



OST

Eastern Switzerland
University of Applied Sciences

Distributed Systems (DSy)

Introduction

Thomas Bocek

24.02.2022

Learning Goals

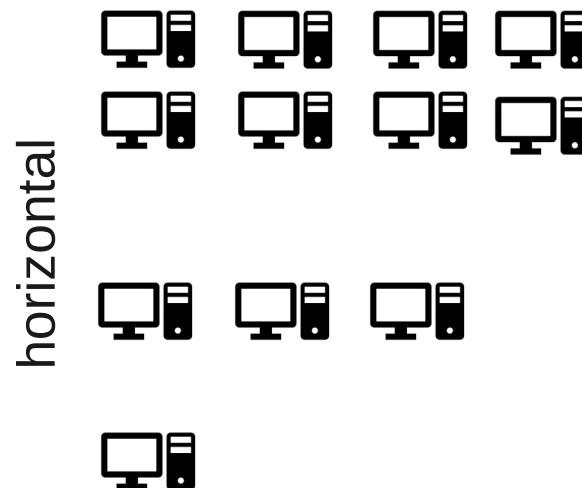
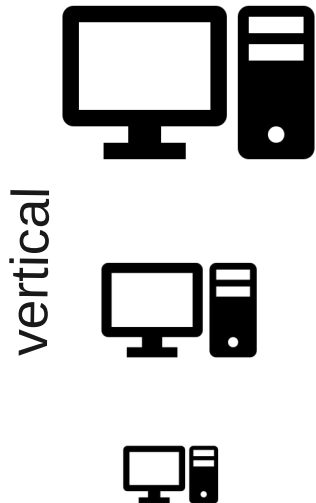
- Distributed systems add complexity. Avoid complexity!
- Why do we need distributed systems?
 - 1) Scaling (if one machine is not enough)
 - 2) Location (to move closer to the user)
 - 3) Fault-tolerance (HW will fail eventually)

Distributed Systems Motivation

- Why Distributed Systems

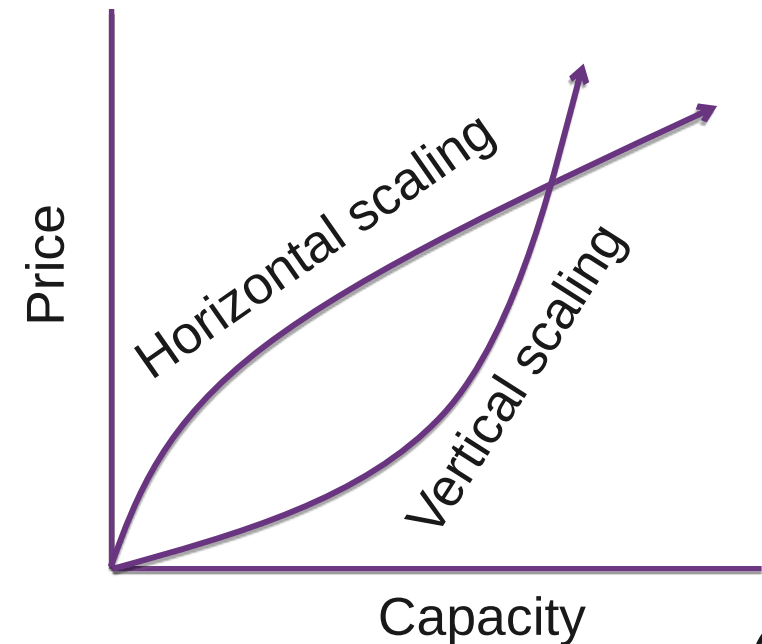
- Scaling

- Vertical (scale up), more memory, faster CPU
 - Horizontal (scale out), more machines
 - Apple has 75'000 Apache Cassandra nodes storing 10 petabytes of data in 2015 [[source](#)]



- Economics

- Initially scaling vertically is cheaper, until you max out HW
 - Current servers are fast: **96cores** ~ 70k TPS



Distributed Systems Motivation

Horizontal Scaling

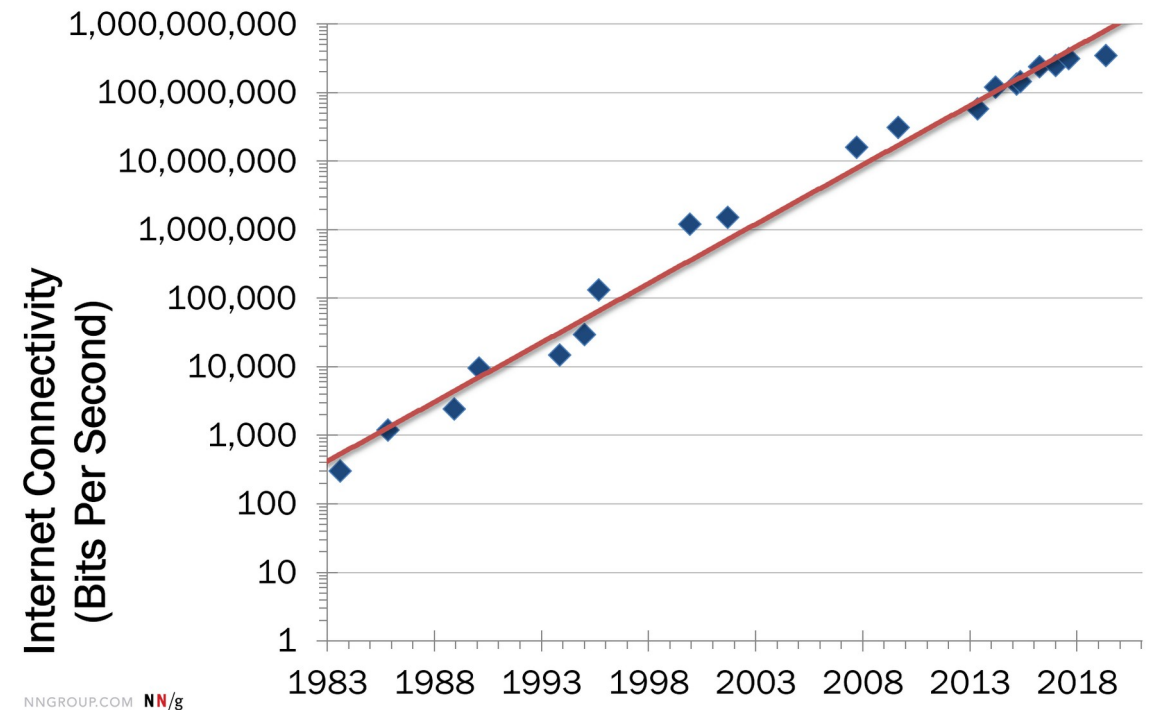
- + Lower cost with massive scale
- + Easier to add fault-tolerance
- + Higher availability
- Adaption of software required
- More complex system, more components involved

Vertical Scaling

- + Lower cost with small scale
- + No adaption of software required
- + Less complexity
- HW limits for scaling
- Risk of HW failure causing outage
- More difficult to add fault-tolerance

Vertical Scaling Performance

- Nielsen's Law: a high-end user's connection speed grows by 50% per year
- **Bandwidth grows slower than computer power**
 - Telecoms companies are conservative
 - Users are reluctant to spend much money on bandwidth
 - The user base is getting broader
- Optimize for bandwidth not for CPU
- **Zmap** complete scan of the IPv4 address space in under 5 minutes
- Init7: **Fiber7-X2** 25/25 Gbit ~65CHF/month



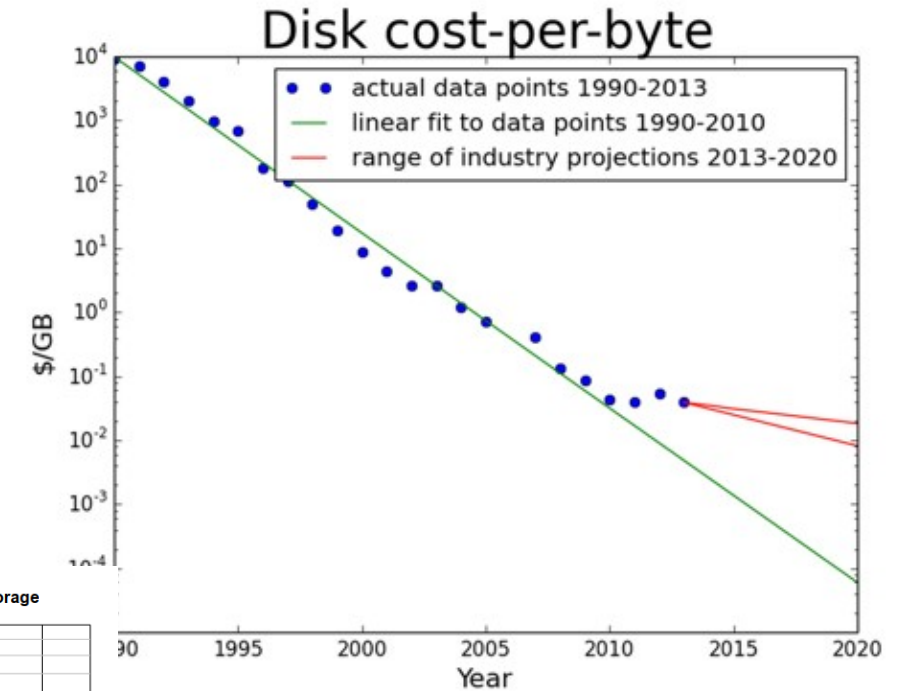
<https://www.nngroup.com/articles/law-of-bandwidth/>

		Annualized Growth Rate	Compound Growth Over 10 Years
Nielsen's law	Internet bandwidth	50%	57×
Moore's law	Computer power	60%	100×

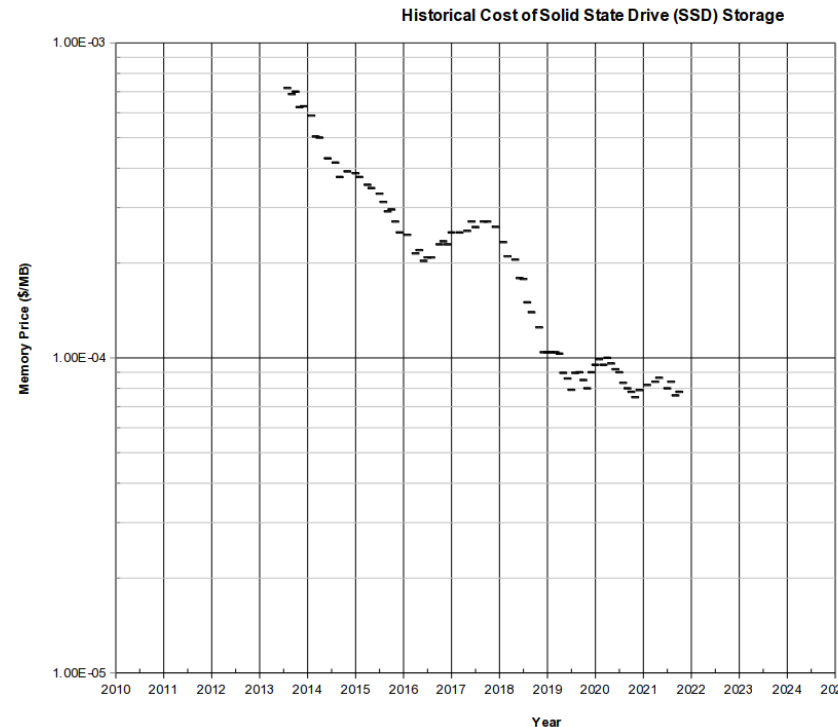
Vertical Scaling Performance

- Kryder's Law: disk density doubling every 13 month
- «Soon hard drives will migrate into phones, still cameras, PDAs, cars and everyday appliances»
<https://www.scientificamerican.com/article/kryders-law/> ,
Aug. 2005

- User behavior changed
 - SSD, speed is important
- Cloud – Dropbox, Spotify
 - Streaming



<http://blog.dshr.org/2016/05/the-future-of-storage.html>

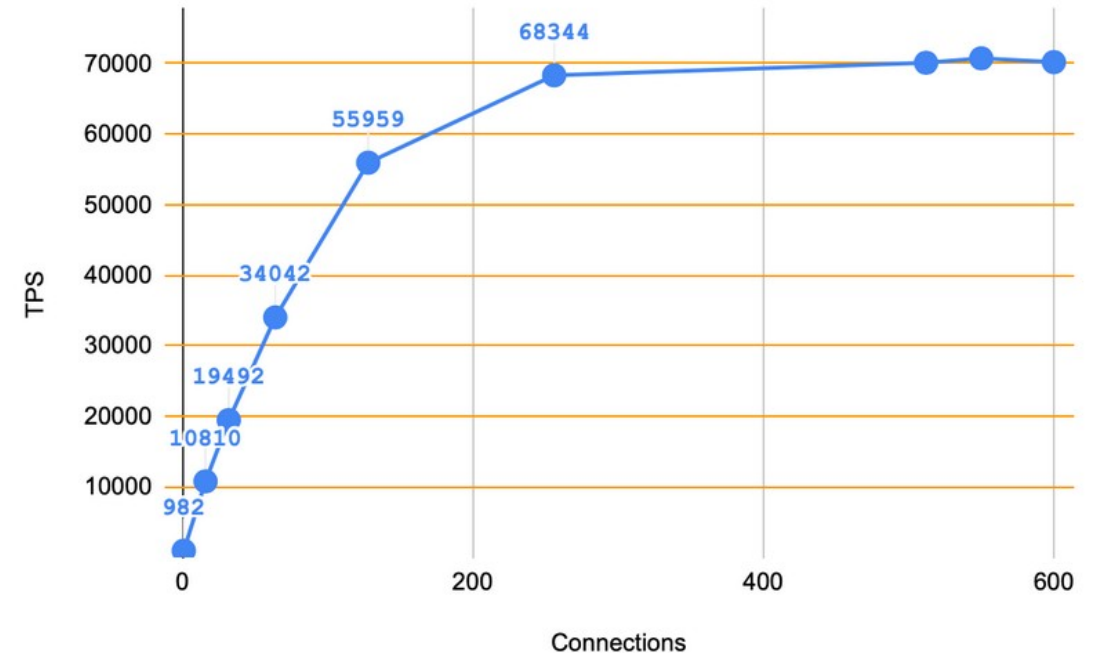


Source: <https://jcmmit.net/flashprice.htm>

Vertical Scaling Performance

- Vertical scaling
 - HW today is fast!
 - Database benchmark with a fast machine in 2020 (96 cores, 384GB RAM, 4 x NVMe SSD)
 - 70k TPS
- Best principle for small and simple applications!
- Simple website with a few DB calls is not HW intensive
 - But: ML, Gaming (**cloud gaming**) are HW intensive

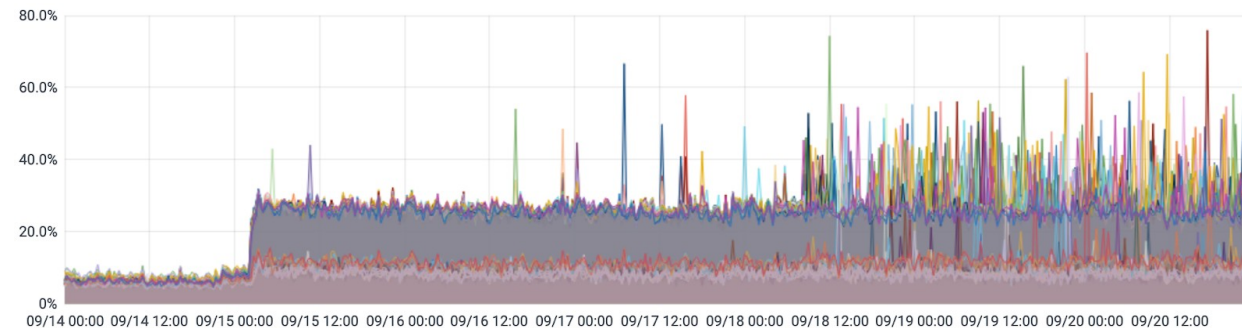
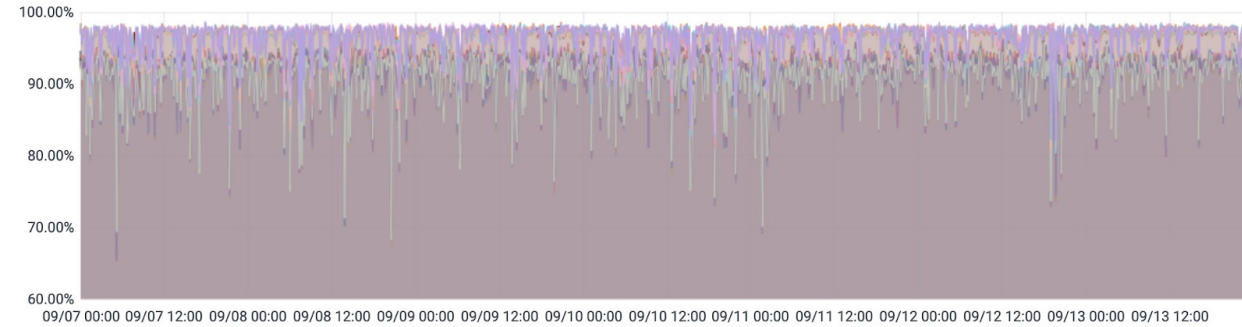
PostgreSQL12: TPS vs. Connections



<https://www.enterprisedb.com/blog/pgbench-performance-benchmark-postgresql-12-and-edb-advanced-server-12>

Vertical Scaling Performance

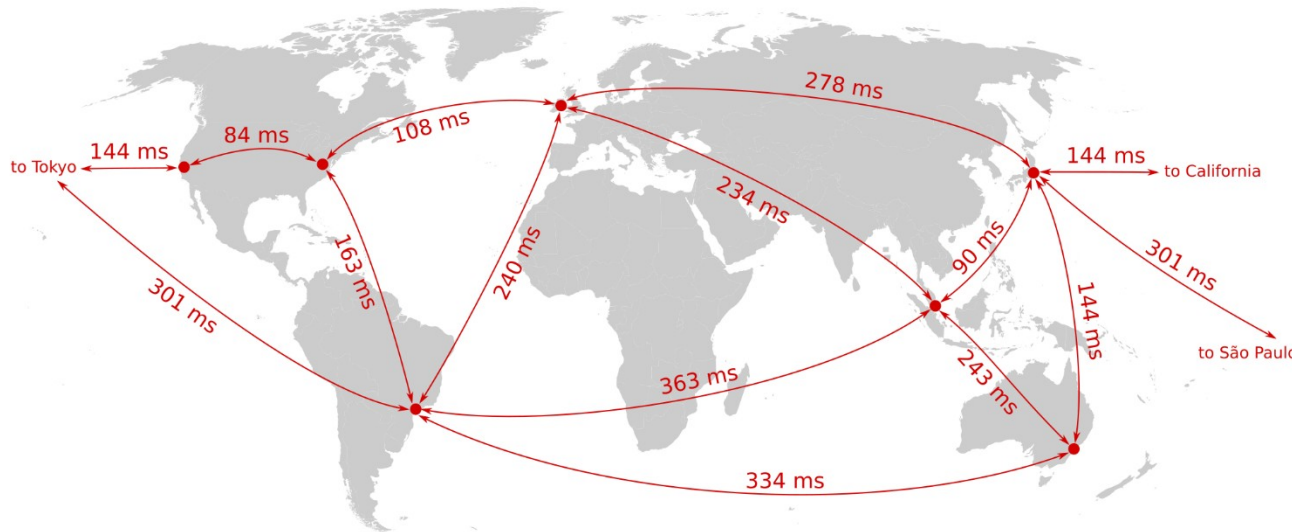
- Example: Let's Encrypt
- 21.01.2021: The Next Gen Database Servers Powering Let's Encrypt [[link](#)]
 - Providing certificates for 275m **websites**
 - “A database is at the heart of how Let’s Encrypt manages certificate issuance” - 1 single MariaDB
 - “We run the CA against a single database in order to minimize complexity” – Some read operations at replicas, one server for writes
 - 2x Xeon 24-cores running at 90%
 - Upgrade to 2x64 Epyc, on 15.09, running at 25%
 - Query 3 times faster
 - SATA → NVMe - IO from 500MB/s to 3 GB/s



Distributed Systems Motivation

- Why Distributed Systems
 - Location
 - Everything gets faster, latency stays
 - Physically bounded by the speed of light

<https://www.inkandswitch.com/local-first.html>



- New protocols can decrease #RT
 - Upcoming lecture
- Place services closer to user
 - Sometimes latency of 310ms is unacceptable
 - ping sydney.edu.au
 - Gaming / **Esports**:
 - Human reaction time 200ms
 - Total **from keypress to display**:
 - Thinkpad 13 ChromeOS: 70ms
 - Lenovo X1 carbon 2016: 150ms
 - TV output lag ~15-30ms (**random TV**)
 - **Keyboard** 15-60ms
- **CDN**: Content delivery network
 - Place your images, sites, scripts close to your users

Distributed Systems Motivation

- Why Distributed Systems
 - Fault-tolerance
 - Any hardware will crash eventually
 - Random bit flips in memory
 - **1990**: “Computers typically experience about one cosmic-ray-induced error per 256 megabytes of RAM per month”
 - **Google study 2009**: more than 8% of DIMMs affected by errors per year
 - **2007**: 44 reported memory errors (41 ECC and 3 double bit) on ~1300 nodes during a period of about 3 months
- Source
 - **Cosmic rays**
 - **Solar flares, Coronal mass ejection, Solar proton events, Background radiation**
- **Cosmic rays** may be blamed for an electronic voting error in Belgium (**2003**)
 - Bit flip in electronic voting machine
 - Added 4096 extra votes to one candidate
 - Candidate more votes than were possible

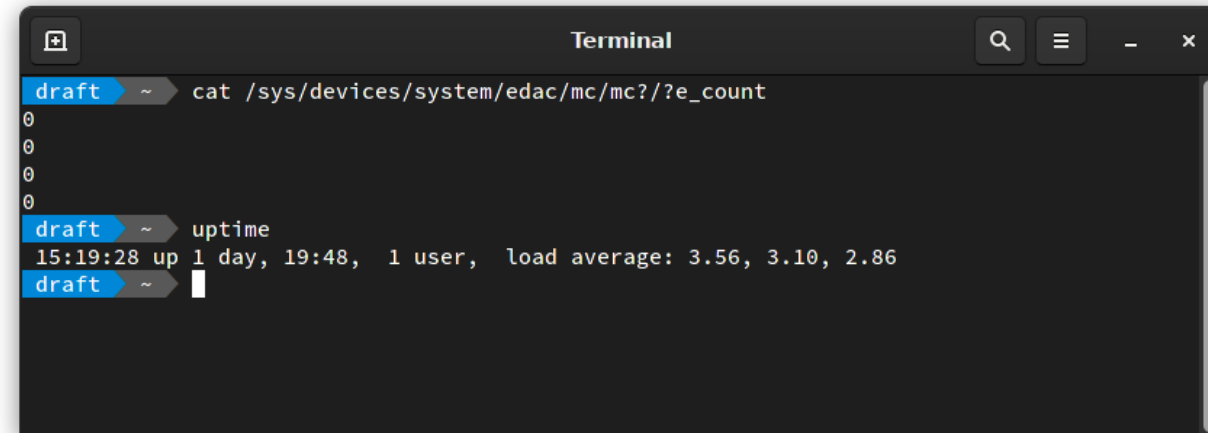


https://en.wikipedia.org/wiki/Solar_flare

Distributed Systems Motivation

- Influencing factors
 - Sensitivity of each transistor, number of transistors on the microchip, altitude
 - Smaller transistors leading to an increased sensitivity per transistor, but smaller cells make smaller targets
- Mars Rover?
 - Cassini reported 280 bitflip/day [link] – max 890 due to solar proton event - TMR with ~300MB RAM
 - Radiation-tolerant FPGAs → TMR
- Error-correcting code memory
 - Uses TMR or Hamming Code, correct 1 bitflip / detect 2 bitflips
 - Used for Servers, not (yet) used for consumer products

- Double bit-flips unlikely?
 - Jaguar super computer with 360TB ECC RAM
 - Double bitflip → happened every 24h
- Check your HW



```
Terminal
draft ~ cat /sys/devices/system/edac/mc/mc??e_count
0
0
0
0
0
draft ~ uptime
15:19:28 up 1 day, 19:48, 1 user, load average: 3.56, 3.10, 2.86
draft ~
```

- What can happen: e.g., expr segfaults

Distributed Systems Motivation

- Random bit flips in memory
 - Bitsquatting: DNS Hijacking without exploitation (2015)
 - Register names with single bit error, e.g,

Bitsquat Domain	Original Domain
ikamai.net	akamai.net
aeazon.com	amazon.com
a-azon.com	amazon.com
amazgn.com	amazon.com
microsmft.com	microsoft.com
micrgsoft.com	microsoft.com

- Idea: if bitflip happens, it may happen for DNS names in your memory
 - Early tests by Artem Dinaburg: “59 unique IPs per day made HTTP requests to my 32 bitsquat domains”
 - 1mio DNS queries every 24h to bitsquatted domains
- Key findings
 - Most users from China (more bitflips on Chinese machines?)
 - 240k session cookies

Fault Tolerance

- Network outages happens **often**
- 22.02.2022: Tonga Cable Successfully Repaired [[link](#)]
 - 38 days broken, see “in the news”
- 26.01.2022: Internet In Yemen Returns After Four Day Outage Caused by Saudi Air Strikes On Telco Facility [[link](#)]
 - Issue lasted 4 days, duet to Air Strike on telecom hub
- 13.01.2022: Fault Reported on Sea-MeWe-4 [[link](#)]
 - Degraded internet performance

- 10.01.2022: Svalbard Suffers Power Fault On Subsea Fiber Cable [[link](#)]
 - 1 out of 2 cables broken (redundancy!)
- [Submarine Cable Map](#)

