# Hadoop Distributed File System(HDFS)

赵敏 老师

**扫描二维码关注微信/微博**
**获取最新IT面试情报及权威解答**

微信: ninechapter
知乎专栏: http://zhuanlan.zhihu.com/jiuzhang
微博: http://www.weibo.com/ninechapter
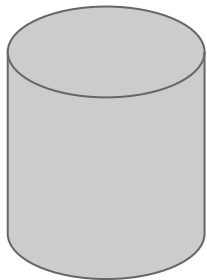官网: www.jiuzhang.com

# Outline

- What is HDFS & Why HDFS

- Demo of HDFS Commands

- Overview of HDFS Architecture

- HDFS System Design.

HDFS

- Hadoop <span style="color:red">Distributed</span> File System

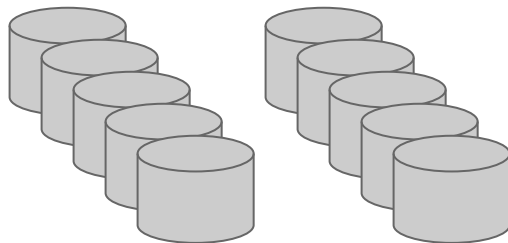- A <span style="color:red">fault-tolerant</span> file system designed to run on inexpensive hardware

Challenge: Read 1TB of data

1 machine
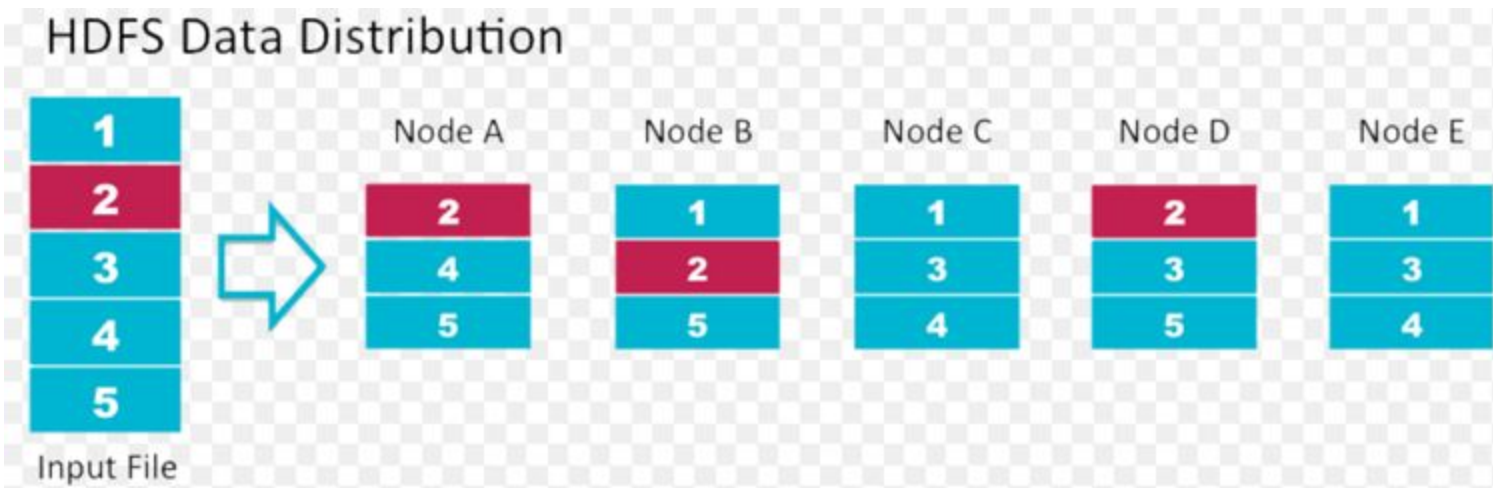
4 I/O channels
Each channel: 100MB/S

45 MIN

10 machines

4 I/O channels
Each channel: 100MB/S

4.5 MIN

Faster

- Split

- Replication

## HDFS Data Distribution

| Input File | Node A | Node B | Node C | Node D | Node E |
|---|---|---|---|---|---|
| 1 | 2 | 1 | 1 | 2 | 1 |
| 2 | 4 | 2 | 3 | 3 | 3 |
| 3 | 5 | 5 | 4 | 5 | 4 |
| 4 | | | | | |
| 5 | | | | | |

Easier

Demo

Why do we learn HDFS architecture?

- In interview: can be applied to all file system design

- Help you understand the whole hadoop ecosystem

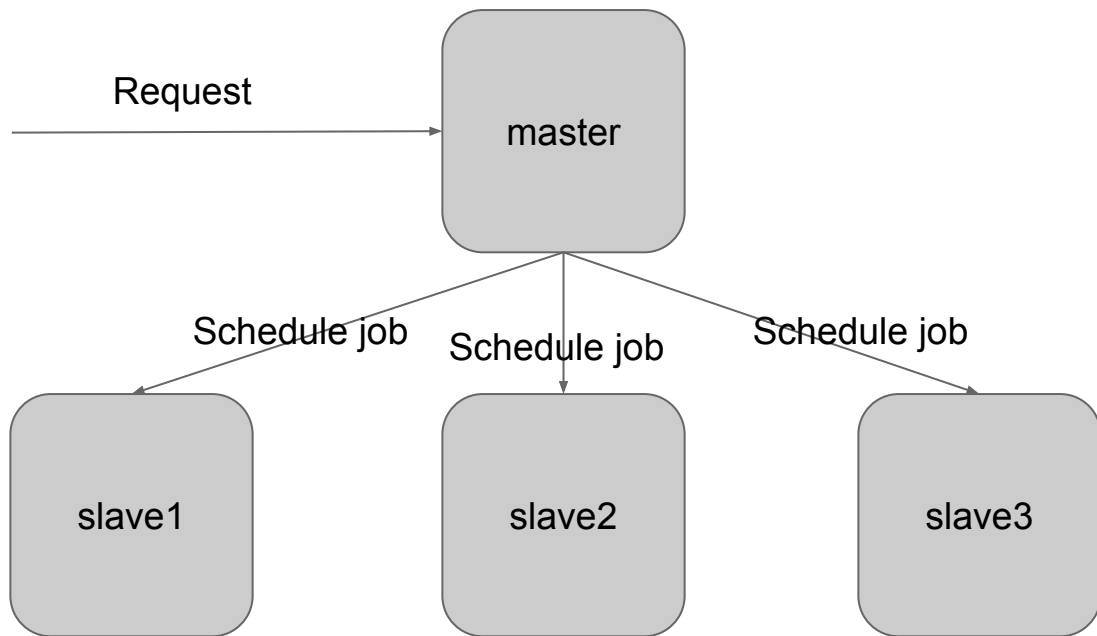Since we have multiple machines, how could they communicate?

● Do they talk to each other?

# Overview of HDFS Architecture

- Master slave model

  - High consistency

  - Simpler design

  - Single master node is not robust

- Peer peer model

  - Distributes read write load

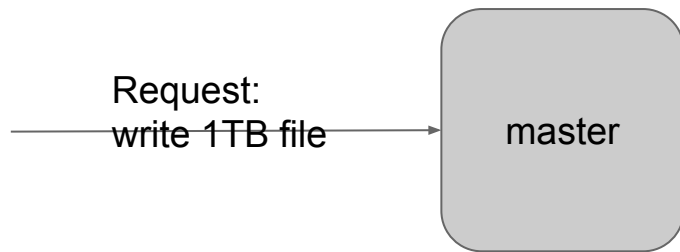  - One node down won't affect others

  - Low consistency

http://stackoverflow.com/questions/3736969/master-master-vs-master-slave-database-architecture

HDFS uses master-slave model.

Request:
write 1TB file

master

Will master transport data?

- No, or it will be a bottleneck.

- Master will decide which slave node to read/write, then client will talk to slave node.

slave1

How to store file?

- Whole big file

- Blocks of small files

- Blocks of small files

- Who will divide file into blocks?

  - Master node?

  - Slave node?

Neither master nor slave….

---> HDFS client

/user/mill/data1.txt      1, 2, 3
/user/mill/data2.txt      4, 5

User → HDFS Client → master

Schedule job → slave1 [ 1 ] [ 3 ]

Schedule job → slave2 [ 2 ] [ 5 ]

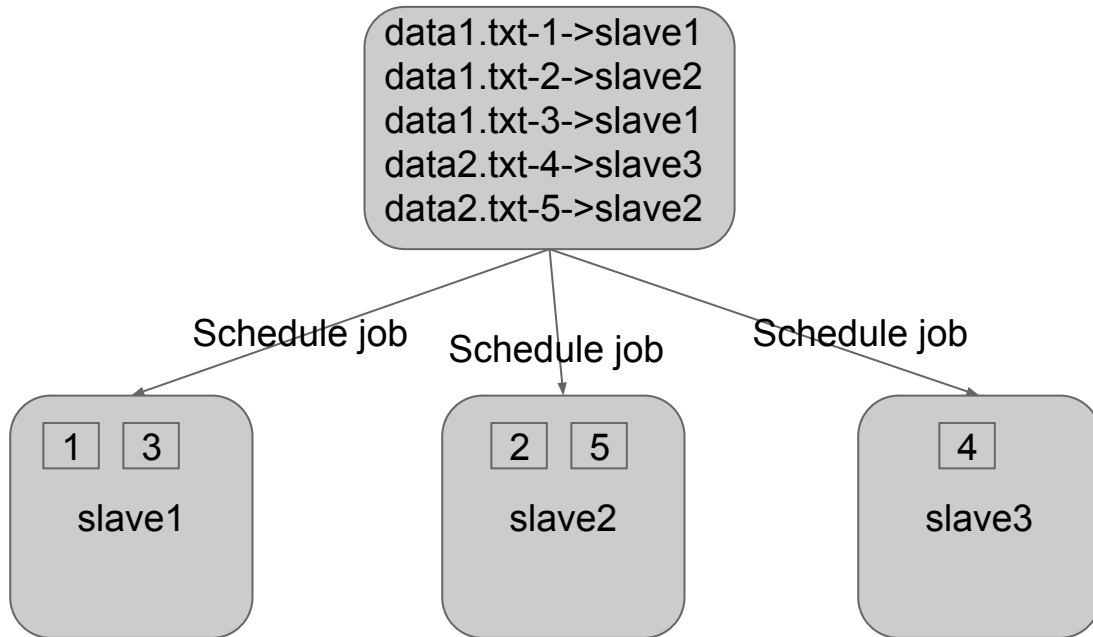Schedule job → slave3 [ 4 ]

master

What do we store on master machines?

Master should know where to write and read

- Metadata
  - File blocks
  - Where are the blocks

# Overview of HDFS Architecture

/user/mill/data1.txt          1, 2, 3
/user/mill/data2.txt          4, 5

data1.txt-1->slave1
data1.txt-2->slave2
data1.txt-3->slave1
data2.txt-4->slave3
data2.txt-5->slave2

Schedule job          Schedule job          Schedule job

| 1 | 3 |              | 2 | 5 |              | 4 |

slave1                slave2                slave3

What if one slave node fail?

# Overview of HDFS Architecture

Lose data!

data1.txt-1->slave1
data1.txt-2->slave2
data1.txt-3->slave1
data2.txt-4->slave3
data2.txt-5->slave2

Schedule job

Schedule job

Schedule job

| 1 | 3 |

slave1

| 2 | 5 |

slave2

| 4 |

slave3

How to avoid data loss?

Data replication

1 > slave1 & 2
2 > slave1 & 2
3 > slave1 & 3
4 > slave2 & 3
5 > slave2 & 3

Schedule job    Schedule job    Schedule job

| 1 | 3 |

slave1

| 2 |

| 2 | 5 |

slave2

| 1 | 4 |

| 4 | 3 |

slave3

| 5 |

What if one master node fail?

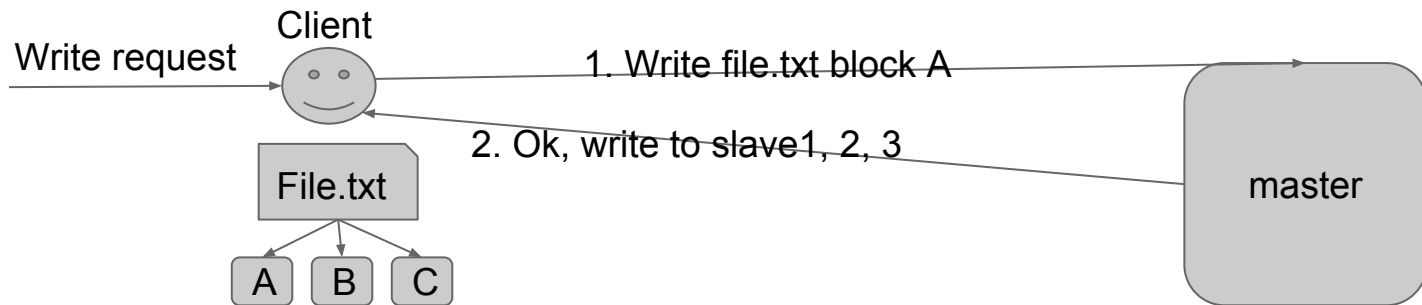We have a checkpoint node to copy all the data from master node per hour.

*Master-slave architecture*

- *Single NameNode -* a master server that manages the file system namespace and regulates access to files by clients.

- *Multiple DataNodes* – typically one per node in the cluster.
  - Manage storage
  - Serving read/write requests from clients
  - Block creation, deletion, replication based on instructions from NameNode

Read & Write

How to do the write operation?

# Overview of HDFS Architecture: Write

Write request → Client

Client:
File.txt
├── A
├── B
└── C

1. Write file.txt block A

2. Ok, write to slave1, 2, 3

master

Who should client write to?

Slave nodes

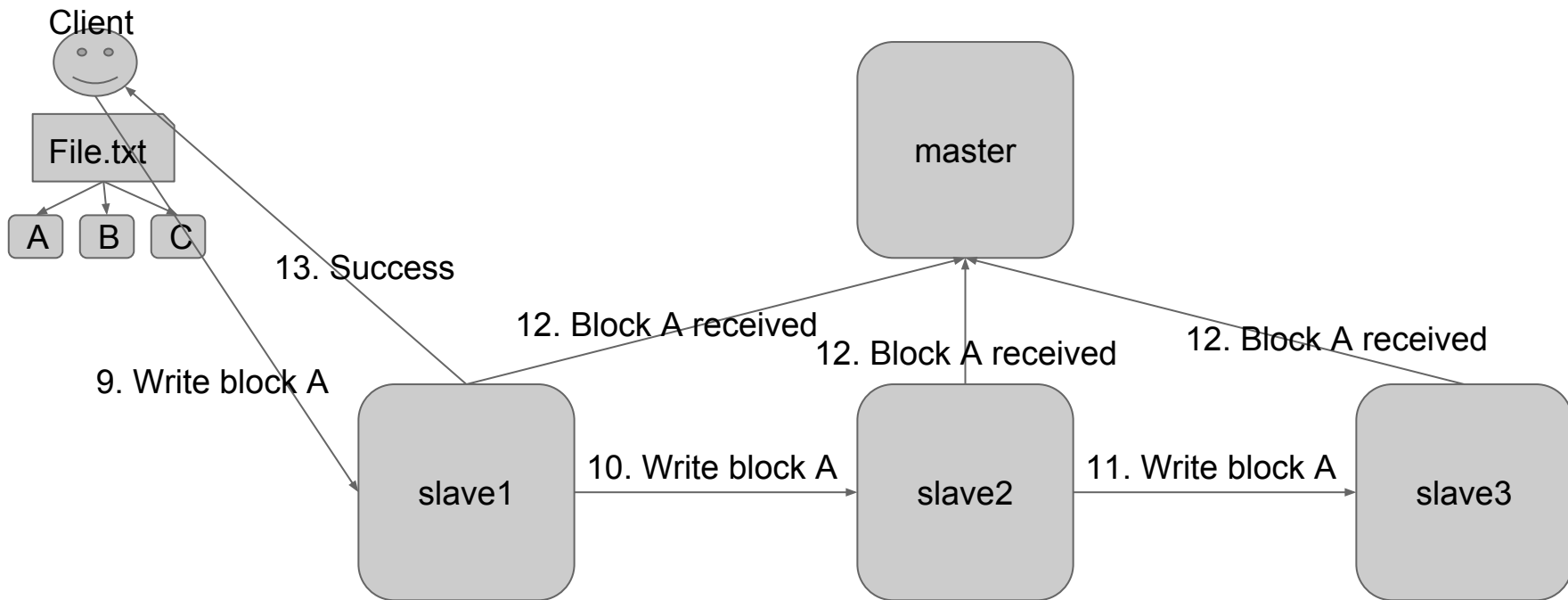Since we have three slave nodes to write, should client writes to three nodes?

No, will result in bottleneck.

Client

8. ACK!

3.Get yourself ready,
I will send you data,
Also get 2,3 ready.

slave1

4. Get 3 ready

slave2

5. Get you ready

slave3

7. ACK!

6. ACK!

# Overview of HDFS Architecture: Write



Client

File.txt

A  B  C

master

13. Success

12. Block A received

12. Block A received

12. Block A received

9. Write block A
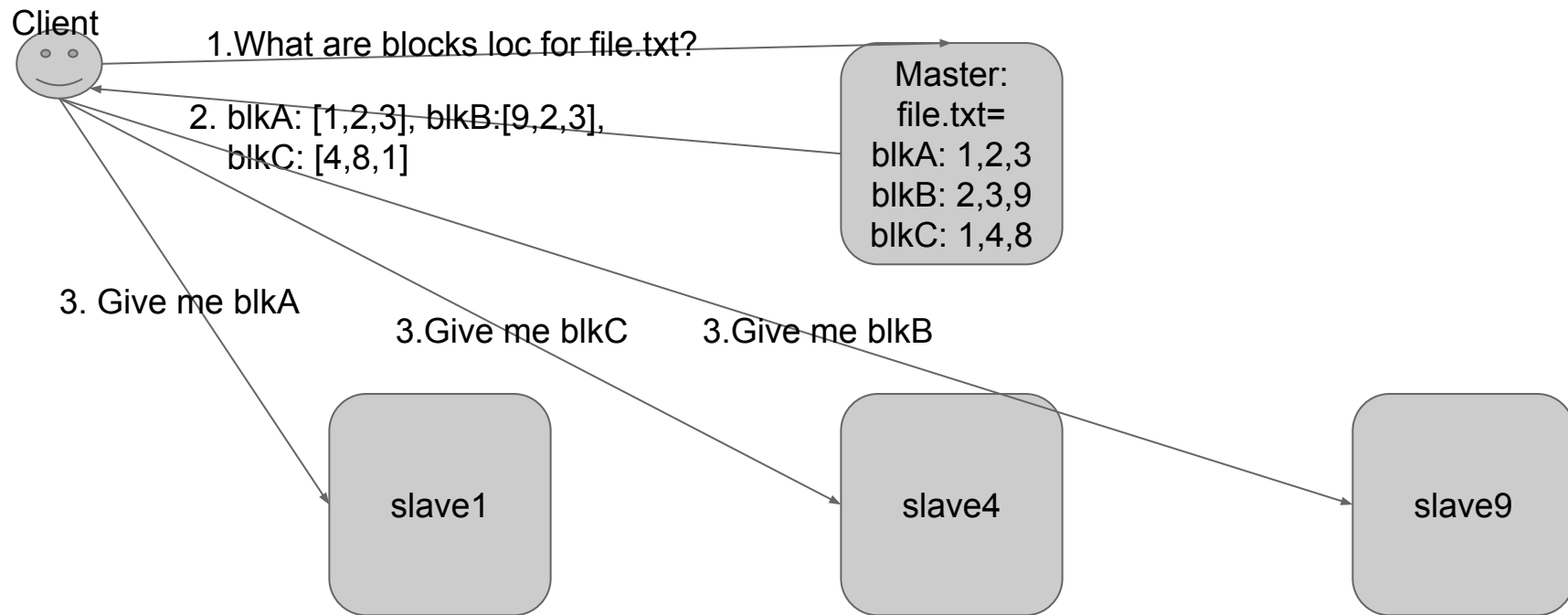
slave1

10. Write block A

slave2

11. Write block A

slave3

- Data is kept in different <span style="color:red">racks</span>. To ensure if one rack fails, we still have another rack to hold the data.
- Keep two blocks in same rack to achieve high throughput while reading data because two machines in same rack have more bandwidth and lower latency.
- Client does not send blocks to all 3 data nodes identified by Name node. The reason is Client will be choked by data transmission at a time.
- Name Node creates metadata from block reports received from data nodes

How to do the read operation?

Client

1.What are blocks loc for file.txt?

2. blkA: [1,2,3], blkB:[9,2,3],
blkC: [4,8,1]

Master:
file.txt=
blkA: 1,2,3
blkB: 2,3,9
blkC: 1,4,8

3. Give me blkA

3.Give me blkC

3.Give me blkB

slave1

slave4

slave9

blkB:[9,2,3]

Since we have 3 replica, which to read?

blkB:[9,2,3]

To visit the closest one.

blkB:[9,2,3]

Why does the master provide this slave nodes order?

- The slave nodes order are decided by the distance between slave node and the client.
- The closer, the faster.

# What you have learned

- Concept about HDFS

- HDFS Commands

- How to design a file system

- How to design the communication in file system