

Infrastructure for Data Analytics and Machine Learning

Part 1 - Distributed Filesystems

Paweł Wiejacha
RTB House

Plan for this lecture

- Filesystem, Distributed Filesystem, Object Storage
 - what they are?
 - and why we need them?
- What can we expect from Distributed Storage?
- How do Distributed Storage work?

File system

File system - a method and data structure used to control how data is stored and retrieved.

- Digital file systems are named and modeled after paper-based filing systems (e.g. file, directory, folder).
- Separating data into pieces (**files**), naming them, and organizing files into **directories** simplifies data storage and retrieval
- Without such separation there would be just unorganized chunks of data written on a storage medium

The file system is responsible for:

- organizing files and directories
- mapping file parts to physical storage blocks
- keeping file and directory metadata (e.g. creation date, owner)



Distributed file system

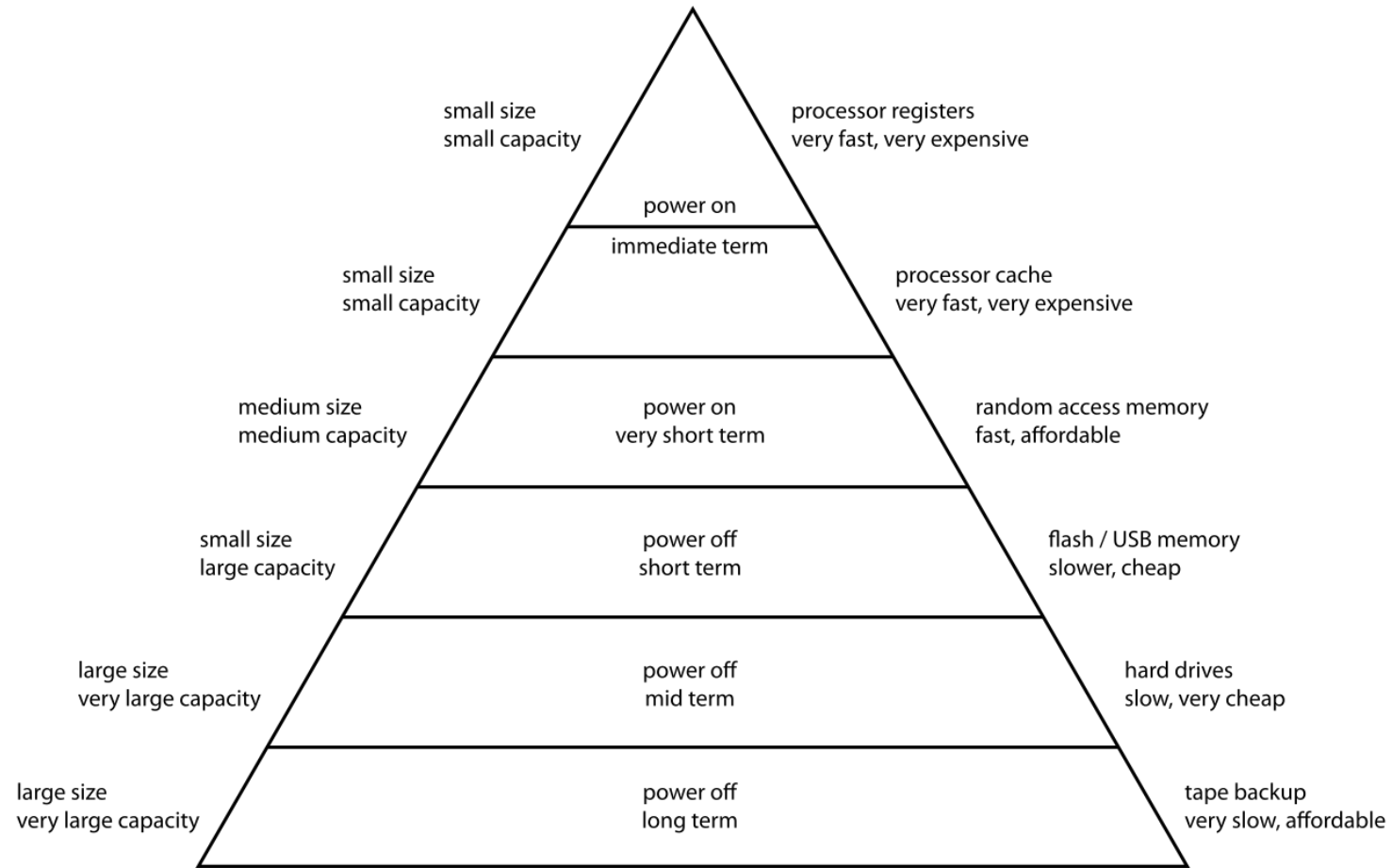
Distributed file system - a file system where data is stored across multiple networked computers which communicate and coordinate their actions by passing messages.

Just like ordinary file system has to map files to disk sectors, distributed file system has to:

- map files to disk sectors stored on different nodes
- but also handle adding and removing machines to a system, handle network partitioning, etc.

Why we need a distributed storage?

Memory hierarchy



Memory hierarchy in 2022

Storage medium	Max. capacity per node	Reading / Writing speed	Seek time / random IOPS	Price per TB
Magnetic tape	20 TB	read: 400 MB/s write: 400 MB/s	10-100 s	drive + \$13
SAS/SATA HDD	20 TB	read: 200 MB/s write: 269 MB/s	4.15 ms (~250 IOPS)	\$22
NVMe SSD	30 TB	read: 7 000 MB/s write: 3 600 MB/s	0.001 ms (1e-6 s) read: 930 000 IOPS write: 190 000 IOPS	\$125
DRAM	4 TB	190 000 MB/s	68-100ns (6.8e-8 s)	\$7 500
HBM3	0.64 TB	600 000 MB/s	107-130ns (1e-7 s)	\$9 000
L3 Cache (SRAM)	0.0002 TB	~1 000 000 MB/s	~12ns (1.2e-8 s)	\$7 000 000

Memory hierarchy in 2024

Storage medium	Max. capacity per node	Reading / Writing speed	Seek time / random IOPS	Price per TB
Magnetic tape	20 TB	read: 400 MB/s write: 400 MB/s	10-100 s	drive + \$13
SAS/SATA HDD	30 TB	read: 200 MB/s write: 269 MB/s	4.15 ms (~250 IOPS)	\$22
NVMe SSD	30 TB	read: 12 400 MB/s write: 11 800 MB/s	0.001 ms (1e-6 s) read: 1 500 000 IOPS write: 1 500 000 IOPS	\$125
DRAM	12 TB	285 000 MB/s	68-100ns (6.8e-8 s)	\$10 000
HBM3	1.28 TB	600 000 MB/s	107-130ns (1e-7 s)	\$9 000
L3 Cache (SRAM)	0.0005 TB	~1 000 000 MB/s	~12ns (1.2e-8 s)	\$7 000 000

Do we need a distributed storage?

- We know what Distributed File Systems are, but **do we really need them?**
 - and if so, when?
- Computers became parallel and very powerful, there are technological advancements in the storage field
 - maybe we can buy an expensive computer, big and fast disks and use a simple non-distributed filesystem?

Problem: non-distributed 10 PB storage

Imagine a "BigData" scenario where:

- we store **10 PB** of data
- and perform multiple batch tasks that read only a part (i.e. **1 PB**) of stored data

10 PB of data is not that much, e.g.:

- 1.5e9 smartphone photos (~4 photos of every product available on Amazon)
- or 1080p movies uploaded to YouTube over 125-250 days
- or bid-logs uploaded to RTB House during 9 hours

Let's try to design and select hardware for a **non-distributed** storage for this scenario.

Is it possible?

Hardware limits and bottlenecks (2022)

- Number of processors in a single server
 - 2 CPUs (2x 48 CPU cores)
- Expansion bus limits
 - 128-160 PCIe gen4 lanes
 - ~2GB/s per PCIe gen4 lane
- Network adapters
 - 200 GbE (~25 GB/s)
 - latency: 1.01 microseconds
 - price ~\$1600 per two port NIC

Storage failure rates

Annualized failure rate (AFR) - estimated probability that a component will fail during a full year of use.

Mean Time Between Failures (MTBF) - $\frac{1}{\lambda}$, where λ is the failure rate.

$$AFR = 1 - e^{-8766/MTBF}$$

Both are population statistic estimated on a sample of given hardware components

Storage medium	<i>Real</i> AFR
Magnetic tape ^[1]	3.44%
SAS/SATA HDD	0.80% - 1.10%
NVMe SSD	0.78% - 1.22%
DRAM (Reg. ECC)	0.07 - 0.11%

[1] assuming unrealistic 100% duty cycle

Failures in real world

- **disk level:** assuming 1.22% AFR, given a system with 1450 disks, we can expect a disk failure every 3 weeks
- **node level:** power supply unit failures, motherboard failures
- **multi-node level:** network switch failures, a rack issues
- **data-center level:** maybe not earthquakes, but DNS/BGP problems
- human errors

With failures so common, they **cannot degrade** storage system performance - they need to be transparent to users.

Hardware setup for a non-distributed storage (1)

Let's try to select hardware for a imaginary non-distributed storage (10 PB of stored data, 1 PB batch task)

Hardware setup for a non-distributed storage (1)

Let's try to select hardware for a imaginary non-distributed storage (10 PB of stored data, 1 PB batch task)

- let's get back to our large non-distributed storage problem
- reading 1 PB using a single HDD: $1 \text{ PB} / 200 \text{ MB/s} = 62 \text{ days}$ - too long

Hardware setup for a non-distributed storage (1)

Let's try to select hardware for a imaginary non-distributed storage (10 PB of stored data, 1 PB batch task)

- let's get back to our large non-distributed storage problem
- reading 1 PB using a single HDD: $1 \text{ PB} / 200 \text{ MB/s} = 62 \text{ days}$ - too long
- another reason to use more disks is storage space:
 - $10 \text{ PB} / 30 \text{ TB disks} = 341 \text{ disks}$
 - we want to handle simultaneous failure of 2 disks: $341 + 25\% = 426 \text{ disks}$
 - we want to have at least 20% of free space: $426 + 25\% = 532 \text{ disks}$
 - our system/company is expecting to grow: $532 + 15\% = 600+ \text{ disks}$

Hardware setup for a non-distributed storage (1)

Let's try to select hardware for a imaginary non-distributed storage (10 PB of stored data, 1 PB batch task)

- let's get back to our large non-distributed storage problem
- reading 1 PB using a single HDD: $1 \text{ PB} / 200 \text{ MB/s} = 62 \text{ days}$ - too long
- another reason to use more disks is storage space:
 - $10 \text{ PB} / 30 \text{ TB disks} = 341 \text{ disks}$
 - we want to handle simultaneous failure of 2 disks: $341 + 25\% = 426 \text{ disks}$
 - we want to have at least 20% of free space: $426 + 25\% = 532 \text{ disks}$
 - our system/company is expecting to grow: $532 + 15\% = 600+ \text{ disks}$
- reading 1 PB using 600 HDDs: $1 \text{ PB} / 600 / 200 \text{ MB/s} = 2.5 \text{ hours}$ - **acceptable**

Hardware setup for a non-distributed storage (2)

Let's try to select hardware for a imaginary non-distributed storage (10 PB of stored data, 1 PB batch task)

- 600 HDDs is too many for a single machine
 - slots for disks - max 72 disks per machine (usually 12-48)
 - 600 HDDs: 420 kg, 12 meters high stack or 90 meters long chain

Hardware setup for a non-distributed storage (3)

Let's try to select hardware for a imaginary non-distributed storage (10 PB of stored data, 1 PB batch task)

- $600 \text{ disks} * 200 \text{ MiB/s} = 117 \text{ GiB/s}$ - slightly too much for a single machine:
 - 59 PCIe lanes just to read data from disks
 - copy to memory (12 cores needed - for a simple `memcpy()`)
 - decrypt, compute checksums (46 cores needed for decryption)
 - read from memory (another 12 cores)
 - compute or send over network
 - sending over network requires 64 PCIe lanes and 5 network adapters
 - Linux kernel would have a problem to handle that many network packets

Hardware setup for a non-distributed storage (4)

Let's try to select hardware for a imaginary non-distributed storage (10 PB of stored data, 1 PB batch task)

Fault tolerance and reliability:

- having spare PSU, motherboard, CPUs, RAM modules is a good idea
- but it's not possible to achieve high availability this way

Evolution of systems:

In real world, it's likely you are scaling up an existing system (e.g. hardware bought 2 years ago), e.g.:

- you have lots of 15TB disks instead of 30 TB ones
- and some PCIe gen3 motherboards instead of new PCIe gen4

Advances in modern hardware allow to create truly high performant storage systems.

But in many cases we still need distributed storage

Distributed storage interfaces for batch processing

Block storage

Block-level storage - storage where interface emulates behaviour of a traditional block device (e.g. HDD).

Data is organized as **fixed-size blocks** which are identified by an arbitrary identifier which can be used to store and retrieve given block. There are no files, directories, no structure - just blocks.

Examples:

- non-cloud: Ceph RADOS Block Device
- cloud: Amazon EBS, Google Cloud Persistent Disks

Verdict: **Too low level to for batch processing** - we need higher level of abstraction.

POSIX file system interface

Portable Operating System Interface (POSIX) - a family of standards specified in the 1988 for maintaining compatibility between (UNIX-like) operating systems.

There is a section of the standard that defines the semantics of the POSIX-compatible file system, e.g.:

- allows hierarchical file names and resolution (e.g. `/dir/dir/file`)
- strong consistency, atomic writes, atomic renames, deleting open files
- implement certain operations, like:
 - random access reads/writes (`pread()`, `pwrite()`)
 - access control (`chown`, `chmod`, etc.)
 - symlinks, hardlinks
 - `ftruncate()`, `fsync()`, `fcntl()`, `mmap()`, `fadvise()`, `fallocate()`

Interfaces for batch processing

POSIX file system interface

Examples^[1]:

- local filesystems: XFS, ext4, Btrfs, ZFS, ...
- distributed filesystems: CephFS, GlusterFS, MooseFS/LizardFS, Lustre, JuiceFS

Verdict:

- **too complex** - implementing a POSIX-compliant distributed filesystem is a challenge
- not all features nor consistency guarantees are needed for batch processing

[1] - with various level of POSIX compliance

Object Storage

Object storage, blob storage - storage architecture that manages data as objects that

- have user-defined globally unique identifier
- contain variable length data
- and optionally a metadata attached to them

Usually, the structure is flat (limited number of buckets, objects inside buckets) - there are no nested file nor directories.

Basically, it is a KVS (**Key-Value Store**) with big "Values" and metadata.

Object Storage

Limited set of operations, usually:

- `get_object(object_name) : Data`
- `put_object(object_name, data)`
- `delete_object(object_name)`
- `set_metadata(object_name, metadata)`
- `list_objects(prefix) : List[ObjectName]` - since there are no directories, we need to use prefix searches like:
 - "list all objects with names starting with `/folder/sub-folder/`" to simulate file directories

Limited functionality:

- sometimes only eventual consistency
- immutable objects
 - `put_object(object_name, data)` creates or overwrites object `object_name`
 - no positioned writes, no `mmap()`
 - in some cases: append-only writes
- sequential reads or limited random reads

Interfaces for batch processing

Object Storage

Examples:

- Cloud: **Amazon S3**, Google Cloud Storage, Azure Blob Storage, Cloudflare R2, Backblaze B2
- Open Source: **Ceph Object Gateway**, MinIO, OpenStack Swift, HDFS ^[1]

Verdict: **just right for batch processing**

- have all operations needed for batch processing
- simplified semantic does not limit storage performance

[1] - not a POSIX filesystem, not a object storage

Interfaces for batch processing

Object Storage - blobs and record sets

Object storage is a kind of KVS but with "big" values, e.g.

1. **binary blobs**: images, videos, or recorded sounds
2. files that contain multiple "**records**" in formats like CSV, JSON, Avro, Parquet

Let's skip the obvious case of batch processing large binary blobs, and focus on the second one.

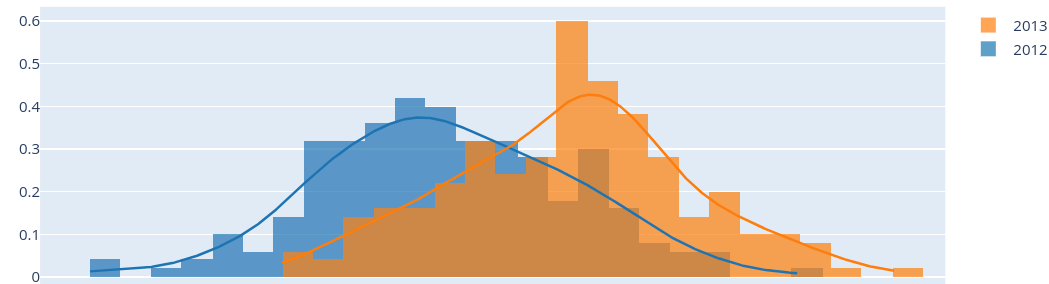
Storing data as JSON records in files? How about using a **database**?

Example scenario: storing and processing data for a huge store

Imagine you own a big e-commerce store that tracks every action user performs on your website (browsing products, performing purchases, writing a product review, etc.)

And that you want to answer **ad-hoc** questions like:

- What was the price of product ABC 200 days ago?
- What how many unexpired promo tokens user XYZ had 183 days ago?
- And what **exactly** was user XYZ's detailed browsing history just **before** request 9742-234237483-56452-554232 was received?
- What is the distribution of basket values for premium users who have more than two delivery addresses?



Database vs Object Storage - DB approach

- Analyze and **design the domain model** carefully
- **Normalize data**
- Design complex database schema and intricate entity relationships
 - with hundreds of tables and relations
- Create many huge **database indices** to speed up queries
- Optimize and maintain indices to match common queries
- Use rigid query language to perform your business logic
- Execute queries using really **huge joins** using huge DB indices

```
SELECT
  event.time,
  event.type,
  product.name,
  product_prices.price,
  user_segment, ...
FROM
  events JOIN
  products JOIN
  product_prices JOIN
  i18n_product_details JOIN
  users JOIN
  order_requests
  ...
```

Database vs Object Storage - denormalization



Database vs Object Storage - filesystem approach

Sometimes it's better to store data records as they are:

- **without restructuring, normalization, deduplication** - e.g. with deeply nested, repeated and redundant fields
- **without indexing**
- using very simple but flexible schema
- just roughly partition them (e.g. by event type and creation date) while writing
- and process them in a **massively parallel** manner when needed

```
/warehouse/products/...  
/warehouse/user-reviews/...  
/warehouse/user-reviews/2022-05-02/reviews-1652658665521_73_QPECDbzZXa.json.snappy  
/warehouse/user-reviews/2022-05-02/reviews-1652658665523_32_fle34scyaV.json.snappy  
/warehouse/user-reviews/...  
/warehouse/browsing_events/...  
/warehouse/browsing_events/2022-05-02/...  
/warehouse/browsing_events/2022-05-03/...  
/warehouse/browsing_events/2022-05-04/region=us/...  
/warehouse/browsing_events/2022-05-04/region=eu/log-1652658664132_40_hXMzbDepoC.json.snappy  
/warehouse/browsing_events/2022-05-04/region=eu/log-1652658665521_70_aXZzbDCEPQ.json.snappy  
/warehouse/browsing_events/2022-05-04/region=eu/log-1652658278519_13_xiTNJreBBT.json.snappy  
/warehouse/browsing_events/2022-05-04/region=eu/...  
/warehouse/subscriptions/...  
/warehouse/purchases/...
```


Database vs Object Storage - why use the filesystem approach?

- because of **performance reasons**
 - creating and maintaining huge indices is a real issue
 - many DB entities - lots of joins, lots of random I/O
 - as opposed to reading data in a massively parallel and **sequential** manner
- the world, your company, your platform and your data **will evolve**
 - maintaining database schema and existing queries would be a nightmare
- sometimes you cannot build an index, because **you cannot predict the future**, e.g.
 - what queries analytic team will need to perform next quarter?
 - what algorithms Machine Learning team will create?
- **sometimes it better to store unstructured data** (e.g. raw requests in a JSON format)
 - it might be easier than reaching consensus about taxonomy, schema, shared indices, priorities

Database vs Object Storage - data locality

- Sometimes the easiest or fastest way is to process unstructured data sequentially
 - this processing might require a lot of computing power
 - and performing computations **close to data** might be a desired property
 - property that might not be possible to achieve when a standard database is used

Distributed storage - features and requirements

Distributed storage - basic requirements

What are our basic requirements for a distributed storage?

- **storing data** and providing reasonable interface
- having **scalable storage capacity**
- having **scalable** data access **throughput**
 - i.e. increasing number nodes/disks increases total bandwidth/IOPS
- being **resilient** to hardware failures
 - both temporary (e.g. HA) and persistent (e.g. broken disks) failures

Distributed storage - non-basic requirements

What else can we expect from a distributed storage used for batch processing?

- Encryption at rest
- Metadata
- Object versioning
- Constant-time file concatenation
- Snapshots
- Compression
- Atomic rename or Compare-and-Set
- Lifecycle and retention policies

Distributed storage

building blocks and common design patterns

Checksums

Computing checksums (error detecting codes) is the **absolute minimum** distributed storage should do.

- the simplest way is parity check (split data into n -bit words and xor them)
- CRC (cyclic redundancy check)
- using cryptographic hash function that uses dedicated CPU instructions (e.g. SHA instruction set)

```
import zlib
zlib.crc32(b"very-long-piece-of-data-abcdefgh") # 2142146772
zlib.crc32(b"very-long-piece-of-data-a6cdefgh") # 2822004080
                                     ^-----
zlib.crc32(b"very-long-piece-of-ERRORabcdefgh") # 1649843363
```

Resiliency - replication

Replication is the simplest way to achieve storage that is resilient to disk or node failures.

Basically, to handle $n - 1$ disk failures, instead of writing each object (file or block) once, we need to **write it n times** on n different disks (and different nodes if we want to handle node failures).

Pros:

- the simplest way to achieve resiliency
- can speed up reading process by reading from different disks
 - if the same file is read simultaneously by multiple readers

Cons:

- you need to buy n times more disks

Resiliency - erasure coding

[optimal] erasure code is a forward error correction:

- that transforms a message of n symbols into a message with $m = n + r$ symbols (n original and r redundant)
- so that the original message can be recovered from **any n -symbol subset** of the m symbol message
- the "symbol" can be any number of bits/bytes (processor word, 512 byte sector, 64 KB block)
- we are assuming "erasure" (disappearing) of symbols instead of "error" (corruption) of symbols
 - that's why checksums are so important

Resiliency - erasure coding - $r = 1$

- the simplest erasure code is xor/parity code ($r = 1$)

```
# n = 3, m = 4

input: 123, 561, 913
parity: 123 ^ 561 ^ 913 == 475

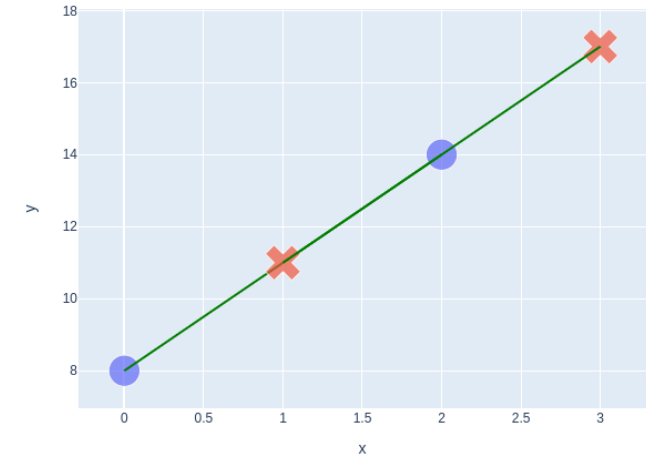
input[0] is missing:
  ???   561   913
(475) ^ 561 ^ 913 # == 123 (recovered 1st chunk)

input[1] is missing:
123 ^ ??? ^ 913
123 ^ (475) ^ 913 # == 561 (recovered 2nd chunk)

input[2] is missing:
123 ^ 561 ^ ???
123 ^ 561 ^ (475) # == 913 (recovered 3rd chunk)
```

Resiliency - erasure coding - $r > 1$

- for general case ($r > 1$), one of methods is polynomial oversampling
 - you need at least $n + 1$ points to determine n -degree polynomial coefficients [over some finite field]
 - so you can:
 - use original message symbols as polynomial $p(x)$ **coefficients**
 - then generate $\{ p(x_i) : i \in 1..m \}$ samples and use them as a erasure-coded message
- e.g. for the input message [5, 3], -we have $p(x) = 5 + 3x$
 - $5 = p(x_0) \implies x_0 = 0$
 - $3 = p(x_1) \implies 3 - 5 = 3x_1 \implies x_1 = -2/3$
 - output 2+2 message: $\langle p(x_0), p(x_1), p(2), p(3) \rangle$
 - [5, 3, 11, 14]
- Reed-Solomon EC - polynomial oversampling in $GF(2^k)$



Resiliency - example comparison

Input: 2000 KiB file, block size = 1000 KiB

Replication (3x):

- split file into 2 blocks: b_1 and b_2 , then for each block
 - write the same block 3 times using 3 different disks

8+2 erasure coding:

- split file into 2 blocks: b_1 and b_2 , then for each block
 - split each block into 8 chunks (125 KiB each): c_1, \dots, c_8
 - compute parity chunks:
 - $c_9 = \text{parity}_\alpha(c_1, \dots, c_8)$,
 - $c_{10} = \text{parity}_\beta(c_1, \dots, c_8)$
 - write each chunk on a different disk

stored data: [0123456789abcdefghijklmnopqrstuw]

	block_1	block_2
replica 1:	[0123456789abcdef]	[ghijklmnopqrstuw]
replica 2:	[0123456789abcdef]	[ghijklmnopqrstuw]
replica 3:	[0123456789abcdef]	[ghijklmnopqrstuw]

	block_1	block_2
original chunk 1:	[01]	[gh]
original chunk 2:	[23]	[ij]
original chunk 3:	[45]	[kl]
original chunk 4:	[67]	[mn]
original chunk 5:	[89]	[op]
original chunk 6:	[ab]	[qr]
original chunk 7:	[cd]	[st]
original chunk 8:	[ef]	[uw]
redundant chunk 1:	[UV]	[PQ]
redundant chunk 2:	[XY]	[RS]

Resiliency - erasure coding (3)

- the more chunks the **smaller storage overhead** (e.g. $1000+2$ vs $10+2$)
 - but need multiple disk accesses (e.g. 1000 seeks) to read a single block
 - and you need system with at least 1002 **different** disks
- computing redundant chunks requires more computing power than a simple copy done by a replication
- you need more disks to get the same resiliency level, e.g.:
 - 3x replication - we need only 3 disks
 - 8+2 EC - we need 10 disks (however we have 8 times more usable space)
- Erasure coding is a much more **complex solution**, e.g.
 - in case of failure of a disk containing original chunk, you need to reconstruct the block on the fly

Sharding

Sharding - partitioning data horizontally (e.g. by rows not columns) so that each data shard is stored or handled by a different instance, process or disk.

Examples:

- storing all blocks with `hash(block_uuid) mod num_machines == i` on machine i
- storing all metadata with `hash(object_name) mod num_metadata_servers == i` on server i

This way we can **distribute load and responsibility** evenly and scale our system.

```
file_names = ['file0', 'file1', 'file2', 'file3', 'file4', 'file5', 'file6', 'file7', 'file8', 'file9']  
  
>>> [hash(file_name) % 3 for file_name in file_names]  
[0, 1, 2, 1, 0, 2, 0, 2, 0, 0]  
  
# server0: 'file0', 'file4', 'file6', 'file8', 'file9'  
# server1: 'file1', 'file3'  
# server3: 'file2', 'file5', 'file7'
```

Tiering and caching

Tiering - arranging something in tiers (layers)

For example:

- Keeping frequently used data (filesystem hierarchy, indexes, file metadata) in RAM to reduce latency
- Keeping hot data (recently read files, blocks prefetched by the read-ahead mechanism) in RAM or SSD
- Keeping cold and ordinary data on slower medium (e.g. HDD)

Data organization

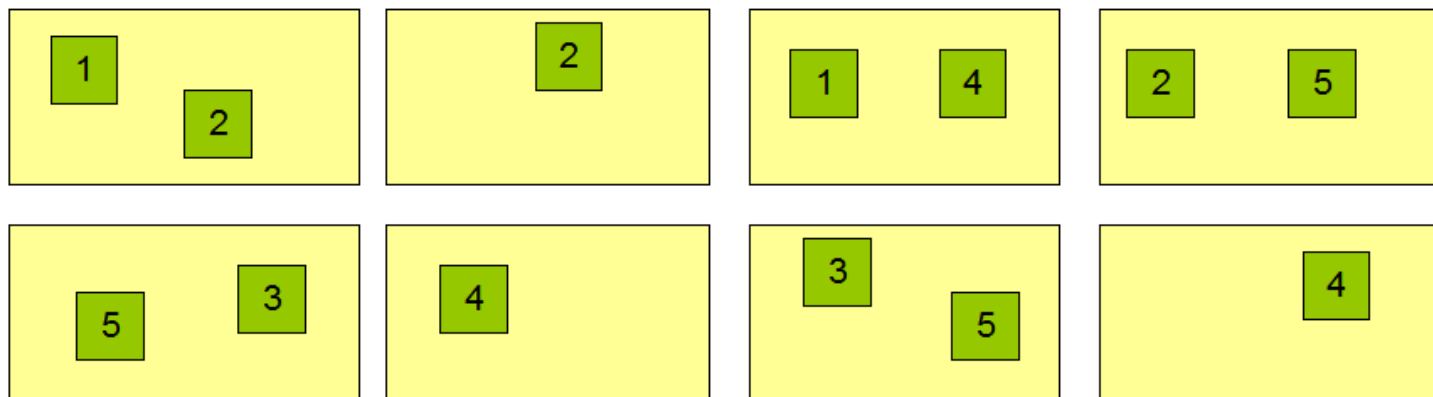
- storing data on multiple machines (*data nodes*)
- splitting files into **blocks** with limited size (e.g. up to 128 MB)
 - simplifies things like reconstruction, balancing, out of space handling
 - allows to parallelize data transfers
- placement awareness (same disk, same node, same rack)
- replication or erasure coding
- data transfers should be done directly to/from a data node
- the simplest way to store blocks is to use a normal filesystem (XFS, EXT4)
 - examples: HDFS, MinIO, OpenStack Swift,
 - this way we can use `fsck.ext4`, `rsync`, `ls`, `cp`, software RAID, etc.
- if we need to squeeze more performance we can create storage layer by ourselves
 - example: BlueStore in CephFS

Data organization

Block Replication

Namenode (Filename, numReplicas, block-ids, ...)
/users/sameerp/data/part-0, r:2, {1,3}, ...
/users/sameerp/data/part-1, r:3, {2,4,5}, ...

Datanodes



Metadata organization (1)

- We stored data blocks in a resilient way on data nodes.
- During read request, client is going to connect directly to data nodes to fetch blocks
- But which data node client should connect?
- We need some place to **store metadata**
 - (e.g. file size, list of file's blocks and their locations, index of stored files)
- We can use deterministic algorithms and placement policies, sharding, broadcast messages, etc. to reduce size of metadata
 - e.g. Ceph uses CRUSH, a scalable pseudorandom data distribution function that efficiently maps data objects to storage devices without relying on a central directory
- And if we want to provide reasonable consistency guarantees we need to coordinate (e.g. serialize) all metadata operations

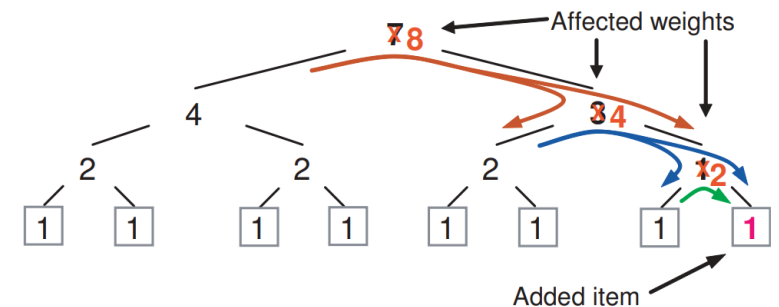
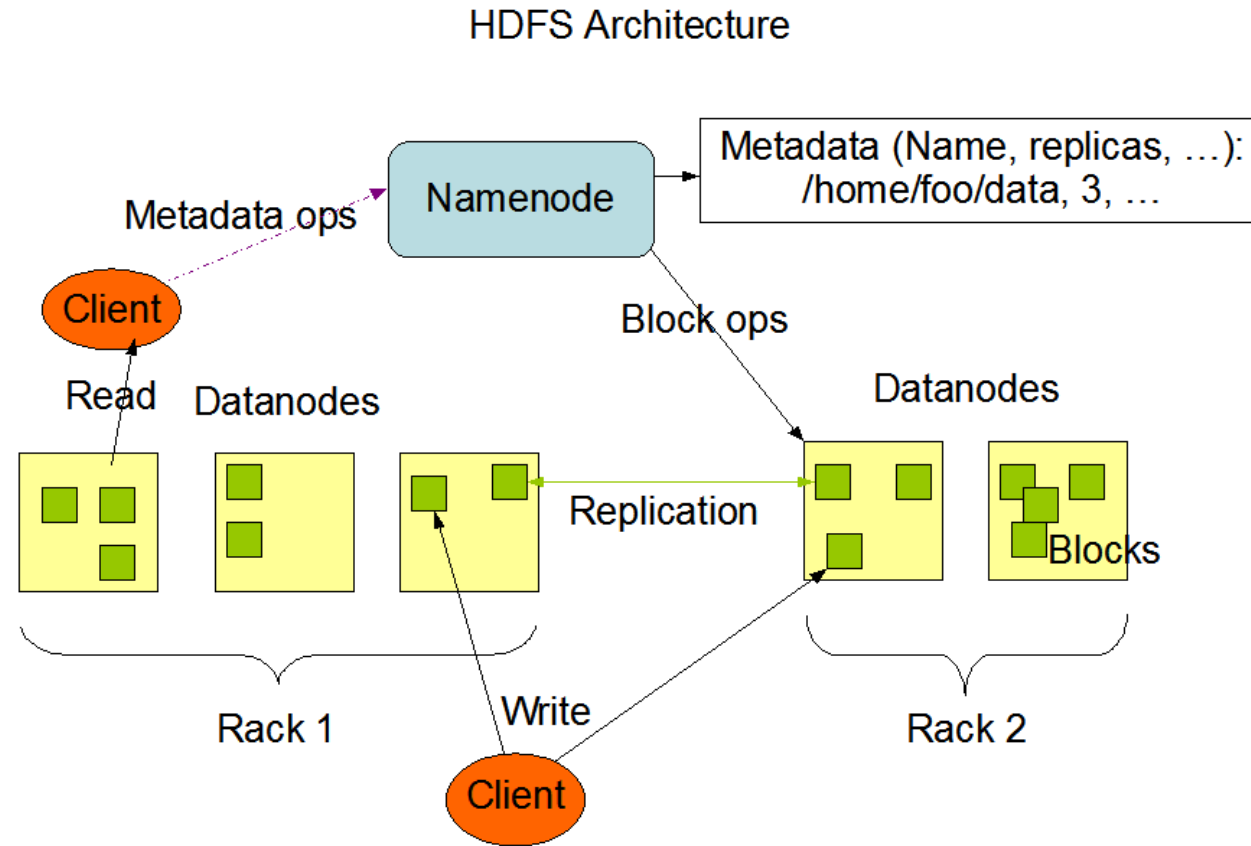


Figure 3: Data movement in a binary hierarchy due to a node addition and the subsequent weight changes.

Metadata organization (2)



Metadata organization (3)

- the simplest way is to keep all metadata on a **single server**
 - powerful enough to handle all metadata operations and keep hot data in RAM
 - forcing every metadata request to go through this server to guarantee consistency
- next level is to use **active/standby** configuration, e.g.
 - writing metadata/log on a secure storage before confirming
 - waiting for replica confirmation storage before confirming to client
- next level: multiple metadata servers that use **distributed consensus** algorithms (Raft, Paxos)
- don't forget about sharding
- how to store metadata?
 - filesystem
 - some basic local KVS/DB (e.g. LevelDB, RocksDB)
 - custom database

Rebalancing

Storage system should handle:

- **extension** - adding more disks and nodes to extend storage and throughput
- **retiring (decommissioning)** - marking old or partially broken nodes for removal or repair

Those actions can result in unbalanced storage cluster state

- e.g. 200 disks are 95% full and 50 new disks are empty

Unbalanced system should automatically start a **rebalancing process** that moves data blocks to new places.

Ideally this process should not degrade user reads/writes.

Reconstruction

In case of (multiple) node or disk failures number of original/redundant chunks falls

- we either lost some data (if we lost too many chunks)
- there are less chunks than required to handle next failure

Storage system can reconstruct blocks on-the-fly during read operations, but it should perform **reconstruction process** to rebuild missing chunks permanently.

- especially for files that are not going to be read in the near future

Implementing it correctly can be a little tricky:

- we should not degrade user experience
- but we also cannot delay/throttle reconstruction because another failure can happen when we are not ready for it
- we also should handle temporary "failures" (e.g. failed switch, machine reboot) gracefully

Scrubbing

Scrubbing - process of reading stored data in the background in order to detect errors.

- it can detect hardware errors (broken sectors, broken drives) and sometimes software errors (incorrect programs)
- read data can be compared to:
 - other replicas
 - reconstructed data from EC
 - stored checksum
- assuming some redundant information is present, data can be reconstructed

It's better to reconstruct data when it's still possible than have a storage system with hidden possibly unrecoverable corrupted data.

Distributed storage - non-basic requirements

Encryption at rest

If someone breaks into a data center and steals our servers/disks he should not be able to read our data

- usually we can encrypt drives and store keys in memory, and store the master key in some secure location (TPM, TEE)

Encryption at rest this might be a legal requirement.

Metadata

- sometimes you don't want to read whole file just to get some basic information
 - e.g. `gdpr_cleanup_id`, `schema_version`
- and you cannot store it in the object name
 - sometimes you want to attach or change metadata long after file was created
- **having those metadata indexed is a nice feature**
 - easier than managing the metadata and consistency in a separate database

Object versioning

- the idea is that when you overwrite an object, storage keeps previous and the current version of it
 - and when you delete it keeps old version and a deletion marker
- this is not something you cannot achieve with sensible naming policy
 - but sometimes it simplifies things and is a nice feature
 - especially if paired with lifecycle policy

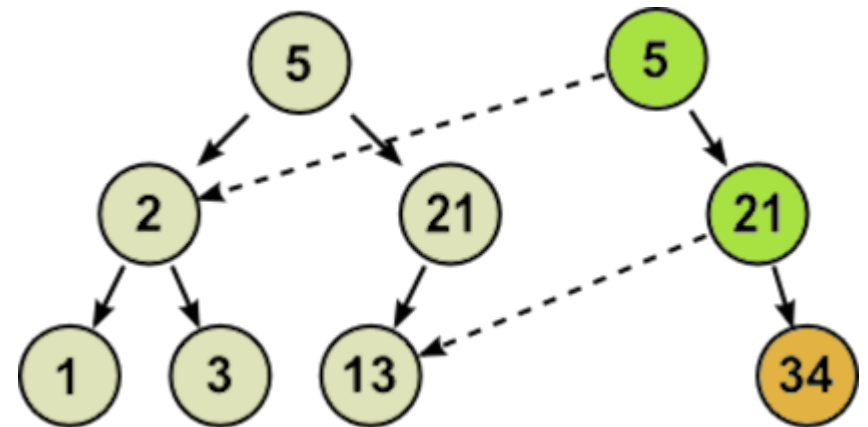
Constant-time file concatenation

- Needed in scenarios when multiple writers create multiple files (e.g. in M/R job) that logically are a single file
 - and when processing concatenated file is easier than processing hundreds of files)
- Also appending header/footer without the need of rewriting source file can be useful

Snapshots

- This is rather file system not object storage feature
- With snapshots, we can make a **instantaneous read-only point-in-time snapshot** of filesystem subtree
 - and later modify the source subtree without changing the snapshot
- We can take multiple snapshots of a given directory
- This is quite simple operation in a filesystem where you cannot modify files

HDFS supports filesystem subtree snapshots.



Compression

- Compression not only **increases storage capacity**
- It also reduces disk utilization
 - and network bandwidth (if we decompress on client)
- If we use compression algorithms that focus on compression/decompression speeds (lz4, zstd, snappy) we can also **increase download and uploads speeds**

Manual compression - compressing data by ourselves when uploading and decompress when downloading

Transparent compression - with compression implemented by the file system we reduce amount of code/logic to write and maintain

Atomic rename or Compare-and-Set

With **Compare-and-Set** operation or at least **atomic rename**, we can implement few useful algorithms and create complex systems on top of a object storage (i.e. without using additional database, distributed locks, etc.)

```
# returns current object data and object sequence number (monotonic version number)
def read_object(key: Key) -> (Value, ObjectSeqNo): ...

# atomically writes object *only* if current object's sequence number is expected_seq_no
def compare_and_set(key: Key, new_value: Value, expected_seq_no: Our) -> Success | Failure: ...
```

```
def deposit_money(amount):
    current_balance, seq_no = read_object("account_123")
    new_balance = current_balance + amount
    result = compare_and_set("account_123", new_balance, seq_no)
    if result != Success:
        deposit_money(amount) # repeat
```

Unfortunately, not many object storages implement this.

- HDFS has atomic rename
- GCS has no rename operation, no CAS
- S3 implemented strong-consistency for read-after-write (for a single PUT/DELETE) 4 years ago

Lifecycle and retention policies

Sometimes storage system can store store object in different ways:

- with more or less resiliency (e.g. handling only one not three disk failures)
- or on a cheaper medium (e.g. tapes) where object retrieval is delayed

It's very convenient if storage system allows to set **retention policies**, e.g.:

- delete objects older than X days
- move objects older than X days to a cold storage

Summary

We have discussed following topics:

- What is a Distributed Storage
- Why batch processing may require Distributed Storage
- What interface for Distributed Storage is best suited for batch processing
- What can we expect from a Distributed Storage
- How (in general) Distributed Stores work

