

Introduction to Cloud COMPUTING



“Data Cloud and Cloud Security”

2025 @ Università di Trieste

Outline



Data Cloud



Cloud Security

Overview

- ❑ Traditional Storage architectures
- ❑ Cloud storage theory and implementations
- ❑ Cloud storage and storage in the Cloud.
- ❑ Cloud storage and Cloud Distributed storage.

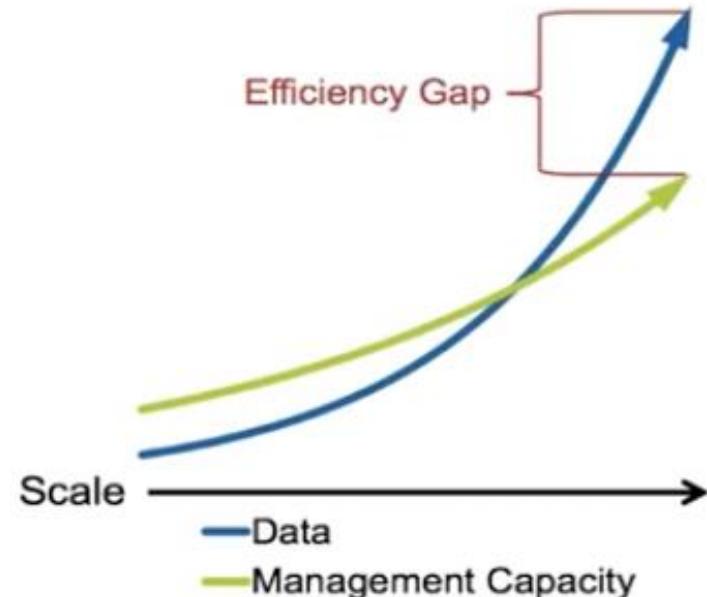


Data Efficiency gap

Cloud storage is a service model in which data is maintained, managed and backed up remotely and made available to users over a network (typically the Internet).

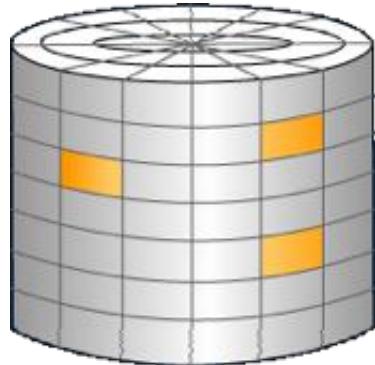
Data explosion and Mobile device growth

- improve scalability (up+out) and security
- improve performance
- simplify storage management
- on demand access
- unstructured data

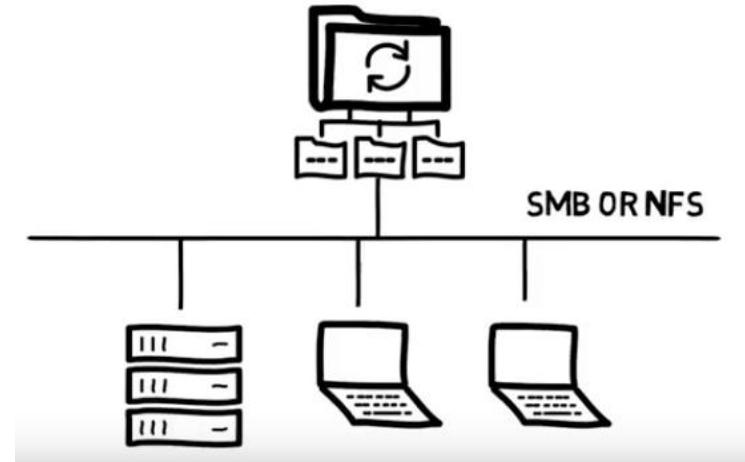


Traditional Storage Architecture

Block Storage



File Storage

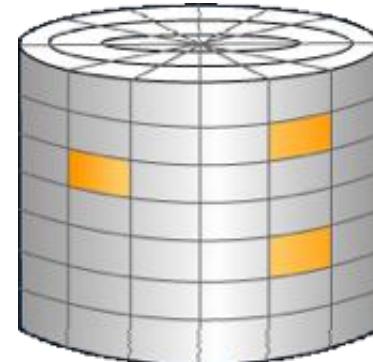


Block Storage

Block storage is similar to an HD

- Volumes
- Blocks (read and Write)
- Fibre Channel or iSCSI protocol
- Local
- Low Latency, high IOPs, low size (<1PB)
- Complex to expand and expensive

IN CASE OF CRASH I CAN
USE ANOTHER NODE TO
ACCESS THE DATA

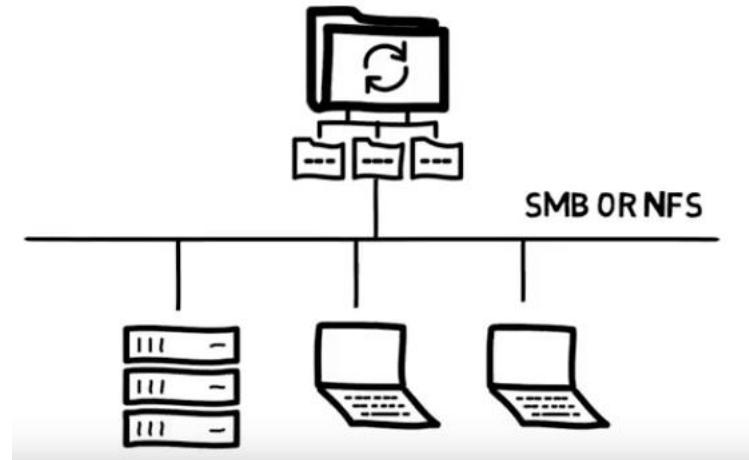


File Storage

File Storage

QUITE THE SAME IDEA
BUT I CAN OFFER THE
SAME FILESYSTEM TO
MORE THAN 1 MACHINE

- Files and directories
- Network Shared (LAN)
- SMB, NFS, OCFS etc
- High throughput, large size (PBs)
- Scale out capabilities
- Multi-tiered architecture
- Expensive



Distributed FileSystem

*“File system that is shared by many distributed clients
The resources (file+dir) on a particular machine are local to
itself. Resources on other machines are remote
Basic layer for many distributed systems (clients) and applications”*

A DFS provides a service for clients. The server interface is the normal set of file operations: create, read, etc. on files.

Servers allow clients to perform operations on resources that resides on servers.

Distributed Filesystem Properties

Transparency:

- Location: a client cannot tell where a file is located
- Migration: a file can transparently move to another server
- Replication: multiple copies of a file may exist
- Concurrency: multiple clients access the same file

Flexibility

- Servers may be added or replaced
- Support for multiple file system types

Distributed Filesystem Properties

Dependability

- Consistency: conflicts with replication & concurrency
- Security: users may have different access rights on clients sharing files & network transmission
- Fault tolerance: server crash, availability of files

Performance

- Requests may be distributed across servers
- Multiple servers allow higher storage capacity

Distributed Filesystem Properties

Caching

- Reduce network traffic by retaining recently accessed disk blocks in a cache, so that repeated accesses to the same information can be handled locally.
- If required data is not already cached, a copy of data is brought from the server to the user.
- Perform accesses on the cached copy.
- Files are identified with one master copy residing at the server machine,
- Copies of (parts of) the file are scattered in different caches.

Cache Consistency Problem -- *Keeping the cached copies consistent with the master file*

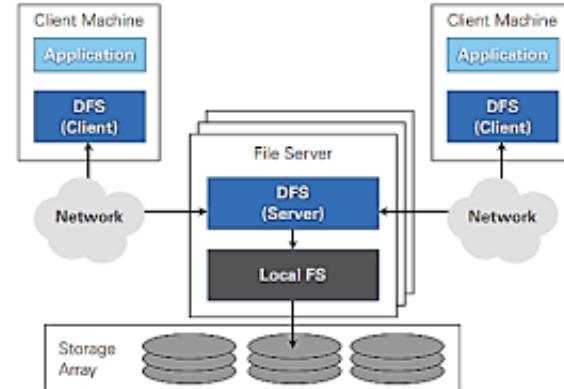
Distributed Filesystem Properties

Ideally, the client would perceive remote files like local ones.

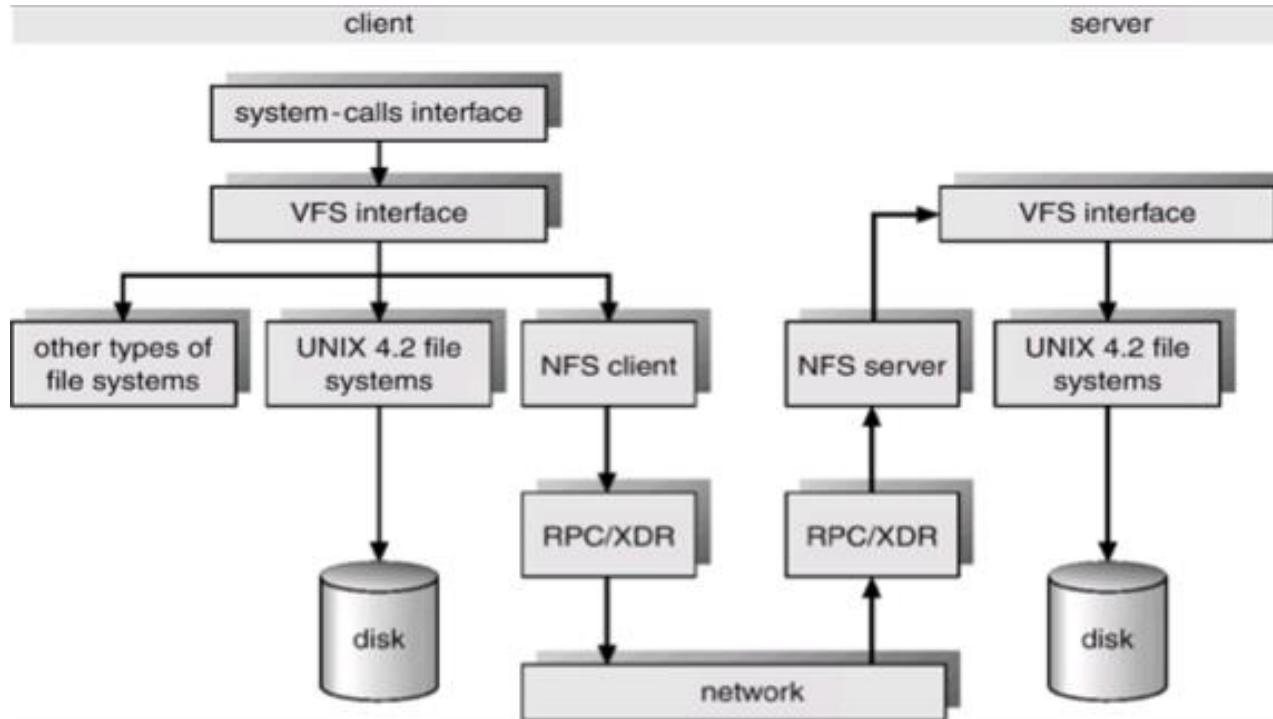
Clients, servers, and storage are dispersed across machines.

Configuration and implementation may vary:

- Servers may run on dedicated machines, OR
- Servers and clients can be on the same machines.



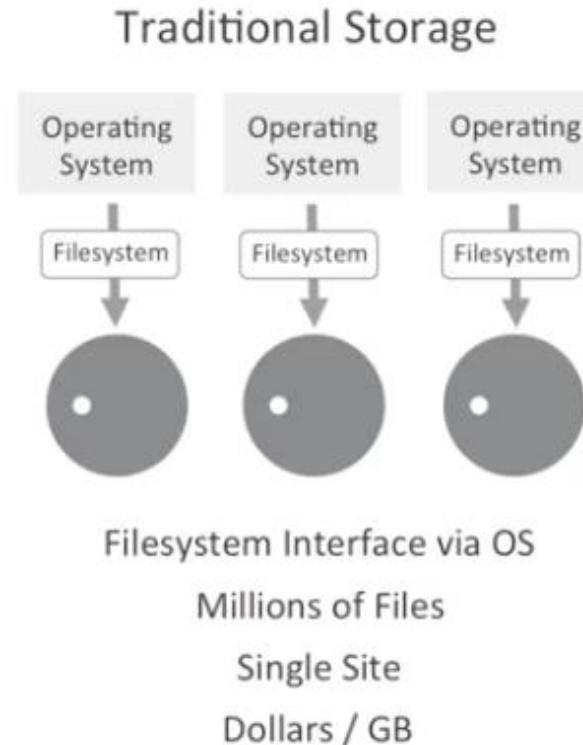
Network filesystem (NFS)



Limitations of traditional Storage

- Handle increasing number of files and users
- Growth over geographic and administrative areas
- Growth of storage space
- No central naming service
- No centralised locking
- No central file store

Applications must be aware of volumes/data location
Data is structured and isolated in Filesystems and volumes



Google File system

Motivation

One single distributed file system

Store big data reliably

Allow parallel processing of big data



Assumptions

Inexpensive components that often fail

Large files (million of files 100+MB)

Large streaming reads and small random reads (500Mb/s
read/write load)

Large sequential writes

Multiple users append to the same file

High bandwidth is more important than low latency.

GFS interface

No common standard like POSIX.

Provides familiar file system interface:

Create, Delete, Open, Close, Read, Write

Snapshot: low cost copy of a whole file with copy-on-write operation

Record append: Atomic append operation

GFS Design Overview

Files split in fixed size chunks of 64 MByte

Chunks stored on chunk servers

BACKUPS ARE UNPRACTICABLE
BUT ↗

Chunks replicated on multiple chunk servers

→ NOT REALLY A BACKUP
BUT IT'S A WAY TO MAKE
THE DATA RELAYABLE

GFS master manages name space

Clients interact with master to get chunk handles

Clients interact with chunk servers for reads and writes

No explicit caching → IMPOSSIBLE

GFS Design Overview

Master server: Single master - Keep metadata - accept requests on metadata - Most management activities

Chunk servers: Multiple - Keep chunks of data- Accept requests on chunk data

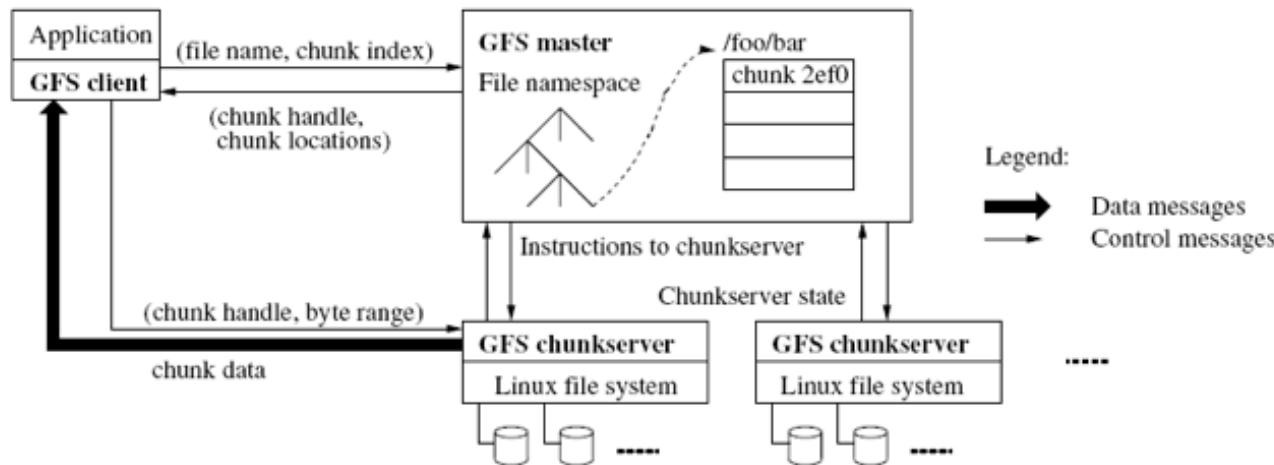
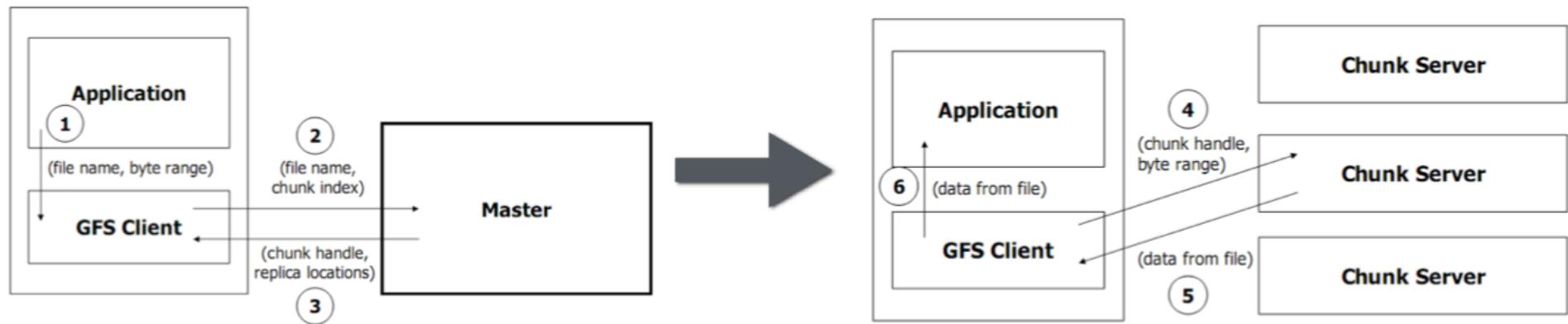


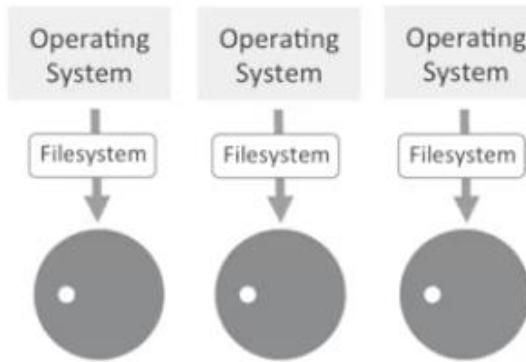
Figure 1: GFS Architecture

GFS example: read files



Object storage

Traditional Storage



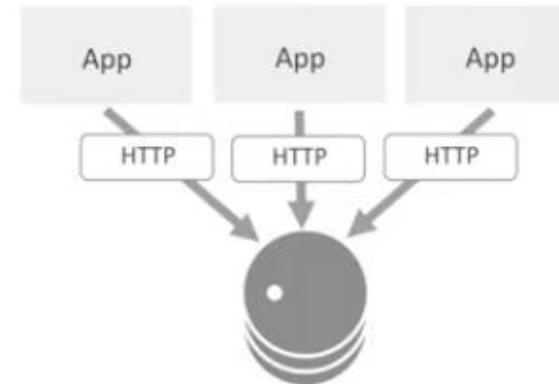
Filesystem Interface via OS

Millions of Files

Single Site

Dollars / GB

Object Storage



HTTP API directly accessible by Applications

Billions of Files

Geo-Distribution

Cents / GB

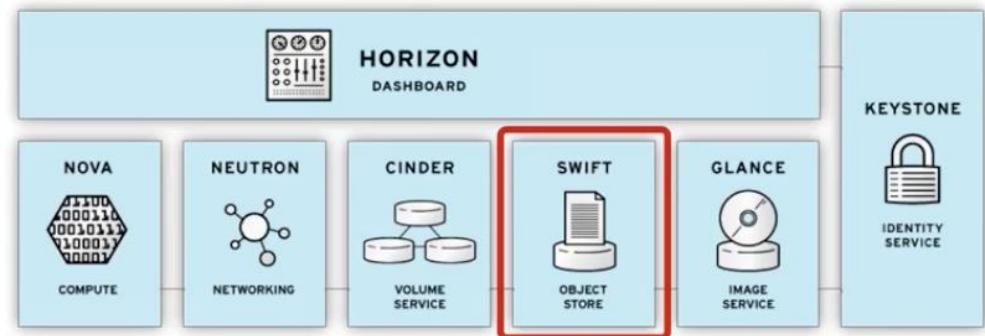
Example of Object storage

Public cloud services

- Amazon S3
- Google Storage (not Google Drive!!!)
- HP Cloud Object Storage
- etc...

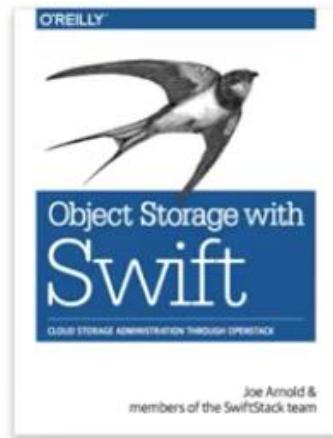
Object Storage Systems

- OpenStack Object Storage System (swift)
- CHEP

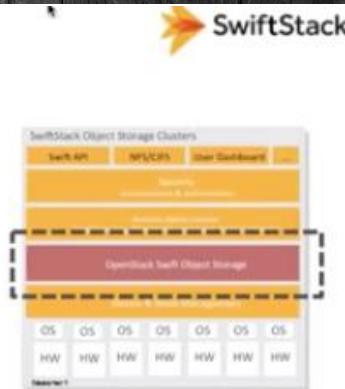


SWIFT

OpenStack Swift



- OpenStack Swift – Worlds most popular Object Storage System
- Powers extremely large storage clouds
- 100% Open-source (Apache 2)
- Rich ecosystem of tools and applications
- SwiftStack is core developer and leads project
- +180 developers today. Contributors include:



What is an Object

An object is a logical unit of storage

- ID (Identification)
- Application data
- Metadata which includes block allocation and length
- Attributes that are accessible by users

Objects have file-like methods
`open`, `close`, `read`, `write`

Object = File + Metadata

What is metadata

Describes the object

- Helps you to find the right one
- Tells you what it is
- The specifications
- Used where and when
- Access permissions

Any and all objects

- Different attributes per object
- Add attributes later

Name
utf8Name
type
size
hashScheme
hash
retention
retentionString
retentionClass
changeTimeString
ingestTime
ingestTimeString
accessTime
accessTimeString
updateTime
updateTimeString
uid
gid
permissions
hold
shred
dpl
index
replicated
customMetadata
acl
namespace
objectPath
owner

Astronomical Data Analisys

Use Case:

Cost effective On-line archiving enriched with metadata.

Challenge:

- Astronomical FITS files are self descriptive
- Content needs to be available for image reduction
- Content must be content of a permanent archival record securely stored

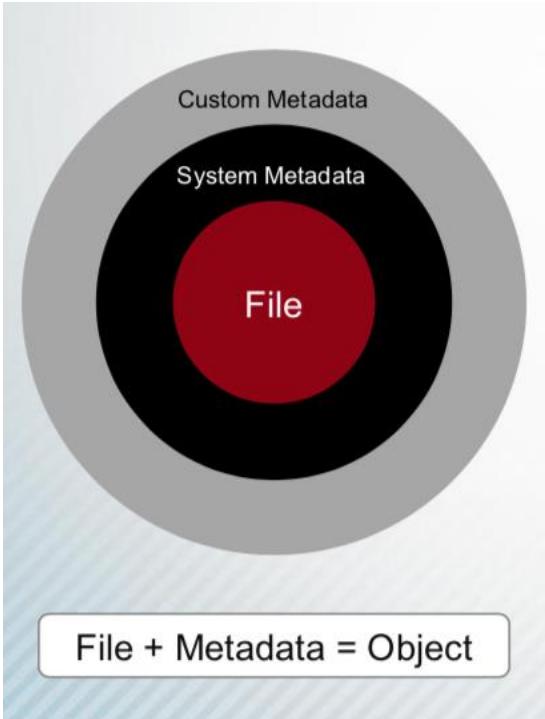
FITS files

FITS. *Flexible Image Transport System (**FITS**) is an open standard defining a digital **file format** useful for storage, transmission and processing of data: formatted as N-dimensional arrays (for example a 2D image), or tables. **FITS** is the most commonly used digital **file format** in astronomy.*

Data Reduction in python example

https://github.com/gtaffoni/Learn-Python/blob/master/Lectures/PythonLecture05-Astronomy_Data_Reduction.ipynb

Objects and Metadata



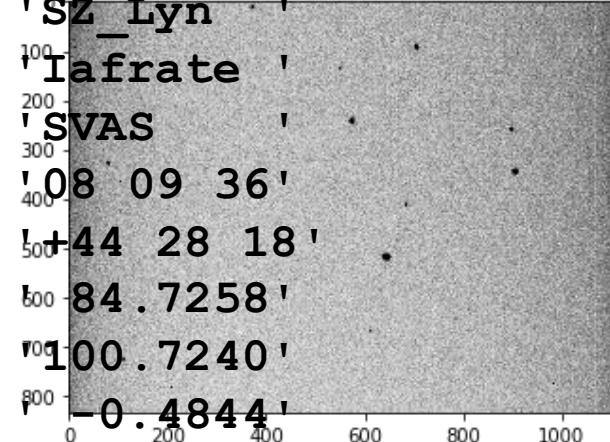
File class = image

Filename = image.FITS

Created = Oct 3, 2018 13:29:59

**Last modified: Oct 10, 2018
11:00:00**

OBJECT = 'SZ_Lyn '
OBSERVER= 'Iafrate '
TELESCOP= 'SVAS '
OBJCTRA = '08 09 36'
OBJCTDEC= '+44 28 18'
OBJCTALT= '84.7258'
OBJCTAZ = '100.7240'
OBJCTHA = '-0.4844'



Mining Metadata

Metadata:

- stored with the object in extend attributes
- can be changed, removed or added (later)
- can be indexed

Identify sets of related objects based on system and custom metadata

Understand the object store – gather object or content metrics on sets of objects based on metadata

- Content discovery

- Return all files owned by "Bob Smith" created after "3/15/2011" with unexpired retention
- Return objects with a specific retention class defined
- Return objects under retention hold, for mitigation purposes

- Metrics gathering

- What is the size distribution of files in the system?
- What percentage of my files are Word documents?
- Which owners have generated the most content?

Access an object

Objects have a 64 bit unique ID

Objects live on a flat namespace

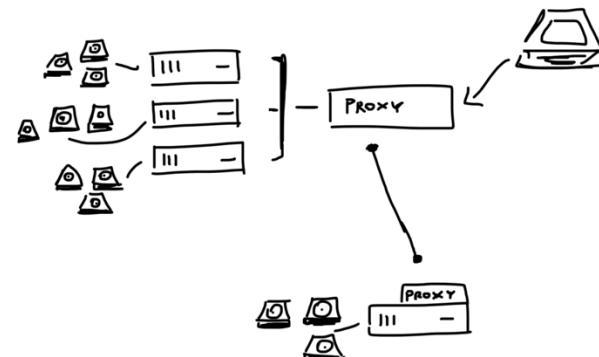
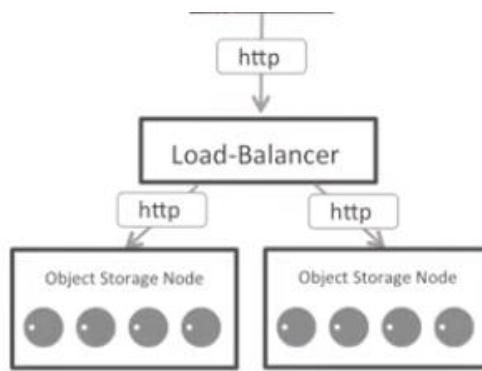
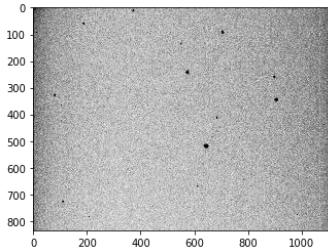
Requests for storage are made with HTTP using RESTfull APIs (or SWIFT client)

Three primary components of the request:

1. HTTP Verb (PUT, GET, DELETE, etc.)
2. Authentication Information
3. SURL
4. Metadata (*Optional*)

HTTP API and SWIFT Architecture

PUT `http://xxx.xxx.xxx.xxx/v1/AUTH_application1/pic/image.fits`



Manage Object with HTTP



- Storage area within the cluster
- NOT a user account
- Each account has its own URL
Swift is multi-tenant

- Namespaces used to group objects within an account
- Containers are unlimited
- Like folders, but can't nest them

- File + Metadata
- Each object is addressed as a URL
- Users name the object
- Objects are not organized based on hierarchy
- Instead, object names may contain "/", so pseudo-nested directories are possible.

Method	Description
GET	Downloads an object with its metadata
PUT	Creates new object with specified data content and metadata
COPY	Copies an object to another object
DELETE	Deletes an object
HEAD	Retrieves object metadata
POST	Creates or updates object metadata

There are also API operations for Accounts and Containers.

| Object Storage Usecase

Video streaming and processing

Map-Reduce

Astronomical data analysis



Web Content Store



Large-File Workflows



Document Sync & Share



Backups and Disaster Recovery



Private Cloud

Data Cloud Examples

Cloud Storage solution based on Nextcloud

The screenshot shows the official website for Nextcloud. At the top, there's a dark blue header with the Nextcloud logo, a navigation bar with links for Products, Solutions, Resources, Community, About, Instant trial, a Download button, and a language selector for EN. Below the header, the main content area has a large blue background. On the left, there's a white sidebar with the text "Regain control over your data" and "Remote collaboration made easy". It lists three benefits with checkmarks: "On-premises or cloud", "Customizable and scalable", and "No data leaks to third parties". Below this are two buttons: "Try Nextcloud now" and "Download". The main content area features a large image of a tablet displaying a video conference with six participants. To the right of the tablet is a screenshot of a desktop interface showing a file browser on the left and a messaging or collaboration tool on the right.

Cloud storage capabilities

Cloud File Storage

- User can upload, download and delete files with full permission on their private directory

User Management

- Groups are needed to manage permissions of multiple users at once
- Login, Logout and Signup must be available

Easy to use Web GUI

- A nice web interface that can be easily used on every device

Security

- Server-Side Encryption, Encrypted communication between client and server, Brute Force attack mitigation Tailored security policies for the OS, Password structure policies

Scalability

- Easily scalable system
- Load Balancing

Multiple access protocols

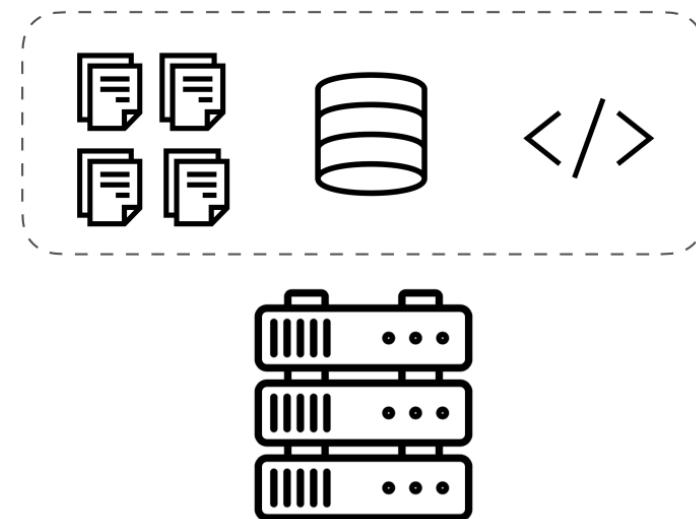
- Native client
- Object Storage

Architecture

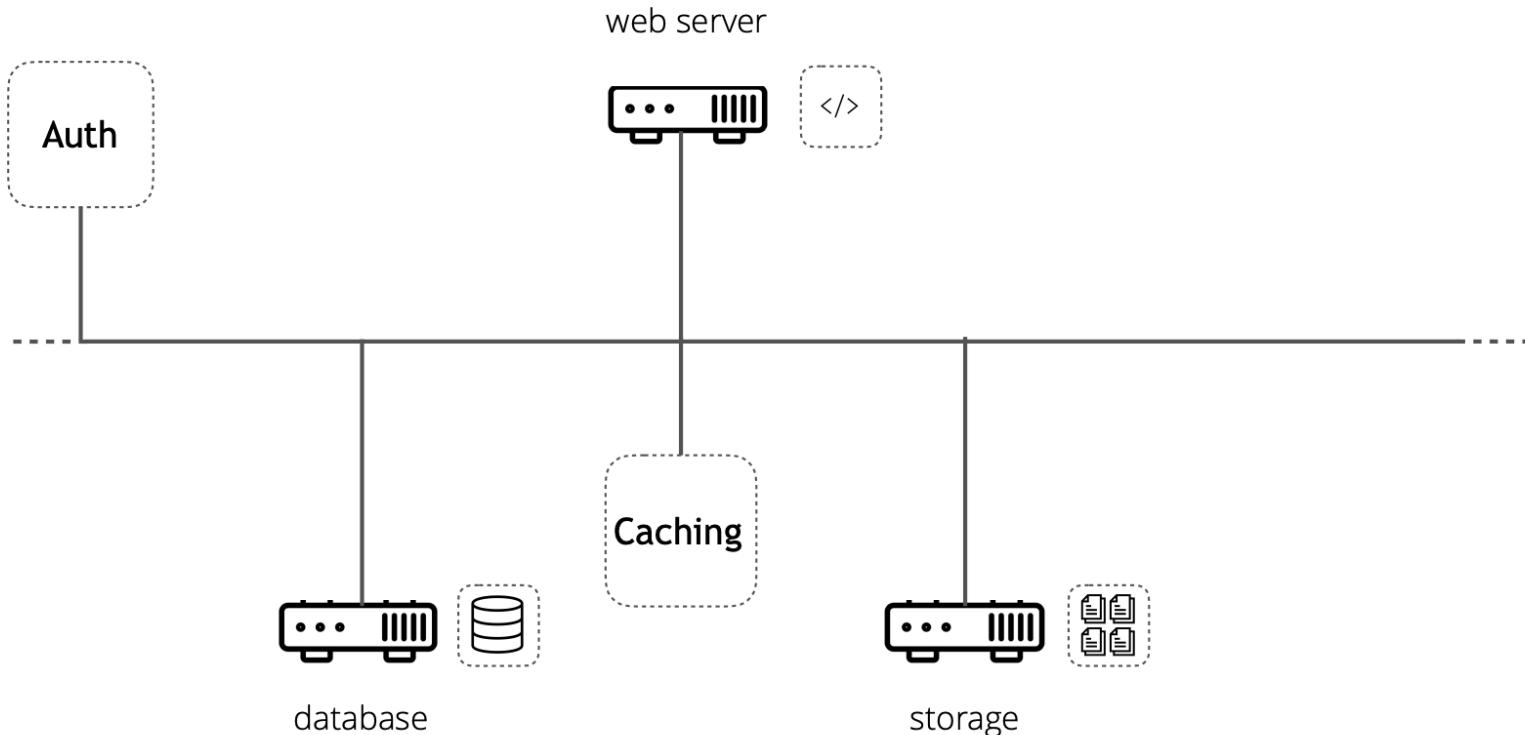
Server & app code • PHP runtime • web server

MySQL/MariaDB or PostgreSQL • stores everything except the binary files • users, shares • settings, ...

user uploaded files • versions • trash bin
• previews

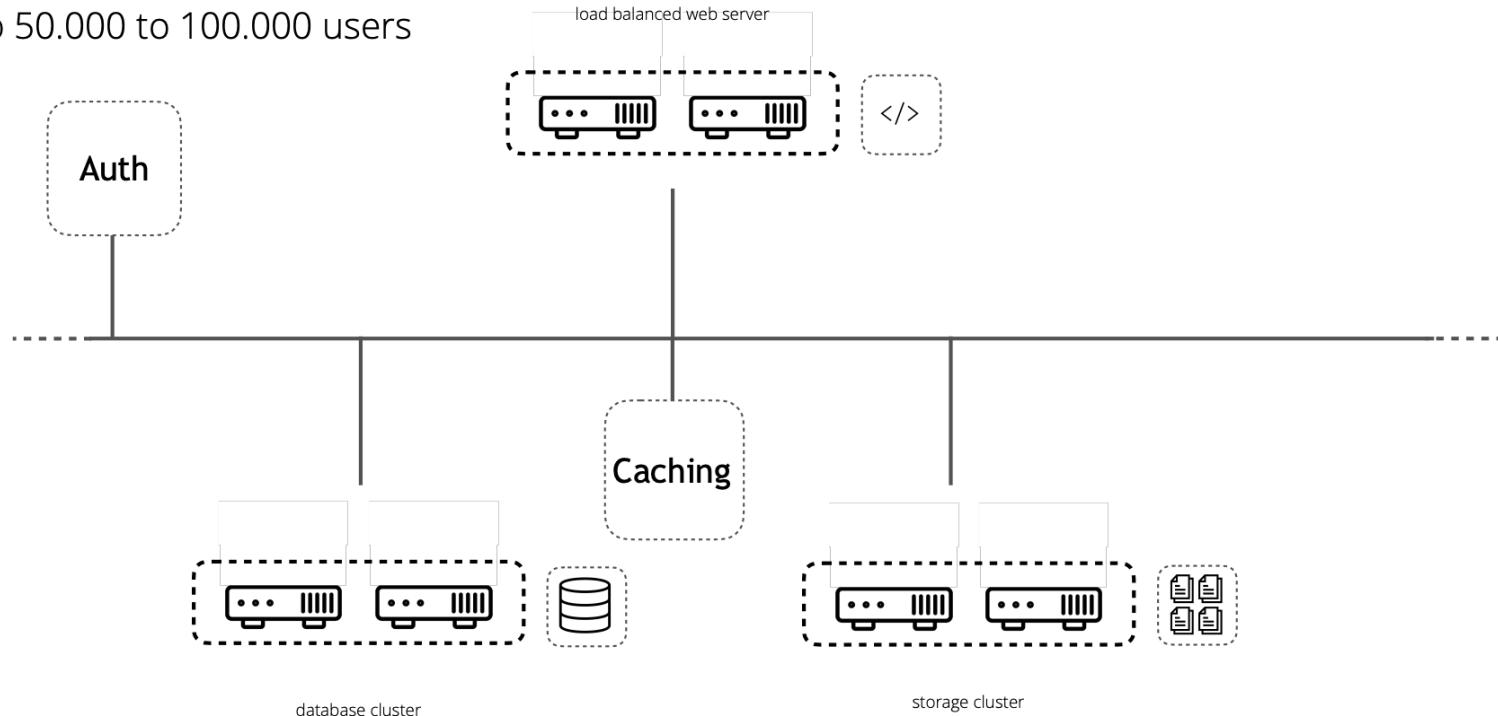


NC Deployment

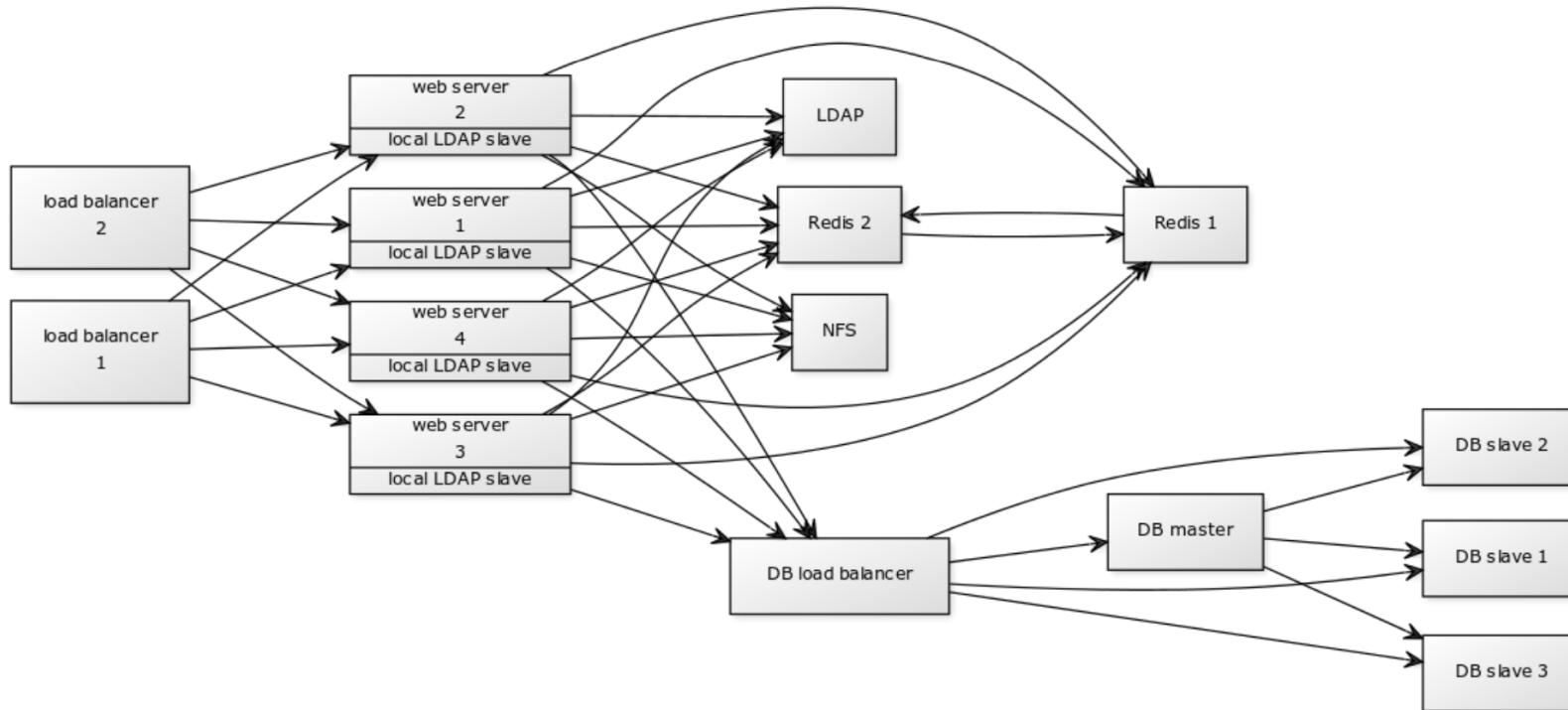


Scaling up...

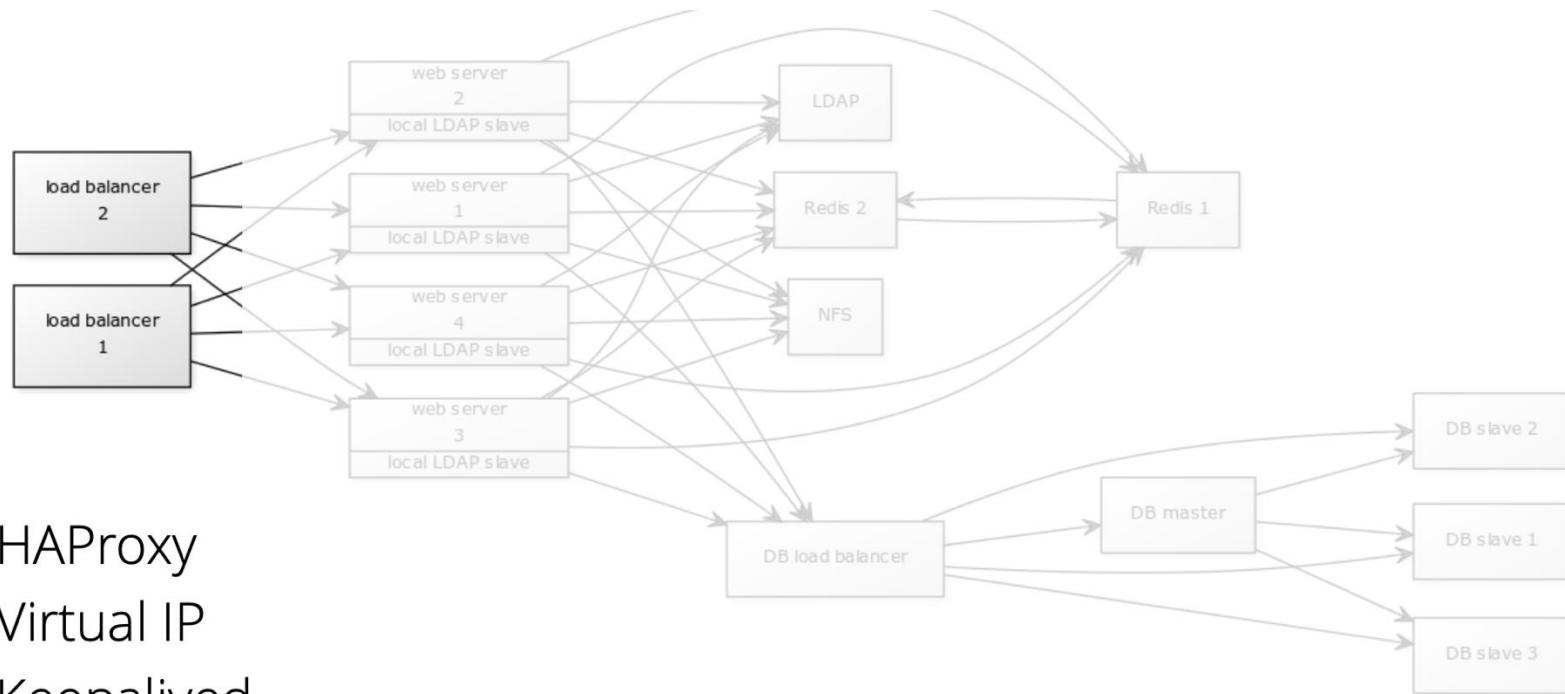
up to 50.000 to 100.000 users



Production deployment

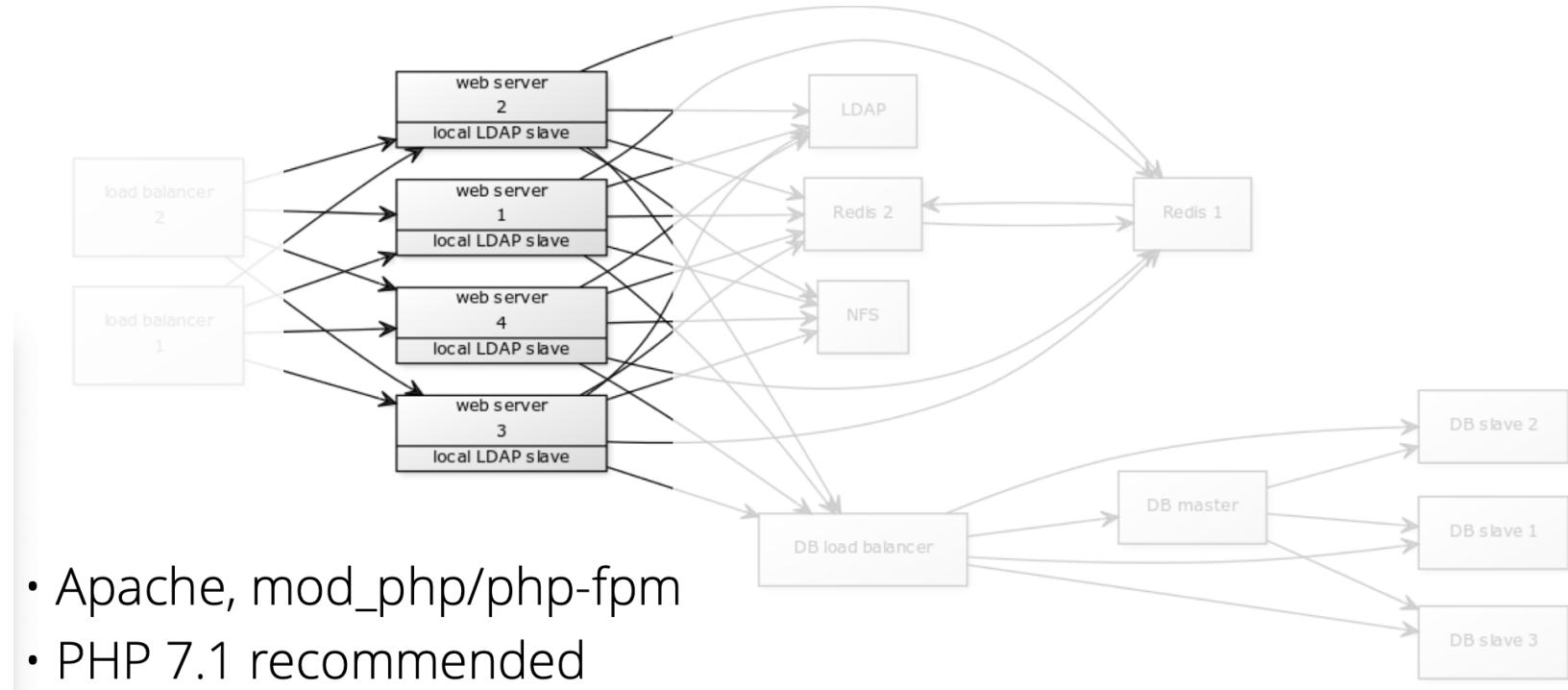


Production deployment



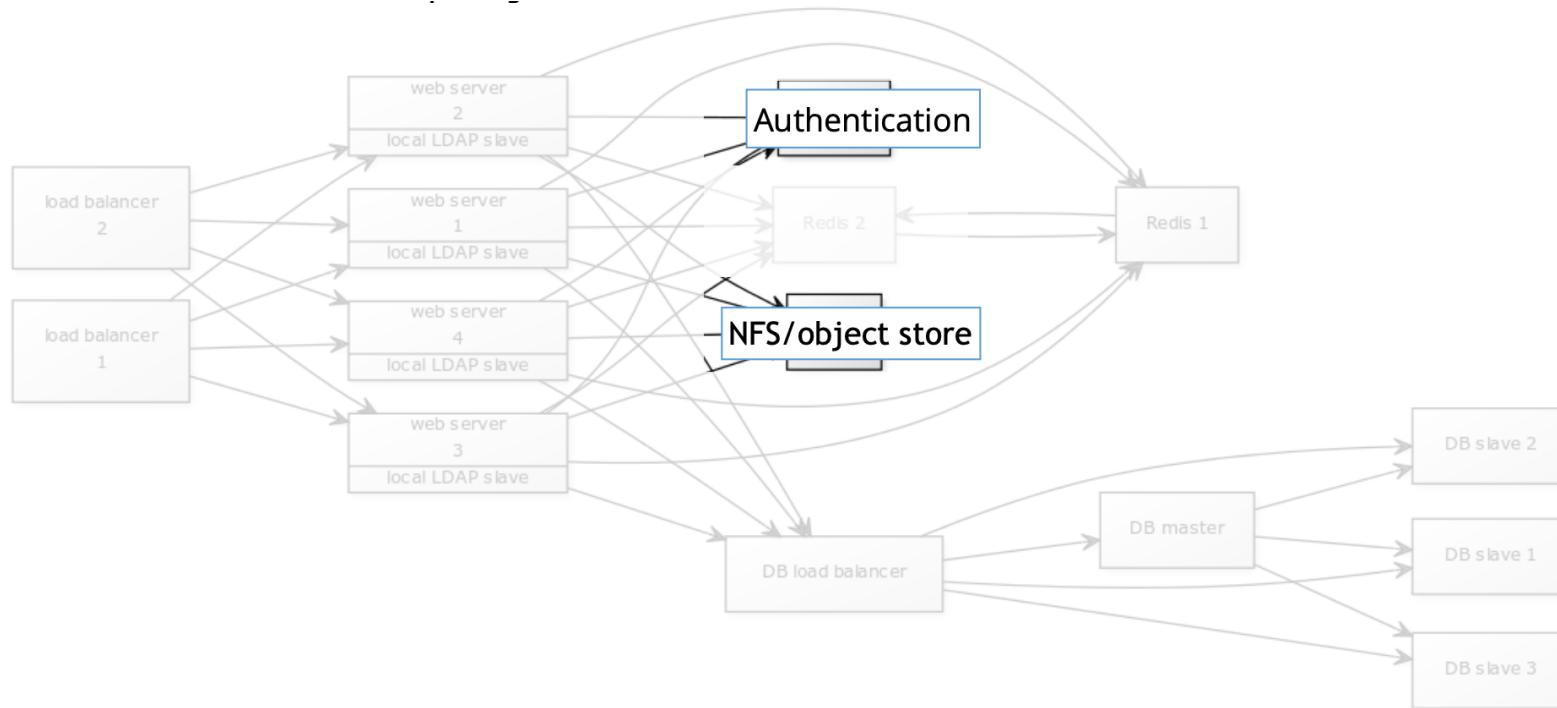
- HAProxy
- Virtual IP
- Keepalived

Production deployment

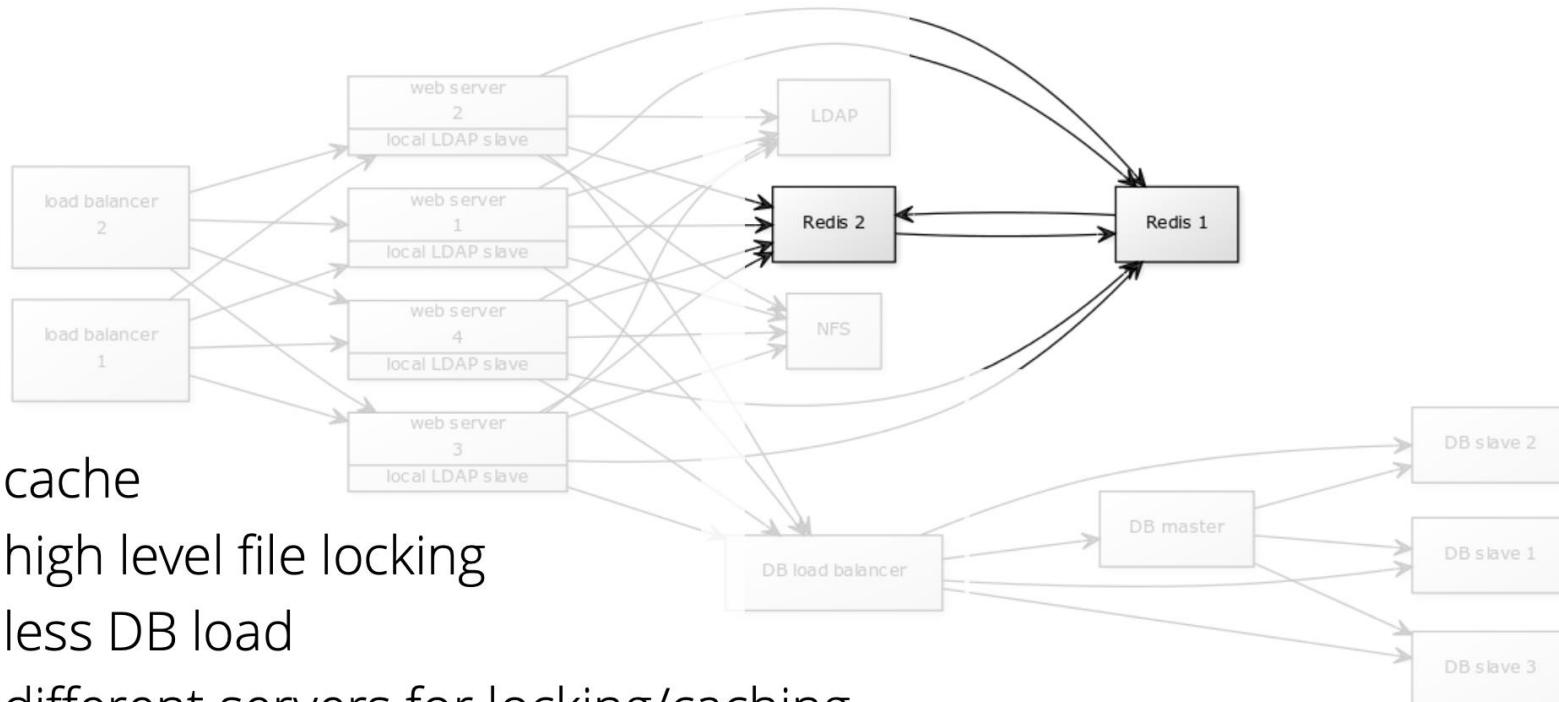


- Apache, mod_php/php-fpm
- PHP 7.1 recommended

Production deployment

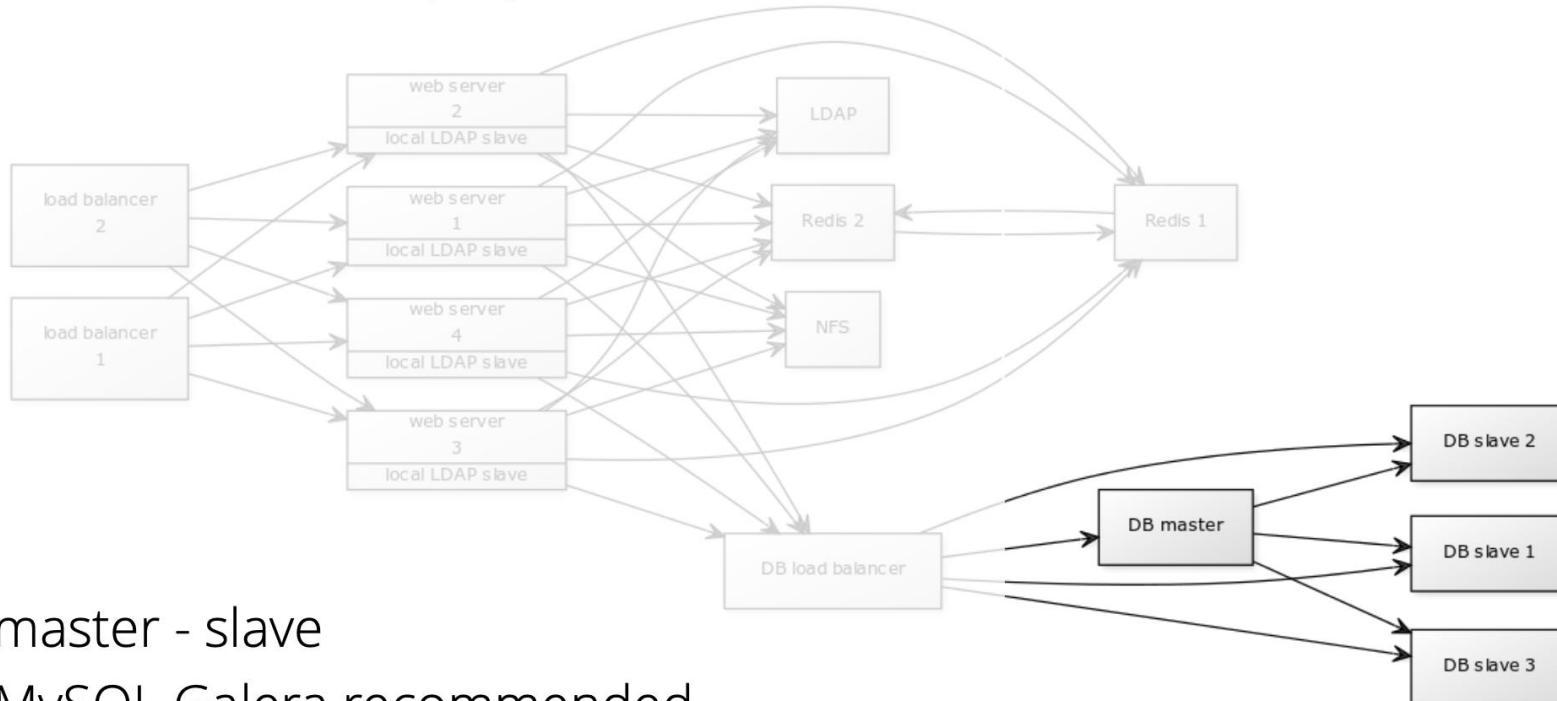


Production deployment



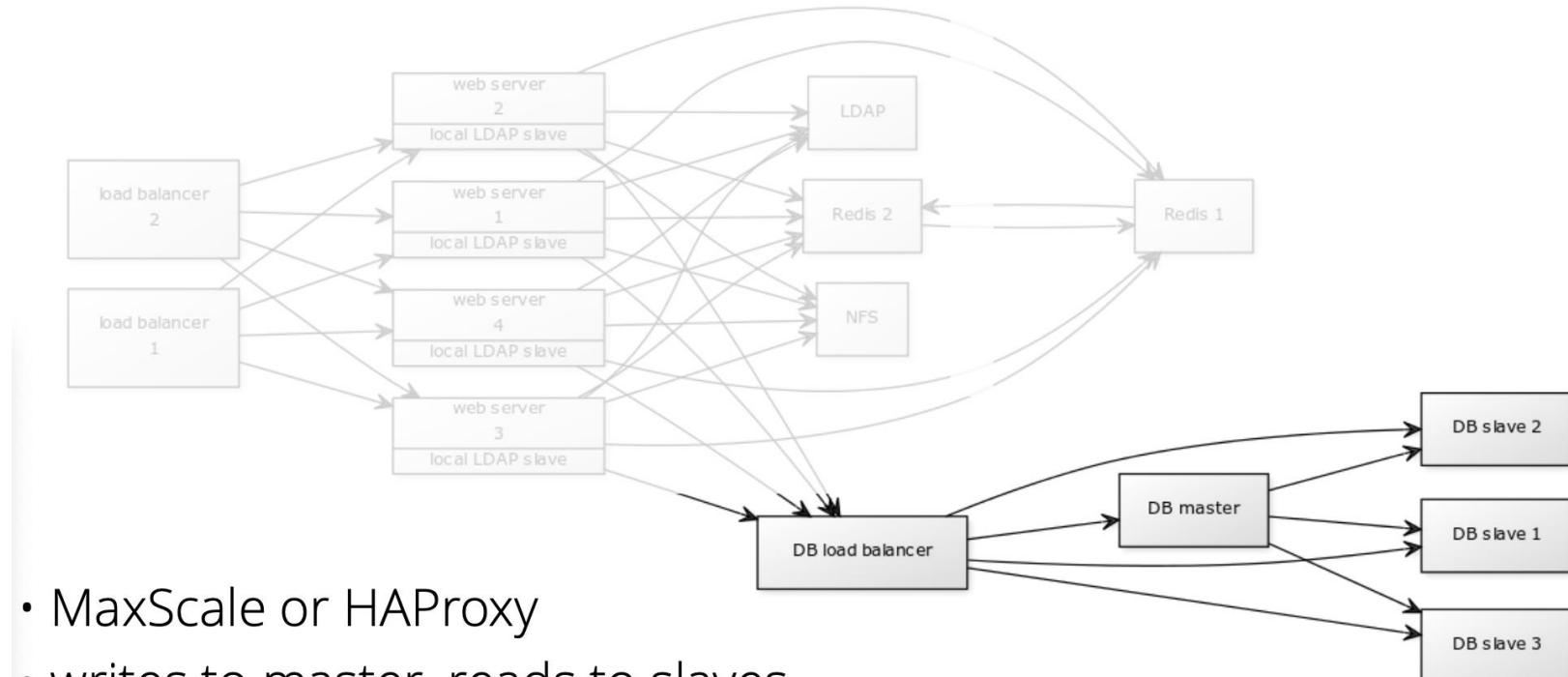
- cache
- high level file locking
- less DB load
- different servers for locking/caching

Production deployment



- master - slave
- MySQL Galera recommended

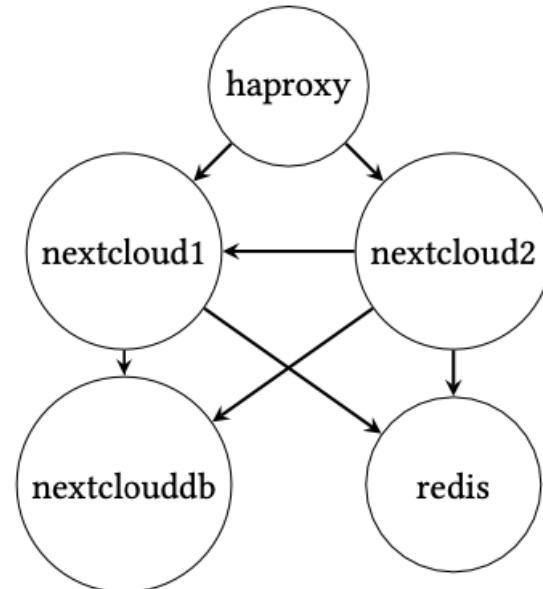
Production deployment



Container deployment

version: '3'

```
services:  
  haproxy:  
    image: haproxy:1ts-alpine3.19  
    restart: unless-stopped  
    networks:  
      - custom_bridge  
    ports:  
      - 80:80  
    volumes:  
      - ./haproxy.cfg:/usr/local/etc/haproxy/haproxy.cfg:ro  
    depends_on:  
      - nextcloud1  
      - nextcloud2  
  deploy:  
    resources:  
      limits:  
        memory: 256MB  
        cpus: "0.5"
```

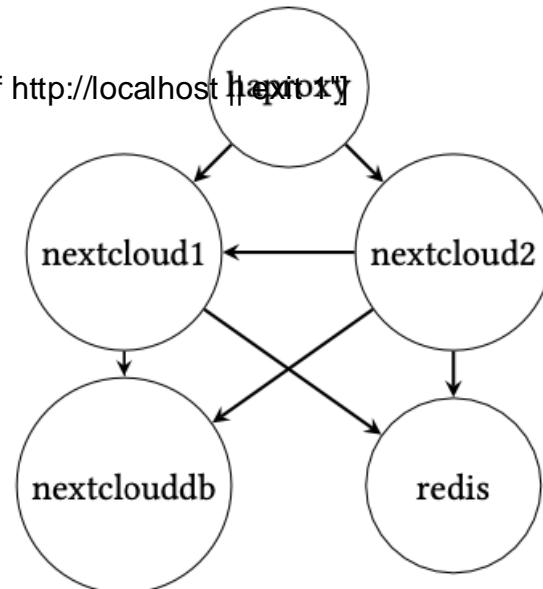


Container deployment

nextcloud1:

```
container_name: nextcloud1
image: nextcloud:latest
restart: unless-stopped
depends_on:
  nextclouddb:
    condition: service_healthy
  redis1:
    condition: service_started
env_file:
  - env_files/.env
  - env_files/nextcloud_1.env
networks:
  - custom_bridge
ports:
  - 8081:80
volumes:
  - ./container_data/nextcloud/html:/var/www/html
  - ./container_data/nextcloud/custom_apps:/var/www/html/custom_apps
  - ./container_data/nextcloud/config:/var/www/html/config
  - ./container_data/nextcloud/data:/var/www/html/data
  - ./install_nextcloud.sh:/install_nextcloud.sh
```

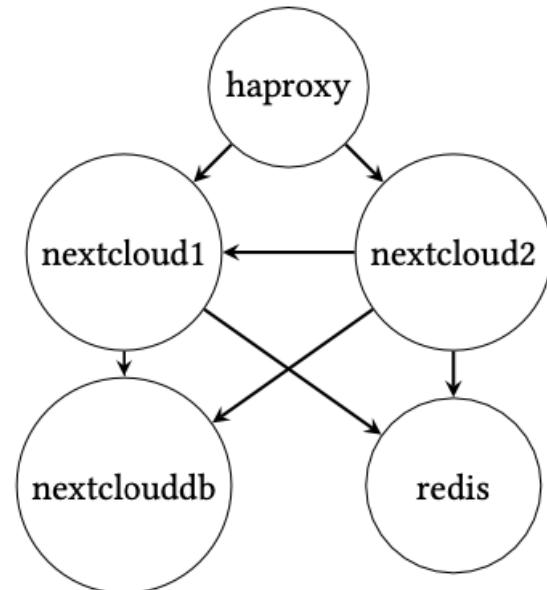
```
healthcheck:
  interval: 20s
  retries: 3
  test: ["CMD-SHELL", "curl -f http://localhost"]
deploy:
  resources:
    limits:
      memory: 2GB
      cpus: "2"
```



Container deployment

nextclouddb:

```
image: mariadb:latest
restart: unless-stopped
networks:
- custom_bridge
volumes:
- ./container_data/nextclouddb:/var/lib/mysql
env_file:
- env_files/.env
- env_files/nextcloud_db.env
healthcheck:
interval: 20s
retries: 3
test:
[
  "CMD",
  "healthcheck.sh",
  "--su-mysql",
  "--connect",
  "--innodb_initialized"
]
```



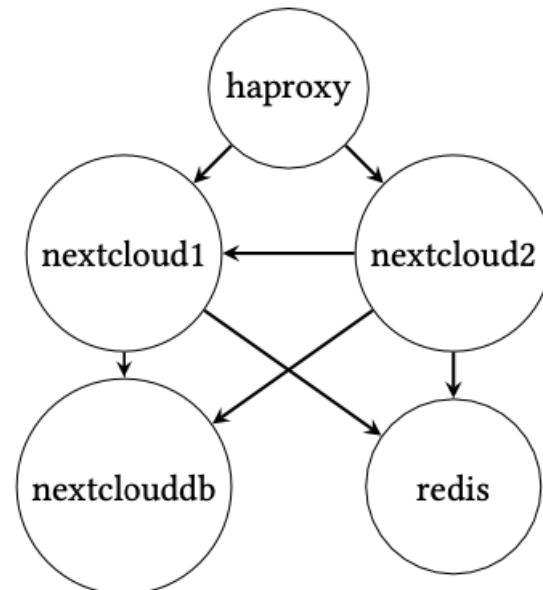
Container deployment

redis1:

```
image: redis:alpine
volumes:
- ./container_data/redis:/data
env_file:
- env_files/.env
networks:
- custom_bridge
deploy:
resources:
limits:
memory: 512MB
cpus: "0.5"
```

networks:

```
custom_bridge:
driver: bridge
```

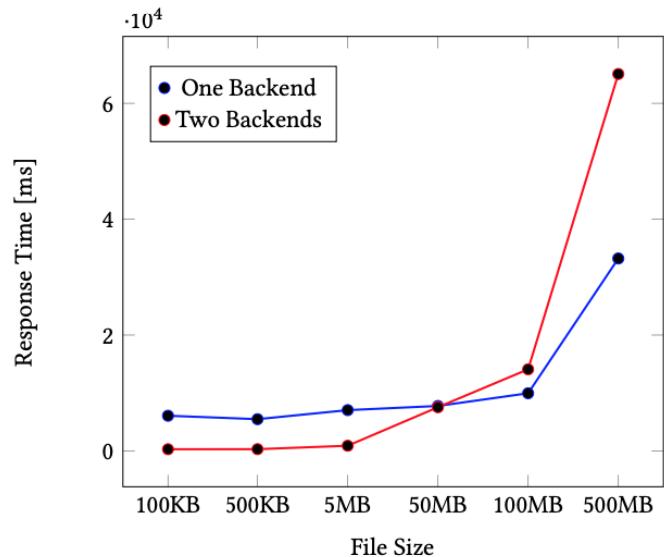
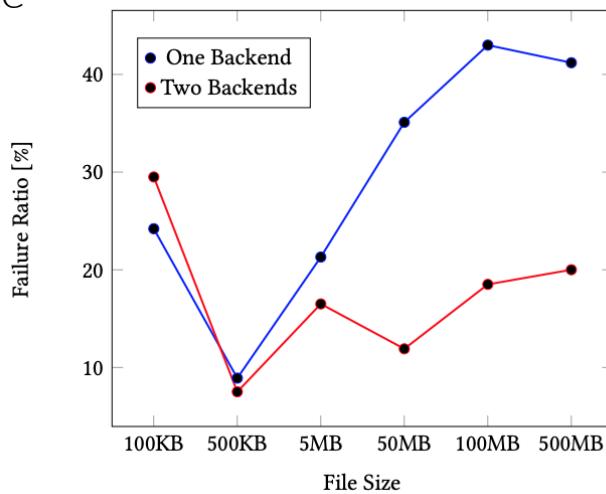


Benchmarking the system

Identify the capabilities of the server(s)

Test service:

- User access/creation
- File upload download
- File read/write



Monitoring

- **Docker stats:** command-line tool that provides real-time resource usage statistics for Docker containers
- **Portainer:** a management tool for Docker environments. Easy to setup
- **Grafana:** analytics and visualization platform that allows users to query, visualize, and monitor data from multiple sources in real-time. Requires configuration.





Cloud Security

Outline

Introduction

Infrastructure security

Data security

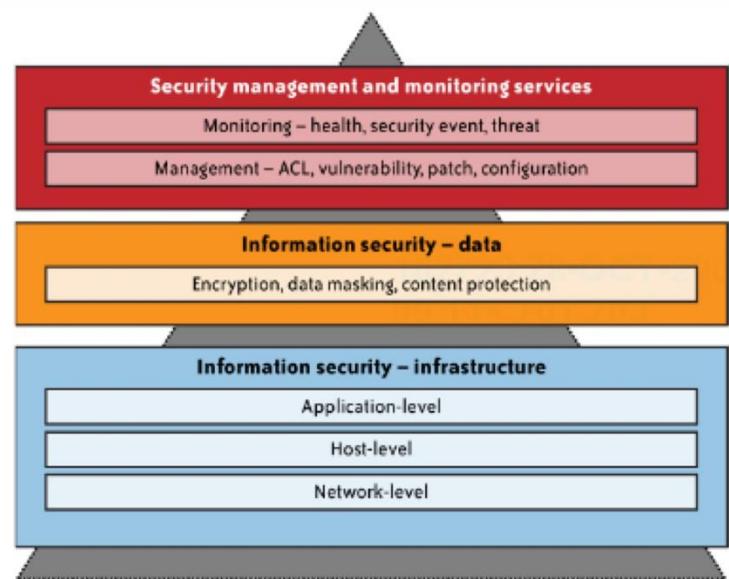
Identity and access management



Introduction

What should I do when I use the cloud, and I want some privacy?

Many security problems in non-cloud environment are still applicable!



General considerations

Security boundary separating the client's and vendor's responsibilities.

Different types of cloud computing service models provide different levels of security

Focus on public clouds

Different levels:

- Network level
- Host level
- Application level

| What is necessary on a cloud

Auditing

Data integrity

e-Discovery for legal compliance Privacy

Recovery

Regulatory compliance

Before approaching the Cloud

- Determine which resources (data, services, or applications) you are planning to move to the cloud.
- Determine the sensitivity of the resource to risk.
- Risks that need to be evaluated are loss of privacy, unauthorized access by others, loss of data, and interruptions in availability.
- Determine the risk associated with the cloud type for a resource
- Consider the cloud service model that you will be using.
- If you have selected a particular cloud service provider, you need to evaluate its system to understand how data is transferred, where it is stored, and how to move data both in and out of the cloud.

Security Boundary

Who has responsibility for a particular security mechanism?

Where is the boundary between the responsibility of the service provider and the responsibility of the customer?

Cloud Security Alliance: CSA is an industry working group that studies security issues in cloud computing and offers recommendations to its members.

| Multi- tenancy security

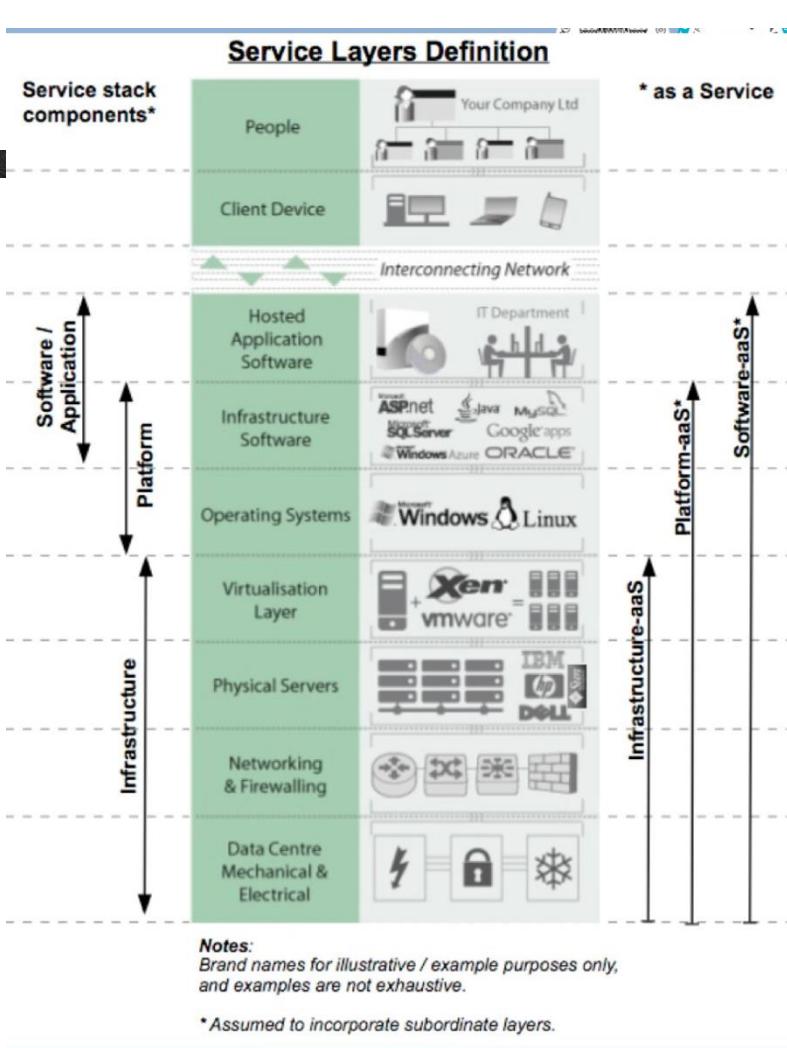
Different customers must be isolated,

Customer data segmented,

Customer services identified and isolated.

Multi project approach (e.g. OpenStack tenants)

XaaS Boundaries



The platform security level is reduced to the security level of the most vulnerable application running on the platform. (isolate with VM and containers)

Network Level Security

Confidentiality and integrity of data-in-transit

Amazon had security bugs with digital signature on SimpleDB, EC2, and SQS accesses (in 2008)

Less or no system logging /monitoring

- Only cloud provider has this capability
- Thus, difficult to trace attacks

Reassigned IP address

- Expose services unexpectedly
- Spammers using EC2 are difficult to identify

Availability of cloud resources

Some factors, such as DNS, controlled by the cloud provider.

Physically separated tiers become logically separated E.g., 3 tier web applications

Host Level

- Hypervisor security
 - “zero-day vulnerability” in VM, if the attacker controls hypervisor
- Virtual machine security
 - ssh private keys (if mode is not appropriately set)
 - VM images (especially private VMs)
 - Vulnerable Services

VM/Containers security advantages

Simple way to implement resource management policies

Improved intrusion prevention and detection.

Secure logging and intrusion protection.

More efficient and flexible software testing

| VM/containers security risks

- Explosion of number of VMs
- Snapshots (roll back to state that can be exploited)
- Shared Images
 - backdoors and leftover credentials,
 - unsolicited connections,
 - Malware

22% of the scanned Linux AMIs contained credentials allowing an intruder to remotely log into the system.

OS Security

Implement minimal security on: Access control, authentication usage, and cryptographic usage policies are all elements of mandatory OS security.

Applications with special privileges that perform security- related functions are called trusted applications. Such applications should only be allowed the lowest level of privileges required to perform their functions.

An OS poorly isolates one application from another, and once an application is compromised, the entire physical platform and all applications running on it can be affected.

Data Security

Identify the **security boundary** separating the client's and vendor's responsibilities
Determine the **sensitivity** of the data to risk

Data should be transferred and stored in an **encrypted format**.

Separate clients from direct access to shared cloud storage.

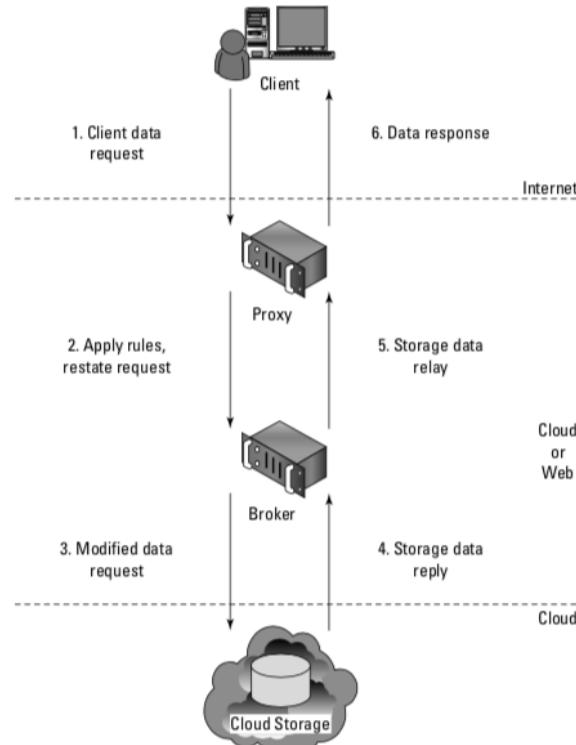
These are the key mechanisms for protecting data mechanisms:

- Access control
- Auditing
- Authentication
- Authorizations

Data Segregation and Isolation

Isolate data from direct client access
creating a layered access to the data.

Data segregation based on tenants



Encryption

Most cloud service providers store data in an encrypted form (e.g. Amazon S3 256-bit Advanced Encryption Standard) on server side or client side.

Some example of java code here:

<https://docs.aws.amazon.com/AmazonS3/latest/dev/UsingClientSideEncryption.html>

Problems:

- a problem with encrypted data may result with data that may not be recoverable.
- it does nothing to prevent data loss: keep your keys!!!!

that's all, have fun

“So long
and thanks
for all the fish”