

Evolution of cloud computing

1. Virtualization (own servers)
2. Hardware on Demand (AWS)
3. "Serverless computing"

Datacenter

- Large collection of commodity components
- High bandwidth networking gear
- Distributed computing software

Why 2-socket CPU

- Amortize cost of other hardware components
- Increase RAM - cores not that closely integrated
- Performance at large scale - No uniform memory access

- Total response time: Latency + Transfer Time

- Speedup: $\frac{\text{Latency old}}{\text{Latency new}}$

- Amdahl's Law: $\frac{1}{1-p+\frac{p}{s}}$ speedup on parallelizable part

- Speedup limited by parallel part
- Constant size

- Gustafson's law: $\cdot 1-p+s-p$

- Constant computing power

- Brawny CPU (Standard: 2 Brawny CPUs)

- High per-core performance

- Good when a great serial part

- Limited by serial part (Amdahl's law)

- Wimpy CPU

- Lower power

- Good when highly parallelizable

- Systems still consume significant power at low load

- Bisection Bandwidth: Available bandwidth between any two nodes
- Scaling bisection bandwidth
 - For hierarchical network topologies
 - Best switches at the top
 - Over-subscription

Disaggregated Resource Pool



- Total Cost of Ownership ← amortize over time
BUT dominates
- CAPEX: capital expenses (Facilities, hardware, ...)
 - OPEX: operative expenses (Power, employees, ...)
 - Efficiency = $\frac{\text{Computation}}{\text{Total Power}}$
 - Power Usage Effectiveness (= PUE) = $\frac{\text{Total facility Power}}{\text{IT Equipment Power}}$

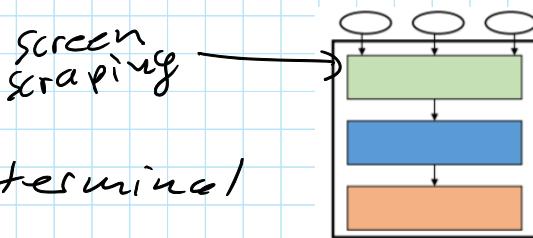
Tiered Software Architecture

- Presentation layer: Client, outside world
- Application logic: Process
- Data layer: Database
- The more boxes, the modular the system is, but more connections need to be kept alive and more context-switches

Enterprise (workload)
data centers → Cloud (workload)
provider

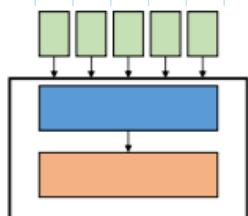
One Tier architecture

- Fully centralized
- No processing, just a terminal
- Everything in one box
- Beginning of computing, no networking yet
- Virtual one tier today: web servers



Two Tier architecture

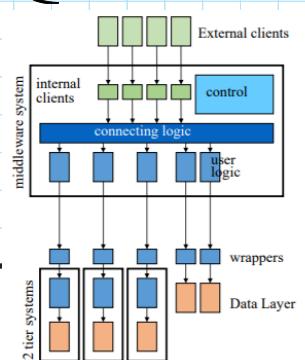
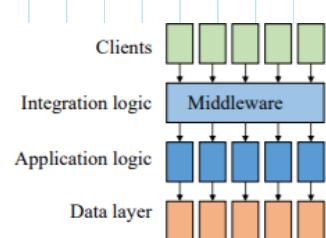
- Only a few clients
- Representation layer moved to client
- Introduces API
- Off-load work to client
- Server has to deal with all possible connections from client
- No load balancing
- Client would need representation layer for each server
- Client needs to know where things are



↳ Solution: Web browsers standardized representation layer

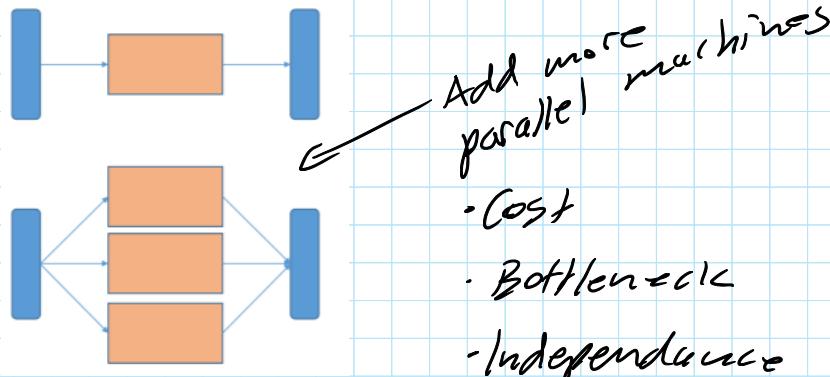
Three Tier architecture

- Completely modular
- Needs to introduce a middleware logic
- Middleware: provides single interface to client

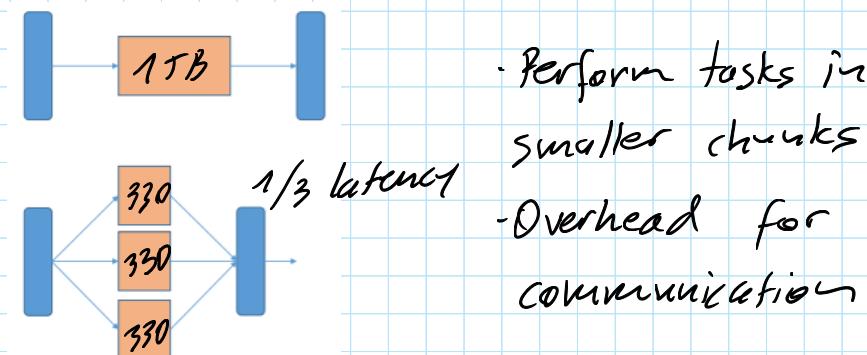


Performance

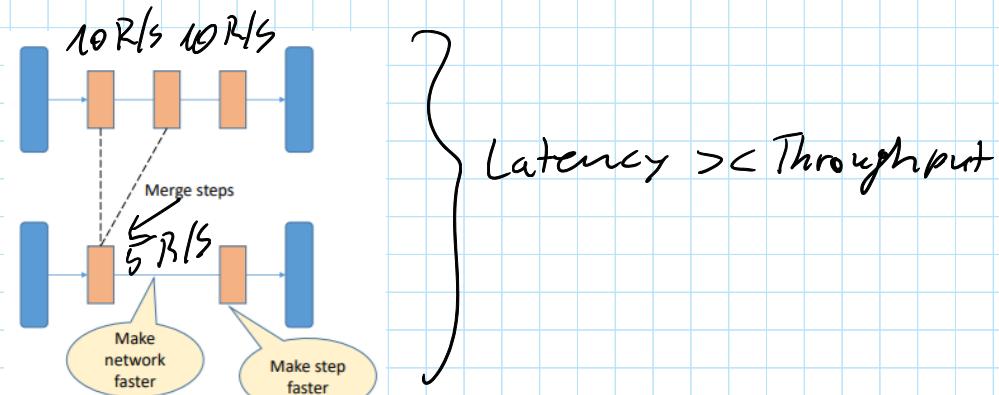
- Requests completed per unit of time
- Throughput : • Can be optimized when parallel



- Latency : • Typically a system measurement
- If tasks take long \Rightarrow Parallelize them



- Reduce length of critical path (consolidation)
- \hookrightarrow Less hops / communication



- Cache it for faster access

- Respond time : perceived by client/user (inverse of throughput)
- The more tiers/components the better
- replication can be done

Open System

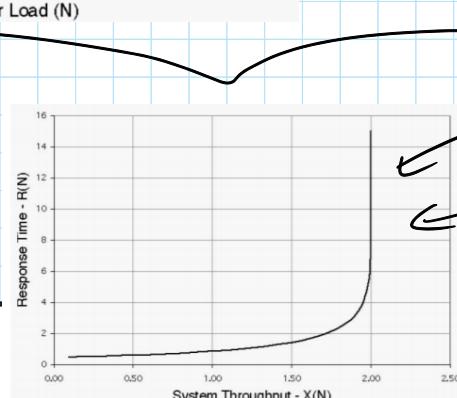
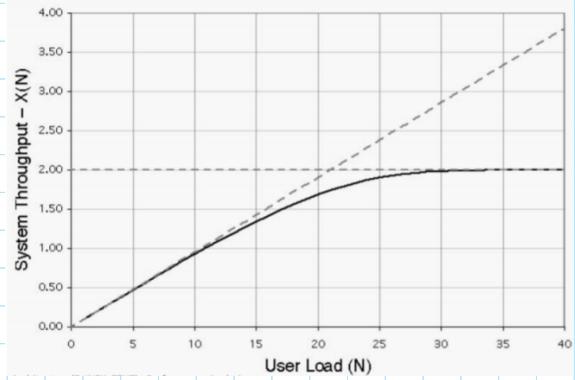
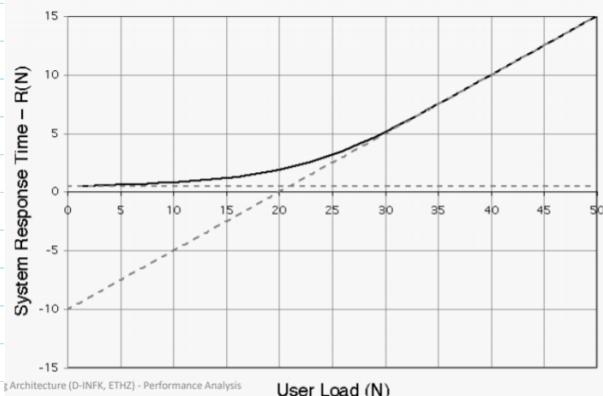
- Load comes from internet
- Not self-adjusting (if server down, more requests)
- e.g. Webserver \hookrightarrow can use flow control
 \hookrightarrow drop packets

Closed System

- Load from limited set of clients
- Client wait for response before sending next request
- e.g. database with local clients

Interactive Response time law

- Response time: R
- Wait for next request: Z
- Total cycle time: $R + Z$
- Number of clients: N
- Latency $X = \frac{\text{Jobs}}{\text{Time}} = \frac{N \cdot T}{R+Z} = \frac{N}{R+Z}$
- $R = \frac{N}{X} - Z$



Why is speed-up sublinear?

- Cost for split/merge operations
- Synchronize servers
- Not splitted into equal chunks

Experiments

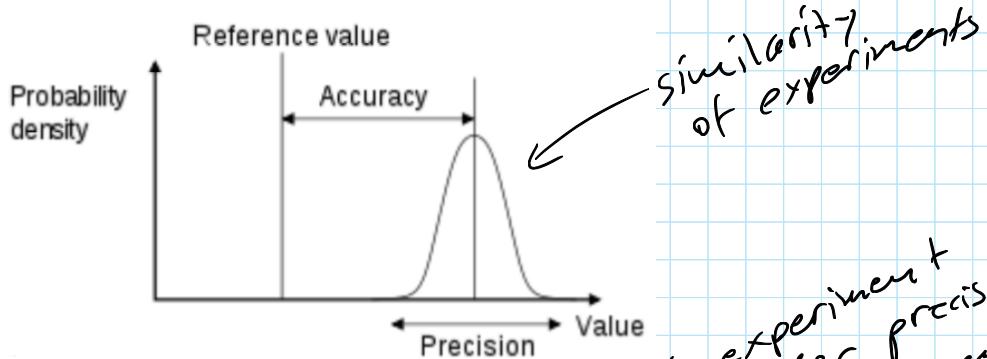
- Done on the machine itself (expensive)
- Benchmark it

Modeling

- Create a model of the system
- Calculate results with model

Simulation

- Implement system that behaves like real one



↗ repeat for higher precision
 ↗ more experiments
 ↗ D
 ↗ n

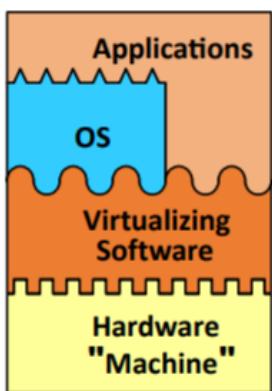
$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

↗ mean

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

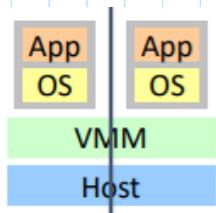
↗ std. deviation (68% of all values)

Virtual Machines



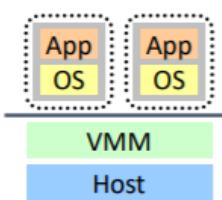
- Guest }
 - Virtual Machine Monitor / Hypervisor }
 - Host }

(local)

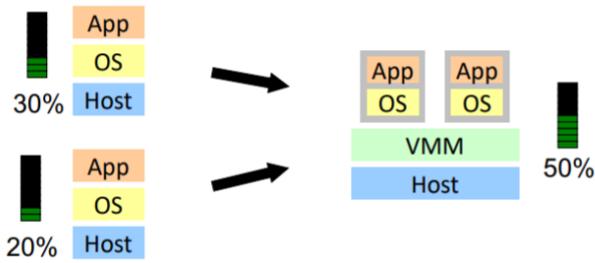


? Partitioning

- Resource sharing
- Context switch

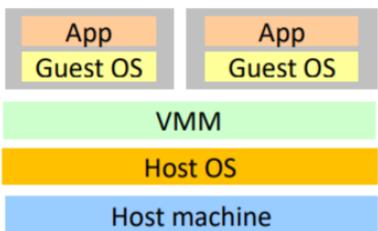


? Encapsulation
 - Security



? Consolidation to improve utilization

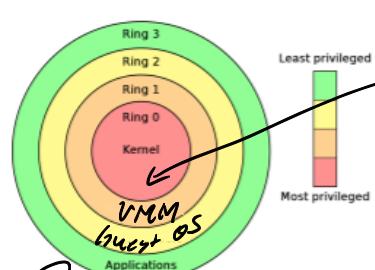
- Hosted interpretation: e.g. Virtual Box



- VMM sits on top of OS
- VMM emulates hardware
- Complete isolation
- Emulation is difficult/slow

- Direct execution with trap and emulate:

- When VM wants to execute privileged instruction \Rightarrow Trap and switch to VMM



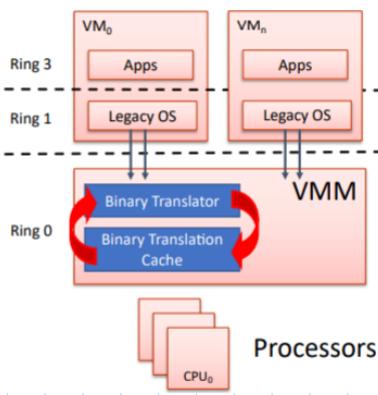
Flush cache,
I/O operations,-

- No host OS is run
- Faster than with host OS

Trap for context switch

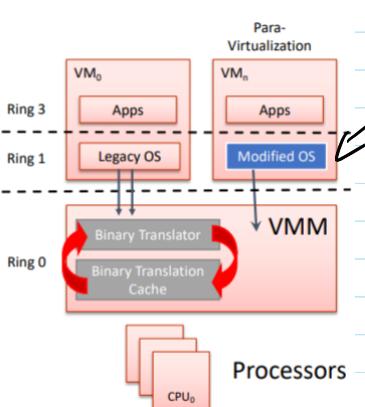
- Trap: context switch
 - System call: privileged, switch to kernel
 - Exception: switch to kernel because of event
 - Interrupt: exterior event, e.g. keyboard input
 - Sensitive Instructions: change hardware configurations
 - Privileged Instructions: cause a trap
- ↳ CPU is virtualizable when also the sensitive instructions trap

Direct execution with binary translation:



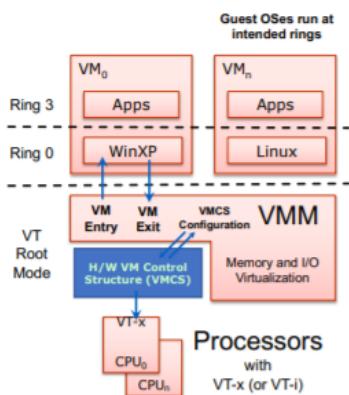
- Can run unmodified guest OS
- Most instructions run directly
- Need to translate virtual memory
- VMM rewrites non-virtualizable instructions

Para-Virtualization



Modify guest-OS to remove sensitive but unprivileged instructions

- Direct execution with hardware support



- Direct execution of VM on processor until privileged instruction is executed
- commonly used today

- Virtualizing Memory

• shadow page-table for os

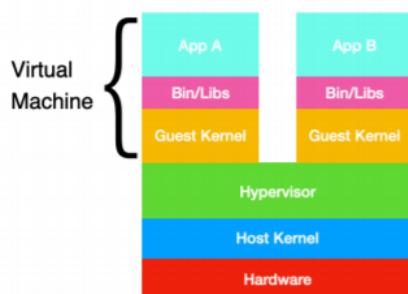
• VMM intercepts

paging operations

• Virtual Page

VMM { ↗ Physical page }
os

↗ Machine page



- Less secure
- No VMM
- Share Host-OS
- Light weight
- Faster start-up
- Bare-metal performance
- Less isolation

- Docker:

• Used for packaging

• Builds on top of Linux containers:

• Image: - creates container
- read only

• Resource isolation
• Manage CPU cores, etc.

- Linux kernel namespaces:

• global system resources

• Restrict what container can see

• e.g. net: network devices

- Linux control groups:
 - Limit resources
 - Groups tell limits
- Why are containers less secure?
 - OS is providing isolation
 - Container runs directly on kernel, which is much bigger \Rightarrow more bugs
 - VMM has a smaller codebase

Kubernetes

- Runs and manages containers
- Run across a cluster
- Coordinates containers
- Pod:
 - bunch of containers with same task
 - used to run a common task together
- Service: group of pods that work together

Cloud Deployments

- On-premises:
 - "enterprise-data center"
 - Often underutilized
- Hybrid:
 - Parts local, parts in the cloud
- Cloud:
 - data center is run by a provider
 - Accessed remotely

Traditional On-Premises IT	Colocation	Hosting	IaaS	PaaS	SaaS
Data	Data	Data	Data	Data	Data
Application	Application	Application	Application	Application	Application
Databases	Databases	Databases	Databases	Databases	Databases
Operating System	Operating System	Operating System	Operating System	Operating System	Operating System
Virtualization	Virtualization	Virtualization	Virtualization	Virtualization	Virtualization
Physical Servers	Physical Servers	Physical Servers	Physical Servers	Physical Servers	Physical Servers
Network & Storage	Network & Storage	Network & Storage	Network & Storage	Network & Storage	Network & Storage
Data Center	Data Center	Data Center	Data Center	Data Center	Data Center

Legend: █ Provider-Supplied █ Self-Managed

FaaS

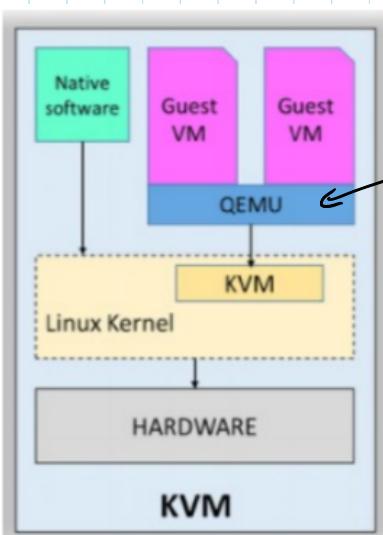
\hookrightarrow serverless/stateless

- small applications
- Executed rarely
- Charged by ms
- short-lived



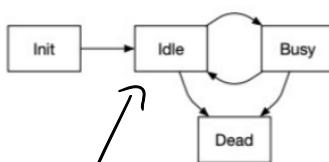
Firecracker

- VMM monitor developed by Amazon
- Uses Linux Kernel-based virtual machine
- Used for lightweight micro VMs
- Used for serverless functions
- Best of both worlds: secure & fast
- Serverless functions only need certain features in the VMM
 - ↳ No USB device support for example
 - ↳ only need networking support
- Linux Kernel VM:
 - Virtualization module
 - Linux as VMM



↳ Emulates the hardware
↳ Replaced by Firecracker

- One Firecracker process per MicroVM
- Written in Rust:
 - Provides memory safety
 - Provides thread safety
 - Fast
- AWS Lambda:
 - Slots for predefined functions



only consumes memory

Resource Assignment

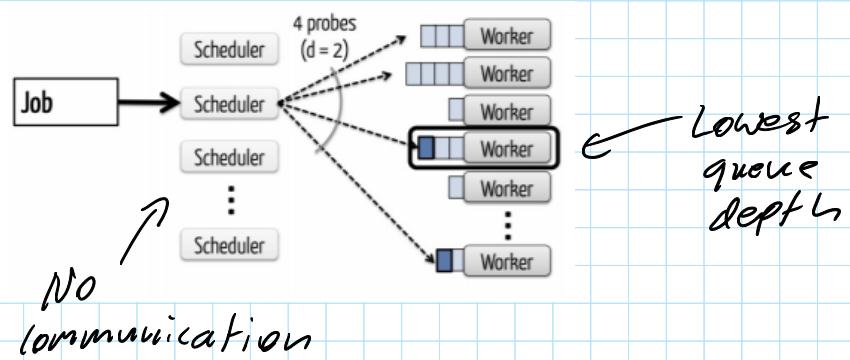
- Private:
 - static set of resources
 - high priority workloads
 - allows hardware specification
- Shared:
 - Needs scheduler
 - Work conserving: unallocated resources get assigned to work

Scheduler

- Single Resource Scheduler:
 - n users
 - give $1/n$ of resource
- Max-Min-Fairness:
 - If user needs more/less
 - min (actual work, max. resource)
 - min (act. work, weight · max resource)
 - ↳ Weighted Max-Min
 - each user gets at least $1/n$
 - Gets less when demand is less
 - Problem: only single resource
- Dominant Resource Fairness:
 - Apply Max-Min to dominant share
 - 9 CPU, 18 GB
 - 1 CPU, 9 GB $\Rightarrow 1/9 < \frac{4}{18}$
 - ↳ x CPU, $9x$ GB
 - 3 CPU, 1 GB $\Rightarrow \frac{3}{9} > \frac{1}{18}$
 - ↳ $3y$ CPU, y GB
 - maximize x/y subject to
 - CPU constraint: $x + 3y \leq 9$
 - Memory constraint: $9x + y \leq 18$
 - Equality: $4x/18 = 3y/9$
 - $\Leftrightarrow 2x/9 = y/3$

Cluster Manager

- Centralized: One single manager responsible for entire cluster
- Two level:
 - One master and multiple frameworks
 - makes resource offers
 - selects resources and provides tasks
- Distributed:
 - multiple scheduling agents
 - Better scalability



- Problem: 70% of tasks use more resources than they actually need
 - Overprovisioning

Quasar

- Observation: resource utilization is low
 - Idea: let user specify performance goals
 - Dynamically adjust based on profiles of service
 - performance can be affected by:
 - Scale up
 - Scale out
 - HW heterogeneity
 - Interference
1. short profiling runs
 2. collaborative filtering to fill missing data
 3. greedy scheduler finds number and type of resources that maximizes performance/utilization

Machine Learning

- Training:
 - Forward pass (loss)
 - Backward pass (gradients)
- Inference:
 - forward pass on small batch of data
- Computation graph:
 - Static: better optimization
 - Dynamic: more intuitive, better for debugging
- Distributed training:
 - Data parallelism:
 - Partition data
 - Train multiple copies of model
 - Synchronize weight updates
 - Model parallelism:
 - More complicated
 - split model across nodes

RPC (= Remote procedure calls)

- Invented when going from 1-Tier \rightarrow 2-Tier
 - Function call is converted into remote procedure
 - ↪ Not seen by program
 - Synchronous
1. Call remote function
 2. Bind: Find remote server with directory service
 3. Marshalling: Convert to a packet
- Tied to TCP/IP
 - Tied to programming languages
 - No standardization

REST (= Representational State Transfer)

- Implemented in HTTP
 - Stateless interaction
 - Still synchronous
 - Get, Put, Delete, Post
- Materialized view: Cache rest-call
- Service aggregations: Make all calls from a centralized place

Data Replication

- why replication:
 - Performance (near me)
 - Fault Tolerance
- Synchronous replication: used for transactions (DBs)
- Asynchronous replication: eventual consistency
- Update everywhere: changes can be done on every replication set
- Primary copy: · only one copy which can be updated
 - databases

Synchronous	<table border="1"><tr><td>Advantages: Updates not coordinated No inconsistencies</td><td>Advantages: No inconsistencies Elegant (symmetrical solution)</td></tr><tr><td>Disadvantages: Longest response time Only useful with few updates Local copies are can only be read</td><td>Disadvantages: Long response times Updates need to be coordinated</td></tr></table>	Advantages: Updates not coordinated No inconsistencies	Advantages: No inconsistencies Elegant (symmetrical solution)	Disadvantages: Longest response time Only useful with few updates Local copies are can only be read	Disadvantages: Long response times Updates need to be coordinated
Advantages: Updates not coordinated No inconsistencies	Advantages: No inconsistencies Elegant (symmetrical solution)				
Disadvantages: Longest response time Only useful with few updates Local copies are can only be read	Disadvantages: Long response times Updates need to be coordinated				
Asynchronous	<table border="1"><tr><td>Advantages: No coordination necessary Short response times</td><td>Advantages: No centralized coordination Shortest response times</td></tr><tr><td>Disadvantages: Local copies are not up to date Inconsistencies</td><td>Disadvantages: Inconsistencies Updates can be lost (reconciliation)</td></tr></table>	Advantages: No coordination necessary Short response times	Advantages: No centralized coordination Shortest response times	Disadvantages: Local copies are not up to date Inconsistencies	Disadvantages: Inconsistencies Updates can be lost (reconciliation)
Advantages: No coordination necessary Short response times	Advantages: No centralized coordination Shortest response times				
Disadvantages: Local copies are not up to date Inconsistencies	Disadvantages: Inconsistencies Updates can be lost (reconciliation)				

Primary copy

Update everywhere

Storage Systems

- Block Storage: virtual raw disk
- Object storage: block blobs / append blobs / page blobs
- Key-Value-Store: distributed hash table in ring
 - ↳ used as cache e.g. hash database queries
- Databases
- File storage: on top of block storage (=NAS)
 - millions of PB files
 - append only

Reliability

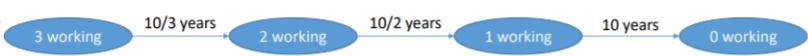
- Errors: operational but produces wrong results
- Faults: operational except some parts
- Failures: system no longer operational
- Performance issues: not usable because too slow
- Frequency of failures: F/τ
- Mean Time To Fail: τ/F ($=MTTF$)
- MTTF multiple components: $MTTF_1/N$
- Failure Rate: $1/MTTF$

↳ When we have a lot of components, the system is more likely to fail

- Mean time between fails: $MTTF + MTTR$
- Availability:
$$\frac{MTTF}{(MTTF + MTTR)}$$
 $\overbrace{\quad\quad\quad}$
mean time
to repair

- Availability of 99% $\hat{=}$ 87,6 hours downtime

- Redundancy:

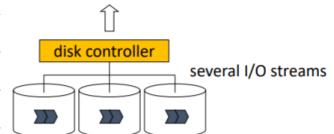


- Countermeasures:

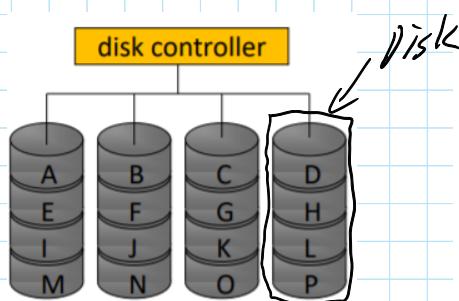
- Replication
- Error correcting codes
- Sharding

Data Striping

- Distribute data among several disks for faster read/writes and backup
- Fine grained:
 - all requests use all disks
 - maximize disk bandwidth
 - only one I/O request serviced at a time { e.g. multimedia }
- Coarse grained:
 - small requests can be serviced in parallel { e.g. DB }
 - large requests still benefit from high transfer rates
- Striping + Parity (detect and correct errors)
 - ↳ RAID

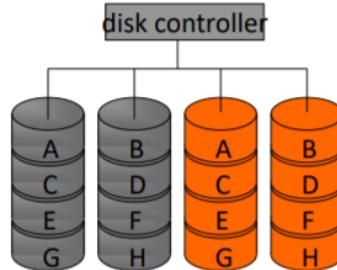


- Raid 0



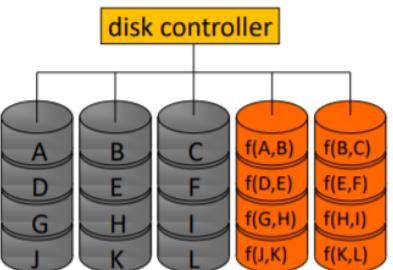
- No redundancy
- Nx faster

- Raid 1



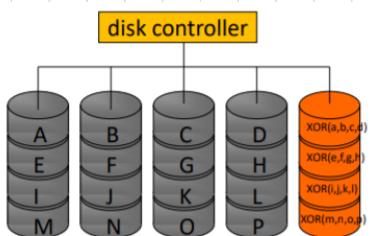
- 50% backup

- Raid 2



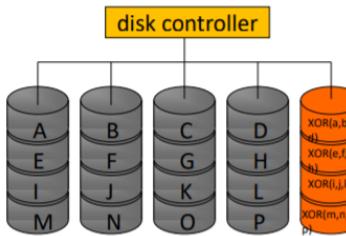
- Hamming code

- Raid 3



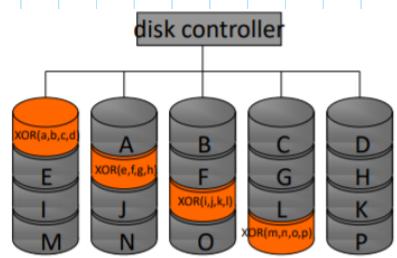
- One disk failure
- parity

- Raid 4

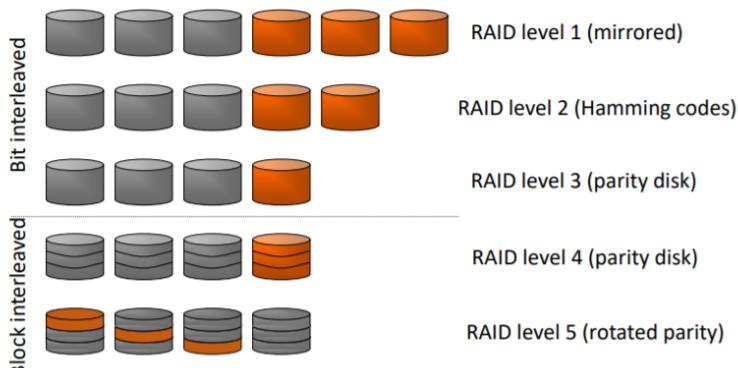


- Blocks instead of bits

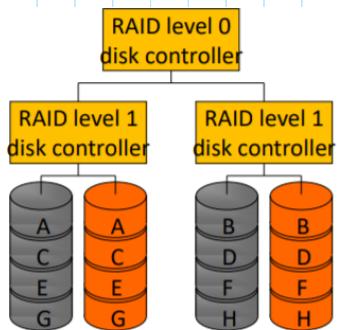
- Raid 5



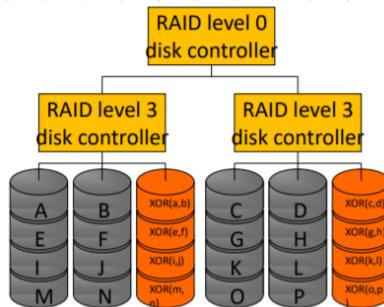
- Read / write in parallel



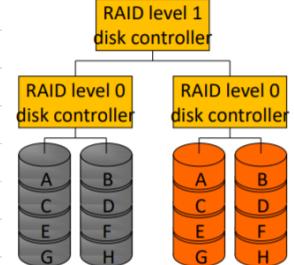
- Raid 10



- Raid 53



- Raid 0+1



- Level 1 controller has less disks to manage

- Sharding

- Split data-sets into smaller pieces
- Place them across multiple machines
- Copyset: set of machines which contain all replicas of one data-block

• If you lose copyset, you lose data

• Random replication eventually

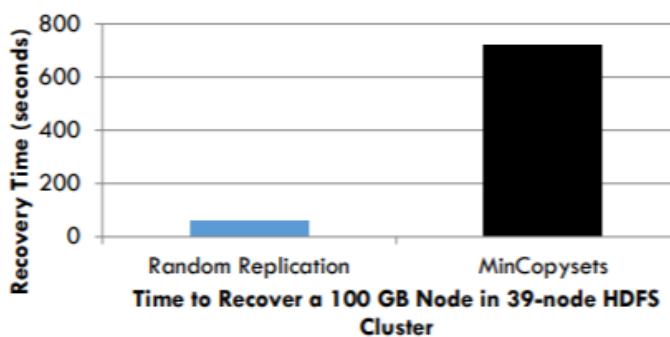
$$\text{leads to } \Pr(\text{data loss}) = \frac{\# \text{current copysets}}{\# \text{max. copysets}}$$

• When $n \rightarrow \infty$: $\Pr(\text{data loss}) = 1$, so if any 3 nodes fail, we have data loss

• MinCopyset: - Minimize #copysets

1. Randomly place first chunk
2. Place two other replicas accordingly to already existing copyset

	MinCopysets	Random Replication
Mean Time to Failure	625 years	1 year
Amount of Data Lost	1 TB	5.5 GB

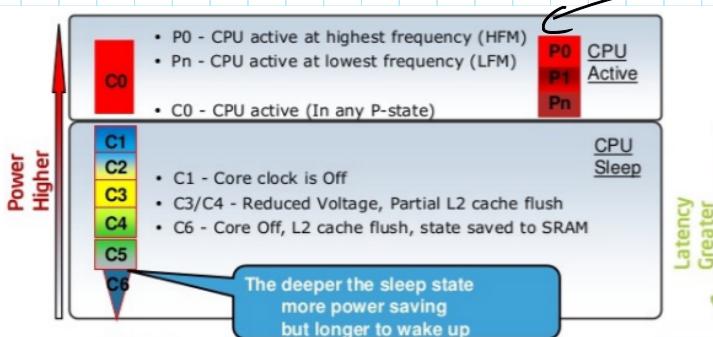


CPU Core Isolation (=Noisy neighbours)

- Intel Cache Allocation Technology: Partition last level cache which is shared

CPU Power Management

- P-states: performance while code being executed
- C-states: CPU is in sleep



Dynamic Voltage/Frequency scaling

- SRAM

- Static Random Access Memory
- Used for caches \Rightarrow Fast
- No refreshing \Rightarrow keeps state
- Bigger than DRAM
- High cost

- DRAM

- Dynamic Random Access Memory
- cell loses charge when read \Rightarrow must periodically refresh
- cell loses charge over time \Rightarrow refresh

- Flash Storage

- Non-volatile
- can only process one operation at a time

\hookrightarrow Flash translation layer (wear count)

Networking

- Memory mapped I/O

- use load/store to read/write locations on device

- Slow

- often used for device configurations

- CPU always involved

- Direct Memory Access

- I/O device to memory transfer

- Without involving CPU

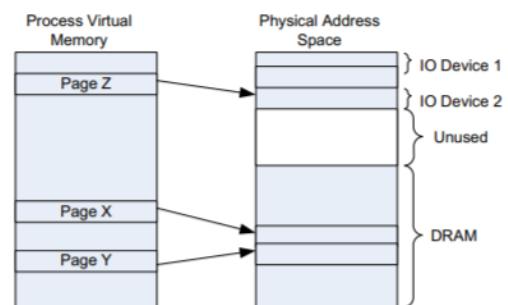
- DMA completes transfer by sending interrupt to CPU

- Inconsistency with CPU caches

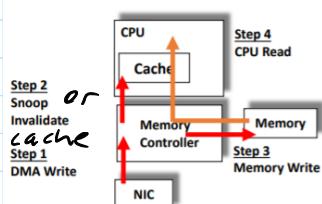
- Map to non-cacheable area

- Cache can snoop DMA bus-line

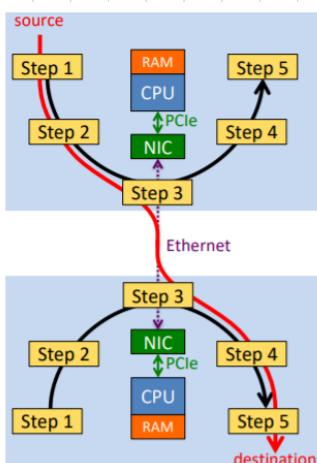
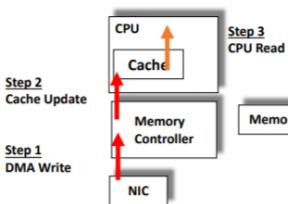
- Direct cache access



Normal DMA Writes



Direct Cache Access



- NIC acceleration techniques
 - checksum
 - Large packets to smaller ones
 - Process multiple packets before interrupt
- Interrupt steering
 - To one core
 - spread across all cores
 - To running app core

Specializations

- CPU: low latency, low parallelization
- GPU: high latency, high parallelization
 - ↳ single instruction
multiple threads
- TPU: matrix multiplication with reduced accuracy
 - Driven by CPU (accelerator)
 - Accelerator

Cloud security

- Goals:
 - Confidentiality: protect data → encrypt it
 - Integrity: Only authorized users can modify
 - Availability: prevent DDoS
 - Access control: Who and what can user do / is
 - Privacy
- Secure boot: check signature of bios / os
- Trusted execution environment (= Enclaves)
 - Run code / memory isolated from system
 - Only trust CPU, all other is encrypted (DRAM)
 - CPU prevents cache access from outside
 - Can only run unprivileged code
 - e.g. for private keys

- Side channels

- communication between components
- Memory access patterns } Examples
- State of branch predictor

- (overt channel)

- Hidden channel to get information out
- Extract patterns from system on shared resources

To send a bit $b \in \{0,1\}$ malware does:

- $b=1$: at 1:00am do CPU intensive calculation
- $b=0$: at 1:00am do nothing

At 1:00am listener does CPU intensive calc. and measures completion time

$$b = 1 \Rightarrow \text{completion-time} > \text{threshold}$$