# Distributed Systems (DSy)

**Introduction, part 2**

Thomas Bocek

25.02.2024

# 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)

OST

# Distributed Systems Motivation

- Why Distributed Systems

  - Location

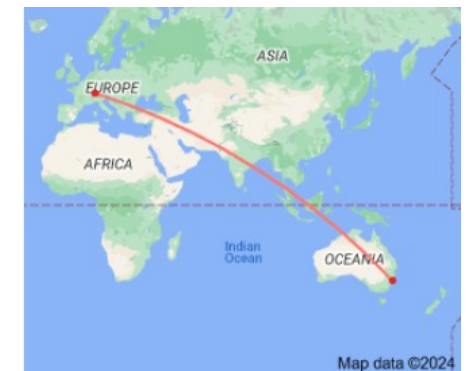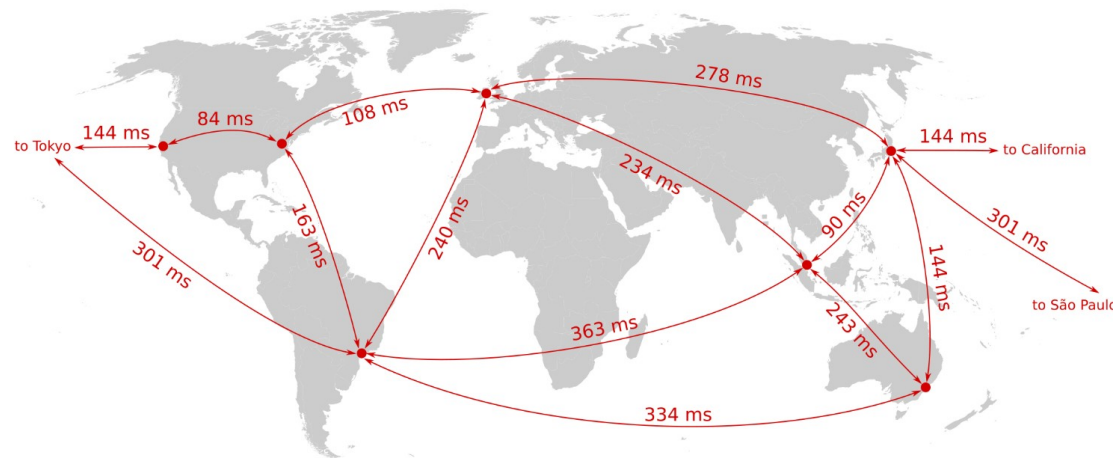    - Everything gets faster, latency stays

- Speed of light in vacuum is ~300'000 km/s

  - Physical limit on how quickly data can travel

- Latency: time for signal to travel from source to destination and back (round-trip time)

  - Perfect vacuum light tube to Sydney: RTT → ~110ms

  - Space? Starlink: 550km

16,540 km

Distance from Rapperswil-Jona to Sydney

# Distributed Systems Motivation

- Copper vs. Fiber

  - Copper propagates faster [link], but not much

  - Depending on the fiber material, latency can change

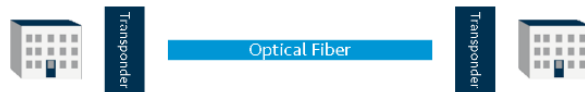  - Reduce latency? (Repeater, Switches, Router)



Figure 1: Low latency options in an example of two data centers 20 km apart

**Example 1: Two data centers 20 km apart**

| | | |
|---|---|---|
| Fiber latency = | 20 x 5 μs = | 100 μs |
| Transponder latency = | 2 x 5 μs = | 10 μs |
| Total latency = | | 110 μs |

Low latency options:

Replace transponders with ultra-low-latency transponders with 4 ns latency per pair.

This effectively removes transponder latency for a 9% savings and a total reduction of 10 μs.

- Importance of latency

  - Amazon: +100ms latency→1% sales loss [link]

  - Google: +500ms latency→20% drop in traffic [link]

  - Bing: +500ms latency → revenue down 1.2% [link]

  - Study: 73% said latency has a critical and direct or an important impact on their revenues [link]

- Gaming



| Sensitivity to latency in online gaming |
|---|
| >300 ms – game unplayable |
| >150 ms – player performance degraded |
| >100 ms – player performance affected |
| 50 ms – target performance |
| 13 ms – lower limit of detectibility |

Source: PubNub

Source: Burger, Thomas, How fast is realtime? Human perception and technology, PubNub, 2015.

OST

# Distributed Systems Motivation

- 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
    - Key travel time!

- Reducing latency

- Faster HW: Repeater, Switch, Router

- New protocols can decrease nr. of RT
  - Upcoming lecture

- Place services closer to user
  - Reduced latency
  - Can increased bandwidth and throughput
  - Can improved reliability and availability
  - Drawback: coordination of data replication and caching

- CDN: Content delivery network – distributed databases, edge computing
  - Place your images, sites, scripts close to your users

OST

# 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 affectedby errors per year

  - 2007: 44 reported memory errors (41 ECC and 3 double bit) on ~1300 nodes during a period of about 3 month

- Source: Bad Pin Connections, Incorrect RAM Timings, Clock Issues, RAM Design Flaws, CPU/RAM/Motherboard Integrated Logic Defects, DRAM Cell Amplification Errors, Cosmic Rays [link]

  - 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
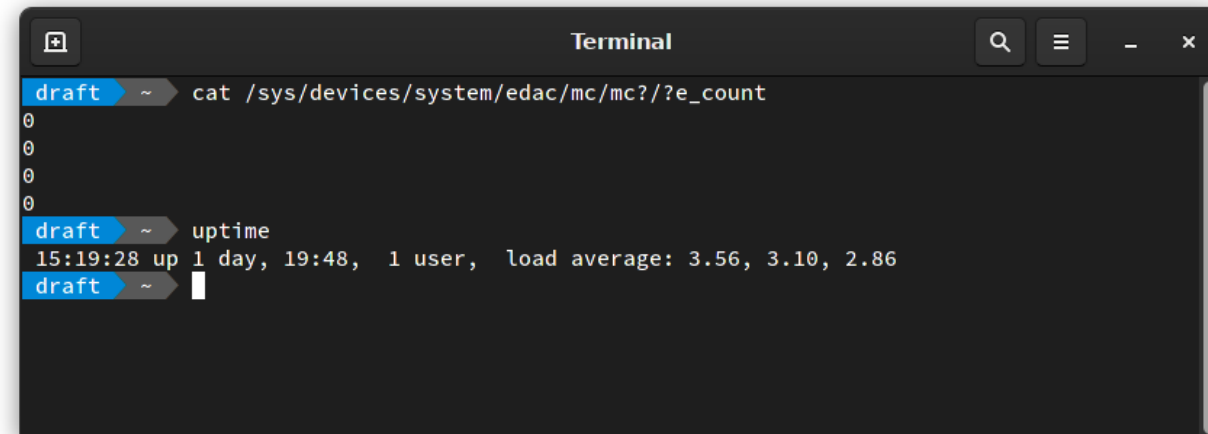
  - Candidate more votes than were possible



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

OST

# Distributed Systems Motivation

- Influencing factors

    - Sensitivity of each transistor, number of transistors on the microchip, altitude, floor level

- 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 – good idea?

- Double bit-flips unlikely?

    - Jaguar super computer with 360TB ECC RAM

    - Double bitflip → happened every 24h

- Check your HW

    - DDR5? on-die vs. traditional ECC [link][link]

```
draft  ~  cat /sys/devices/system/edac/mc/mc?/?e_count
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

OST

# Distributed Systems Motivation

- HDD break [link], SSDs wear out

  - SSDs consist of NAND cells with a limited lifetime

  - An SSD disk has spare NAND that are used when cells break

  - `smartctl –a /dev/xyz`

- SLC, MLC, TLC, QLC

  - SLC: 10'000 – 100'000 write/erase cycles

  - MLC: 10'000 – write/erase cycles

  - TLC: 1'000 – write/erase cycles

  - QLC < 1'000 – write/erase cycles

- 100% → no spare used, My old laptop was at 92%

  - When value is down at 0% disk capacity degrades

- E.g., Samsung 4TB drive uses QLC [link]

  - Write 100 times the same 4kb file, and cells are broken?

- Wear leveling: distribute write and erase operations across all memory cells

  - 4TB → 1b cells, write each 100 → after 100b writes, then cells are broken (TBW)

  - If wear leveling goes wrong: Samsung 990 Pro [link]

- Caching with SLC → files / cells that are frequently changed, store on SLC, once they don't change that often move to MLC/TLC/QLC

```
SMART Attributes Data Structure revision number: 1
Vendor Specific SMART Attributes with Thresholds:
ID# ATTRIBUTE_NAME          FLAG     VALUE WORST THRESH TYPE      UPDATED   WHEN_FAILED RAW_VALUE
  5 Reallocated_Sector_Ct   0x0033   100   100   010    Pre-fail  Always       -        0
  9 Power_On_Hours          0x0032   096   096   000    Old_age   Always       -        18227
 12 Power_Cycle_Count       0x0032   097   097   000    Old_age   Always       -        2430
177 Wear_Leveling_Count     0x0013   092   092   000    Pre-fail  Always       -        288
179 Used_Rsvd_Blk_Cnt_Tot   0x0013   100   100   010    Pre-fail  Always       -        0
181 Program_Fail_Cnt_Total  0x0032   100   100   010    Old_age   Always       -        0
182 Erase_Fail_Count_Total  0x0032   100   100   010    Old_age   Always       -        0
183 Runtime_Bad_Block       0x0013   100   100   010    Pre-fail  Always       -        0
187 Uncorrectable_Error_Cnt 0x0032   100   100   000    Old_age   Always       -        0
190 Airflow_Temperature_Cel 0x0032   071   036   000    Old_age   Always       -        29
195 ECC_Error_Rate          0x001a   200   200   000    Old_age   Always       -        0
199 CRC_Error_Count         0x003e   099   099   000    Old_age   Always       -        15
235 POR_Recovery_Count      0x0012   099   099   000    Old_age   Always       -        682
241 Total_LBAs_Written      0x0032   099   099   000    Old_age   Always       -        3205032857
```

OST

# 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"

- Key findings

  - Most users from China (more bitflips on Chinese machines?)

OST

# Fault Tolerance

- Network outages happens often

- 07.08.2023: Slow Internet Speeds in South Africa — Break in Undersea Cables [link]

- 19.07.2023: Cable Breaks Plague Asian Subsea Cable Operators [link]

- 25.05.2023: Owners of Ship That Damaged Solomon Islands' Coral Sea Cable to Face Charges [link]

- [Submarine Cable Map](#)