Agenda.
→ Design   (S3) → Aws.
→ file storage (for storing big files)
↳ Blob Storage
→ S3
→ HDFS.

```
┌─────────────┐
│   Client    │
└─────────────┘
       │
       ▼
┌─────────────┐
│  API GW + LB │
└─────────────┘
```

```
┌──────┐  ┌──────┐  ┌──────┐  ┌──────┐  ┌──────┐
│      │  │      │  │      │  │      │  │      │
└──────┘  └──────┘  └──────┘  └──────┘  └──────┘
```

Posts.

```
┌───────────────────────────────────────┐
│   ▭      ▭      ▭                       │
│  │ │    │ │    │ │                      │
│  │ │    │ │    │ │                      │
└───────────────────────────────────────┘
```

S3
file
Storage

```
┌───────────────────────────────┐
│   ▭      ▭      ▭              │
│  │ │    │ │    │ │             │
│  │ │    │ │    │ │             │
└───────────────────────────────┘
```
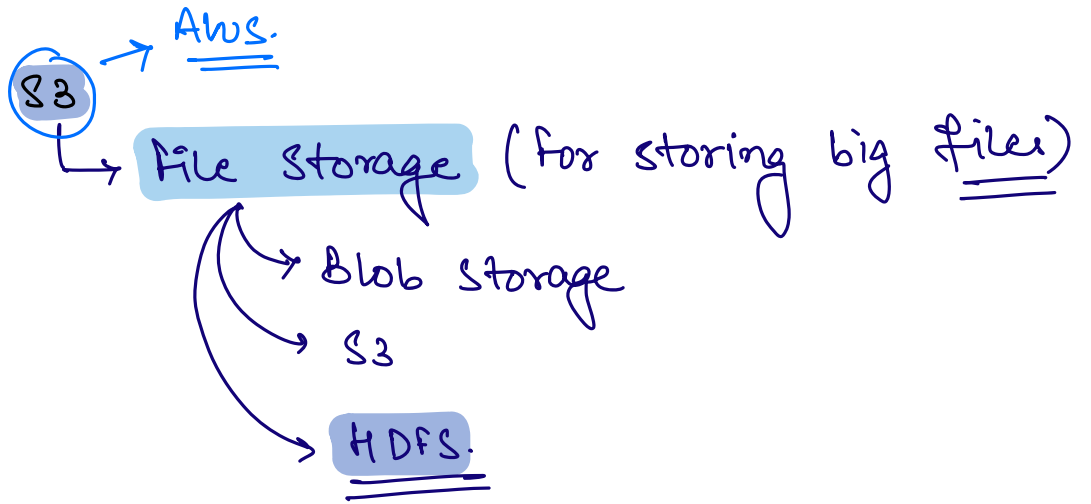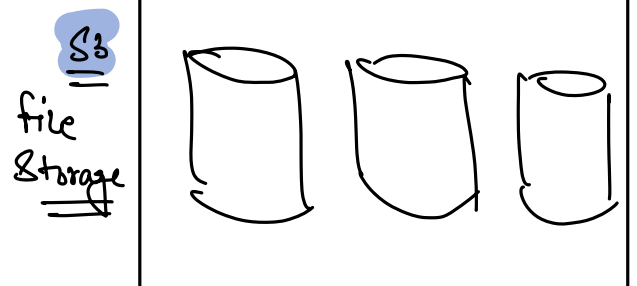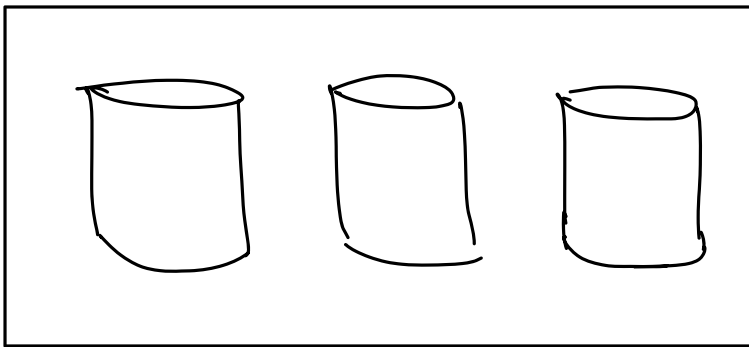
⇒ File Storage.

    → Media Content
       (Image/Video/file etc)

    → In main DB (SQL or NOSQL), we just
       store the S3 URL to refer the content
       from S3.

→ File Storage System should be able to store
big files.

       ↳ 1GB | 10GB | 1TB | 10TB ...

→ Durable

→ Performance of uploading (or) downloading.

       ↳ What happens if our connection drops
         in b/w ?

**Option 1**: File stored as 1 unit in 1 machines.

**Pros**: 1) No need to maintain the chunk info.
(Cost of entries)

2) We need not to collate the chunks.

**Cons**: 1) File size is limited by machine size.

2) Parallelism NOT possible.

**Option 2**: Divide file into smaller chunks and store each chunk into a different m/c.

**Pros**:

**Cons**:

Note: We are going to divide a file into multiple chunks so that the file size isn't limited by the m/c size & parallelism would be possible.

But we shouldn't be dividing file into very small chunks otherwise managing these many chunks would become an overhead.

⇒ **HDFS.**

  ↳ Hadoop Distributed file System.

→ Replication

1) **Data Nodes** : Machines where we store file Chunks.
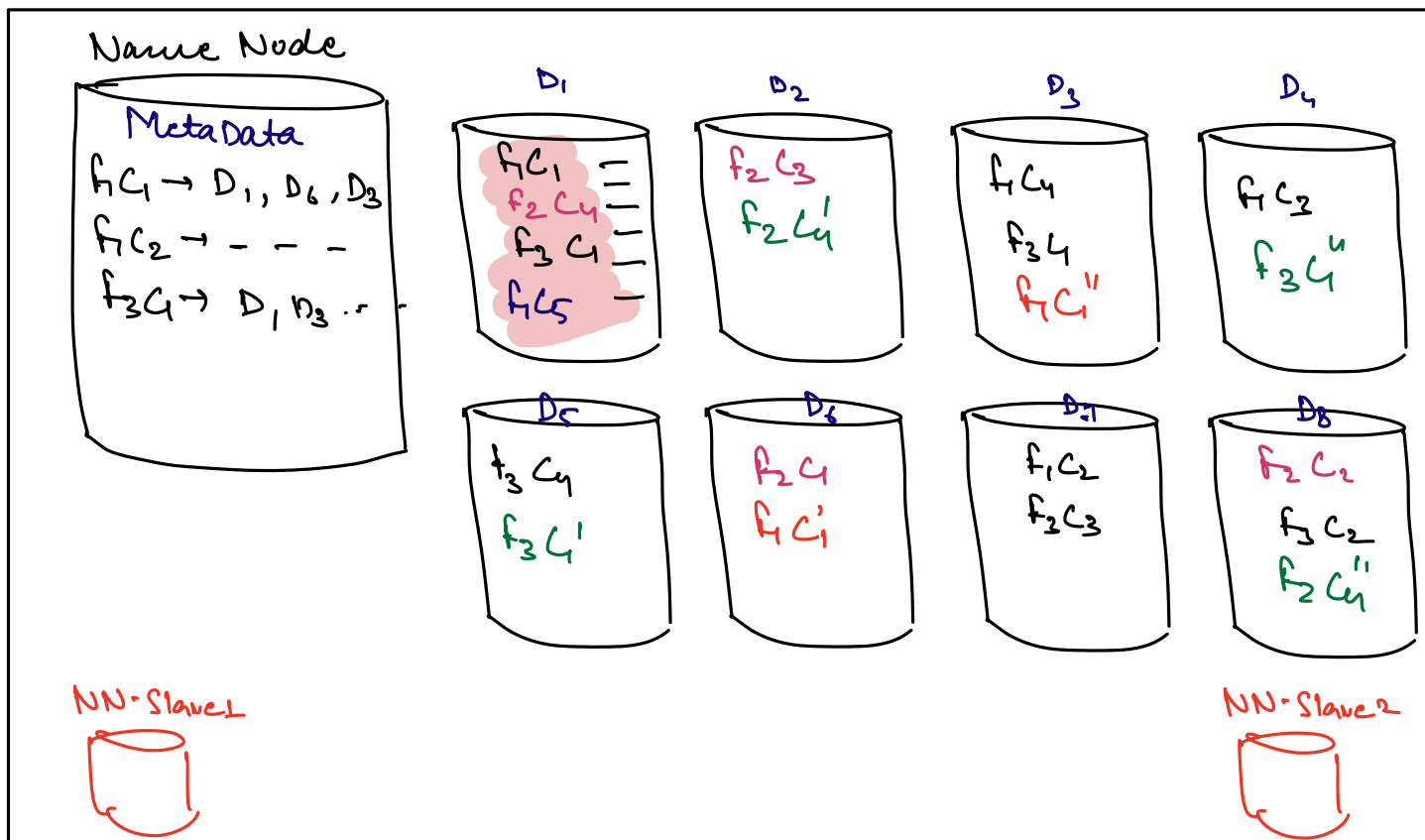
2) **Name Nodes** : Store metadata about the files.

  → Maintains the mapping of file chunks and in which m/c these chunks are present.
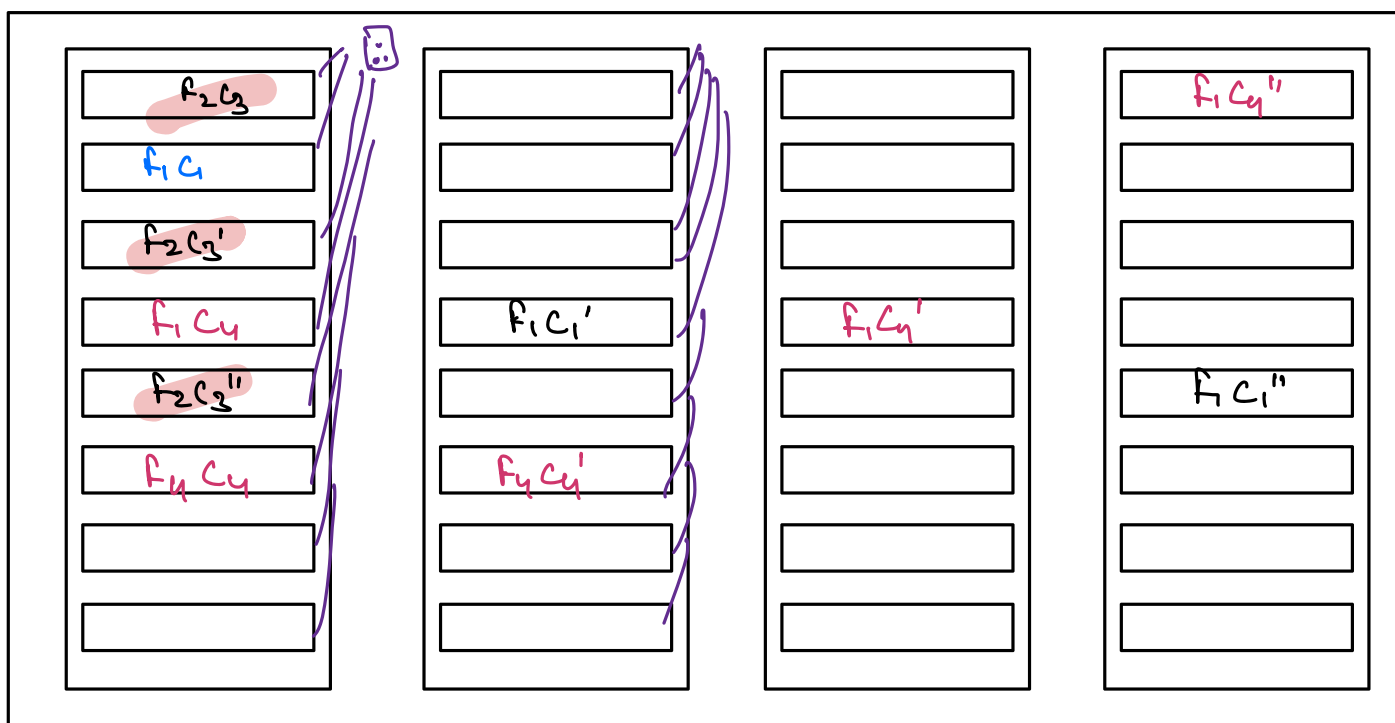
  → Replication

---

F1 → 200 Chunks

$f_1 C_1$ → $D_7$

$f_1 C_2$ → $D_{10}$

$f_1 C_3$ → $D_4$

$f_1 C_4$ → $D_8$

    ⋮    ∴

---

⇒ In one HDFS cluster, we don't shard a Name Node server but replication would be there.

⇒ HDFS Cluster.

Name Node

MetaData
$f_1 C_1 \rightarrow D_1, D_6, D_3$
$f_1 C_2 \rightarrow - - -$
$f_3 C_1 \rightarrow D_1, D_3 \cdots$

**$D_1$**
$f_1 C_1$
$f_2 C_4$
$f_3 C_1$
$f_1 C_5$

**$D_2$**
$f_2 C_3$
$f_2 C_4'$

**$D_3$**
$f_1 C_4$
$f_3 C_1$
$f_1 C_1''$

**$D_4$**
$f_1 C_2$
$f_3 C_1''$

**$D_5$**
$f_3 C_4$
$f_3 C_1'$

**$D_6$**
$f_2 C_1$
$f_1 C_1'$

**$D_7$**
$f_1 C_2$
$f_2 C_3$

**$D_8$**
$f_2 C_2$
$f_3 C_2$
$f_2 C_4''$

NN·Slave1

NN·Slave2

# Rack Aware Algo

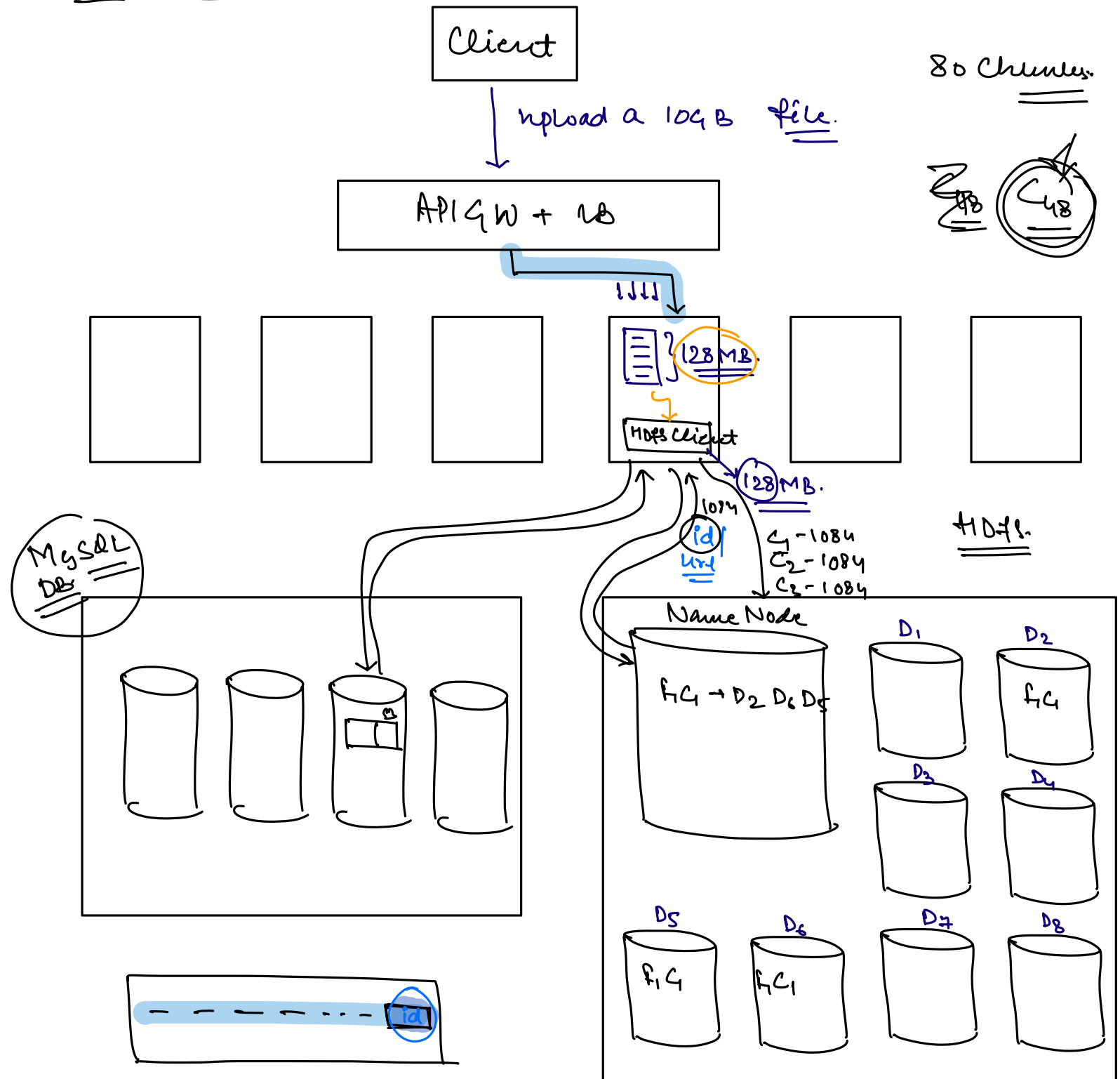| | | | |
|---|---|---|---|
| $f_2 C_3$ | | | $f_1 C_4''$ |
| $f_1 C_1$ | | | |
| $f_2 C_3'$ | | | |
| $f_1 C_4$ | $f_1 C_1'$ | $f_1 C_4'$ | |
| $f_2 C_3''$ | | | $f_1 C_1''$ |
| $f_4 C_4$ | $f_4 C_4'$ | | |
| | | | |
| | | | |

for more *reliability*, HDFS keeps data on different racks so that we do not loose our data even if a rack goes down.
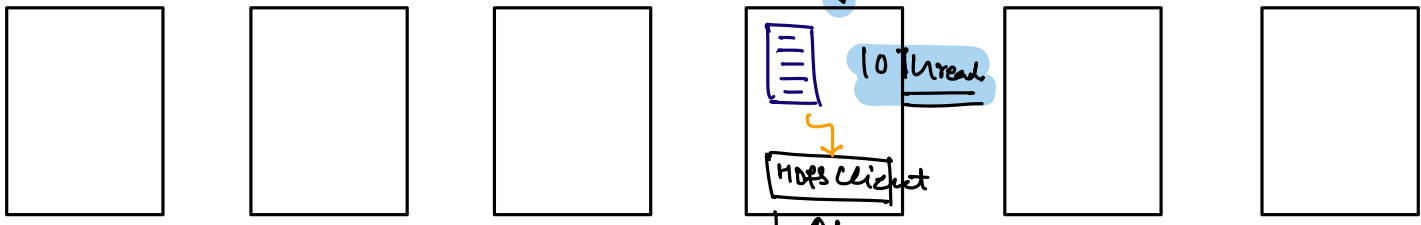
→ Avoid replicating same chunks within the same rack.

## Upload flow.

Client

upload a 10GB file.

API GW + LB

128 MB

HDFS Client

128 MB.

80 Chunks.

80 GB → $C_{48}$

1024

id
url

$C_1 - 1084$
$C_2 - 1084$
$C_3 - 1084$

HDFS.

MySQL DB

Name Node

$f_1 C_1 → D_2 \ D_6 \ D_5$

$D_1$

$D_2$

$f_1 C_1$

$D_3$

$D_4$

$D_5$

$f_1 C_1$

$D_6$

$f_1 C_1$

$D_7$

$D_8$

id