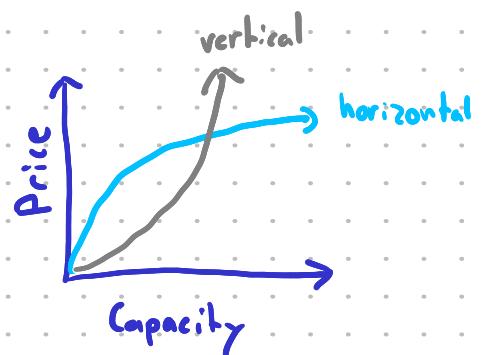


# Distributed System 25

## 1) Scaling

Vertical Scaling: more memory, faster CPUs

- + lower cost at massive scale
- + easier fault tolerance
- + higher availability
- adoption of code needed
- more complex  $\Rightarrow$  more components



Horizontal Scaling: more machines

- + current trend
- + lower cost at smaller scale
- + no SW adaption needed
- HW limits for scaling
- More HW  $\Rightarrow$  bigger risk of failure

## 2) Location

Move the server closer to the user  $\rightarrow$  e.g. CDN

Latency: RTT from src to dst and back.

Bound to speed of light  $\rightarrow$  there is always latency

Fault Tolerance  $\rightarrow$  how a (computer) system maintains proper operation in case of errors/faults

All Systems will fail eventually! Not if but when.

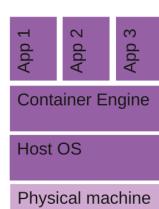
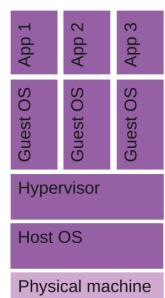
Bitflips can happen and will cause wrong results

- $\hookrightarrow$  e.g. from solar rays  $\rightarrow$  happens on earth too, rover on Mars  $\phi$  280 bflips/day
- $\hookrightarrow$  error correcting memory exists

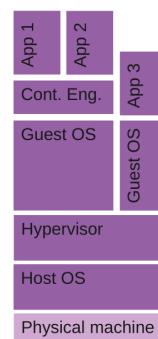
Distributed Systems provide a solution by adding redundancy, but add complexity.

## 3) Containers / Virtual Machines

Virtualization - Visualization



• Virtual machines



• Container

(Docker) container

- + small size
- + quick spinning up
- +/- shared memory
- +/- process based isolation (same CPU)

VMs

- + app can access all OS resources
- + live migrations
- +/- Full isolation

Monorepo / Multirepo → skipped

#### 4) Docker & Docker Debugging.

→ build your own container

• Every Line of the dockerfile is one Layer

- a Docker image contains all layers, they can be unpacked individually
- a layer is cached and not rebuilt unless changed (using docker build)

```
Dockerfile 499 B
1 ## 
2 ## Build
3 ##
4 FROM golang:1.24-alpine AS build
5 ARG VERSION
6 WORKDIR /app
7
8 COPY go.mod go.sum ./
9 RUN go mod download
10
11 COPY ./
12 RUN CGO_ENABLED=0 go build \
13     -ldflags "-s -w -X 'ost.ch/edudn/cmd.version=${VERSION}'" \
14     -trimpath \
15     -o ./bin/edudn \
16     main.go
17
18 RUN apk add --no-cache ca-certificates
19
20 ##
21 ## Prepare Container
22 ##
23 FROM scratch
24
25 COPY --from=build /app/bin/edudn /edudn
26 COPY --from=build /etc/ssl/certs /etc/ssl/certs
27
28 EXPOSE 8080
29 ENTRYPOINT ["/edudn"]
```

Overlay FS does the management of files across docker layers. Data is kept in lower layers but only marked as deleted.

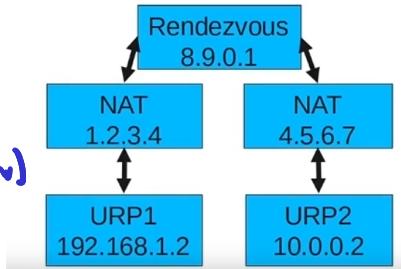
Cgroups is a linux tool to limit CPU usage/memory/disk/... can be used in docker run with flags e.g.: docker run --cpuset-cpus=0 --cpu-shares=20  
↳ kernel 0 ↳ max 20% usage

Network namespaces: provides network isolation for docker containers. Use NAT to communicate out of the network  
↳ abstracted from Docker

Hole punching: Method of connecting two unreachable peers via Nat and rendezvous point.

- Entries in NAT tables make connection work

Docker compose: deploy multiple servers (front-end, back-end etc.)



Docker Security:

- keep images small and up to date
- use tiny images → alpine
- don't expose Docker daemon socket
- don't use root
- limit resources, a garbage collector might not be aggressive enough, so when process gets too big it is killed by docker  
↳ limit resources to make garbage collector docker aware

- Filesystem read only if possible

Suitable log level: NOTICE → production

INFO:

DEBUG:

TRACE:

WARN:

ERROR:

FATAL: → doomsday

log events into correct levels

## 5) Load Balancing

Distribution of workloads across multiple computing resources (horizontal scaling)

- + high availability
- + add / remove containers

Hardware LB: F5 / Cisco → only used if you have your Datacenter \$\$\$

Software LB: L2/L3 : Seesaw

L4: Envoy, trino, Nginx, Traefik, ...

L7: Neutrino, traefik, caddy (use for next project)

Cloud based:

- : pay for use
- : all big internet companies have a variation of this
- : useful if one has no control over datacenter

Types of load balancing:

- DNS load balancing:
  - : Multiple IP addresses returned, client can choose what IP to use
    - : Depending on location, different IPs returned
    - : split horizon DNS
- Anycast loadbalancing: same IP for multiple servers, AS needed, bgp returns fastest path

Load balancing algorithms:

- Round robin
- weighted round robin → some servers prioritized
- least connections
- random (Traefik)

Health checks → load balancer realizes when a server is down

L7 load balancing more resource intense as TLS / HTTPS is terminated

Stickyiness: can be achieved with cookie, session is connected to same container

Normally session lb is done. Packet LB is tricky

## 6) Web architecture

Server side rendering:

- html / css / js rendered by server and sent to client
- rendered by server side code in Java, Go, ...
- + Search engine optimization as crawler gets the whole page

Single Page application:

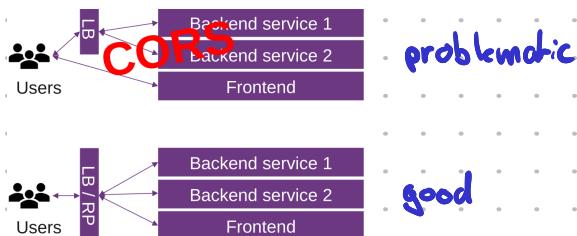
- interaction with single web page
- site updates dynamically with js
- data received via API communication

clicks caught by js and site reloaded  $\rightarrow$  site not loaded again  
Vue, React, ...

CORS: Cross Origin Resource Sharing

For security reasons, browsers restrict access from different origins

$\Rightarrow$  use reverse proxy



State of the art: initial HTML with SSR, then dynamically updated with APIs / js

7) Categorization What is a distributed system?

A group of computers that appear to be working as one to the user

Tightly coupled: access to a common memory

loosely coupled: no access to a common memory

homogeneous: all processors are the same type (cluster of identical servers)

heterogeneous: contains processors of different types

OS  
Rechenleistung

small scale: web app + DB

large scale: >2 machines

Decentralized  $\rightarrow$  not owned by one actor  $\rightarrow$  chains

$\curvearrowright$  can't be all three at the same time neither CP or AP

CAP theorem: Consistency  $\rightarrow$  every node has the same state

Availability  $\rightarrow$  all nodes respond

Partition tolerant  $\rightarrow$  system continues to be consistent even when the network partitions

More ways of classifying Distributed systems:

Architecture: client-server, peer to peer, ...

Communication model: synchronous or asynchronous (pub/sub events)

! No single universally applicable categorization of distributed system !

Distributed System	Decentralized System
<ul style="list-style-type: none"><li>• so many</li><li>• 1 responsible org</li></ul> <p>Secure environment</p> <p>Consistent Hashing</p> <p>Master node</p> <p>No NAT issues as network under control</p> <p>Consistency managed by leader</p>	<ul style="list-style-type: none"><li>• BitTorrent / Blockchain</li><li>• N responsible orgs</li></ul> <p>Hostile environment</p> <p>Consistent Hashing (distributed hash table)</p> <p>Flooding Broadcasting</p> <p>NAT issues → peer to peer connection difficult</p> <p>Weak consistency → proof of work / state ↳ longest chain wins</p>

Transparency in Dsys: user should not be aware of / care about →

- location
- access → resource in single/uniform way
- replication
- concurrent → don't be aware of other users
- Failure → not be aware of recovery mechanisms

8 Fallacies (wrong assumptions) programmers like to make

1) network is reliable

2) Latency is zero

3) Bandwidth is infinite

4) network is secure

5) topology doesn't change (request can take different route than reply)

6) There is one admin

7) Transport cost is zero

8) network is homogeneous

## 8) Authentication

Important concepts of authentication:

- CIA triad : Confidentiality → message can't be eavesdropped → use encryption  
Integrity → Message hasn't been changed → Hash functions, digital signatures  
Availability
- Non-repudiation: So both parties involved can't claim a tx (e.g. transaction) hasn't happened
- Identification: claiming to be whom you say you are
- Authentication: using pw or other factors to prove you are
- Authorisation: Permission applied to an account after successful login

Single Factor: Password only

Multi Factor: something you have, something you own, something you are

TOTP: time-based one-time password

### Basic Auth

- Basic auth can happen at the load balancer or in the server itself.
- Basic auth should only be used with https  
↳ base64 is no encryption!

Format username:password → can be encoded as a link: <https://username:pw@url.ch>

Suitable for simple, intern apps but if you can don't use it!

### Digest Auth

- Fires pw in plain text weakness in basic auth (among others)
- Uses Nonce and Hash to prevent replay attacks
  - Pw encrypted using MD5 (or Sha256)
- Browser looks and feels bad and quite outdated too

### Certificate based auth

→ CA must be trusted!

In the lecture, Thomas created a CA and created certs for client and user. After some configuration, he was able to authenticate between client and loadbalancer. Not really used in practice. Used more often to use TLS between LB and back-end/frontend.

Let's encrypt automated cert / https renewal process.

- Domain ownership proven by http challenge to ./well-known/acme-challenge
- if site not public → DNS challenge using TXT record implemented in Caddy
- cert invalid after 90 days → more automatic renewals → can be automated

## Stateful authentication Session based

Login with upn + pw → server creates and stores the session

Session is also returned to client and can be used to authenticate

If server changes → login not possible. Use sticky sessions on LB to always connect to the same server

## Stateless session based authentication used by OAuth

Better suited for distributed systems, as every server can verify the sess

Instead of only the session, the auth server replies with JWT Token.

- JSON-based access tokens

- Header: {"alg": "HS256"}
- Payload: {"sub": "tom", "role": "admin", "exp": 1422779638}

more fields can be added  
if needed

base64 encoded

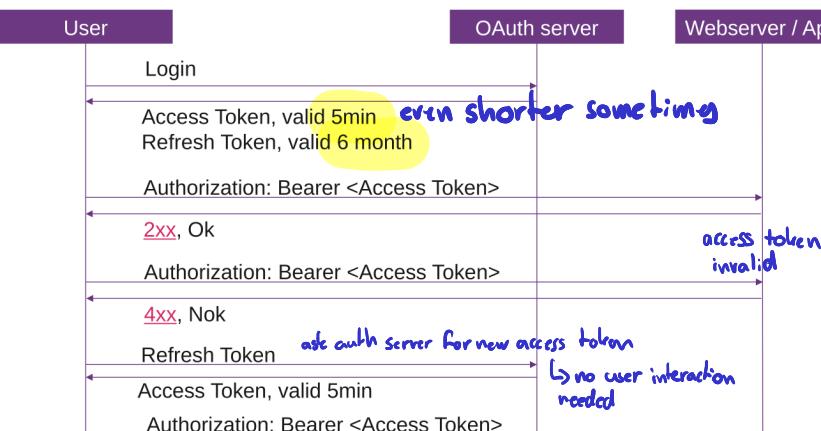
<https://jwt.io/introduction>

not entirely sure how this works

The signature is created with a public/private key pair or a secret (using HMAC)

a webserver can now verify the JWT token by using header, payload and public key.  
If the signature is equal → verification ✓

## Access Token / Refresh Token



- Refresh token only used against bsp
- Longer lifetime

- Access token short lifetime
- usable against all kinds of systems

This enables SSO into "all" apps.

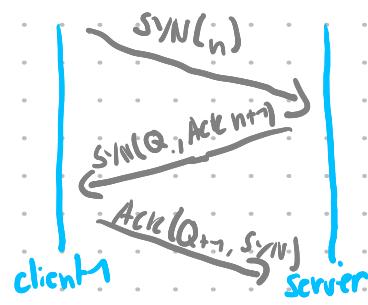
If account disabled → Refresh token disabled and no new access token granted  
But access token still valid until timestamp abgelaufen

↳ can't be changed as it is in header.

## 3) Protocols TCP & Quic

Goes into OSI model, but skipped here

TCP handshake explained  
random numbers used.



Sequence number to order packages and detect missing ones.

TCP does head-of-line blocking.  
One stream does not affect the other. QUIC



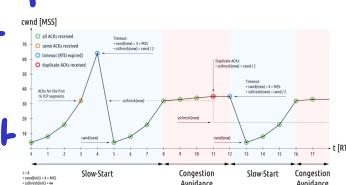
Flow control: Recipient tells the sender how many packages he can handle  
↳ sliding window, adjusted after each packet

congestion control: sender tests the limit of before the traffic overflows

how many packets can be sent in the network

increasing

Starts exponentially → then linearly → repeat



Due to latency, it can take up to a second until data can be sent. DNS query, TCP Handshake + Security = 3 Round Trips → newer protocols HTTP3

↳ Key Exchange

Normally, TLS is also used (basis for https). TLS < 1.3 requires two round trips, with TLS 1.3 it is only one.

While developing locally, RTT is almost 0 it can have a big impact on distributed systems. Thus it is important to take this into consideration.

- Resolve → DNS, TCP Connect → TCP Handshake, HTTPS Handshake → TLS/SSL



From LA to dsl.i.ast.ch took almost 1 sec for initial connection, loading the content was quicker 0.7 sec

## QUIC / HTTP3

Known connections → 0 RTT ✓, UDP based

New connection: 1 RTT for connection and security

Security built in by default  
read blog post about Quic, interesting  
and proto relevant for exam.

## 10) Application protocol

In this lecture it was difficult to figure out what is relevant.

HTTP Hypertext Transfer Protocol

Version 1 introduced in 1989, currently using HTTP2 or HTTP3<sup>↗ with QUIC</sup>  
stateless → server does not keep track of sessions

HTTP stores data in header, most important Status code 200 OK 404 not found  
403 Forbidden ...  
Uses GET, PUT, PATCH, DELETE, ...

Websockets

bi-directional

Websockets are a two way communication channel between client/server

HTTP handshake first, then elevated to websocket

Server side Events (SSE)

One way communication channel from server to client

Implemented in challenge Task :)

automatic reconnection if site is reloaded, uses standard HTTP connections

JSON / REST API

De facto standard, used in many web applications

Human readable json used to transmit data

This has overhead and json parsing is slower compared to other methods

↳ This is why one might develop its own layer 7 protocol.

## Custom protocol

- + control every aspect of it
- more time developing and troubleshooting

## Avro IDL (interface description language)

↳ define the protocol / interface and can be converted into many different languages

You can use binary encoding, which is usually faster

Examples are:

- Avro
- Protobuf → smaller and faster than XML, used by pretty much all of Google internal communication.

## Application Part II

(Binary) Encoding → binary encoding and there are A LOT of options

↳ Eg: Bencoding (used in BitTorrent), UBJSON, Cap'n Proto

- Which encoding to choose depends on the implementation → do Benchmarks

## WebRTC    RTC = Real time communication

A protocol mostly used to enable video calls across browsers. Peer to peer for the most part.

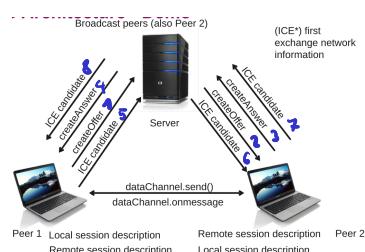
↳ Eliminates Plugins and Flash Applications as supported by all browsers.

- First cross browser call 2013 → standardized in 2021.

The only real downside is that browsers get bloated

WebRTC can have problems with NAT. As a developer there is no need to worry about it, as there are protocols implemented to deal with NAT.

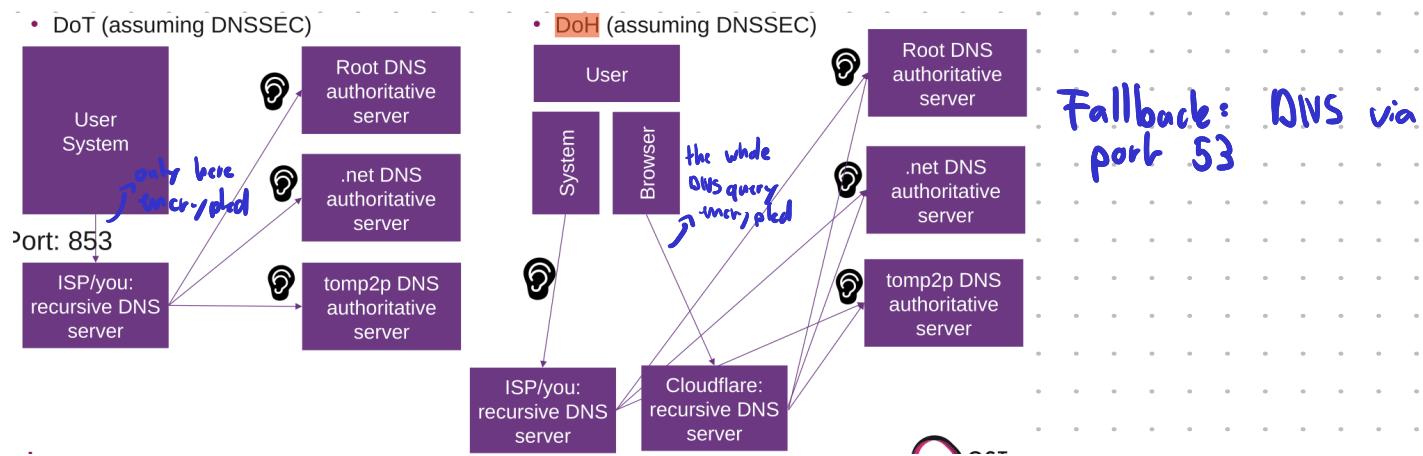
- STUN: help a client find its public IP
- TURN (Traverses using relays around NAT): When P2P isn't possible, a relay server is used where both clients connect to.
- ICE (Internet connectivity establishment): Test multiple ways of establishing a WebRTC connection (p2p preferred) and helps find an optimal way that works.



## DNS

Goes into details that I'll skip, because it is like the 4th time looking at this

DoT = DNS over HTTPS , DoT = DNS over TLS



11) Exam prep → Todo

12) Deployment

**Rolling Deployment:** new version gradually deployed to replace old system  
the whole system is not taken down at once

**Blue - Green Deployment:** one live, develop / test the other then switch

**Canary Release:** new features to a few people, if all goes well extend

**Feature toggle:** finer canary release, enable features for a group of users

**Big Bang:** deploy everything at once

Technical Deployment : Dockerization as the basis

- Ansible → push based via ssh, ssh device list needed
- Plain docker / podman: as in seProj → docker context
- Docker Swarm → has lost the battle against Kubernetes
- Kubernetes - Container orchestration → industry standard
  - define goal (how many containers...) in Yaml and Kubernetes makes sure to deploy it

12) Performance

Idea is to get a feeling how long different operations take.

Network RTT remained stable across the last 20 years, Hardware performance got faster

Branch prediction CPU → example showed that JS code runs faster if the if-clause is predictable. As in always true because array is sorted. → no real use in real life

Latency: time per operation  
Throughput: amount of operations

more compression leads to less data to be sent

Best practice → enable compression, use the default value (lossless compression)

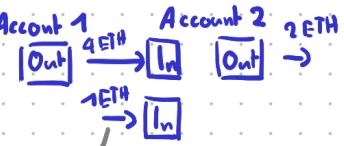
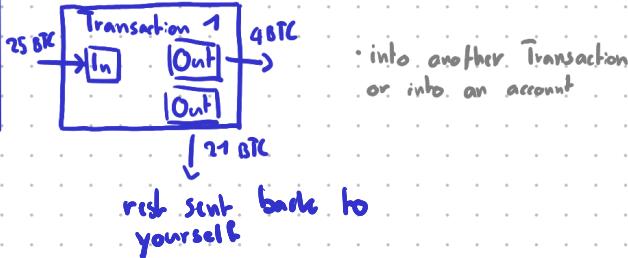
Image compression → lossful compression

### 13) Blockchain

No central entity. Block validation by proof of work

Wallet has private and public key

Transaction broadcasted to all peers and placed in block → verified in around 10 min

Account-based (Ethereum)	UTXO-based (Bitcoin)
<ul style="list-style-type: none"><li>Global state has a list with balance on each account.</li><li>Transaction is valid if sender has enough.</li></ul>  <p>Transactions (from another account)</p>	<p>unspent transaction output</p> <ul style="list-style-type: none"><li>Every referenced input must be valid and not yet spent</li><li>All outputs are always spent</li></ul>  <p>into another Transaction or into an account</p> <p>rest sent back to yourself</p>

Mining: Block SHA256 hashed → miners try to find hashes starting with 0000... by changing Nonce

- Each block has a pointer to the previous one → blockchain
- 51% can take over chain attack