updates

2025-11-02 10:37:56,296 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, overwrite=false, append=false, useDiff=false, userDiff=false, fromSnapshot=null, toSnapshot=null, skipCRC=false, blocking=true, numListStatusThreads=0, maxMaps=20, mapBandwidth=0.0, copyStrategy='uniformSize', preserveStatus=[FILELOCKSIZE], atomicWorkPath=null, logPath=null, sourceFilelisting=null, sourcePaths=[hdfs://namenode-c1:8020/user/madeel/updates], targetPath=hdfs://namenode-c2:8020/user/madeel/updates, filterFile='null', blocksPerChunk=6, copyBufferSize=8192, verboseLog=false, directWrite=false], sourcePaths=[hdfs://namenode-c1:8020/user/madeel/updates], targetPathExists=false, preserveRawAttributes=false
2025-11-02 10:37:56,352 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-11-02 10:37:56,403 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-11-02 10:37:56,403 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-11-02 10:37:56,575 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 2; dirCnt = 1
2025-11-02 10:37:56,577 INFO tools.SimpleCopyListing: Build file listing completed.
2025-11-02 10:37:56,577 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
2025-11-02 10:37:56,592 INFO tools.DistCp: Number of paths in the copy list: 2
2025-11-02 10:37:56,603 INFO tools.DistCp: Number of paths in the copy list: 2
2025-11-02 10:37:56,603 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2025-11-02 10:37:56,707 INFO mapreduce.JobSubmitter: number of splits:1
2025-11-02 10:37:57,039 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local234076817_0001
2025-11-02 10:37:57,039 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-11-02 10:37:57,117 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-11-02 10:37:57,117 INFO tools.DistCp: DistCp job-id: job_local234076817_0001
2025-11-02 10:37:57,119 INFO mapreduce.Job: Running job: job_local234076817_0001
2025-11-02 10:37:57,124 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-11-02 10:37:57,144 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-11-02 10:37:57,144 INFO output.FileOutputCommitter: File Output Committer skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-11-02 10:37:57,145 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.tools.mapred.CopyCommitter
2025-11-02 10:37:57,205 INFO mapred.LocalJobRunner: Waiting for map tasks
2025-11-02 10:37:57,208 INFO mapred.LocalJobRunner: Starting task: attempt_local234076817_0001_m_000000_0
2025-11-02 10:37:57,243 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-11-02 10:37:57,243 INFO output.FileOutputCommitter: File Output Committer skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-11-02 10:37:57,257 INFO mapred.Task: Using ResourceCalculatorProcessTree : [ ]
2025-11-02 10:37:57,261 INFO mapred.MapTask: Processing split: file:/tmp/hadoop/mapred/staging/root1764787949/.staging/_distcp871506186/filelist.seq:0+402
2025-11-02 10:37:57,267 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-11-02 10:37:57,267 INFO output.FileOutputCommitter: File Output Committer skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-11-02 10:37:57,286 INFO mapred.CopyMapper: Copying hdfs://namenode-c1:8020/user/madeel/updates to hdfs://namenode-c2:8020/user/madeel/updates
base-11-02 10:37:57 2025 INFO mapred.CopyMapper: Cancelling hdfs://namenode-c1:8020/user/madeel/updates from update but to hdfs://namenode-c2:8020/user/madeel/updates but

2025-11-02 10:37:57,300 INFO mapred.CopyMapper: Copying hdfs://namenode-c1:8020/user/madeel/updates/new_update.txt to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt
2025-11-02 10:37:57,304 INFO mapred.RetrifiableFileCopyCommand: Copying hdfs://namenode-c1:8020/user/madeel/updates/new_update.txt to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt
2025-11-02 10:37:57,305 INFO mapred.RetrifiableFileCopyCommand: Creating temp file: hdfs://namenode-c2:8020/user/madeel/updates/_distcp.tmp.attempt_local234076817_0001_m_000000_0
2025-11-02 10:37:57,305 INFO mapred.RetrifiableFileCopyCommand: Writing to temporary target file path hdfs://namenode-c2:8020/user/madeel/updates/_distcp.tmp.attempt_local234076817_0001_m_000000_0
2025-11-02 10:37:57,374 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2025-11-02 10:37:57,467 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2025-11-02 10:37:57,565 INFO mapred.RetrifiableFileCopyCommand: Renaming temporary target file path hdfs://namenode-c2:8020/user/madeel/updates/_distcp.tmp.attempt_local234076817_0001_m_000000_0 to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt
2025-11-02 10:37:57,578 INFO mapred.RetrifiableFileCopyCommand: Completed writing hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt (46 bytes)
2025-11-02 10:37:57,588 INFO mapred.LocalJobRunner: 
2025-11-02 10:37:57,610 INFO mapred.Task: Task:attempt_local234076817_0001_m_000000_0 is done. And is in the process of committing
2025-11-02 10:37:57,612 INFO mapred.LocalJobRunner: 
2025-11-02 10:37:57,612 INFO mapred.Task: Task attempt_local234076817_0001_m_000000_0 is allowed to commit now
2025-11-02 10:37:57,614 INFO output.FileOutputCommitter: Saved output of task 'attempt_local234076817_0001_m_000000_0' to file:/tmp/hadoop/mapred/staging/root1764787949/.staging/_distcp871506186/_logs
2025-11-02 10:37:57,615 INFO mapred.LocalJobRunner: 100.0% Copying hdfs://namenode-c1:8020/user/madeel/updates/new_update.txt to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt [-6.09/46.08]
2025-11-02 10:37:57,615 INFO mapred.Task: Task 'attempt_local234076817_0001_m_000000_0' done.
2025-11-02 10:37:57,624 INFO mapred.Task: Final Counters for attempt_local234076817_0001_m_000000_0: Counters: 26
  File System Counters
    FILE: Number of bytes read=187932
    FILE: Number of bytes written=716982
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=46
    HDFS: Number of bytes written=46
    HDFS: Number of read operations=19
    HDFS: Number of large read operations=4
    HDFS: Number of write operations=4
    HDFS: Number of bytes read erasure-coded=0
  Map-Reduce Framework
    Map input records=2
    Map output records=0
    Input split bytes=150
    Spilled Records=0
    Failed Shuffles=0

```

```

racted priorities--v
Merged Map outputs=0
GC time elapsed (ms)=0
Total committed heap usage (bytes)=188743680
File Input Format Counters
    Bytes Read=434
File Output Format Counters
    Bytes Written=8
DistCp Counters
    Bandwidth in Btyes=46
    Bytes Copied=46
    Bytes Expected=46
    Files Copied=1
    DIR_COPY=1
2025-11-02 10:37:57,626 INFO mapred.LocalJobRunner: Finishing task: attempt_local234076817_0001_m_000000_0
2025-11-02 10:37:57,627 INFO mapred.LocalJobRunner: map task executor complete.
2025-11-02 10:37:57,644 INFO mapred.CopyCommitter: About to preserve attributes: B
2025-11-02 10:37:57,647 INFO mapred.CopyCommitter: Preserved status on 0 dir entries on target
2025-11-02 10:37:57,647 INFO mapred.CopyCommitter: Cleaning up temporary work folder: file:/tmp/hadoop/mapred/staging/root1764787949/.staging/_distcp871506186
2025-11-02 10:37:58,130 INFO mapreduce.Job: Job job_local234076817_0001 running in uber mode : false
2025-11-02 10:37:58,133 INFO mapreduce.Job: map 100% reduce 0%
2025-11-02 10:37:58,136 INFO mapreduce.Job: Job job_local234076817_0001 completed successfully
2025-11-02 10:37:58,141 INFO mapreduce.Job: Counters: 26
File System Counters
    FILE: Number of bytes read=187932
    FILE: Number of bytes written=716982
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=46
    HDFS: Number of bytes written=46
    HDFS: Number of read operations=19
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=4
    HDFS: Number of bytes read erasure-coded=0
Map-Reduce Framework
    Map input records=2
    Map output records=0

```

```

Merged Map outputs=0
GC time elapsed (ms)=0
Total committed heap usage (bytes)=188743680
File Input Format Counters
    Bytes Read=434
File Output Format Counters
    Bytes Written=8
DistCp Counters
    Bandwidth in Btyes=46
    Bytes Copied=46
    Bytes Expected=46
    Files Copied=1
    DIR_COPY=1
root@namenode-c1:/# 
root@namenode-c1:/# S█

```

NOW VERIFY ON CLUSTER2

```
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab
File Edit View Search Terminal Help
madeel@bdacourse:~/mahnoor-hadoop-multicluster-lab$ docker exec -it namenode-c2
bash
root@namenode-c2:/# hdfs dfs -ls -R $DST_URI/user/madeel/updates/
-rw-r--r-- 3 root supergroup 46 2025-11-02 10:37 /user/madeel/updates/
new_update.txt
root@namenode-c2:/#
root@namenode-c2:/# s
root@namenode-c2:/# hdfs dfs -stat "%n | %b bytes | %y" $DST_URI/user/madeel/updates/
new_update.txt
new_update.txt | 46 bytes | 2025-11-02 10:37:57
root@namenode-c2:/#
```

**NOW LETS ONLY SYNC ONE FILE
ON CLUSTER 1, LETS ADD A SECOND FILE**

```
root@namenode-c1:/# echo "Second file, should be the only new one copied." > second.txt
root@namenode-c1:/# hdfs dfs -put -f second.txt /user/madeel/updates/
2025-11-02 10:45:55,519 INFO sasl.SaslDataTransferClient: SASL encryption trust
check: localHostTrusted = false, remoteHostTrusted = false
root@namenode-c1:/#
root@namenode-c1:/# hdfs dfs -ls $SRC_URI/user/madeel/updates/
Found 2 items
-rw-r--r-- 3 root supergroup 46 2025-11-02 10:36 hdfs://namenode-c1:80
20/user/madeel/updates/new_update.txt
-rw-r--r-- 3 root supergroup 48 2025-11-02 10:45 hdfs://namenode-c1:80
20/user/madeel/updates/second.txt
root@namenode-c1:/#
root@namenode-c1:/#
```

AND ON CLUSTER 2 LETS ONLY SYNC THIS SECOND FILE

```

root@namenode-c1:/# hadoop distcp -update -pb \
>   $SRC_URI/user/madeel/updates \
>   $DST_URI/user/madeel/updates

2025-11-02 10:46:36,243 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=true, deleteMissing=false, ignoreFailures=false, overwrite=false, append=false, useDiff=false, userDiff=false, fromSnapshot=null, toSnapshot=null, skipCRC=false, blocking=true, numListStatusThreads=0, maxMaps=20, mapBandwidth=0.0, copyStrategy='uniformsize', preserveStatus=[BLOCKSIZE], atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[hdfs://namenode-c1:8020/user/madeel/updates], targetPath=hdfs://namenode-c2:8020/user/madeel/updates, filterFile='null', blocksPerChunk=0, copyBufferSize=8192, verboseLog=false, directWrite=false], sourcePaths=[hdfs://namenode-c1:8020/user/madeel/updates], targetPathExists=true, preserveRawXattr:false
2025-11-02 10:46:36,304 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-11-02 10:46:36,319 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-11-02 10:46:36,319 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-11-02 10:46:36,476 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 2; dirCnt = 0
2025-11-02 10:46:36,476 INFO tools.SimpleCopyListing: Build file listing completed.
2025-11-02 10:46:36,479 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
2025-11-02 10:46:36,479 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
2025-11-02 10:46:36,501 INFO tools.DistCp: Number of paths in the copy list: 2
2025-11-02 10:46:36,504 INFO tools.DistCp: Number of paths in the copy list: 2
2025-11-02 10:46:36,508 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2025-11-02 10:46:36,557 INFO mapreduce.JobSubmitter: number of splits:1
2025-11-02 10:46:36,630 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local379944095_0001
2025-11-02 10:46:36,630 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-11-02 10:46:36,709 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-11-02 10:46:36,709 INFO tools.DistCp: DistCp job-id: job_local379944095_0001
2025-11-02 10:46:36,710 INFO mapreduce.Job: Running job: job_local379944095_0001
2025-11-02 10:46:36,711 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-11-02 10:46:36,716 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-11-02 10:46:36,716 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-11-02 10:46:36,717 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.tools.mapred.CopyCommitter
2025-11-02 10:46:36,718 INFO mapreduce.Job: Waiting for map tasks
2025-11-02 10:46:36,741 INFO mapred.LocalJobRunner: Starting task: attempt_local379944095_0001_m_000000_0
2025-11-02 10:46:36,758 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-11-02 10:46:36,758 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-11-02 10:46:36,775 INFO mapred.Task: Using ResourceCalculatorProcessTree : [ ]
2025-11-02 10:46:36,779 INFO mapred.MapTask: Processing split: file:/tmp/hadoop/mapred/staging/root306408745/.staging/_distcp677165651/fileList.seq:0+424
2025-11-02 10:46:36,787 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-11-02 10:46:36,787 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-11-02 10:46:36,794 INFO mapred.CopyMapper: Copying hdfs://namenode-c1:8020/user/madeel/updates/second.txt to hdfs://namenode-c2:8020/user/madeel/updates/second.txt
2025-11-02 10:46:36,801 INFO mapred.RetriableFileCopyCommand: Writing to temporary target file path hdfs://namenode-c2:8020/user/madeel/updates/_distcp.tmp.attempt_local379944095_0001_m_000000_0
2025-11-02 10:46:36,863 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2025-11-02 10:46:36,939 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2025-11-02 10:46:37,037 INFO mapred.RetriableFileCopyCommand: Renaming temporary target file path hdfs://namenode-c2:8020/user/madeel/updates/_distcp.tmp.attempt_local379944095_0001_m_000000_0 to hdfs://namenode-c2:8020/user/madeel/updates/second.txt
2025-11-02 10:46:37,047 INFO mapred.RetriableFileCopyCommand: Completed writing hdfs://namenode-c2:8020/user/madeel/updates/second.txt (48 bytes)
2025-11-02 10:46:37,056 INFO mapred.CopyMapper: Copying hdfs://namenode-c1:8020/user/madeel/updates/new_update.txt to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt
2025-11-02 10:46:37,067 INFO mapred.CopyMapper: Skipping copy of hdfs://namenode-c1:8020/user/madeel/updates/new_update.txt to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt
2025-11-02 10:46:37,071 INFO mapred.LocalJobRunner:
2025-11-02 10:46:37,078 INFO mapred.Task: Task@attempt_local379944095_0001_m_000000_0 is done. And is in the process of committing
2025-11-02 10:46:37,079 INFO mapred.Task: Task@attempt_local379944095_0001_m_000000_0 is allowed to commit now
2025-11-02 10:46:37,081 INFO output.FileOutputCommitter: Saved output of task 'attempt_local379944095_0001_m_000000_0' to file:/tmp/hadoop/mapred/staging/root306408745/.staging/_distcp677165651/logs
2025-11-02 10:46:37,082 INFO mapred.LocalJobRunner: Copying hdfs://namenode-c1:8020/user/madeel/updates/new_update.txt to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt
2025-11-02 10:46:37,082 INFO mapred.Task: Task '@attempt_local379944095_0001_m_000000_0' done
2025-11-02 10:46:37,088 INFO mapred.Task: Final Counters for attempt_local379944095_0001_m_000000_0: Counters: 27
  File System Counters
    FILE: Number of bytes read=188063
    FILE: Number of bytes written=717106
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=48
    HDFS: Number of bytes written=48
    HDFS: Number of read operations=22
    HDFS: Number of large read operations=3
    HDFS: Number of write operations=3
    HDFS: Number of bytes read erasure-coded=0
  Map-Reduce Framework
    Map input records=2
    Map output records=1
    Input split bytes=149
    Spilled Records=0
    Failed Shuffles=0
    Merged Map outputs=0
    GC time elapsed (ms)=0

```

```

mapred map output
GC time elapsed (ms)=0
Total committed heap usage (bytes)=190840832
File Input Format Counters
Bytes Read=456
File Output Format Counters
Bytes Written=77
DistCp Counters
Bandwidth in Bbytes=48
Bytes Copied=48
Bytes Expected=48
Bytes Skipped=46
Files Copied=1
Files Skipped=1
2025-11-02 10:46:37,089 INFO mapred.LocalJobRunner: Finishing task: attempt_local379944095_0001_m_000000_0
2025-11-02 10:46:37,089 INFO mapred.LocalJobRunner: map task executor complete.
2025-11-02 10:46:37,100 INFO mapred.CopyCommitter: About to preserve attributes: B
2025-11-02 10:46:37,101 INFO mapred.CopyCommitter: Preserved status on 0 dir entries on target
2025-11-02 10:46:37,102 INFO mapred.CopyCommitter: Cleaning up temporary work folder: file:/tmp/hadoop/mapred/staging/root306408745/.staging/_distcp677165651
2025-11-02 10:46:37,714 INFO mapreduce.Job: Job job_local379944095_0001 running in uber mode : false
2025-11-02 10:46:37,715 INFO mapreduce.Job: map 100% reduce 0%
2025-11-02 10:46:37,717 INFO mapreduce.Job: Job job_local379944095_0001 completed successfully
2025-11-02 10:46:37,721 INFO mapreduce.Job: Counters: 27
File System Counters
FILE: Number of bytes read=188063
FILE: Number of bytes written=717106
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=48
HDFS: Number of bytes written=48
HDFS: Number of read operations=22
HDFS: Number of large read operations=0
HDFS: Number of write operations=3
HDFS: Number of bytes read erasure-coded=0
Map-Reduce Framework
Map input records=2
Map output records=1
Input split bytes=149
Spilled Records=0

```

```

Spilled records=0
Failed Shuffles=0
Merged Map outputs=0
GC time elapsed (ms)=0
Total committed heap usage (bytes)=190840832

```

File Input Format Counters

Bytes Read=456

File Output Format Counters

Bytes Written=77

DistCp Counters

Bandwidth in Bbytes=48

Bytes Copied=48

Bytes Expected=48

Bytes Skipped=46

Files Copied=1

Files Skipped=1

root@namenode-c1:/#

root@namenode-c1:/# █

NOW LETS VERIFY ON CLUSTER2

```
root@namenode-c2:/# hdfs dfs -ls $DST_URI/user/madeel/updates/
Found 2 items
-rw-r--r-- 3 root supergroup          46 2025-11-02 10:37 /user/madeel/updates/
new_update.txt
-rw-r--r-- 3 root supergroup          48 2025-11-02 10:46 /user/madeel/updates/
second.txt
root@namenode-c2:/#
root@namenode-c2:/#
```

RECEIVED. NOW LETS EDIT A FILE ON C1 AND RESYNC

```
root@namenode-c1:/# echo " + appended line v2" | hdfs dfs -appendToFile - /user/
madeel/updates/new_update.txt

2025-11-02 10:50:49,489 INFO sasl.SaslDataTransferClient: SASL encryption trust
check: localHostTrusted = false, remoteHostTrusted = false
root@namenode-c1:/#
root@namenode-c1:/# hdfs dfs -stat "%n | %b bytes | %y" $SRC_URI/user/madeel/upd
ates/new_update.txt

new_update.txt | 66 bytes | 2025-11-02 10:50:49
root@namenode-c1:/#
root@namenode-c1:/#
```

NOW LETS DO DISTCP

```

root@namenode-c1:~# hadoop distcp -update -pb \
>   $SRC_URI/user/madeel/updates \
>   $DST_URI/user/madeel/updates

2025-11-02 10:51:37,541 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=true, deleteMissing=false, ignoreFailures=false, overwrite=false, append=false, useDiff=false, useRdiff=false, fromSnapshot=null, toSnapshot=null, skipCRC=false, blocking=true, numListStatusThreads=0, maxMaps=20, mapBandwidth=0.0, copyStrategy='uniformsize', preserveStatus=[BLOCKSIZE], atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[hdfs://namenode-c1:8020/user/madeel/updates], targetPath=hdfs://namenode-c2:8020/user/madeel/updates, filterFile='null', blocksPerChunk=0, copyBufferSize=8192, verboseLog=false, directWrite=false}, sourcePaths=[hdfs://namenode-c1:8020/user/madeel/updates], targetPathExists=true, preserveRawXattr=false
2025-11-02 10:51:37,611 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-11-02 10:51:37,705 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-11-02 10:51:37,705 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-11-02 10:51:37,789 INFO tools.SimpleCopyListing: Paths (files+dirs) Cnt = 2; dirCnt = 0
2025-11-02 10:51:37,789 INFO tools.SimpleCopyListing: Build file listing completed.
2025-11-02 10:51:37,799 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
2025-11-02 10:51:37,799 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
2025-11-02 10:51:37,813 INFO tools.DistCp: Number of paths in the copy list: 2
2025-11-02 10:51:37,813 INFO tools.DistCp: Number of paths in the copy list: 2
2025-11-02 10:51:37,823 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2025-11-02 10:51:37,823 INFO mapred.JobSubmitter: number of splits:1
2025-11-02 10:51:37,971 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local710483887_0001
2025-11-02 10:51:37,971 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-11-02 10:51:38,055 INFO tools.DistCp: DistCp job-id: job_local710483887_0001
2025-11-02 10:51:38,055 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2025-11-02 10:51:39,072 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-11-02 10:51:39,072 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-11-02 10:51:39,072 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.tools.mapred.CopyCommitter
2025-11-02 10:51:39,101 INFO mapred.LocalJobRunner: Waiting for map tasks
2025-11-02 10:51:39,101 INFO mapred.LocalJobRunner: Starting task attempt_local710483887_0001_m_000000_0
2025-11-02 10:51:39,118 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-11-02 10:51:39,118 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-11-02 10:51:39,133 INFO mapred.Task: Using ResourceCalculatorProcessTree: []
2025-11-02 10:51:39,137 INFO mapred.MapTask: Processing split: file:///tmp/hadoop/mapred/staging/root1347886152/.staging/_distcp18117487/fileList.seq:0+424
2025-11-02 10:51:39,144 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-11-02 10:51:39,144 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2025-11-02 10:51:39,153 INFO mapred.CopyHapper: Copying hdfs://namenode-c1:8020/user/madeel/updates/new_update.txt to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt
2025-11-02 10:51:39,161 INFO mapred.RetrivableFileCopyCommand: Copying hdfs://namenode-c1:8020/user/madeel/updates/new_update.txt to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt

2025-11-02 10:51:38,161 INFO mapred.RetrivableFileCopyCommand: Copying hdfs://namenode-c1:8020/user/madeel/updates/new_update.txt to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt
2025-11-02 10:51:38,161 INFO mapred.RetrivableFileCopyCommand: Creating temp file: hdfs://namenode-c2:8020/user/madeel/updates/_distcp.tmp.attempt_local710483887_0001_m_000000_0
2025-11-02 10:51:38,161 INFO mapred.RetrivableFileCopyCommand: Writing to temporary target file path hdfs://namenode-c2:8020/user/madeel/updates/_distcp.tmp.attempt_local710483887_0001_m_000000_0
2025-11-02 10:51:38,239 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localhostTrusted = false, remoteHostTrusted = false
2025-11-02 10:51:38,305 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localhostTrusted = false, remoteHostTrusted = false
2025-11-02 10:51:38,384 INFO mapred.RetrivableFileCopyCommand: Renaming temporary target file path hdfs://namenode-c2:8020/user/madeel/updates/_distcp.tmp.attempt_local710483887_0001_m_000000_0 to hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt
2025-11-02 10:51:38,410 INFO mapred.RetrivableFileCopyCommand: Completed writing hdfs://namenode-c2:8020/user/madeel/updates/new_update.txt (66 bytes)
2025-11-02 10:51:38,410 INFO mapred.CopyHapper: Copying hdfs://namenode-c1:8020/user/madeel/updates/second.txt to hdfs://namenode-c2:8020/user/madeel/updates/second.txt
2025-11-02 10:51:38,442 INFO mapred.CopyHapper: Skipping copy of hdfs://namenode-c1:8020/user/madeel/updates/second.txt to hdfs://namenode-c2:8020/user/madeel/updates/second.txt
2025-11-02 10:51:38,442 INFO mapred.CopyHapper: LocalJobRunner:
2025-11-02 10:51:38,456 INFO mapred.Task: Task attempt_local710483887_0001_m_000000_0 is done. And is in the process of committing
2025-11-02 10:51:38,456 INFO mapred.Task: Task attempt_local710483887_0001_m_000000_0 is allowed to commit now
2025-11-02 10:51:38,459 INFO output.FileOutputCommitter: Saved output of task 'attempt_local710483887_0001_m_000000_0' to file:/tmp/hadoop/mapred/staging/root1347886152/.staging/_distcp18117487/logs
2025-11-02 10:51:38,468 INFO mapred.LocalJobRunner: Copying hdfs://namenode-c1:8020/user/madeel/updates/second.txt to hdfs://namenode-c2:8020/user/madeel/updates/second.txt
2025-11-02 10:51:38,468 INFO mapred.Task: Task attempt_local710483887_0001_m_000000_0 done.
2025-11-02 10:51:38,466 INFO mapred.Task: Final Counters for attempt_local710483887_0001_m_000000_0: Counters: 27
    File System Counters
      FILE: Number of bytes read=188063
      FILE: Number of bytes written=717104
      FILE: Number of read operations=0
      FILE: Number of large read operations=0
      FILE: Number of write operations=0
      HDFS: Number of bytes read=66
      HDFS: Number of bytes written=66
      HDFS: Number of read operations=22
      HDFS: Number of large read operations=0
      HDFS: Number of write operations=4
      HDFS: Number of bytes read erasure-coded=0
    Map-Reduce Framework
      Map input records=2
      Map output records=1
      Input split bytes=149
      Spilled Records=0
      Failed Shuffles=0
      Merged Map outputs=0
      GC time elapsed (ms)=0
      Total committed heap usage (bytes)=190840832
    File Input Format Counters
      Bytes Read=ACK

```

```

File Input Format Counters
    Bytes Read=456
File Output Format Counters
    Bytes Written=73
DistCp Counters
    Bandwidth in Bbytes=66
    Bytes Copied=66
    Bytes Expected=66
    Bytes Skipped=48
    Files Copied=1
    Files Skipped=1
2025-11-02 10:51:38,467 INFO mapred.LocalJobRunner: Finishing task: attempt_local710483887_0001_m_000000_0
2025-11-02 10:51:38,467 INFO mapred.LocalJobRunner: map task executor complete.
2025-11-02 10:51:38,481 INFO mapred.CopyCommitter: About to preserve attributes: B
2025-11-02 10:51:38,484 INFO mapred.CopyCommitter: Preserved status on 0 dir entries on target
2025-11-02 10:51:38,484 INFO mapred.CopyCommitter: Cleaning up temporary work folder: file:/tmp/hadoop/mapred/staging/root1347886152/.staging/_distcp18117487
2025-11-02 10:51:39,058 INFO mapreduce.Job: Job job_local710483887_0001 running in uber mode : false
2025-11-02 10:51:39,058 INFO mapreduce.Job: map 100% reduce 0%
2025-11-02 10:51:39,060 INFO mapreduce.Job: Job job_local710483887_0001 completed successfully
2025-11-02 10:51:39,065 INFO mapreduce.Job: Counters: 27
    File System Counters
        FILE: Number of bytes read=188663
        FILE: Number of bytes written=717104
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=66
        HDFS: Number of bytes written=66
        HDFS: Number of read operations=22
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=4
        HDFS: Number of bytes read erasure-coded=0
    Map-Reduce Framework
        Map input records=2
        Map output records=1
        Input split bytes=149
        Spilled Records=0
        Failed Shuffles=0

```

Spilled Records=0
Failed Shuffles=0
Merged Map outputs=0
GC time elapsed (ms)=0
Total committed heap usage (bytes)=190840832

File Input Format Counters
 Bytes Read=456
File Output Format Counters
 Bytes Written=73
DistCp Counters
 Bandwidth in Bbytes=66
 Bytes Copied=66
 Bytes Expected=66
 Bytes Skipped=48
 Files Copied=1
 Files Skipped=1

```

root@namenode-c1:#
root@namenode-c1:#

```

NOW LETS VERIFY

```
root@namenode-c2:/# hdfs dfs -stat "%n | %b bytes | %y" $DST_URI/user/madeel/updates/new_update.txt  
new_update.txt | 66 bytes | 2025-11-02 10:51:38  
root@namenode-c2:/#  
root@namenode-c2:/# hdfs dfs -cat $DST_URI/user/madeel/updates/new_update.txt |  
tail -n +1  
  
2025-11-02 10:55:03,525 INFO sasl.SaslDataTransferClient: SASL encryption trust  
check: localHostTrusted = false, remoteHostTrusted = false  
This is a new update file for sync test (v1).  
+ appended line v2  
root@namenode-c2:/#  
root@namenode-c2:/#
```

RECEIVED UPDATED

Highest Weightage.

Exercise 13: MapReduce on multi-cluster and single-cluster

Task: Run a MR task on the multi-cluster and the single-cluster environment. Compare the performance using any of the information above and other (eg timing consumed in each task). Use the dashboards' data to support your answer if you want. You can paste as many snapshots as you want.

So using the same data as TASK 4

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it namenode-c1 hdfs dfs -ls /user/madeel/sizetest/
Found 3 items
-rw-r--r--  3 madeel supergroup  208894 2025-11-01 22:45 /user/madeel/sizetest/large.txt
-rw-r--r--  3 madeel supergroup   48893 2025-11-01 22:44 /user/madeel/sizetest/medium.txt
-rw-r--r--  3 madeel supergroup     792 2025-11-01 22:44 /user/madeel/sizetest/small.txt
```

Single Cluster:

```
madeel@bdacourse:~/mahnoor-hadoop-multiclus...$ docker exec -it resourcemanager-c1 bash
root@resourcemanager-c1:/# time hadoop jar /opt/hadoop-3.2.1/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.2.1.jar wordcount \
> /user/madeel/sizetest/large.txt /user/madeel/wc-output-c1
```

Bytes Read=200000
File Output Format Counters
Bytes Written=148905

```
real    0m3.130s
user    0m5.042s
sys     0m0.409s
root@resourcemanager-c1:/#
root@resourcemanager-c1:/#
```

Multi-Cluster-Time = 3.13

Multi cluster: splitted the data of large.txt

```
root@namenode-c3:/# split -l 10000 /tmp/large.txt part_
root@namenode-c3:/# ls
KEYS  dev         .hadoop      lib      mnt      part_ab  run      srv  usr
bin   entrypoint.sh.hadoop-data lib64    opt      proc     run.sh  sys  var
boot  etc          home        media    part_aa root     sbin  tmp
root@namenode-c3:/#
```

Then uploaded the data to each cluster

```
root@namenode-c3:/# hdfs dfs -mkdir -p /user/madeel/sizetest/
root@namenode-c3:/# hdfs dfs -put part_aa /user/madeel/sizetest/
```

```
2025-11-02 17:10:34,846 INFO sasl.SaslDataTransferClient: SASL encryption trust
check: localHostTrusted = false, remoteHostTrusted = false
root@namenode-c3:/#
root@namenode-c3:/#
```

```
root@namenode-c4:/# hdfs dfs -put part_ab /user/madeel/sizetest/
2025-11-02 17:18:18,829 INFO sasl.SaslDataTransferClient: SASL encryption trust
check: localHostTrusted = false, remoteHostTrusted = false
root@namenode-c4:/#
```

Now we run wordcount on parallel on each namenode cluster

```
madel@bdacourse:~/mahnoor-hadoop-multiclusler-lab
File Edit View Search Terminal Help
Reduce input records=10002
Reduce output records=10002
Spilled Records=20004
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=9
Total committed heap usage (bytes)=685768704
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=178894
File Output Format Counters
Bytes Written=68919
real    0m3.047s
user    0m4.803s
sys     0m0.390s
root@namenode-c3:/#
```



```
madel@bdacourse:~/mahnoor-hadoop-multiclusler-lab
File Edit View Search Terminal Help
Reduce output records=10002
Spilled Records=20004
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=11
Total committed heap usage (bytes)=691535872
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=190000
File Output Format Counters
Bytes Written=80025
real    0m3.076s
user    0m4.767s
sys     0m0.411s
root@namenode-c4:/#
root@namenode-c4:/#
```

Now we will copy the cluster 4 output to cluster 3 using distcp

```
root@namenode-c3:/# time hadoop distcp hdfs://namenode-c4:8020/user/madeel/wc-ou
tput-c4 \
>           hdfs://namenode-c3:8020/user/madeel/wc-output-merged/
2025-11-02 17:26:27,782 INFO tools.DistCp: Input Options: DistCpOptions{atomicCo
mmit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, overwri
te=false, append=false, useDiff=false, useRdiff=false, fromSnapshot=null, toSnap
shot=null, skipCRC=false, blocking=true, numListstatusThreads=0, maxMaps=20, map
Bandwidth=0.0, copyStrategy='uniformsize', preserveStatus=[BLOCKSIZE], atomicWor
kPath=null, logPath=null, sourceFileListing=null, sourcePaths=[hdfs://namenode-c
4:8020/user/madeel/wc-output-c4], targetPath=hdfs://namenode-c3:8020/user/madeel
/wc-output-merged, filtersFile='null', blocksPerChunk=0, copyBufferSize=8192, ve
rboseLog=false, directWrite=false}, sourcePaths=[hdfs://namenode-c4:8020/user/m
adeel/wc-output-c4], targetPathExists=false, preserveRawXattrs=false
2025-11-02 17:26:27,838 INFO impl.MetricsConfig: Loaded properties from hadoop-m
----- -----
Bytes Expected=80025
Files Copied=2
DIR_COPY=1

real    0m3.050s
user    0m3.250s
sys     0m0.292s
root@namenode-c3:/#
```

Now we merge the two outputs together

```
root@namenode-c3:/# hdfs dfs -mkdir -p /user/madeel/wc-output-final
root@namenode-c3:/# time hdfs dfs -cat /user/madeel/wc-output-c3/part-r-* \
>           /user/madeel/wc-output-merged/part-r-* \
>           | hdfs dfs -put - /user/madeel/wc-output-final/combined.txt
2025-11-02 17:28:56,271 INFO sasl.SaslDataTransferClient: SASL encryption trust
check: localHostTrusted = false, remoteHostTrusted = false
2025-11-02 17:28:56,382 INFO sasl.SaslDataTransferClient: SASL encryption trust
check: localHostTrusted = false, remoteHostTrusted = false

real    0m2.676s
user    0m4.760s
sys     0m0.409s
root@namenode-c3:/#
```

Now we verify that one output file has been made

```
root@namenode-c3:/# hdfs dfs -ls /user/madeel/wc-output-final/
Found 1 items
-rw-r--r--  3 root supergroup      148944 2025-11-02 17:28 /user/madeel/wc-outpu
t-final/combined.txt
root@namenode-c3:/# hdfs dfs -cat /user/madeel/wc-output-final/combined.txt | he
ad
2025-11-02 17:30:20,869 INFO sasl.SaslDataTransferClient: SASL encryption trust
check: localHostTrusted = false, remoteHostTrusted = false
1      1
10     1
100    1
1000   1
10000  1
1001   1
1002   1
1003   1
1004   1
1005   1
cat: Unable to write to output stream.
root@namenode-c3:/#
```

Done.

The total multi-cluster time is the time taken by the slower of the two parallel jobs plus the time to copy Cluster 4's output to Cluster 3 (DistCp) and merge the results. This accounts for parallel execution and the overhead of combining outputs into a single final dataset.

```
Multi-Cluster-Time = max(T_cluster3_job, T_cluster4_job) + T_DistCp +
T_merge
Multi-Cluster-Time = max(3.047, 3.076) + 3.050 + 2.676
Multi-Cluster-Time = 8.802
```

For this experiment, the single-cluster job completed in about 3 seconds, while the multi-cluster setup took longer (~8 seconds) even though each cluster processed only half the data. This is because the dataset was small (around 20,000 lines), so the Map tasks themselves were quick, but the multi-cluster workflow incurred overhead: 3 seconds for running the Map tasks, 3 seconds to copy the output from one cluster to the other using DistCp, and 2 seconds for merging the results. This demonstrates that for small datasets, the coordination and merge overhead in a multi-cluster setup can outweigh the benefits of parallelism. For larger files, the Map task time per cluster would decrease because the workload is split across multiple clusters, and the merge overhead becomes proportionally smaller relative to the total processing time. This allows the system to scale efficiently, making multi-cluster setups advantageous for handling big datasets.

Cleanup

Stop Both Clusters

```
cd ~/hadoop-multicloud-lab
```

```
# Stop Cluster 1
docker compose -f docker-compose-cluster1.yml down
```

```
# Stop Cluster 2
docker compose -f docker-compose-cluster2.yml down
```

```
# Remove shared network
docker network rm shared-net
```

Show proof

```
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ docker compose -f docker-
compose-cluster1.yml down
WARN[0000] /home/madeel/mahnoor-hadoop-multicloud-lab/docker-compose-cluste
r1.yml: the attribute 'version' is obsolete, it will be ignored, please remov
e it to avoid potential confusion
[+] Running 8/8
✓ Container datanode1-c1                                Removed    11.6s
✓ Container nodemanager1-c1                            Removed    11.7s
✓ Container datanode3-c1                                Removed    11.8s
✓ Container nodemanager2-c1                            Removed    11.6s
✓ Container datanode2-c1                                Removed    11.7s
✓ Container resourcemanager-c1                         Removed    10.4s
✓ Container namenode-c1                                Removed    10.7s
✓ Network mahnoor-hadoop-multicloud-lab_cluster1-net   Removed     0.1s
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ docker compose -f docker-
compose-cluster2.yml down
WARN[0000] /home/madeel/mahnoor-hadoop-multicloud-lab/docker-compose-cluste
r2.yml: the attribute 'version' is obsolete, it will be ignored, please remov
e it to avoid potential confusion
[+] Running 8/8
✓ Container datanode3-c2                                Removed    11.4s
✓ Container nodemanager2-c2                            Removed    11.4s
✓ Container datanode1-c2                                Removed    11.2s
✓ Container nodemanager1-c2                            Removed    11.2s
✓ Container datanode2-c2                                Removed    11.4s
✓ Container resourcemanager-c2                         Removed    10.4s
✓ Container namenode-c2                                Removed    10.6s
✓ Network mahnoor-hadoop-multicloud-lab_cluster2-net   Removed     0.1s
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ docker network rm shared-
net
shared-net
madeel@bdacourse:~/mahnoor-hadoop-multicloud-lab$ |
```