


ETH zürich spci.inf.ethz.ch  
@spci\_eth

**ADRIAN PERRIG & TORSTEN HOEFLER**

## Networks and Operating Systems Chapter 12: Reliable Storage, NUMA & The Future

**An insecure dishwasher has entered the IoT war against humanity** 27 Nov 2017

We are doomed, but look how shiny those plates are




**BAD NEWS FOR THE STEPFORD WIVES**, the dream of the internet-connected dishwasher has been busted wide apart by a security bod who found a vulnerability with one.


"An issue was discovered on Miele Professional PG 8528 PST10 devices. The corresponding embedded web server PST10 WebServer typically listens to port 80 and is prone to a directory traversal attack; therefore, an unauthenticated attacker may be able to exploit this issue to access sensitive information to aid in subsequent attacks," it says.

AS A PROJECT WEARS ON, STANDARDS FOR SUCCESS SLIP LOWER AND LOWER.


0 HOURS



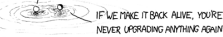
6 HOURS



10 HOURS



24 HOURS



Source: xkcd

ETH zürich spci.inf.ethz.ch  
@spci\_eth

## Last lecture -- basic exam tips

- **First of all, read the instructions**
- **Then, read the whole exam paper through**
- **Look at the number of points for each question**
  - This shows how long we think it will take to answer!
- **Find one you know you can answer, and answer it**
  - This will make you feel better early on.
- **Watch the clock!**
  - If you are taking too long on a question, consider dropping it and moving on to another one.
- **Always show your working**
- **You should be able to explain each summary slide**
  - Tip: form learning groups and present the slides to each other
- **Do NOT overly focus on the quiz questions!**
- **Ask TAs if there are questions**

ETH zürich spci.inf.ethz.ch  
@spci\_eth

## Our Small Quiz

- **True or false (raise hand)**
  - Receiver side scaling randomizes on a per-packet basis
  - Virtual machines can be used to improve application performance
  - Virtual machines can be used to consolidate servers
  - A hypervisor implements functions similar to a normal OS
  - If a CPU is strictly virtualizable, then OS code execution causes nearly no overheads
  - x86 is not strictly virtualizable because some instructions fail when executed in ring 1
  - x86 can be virtualized by binary rewriting
  - A virtualized host operating system can write the hardware PTBR directly
  - Paravirtualization does not require changes to the guest OS
  - A page fault with shadow page tables is faster than nested page tables
  - A page fault with writeable page tables is faster than shadow page tables
  - Shadow page tables are safer than writable page tables
  - Shadow page tables require paravirtualization

ETH zürich spci.inf.ethz.ch  
@spci\_eth

## Virtualizing Devices

- **Familiar by now: trap-and-emulate**
  - I/O space traps
  - Protect memory and trap
  - "Device model": software model of device in VMM
- **Interrupts → upcalls to Guest OS**
  - Emulate interrupt controller (APIC) in Guest
  - Emulate DMA with copy into Guest PAS
- **Significant performance overhead!**

ETH zürich spci.inf.ethz.ch  
@spci\_eth

## Paravirtualized devices

- **"Fake" device drivers which communicate efficiently with VMM via hypercalls**
  - Used for block devices like disk controllers
  - Network interfaces
  - "VMware tools" is mostly about these
- **Dramatically better performance!**

ETH zürich spci.inf.ethz.ch  
@spci\_eth

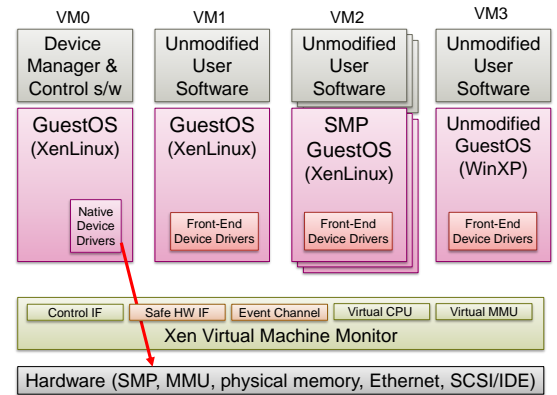
## Networking

- **Virtual network device in the Guest VM**
- **Hypervisor implements a "soft switch"**
  - Entire virtual IP/Ethernet network on a machine
- **Many different addressing options**
  - Separate IP addresses
  - Separate MAC addresses
  - Network Address Translation (NAT)
- **Etc.**

## Where are the real drivers?

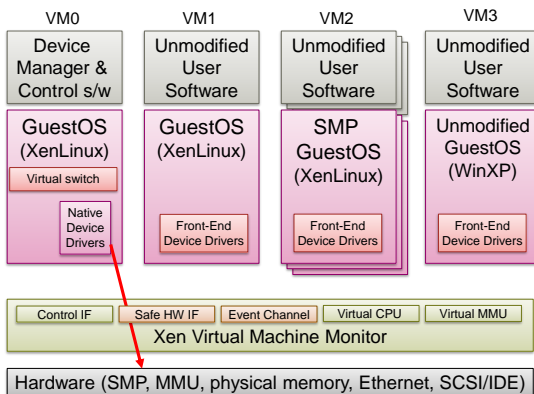
- In the Hypervisor**
  - E.g., VMware ESX
  - Problem: need to rewrite device drivers (new OS)
- In the console OS**
  - Export virtual devices to other VMs
- In “driver domains”**
  - Map hardware directly into a “trusted” VM
  - Device Passthrough**
  - Run your favorite OS just for the device driver
  - Use IOMMU hardware to protect other memory from driver VM
- Use “self-virtualizing devices”**

## Xen 3.x Architecture



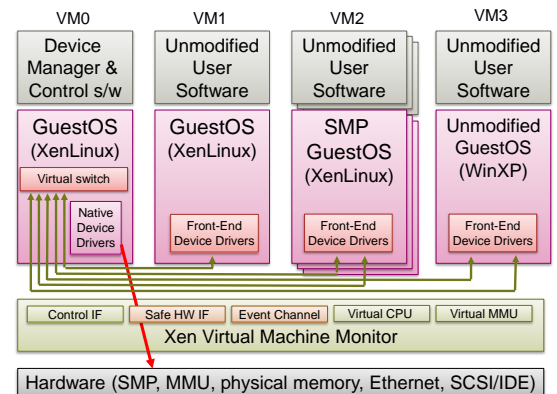
Thanks to Steve Hand for some of these diagrams

## Xen 3.x Architecture



Thanks to Steve Hand for some of these diagrams

## Xen 3.x Architecture



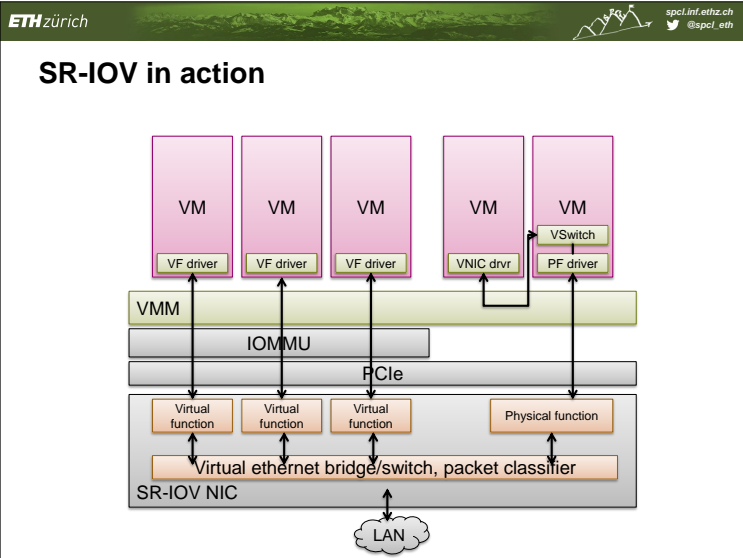
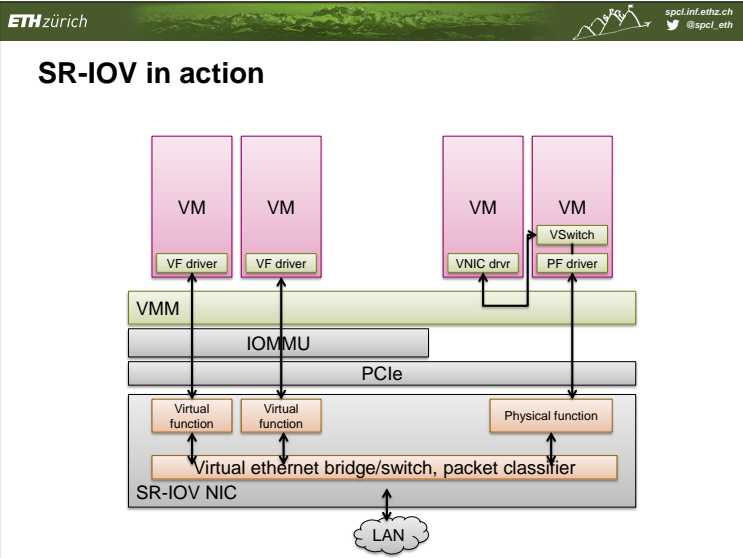
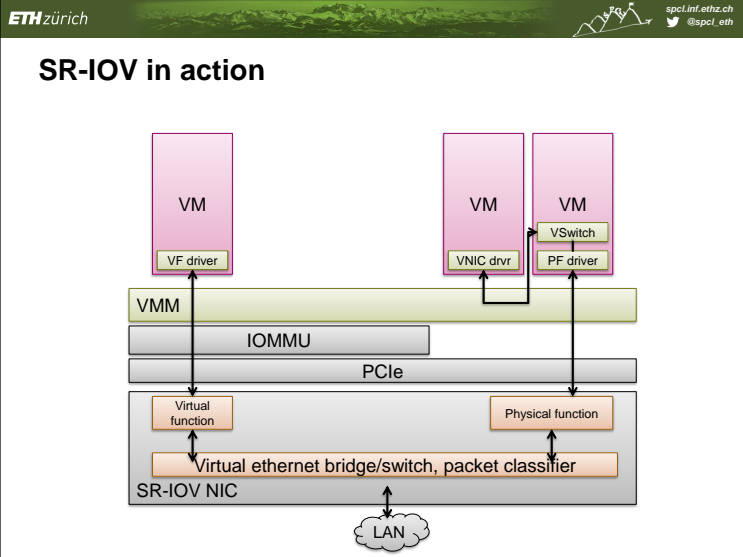
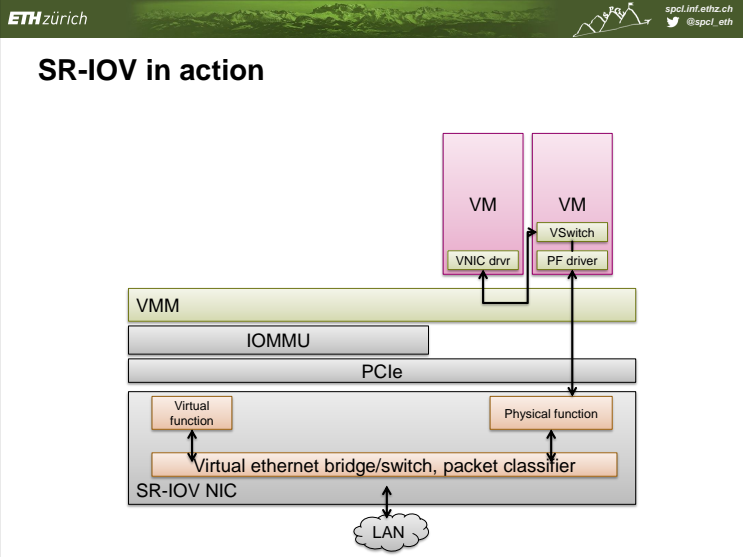
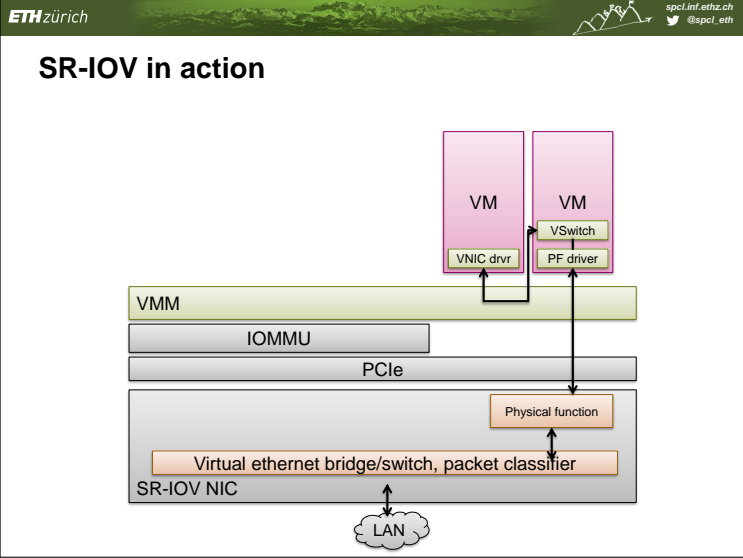
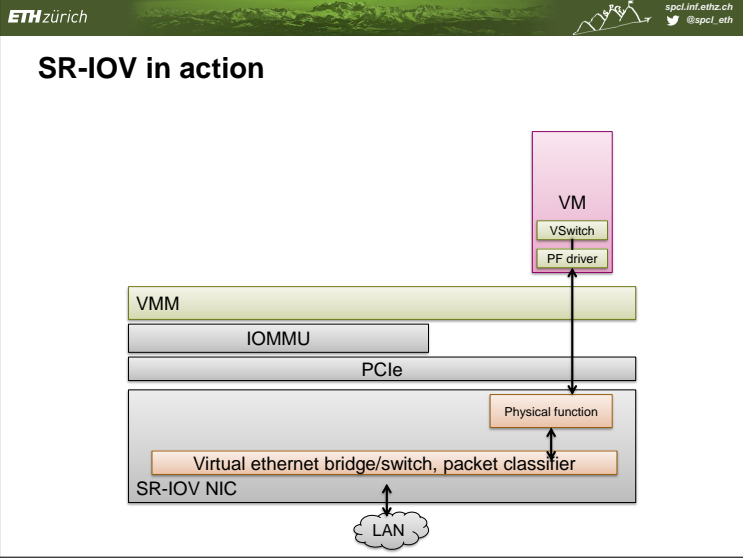
Thanks to Steve Hand for some of these diagrams

## Remember this card?



## SR-IOV

- Single-Root I/O Virtualization**
- Key idea: dynamically create new “PCIe devices”**
  - Physical Function (PF): original device, full functionality
  - Virtual Function (VF): extra “device”, limited functionality
  - VFs created/destroyed via PF registers
- For networking:**
  - Partitions a network card’s resources
  - With direct assignment can implement passthrough





## Self-virtualizing devices

- Can dynamically create up to 2048 distinct *PCI devices* on demand!
  - Hypervisor can create a virtual NIC for each VM
  - Softswitch driver programs “master” NIC to demux packets to each virtual NIC
  - PCI bus is virtualized in each VM
  - Each Guest OS appears to have a “real” NIC, talks directly to the real hardware



## Reliable Storage

OSPP Chapter 14



## Reliability and Availability

A storage system is:

- **Reliable** if it continues to store data and can read and write it.  
 ⇒ **Reliability**: probability it will be reliable for some period of time
- **Available** if it responds to requests  
 ⇒ **Availability**: probability it is available at any given time



## What goes wrong?

1. **Operating interruption: Crash, power failure**
  - Approach: use **transactions** to ensure data is consistent
  - Covered in the databases course
  - See book for additional material



## File system transactions

- Not widely supported
- Only one atomic operation in POSIX:
  - Rename
- Careful design of file system data structures
- Recovery using fsck
- Superseded by transactions
  - Internal to the file system
  - Exposed to applications



## What goes wrong?

1. **Operating interruption: Crash, power failure**
  - Approach: use **transactions** to ensure data is consistent
  - Covered in the databases course
  - See book for additional material
2. **Loss of data: Media failure**
  - Approach: use **redundancy** to tolerate loss of media
  - E.g., RAID storage
  - Topic for today

## Media failures 1: Sector and page failures

### Disk keeps working, but a sector doesn't

- Sector writes don't work, reads are corrupted
- Page failure: the same for Flash memory

### Approaches:

1. **Error correcting codes:**
  - Encode data with redundancy to recover from errors
  - Internally in the drive
2. **Remapping: identify bad sectors and avoid them**
  - Internally in the disk drive
  - Externally in the OS / file system

## Caveats

- **Nonrecoverable error rates are significant**
  - And getting more so!
- **Nonrecoverable error rates are not constant**
  - Affected by age, workload, etc.
- **Failures are not independent**
  - Correlation in time and space
- **Error rates are not uniform**
  - Different models of disk have different behavior over time

## A well-respected disk available now from pcp.ch

Seagate Barracuda 3TB,  
7200rpm, 64MB, 3TB, SATA-3

Price today: CHF 119,95  
(last year: EUR 93,50 (only amazon))  
(in 2015 CHF 119,-)  
(in 2014 CHF 105,-)  
(in 2013 CHF 150,-)



## Specifications (from manufacturer's website)



Persistent errors that are not masked by coding inside the drive

Specifications	3TB'	2TB'
Model Number	ST33000651AS	ST32000641AS
Interface Options	SATA 6Gb/s NCO	SATA 6Gb/s NCO
<b>Performance</b>		
Transfer Rate, Max Ext (MB/s)	600	600
Max Sustained Data Rate OD (MB/s)	149	138
Cache (MB)	64	64
Average Latency (ms)	4.16	4.16
Spindle Speed (RPM)	7200	7200
<b>Configuration/Organization</b>		
Heads/Disks	10/5	8/4
Bytes per Sector	512	512
<b>Reliability/Data Integrity</b>		
Load/Unload Cycles	300K	300K
Nonrecoverable Read Errors per Bits Read, Max	1 per 10E14	1 per 10E14
Annualized Failure Rate (AFR)	0.34%	0.34%
Mean Time Between Failures (hours)	750,000	750,000
Limited Warranty (years)	5	5
<b>Power Management</b>		
Startup Current @ +12V Peak (A) <math>\pm 10\%</math>	2.8	2.8

## Unrecoverable read errors

- What's the chance we could read a *full* 3TB disk without errors?

- For each bit:

$$\Pr(\text{success}) = 1 - 10^{-14}$$

- Whole disk:

$$\Pr(\text{success}) = (1 - 10^{-14})^{8 \times 3 \times 10^{12}} \\ \approx \mathbf{0.7868}$$

- Feeling lucky?

Lots of assumptions:  
Independent errors,  
etc.

## Media failures 2: Device failure

- **Entire disk (or SSD) just stops working**
  - Note: always detected by the OS
  - Explicit failure  $\Rightarrow$  less redundancy required
- **Expressed as:**
  - Mean Time to Failure (MTTF)  
(expected time before disk fails)
  - Annual Failure Rate =  $1/\text{MTTF}$   
(fraction of disks failing in a year)



## Specifications (from manufacturer's website)



Specifications	3TB'	2TB'
Model Number	ST33000651AS	ST32000641AS
Interface Options	SATA 6Gb/s NCQ	SATA 6Gb/s NCQ
<b>Performance</b>		
Transfer Rate, Max Ext (MB/s)	600	600
Max Sustained Data Rate OD (MB/s)	149	138
Cache (MB)	64	64
Average Latency (ms)	4.16	4.16
Spindle Speed (RPM)	7200	7200
<b>Configuration/Organization</b>		
Heads/Disks	10/5	8/4
Bytes per Sector	512	512
<b>Reliability/Data Integrity</b>		
Load/Unload Cycles	300K	300K
Nonrecoverable Read Errors per Bits Read, Max	1 per 10E14	1 per 10E14
Annualized Failure Rate (AFR)	0.34%	0.34%
Mean Time Between Failures (hours)	796,000	750,000
Limited Warranty (years)	5	5
<b>Power Management</b>		
Start-up Current ±1% Peak (A, +10%)	0.9	0.8

86 years!?

## Caveats

- Advertised failure rates can be misleading
  - Depend on conditions, tests, definitions of failure...
- Failures are **not uncorrelated**
  - Disks of similar age, close together in a rack, etc.
- MTTF is not useful life!
  - Annual failure rate only applies during design life!
- Failure rates are **not constant**
  - Devices fail very quickly or last a long time

## ... and reality?

Appears in the Proceedings of the 5th USENIX Conference on File and Storage Technologies (FAST'07), February 2007

### Failure Trends in a Large Disk Drive Population

Eduardo Pinheiro, Wolf-Dietrich Weber and Luiz André Barroso  
Google Inc.

1600 Amphitheatre Pkwy  
Mountain View, CA 94043  
{edpin,wolf,lui}@google.com

(S.M.A.R.T – Self-Monitoring, Analysis, and Reporting Technology)

#### Abstract

It is estimated that over 90% of all new information produced in the world is being stored on magnetic media, most of it on hard disk drives. Despite their importance, there is relatively little published work on the failure patterns of disk drives, and the key factors that affect their lifetime. Most available data are either based on extrapolation from accelerated aging experiments or from relatively modest sized field studies. Moreover, larger population studies rarely have the infrastructure in place to collect health signals from components in operation, which is critical information for detailed failure analysis. We present data collected from detailed observations of a large disk drive population in a production Internet services deployment. The population observed is many times larger than that of previous studies. In addition to presenting failure statistics, we analyze the correlation between failures and several parameters generally believed to impact longevity.

for guiding the design of

vising deployment and m  
Despite the importance  
few published studies or  
drives. Most of the avai  
the disk manufacturers t  
typically based on extra  
test data of small popu  
databases. Accelerated li  
viding insight into how s  
affect disk drive lifetime  
predictors of actual fail  
in the field [7]. Statistic  
cally based on much larger  
populations, but since there  
is little or no visibility into  
the deployment characteris  
tics, the analysis lacks valuable  
insight into what actually  
happened to the drive during  
operation. In addition,

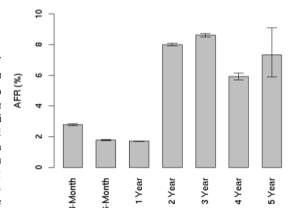
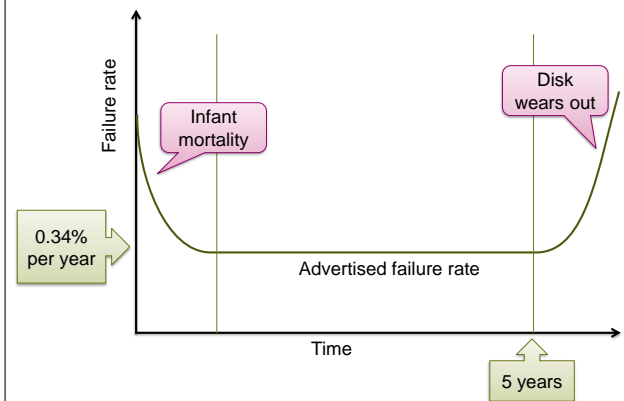
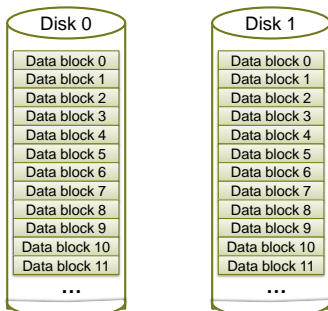


Figure 2: Annualized failure rates broken down by age groups

## Bathtub curve



## RAID 1: simple mirroring

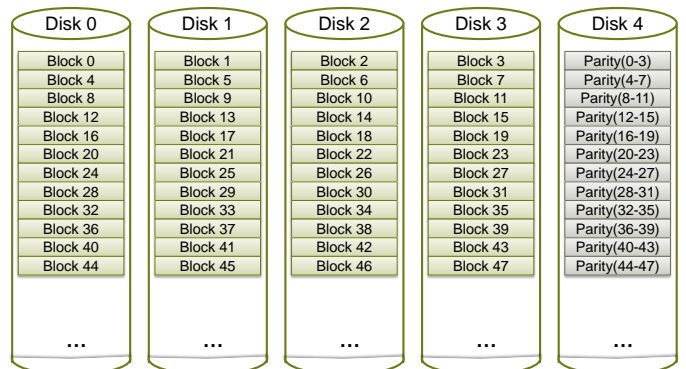


Writes go to both disks

Reads from either disk (may be faster)

Sector or whole disk failure => data can still be recovered

## Parity disks and striping

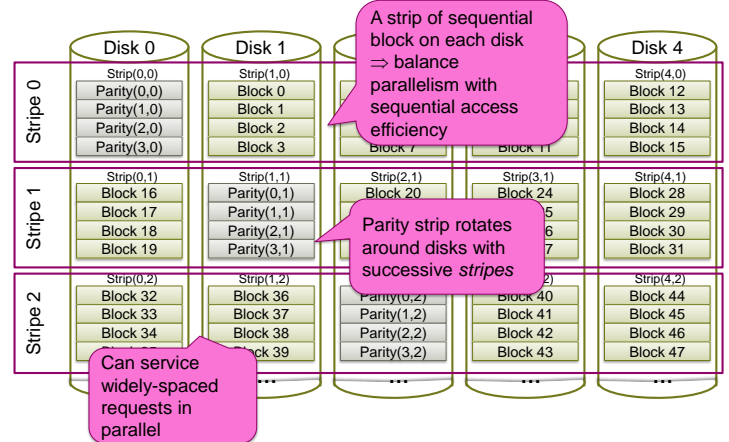


## Parity disks

- Note: errors are always detected  
⇒ Parity allows errors to be corrected
- Write  $d'$  to block  $n$  ⇒ must also update parity, e.g.
  - Read  $d$  from block, parity block, then:  
$$\text{parity}' = \text{parity} \oplus n' \oplus n$$
  - Write  $d'$  to block  $n$ ,  $\text{parity}'$  to parity block
- Problem: with 5 disks, parity disk is accessed 4 times as often on average!

High overhead for small writes

## RAID5: Rotating parity



## Atomic update of data and parity

What if system crashes in the middle?

- Use non-volatile write buffer
- Transactional update to blocks
- Recovery scan
  - And hope nothing goes wrong during the scan
- Do nothing (seriously)

## Recovery

- Unrecoverable read error on a sector:
  - Remap bad sector
  - Reconstruct contents from stripe and parity
- Whole disk failure:
  - Replace disk
  - Reconstruct data from the other disks
  - Hope nothing else goes wrong...

## Mean time to repair (MTTR)

RAID-5 can lose data in three ways:

- Two full disk failures (second while the first is recovering)
- Full disk failure and sector failure on another disk
- Overlapping sector failures on two disks

- MTTR: Mean time to repair
  - Expected time from disk failure to when new disk is fully rewritten, often hours
- MTTDL: Mean time to data loss
  - Expected time until 1, 2 or 3 happens

## Analysis

See the book for *independent failures*

- Key result: most likely scenario is #2.

Solutions:

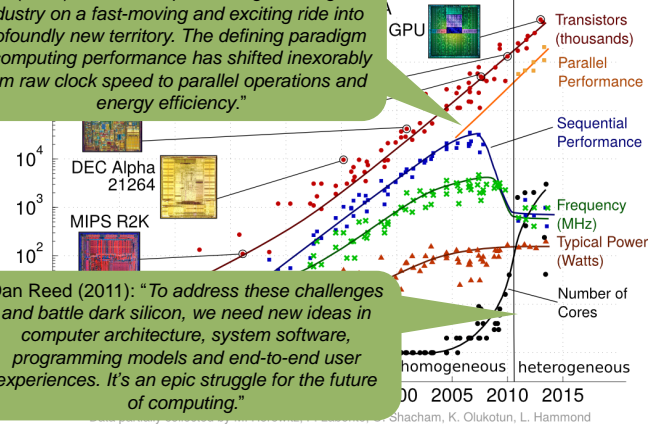
- More redundant disks, erasure coding
- Scrubbing
  - Regularly read the whole disk to catch UREs early
- Buy more expensive disks.
  - I.e., disks with much lower error rates
- Hot spares
  - Reduce time to plug/unplug disk

## Hardware Trends

## The future is exciting!

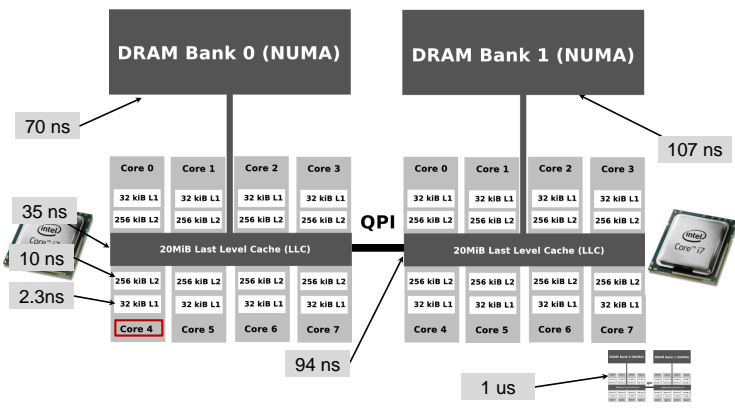
Intel (2006): "Multi-core processing is taking the industry on a fast-moving and exciting ride into profoundly new territory. The defining paradigm in computing performance has shifted inexorably from raw clock speed to parallel operations and energy efficiency."

Dan Reed (2011): "To address these challenges and battle dark silicon, we need new ideas in computer architecture, system software, programming models and end-to-end user experiences. It's an epic struggle for the future of computing."



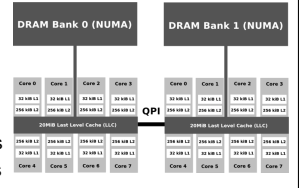
## More and more cores ...

- Like this dual-socket Sandy Bridge system:

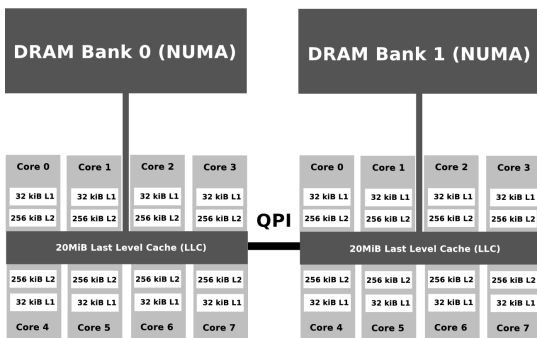


## What does that mean, a nanosecond is short!!

- How fast can you add two (double precision FP) numbers?
  - You're smart, so let's say 1s ☺
- One core performs 8 floating point operations per cycle
  - A cycle takes 0.45ns
- Then ....
  - A L1 cache access (2.3ns) takes 5s
  - A L2 cache access (10ns) takes 22s
  - A L3 cache access (35ns) takes 78s
  - A local DRAM access (70ns) takes 2.5 mins
  - A remote chip access (94ns) takes 3.5 mins
  - A remote DRAM access (107ns) takes 4 mins
  - A remote node memory access (1us) takes 37 mins



## Non-Uniform Memory Access (NUMA)



## NUMA in Operating Systems

- Classify memory into NUMA nodes
  - Affinity to processors and devices
  - Node-local accesses are fastest
- Memory allocator and scheduler should cooperate!
  - Schedule processes close to the NUMA node with their memory
- State of the art:
  - Ignore it (no semantic difference)
  - Striping in hardware (consecutive CLs come from different NUMA nodes)
    - Homogeneous performance, no support in OS needed
  - Heuristics in NUMA-aware OS
  - Special NUMA control in OS
  - Application control



## Heuristics in NUMA-aware OS

- **“First touch” allocation policy**
  - Allocate memory in the node where the process is running
  - Can create big problems for parallel applications (see DPHPC class)
- **NUMA-aware scheduling**
  - Prefer CPUs in NUMA nodes where a process has memory
- **Replicate “hot” OS data structures**
  - One copy per NUMA node
- **Some do page striping in software**
  - Allocate pages round robin
  - Unclear benefits

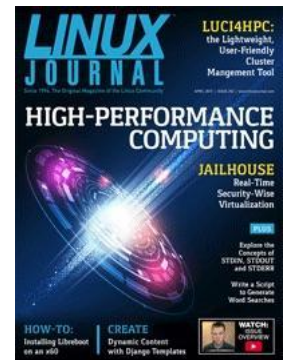
## Special configurations

- **Administrator/command line configurations**
  - Special tools (e.g., Linux)
    - taskset: set a process' CPU affinity*
    - numactl: set NUMA policies*
- **Application configuration**
  - Syscalls to control NUMA (e.g., Linux)
    - cpuset and friends, see “man 7 numa”*

## Non-local system times ☺

- **One core performs 8 floating point operations per cycle**
  - A cycle takes 0.45ns
- **Then ....**
  - A L1 cache access (2.3ns) takes 5s
  - A L2 cache access (10ns) takes 22s
  - A L3 cache access (35ns) takes 78s
  - A local DRAM access (70ns) takes 2.5 mins
  - A remote chip access (94ns) takes 3.5 mins
  - A remote DRAM access (107ns) takes 4 mins
  - A remote node memory access (1us) takes 37 mins
  - Solid state disk access (100us) takes 2.6 days
  - Magnetic disk access (5ms) takes 8.3 months
  - Internet Zurich to Chicago (150ms) takes 10.3 years
  - VMM OS reboot (4s) takes 277 years
  - Physical machine reboot (30s) 2 millennia

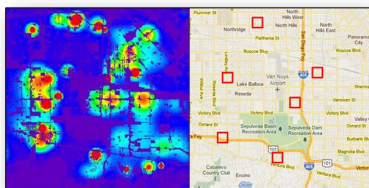
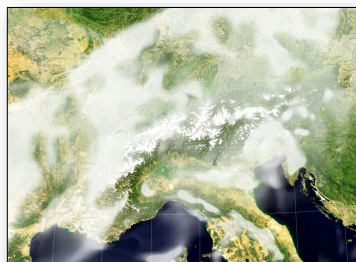
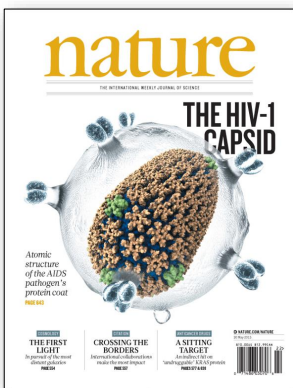
## How to compute fast?



March 2015

## Why computing fast?

- **Computation is the third pillar of science**



## 1 Teraflop in 1997




\$67 Million

ETH zürich spcl.inf.ethz.ch @spcl\_eth

## 1 TeraFlop 18 years later (2015)

Want to play with any of these?



2.9TF  
1 TF  
3 TF

“Amazon.com by Intel even has the co-processor selling for just \$142 (plus \$12 shipping) though they seem to be now out of stock until early December.” (Nov. 11, 2014)

ETH zürich spcl.inf.ethz.ch @spcl\_eth

## 1 TeraFlop 20 years later (2017)



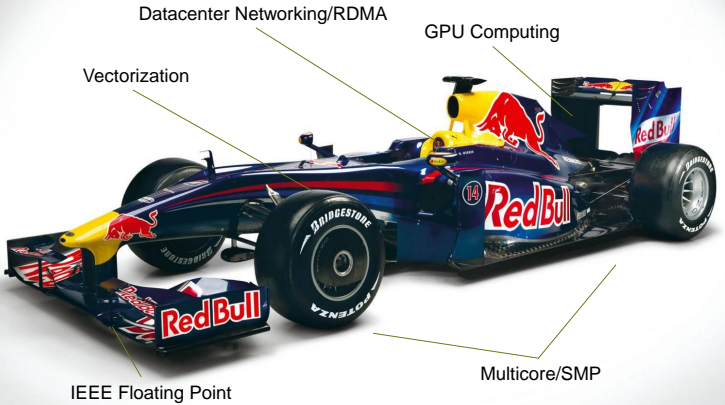
ETH zürich spcl.inf.ethz.ch @spcl\_eth

## 1 TeraFlop 25 years later (2022)



ETH zürich spcl.inf.ethz.ch @spcl\_eth

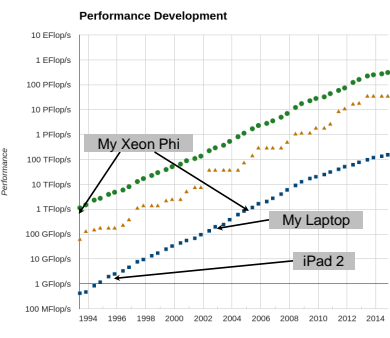
## High-performance Computing (Supercomputing)



ETH zürich spcl.inf.ethz.ch @spcl\_eth

## Top 500

- A benchmark, solve  $Ax=b$ 
  - As fast as possible! → as big as possible ☺
  - Reflects some applications, not all, not even many
  - Very good historic data!
- Speed comparison for computing centers, states, countries, nations, continents ☺
  - Politicized (sometimes good, sometimes bad)
  - Yet, fun to watch



**TOP 500**  
The List.

ETH zürich spcl.inf.ethz.ch @spcl\_eth

## The November 2016 List

Rank	Site	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	National Supercomputing Center in Wuxi, China	Sunway TaihuLight - Sunway MPP; Sunway SW26010 260C 1.45GHz, Sunway NRCCPC	10,649,600	93,014.6	126,435.9	15,371
2	National Super Computer Center in Guangzhou, China	Tianhe-2 (MilkyWay-2) - TH-1TB-FEP Cluster; Intel Xeon E5-2692 V2 @ 2.200GHz, TH Express-2; Intel Xeon Phi 3151P NUDT	3,120,000	33,862.7	54,902.4	17,808
3	DOE/SC/Oak Ridge National Laboratory, United States	Titan - Cray XK7, Opteron 6274 16C 2.200GHz; Cray Gemini interconnect; NVIDIA K20x Cray Inc.	566,640	17,590.0	27,112.5	8,209
4	DOE/NNSA/LLNL, United States	Sequoia - BlueGene/Q, Power BOC 16C 1.60 GHz, Custom IBM	1,572,864	17,173.2	20,132.7	7,890
5	DOE/SC/LLNL/NERSC, United States	Cori - Cray XC40, Intel Xeon Phi 7250 8C 1.4GHz, Aries interconnect; Cray Inc.	622,336	14,014.7	27,880.7	3,939
6	Joint Center for Advanced High Performance Computing, Japan	Oakforest-PACS - PRIMERGY CX1640 M1, Intel Xeon Phi 7250 8C 1.4GHz, Intel Omni-Path; Fujitsu	556,104	13,556.6	24,913.5	2,719
7	RIKEN Advanced Institute for Computational Science (AICS), Japan	K computer; SPARC64 VIIIfx 2.00Hz, Tofu interconnect; Fujitsu	705,024	10,510.0	11,280.4	12,640
8	Swiss National Supercomputing Centre (CS3), Switzerland	Piz Daint - Cray XC30, Jesh ES-2990v3 12C 2.6GHz, Aries interconnect, NVIDIA Tesla P100	206,720	9,779.0	15,968.0	1,312

Want to run on that system?

www.top500.org



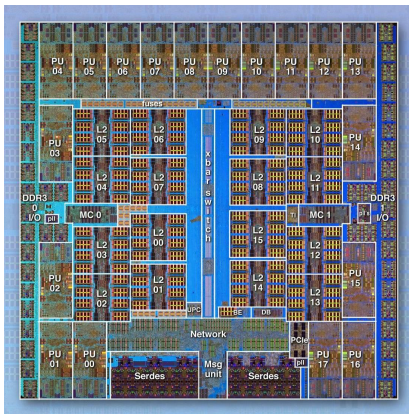
### Case study: OS for High-Performance Computing

- Remember the OS design goals?
  - What if performance is #1?
- Different environment
  - Clusters, special architectures, datacenters
  - Tens of thousands of nodes
  - Hundreds of thousands of cores
  - Millions of CHF
  - Unlimited fun ☺

### Case Study: IBM Blue Gene



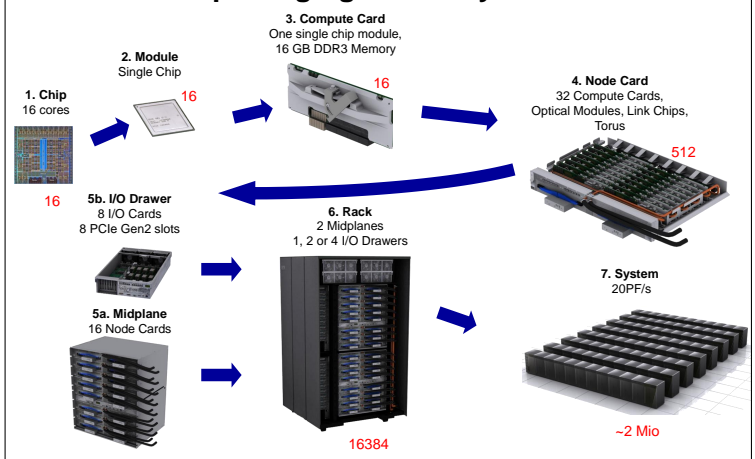
### BlueGene/Q Compute chip



- 360 mm<sup>2</sup> Cu-45 technology (SOI)
  - ~ 1.47 B transistors
- 16 user + 1 service processors
  - plus 1 redundant processor
  - all processors are symmetric
  - each 4-way multi-threaded
  - 64 bits PowerISA™
  - 1.6 GHz
  - L1 I/D cache = 16kB/16kB
  - L1 prefetch engines
  - each processor has Quad FPU (4-wide double precision, SIMD)
- peak performance 204.8 GFLOPS@55W
- Central shared L2 cache: 32 MB
  - eDRAM
  - multiversioned cache will support transactional memory, speculative execution.
  - supports atomic ops
- Dual memory controller
  - 16 GB external DDR3 memory
  - 1.33 Gb/s
  - 2 \* 16 byte-wide interface (+ECC)
- Chip-to-chip networking
  - Router logic integrated into BQC chip.

Ref: SC2010, IBM

### Blue Gene/Q packaging hierarchy

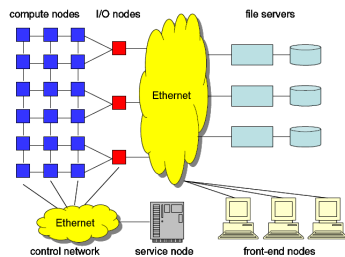


Ref: SC2010, IBM

### Blue Gene/L System Organization

#### Heterogeneous nodes:

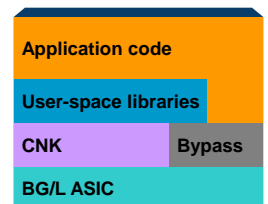
- Compute (BG/L specific)
  - Run specialized OS supporting computations efficiently
- I/O (BG/L specific)
  - Use OS flexibly supporting various forms of I/O
- Service (generic)
  - Uses conventional off-the-shelf OS
  - Provides support for the execution of compute and I/O node operating systems
- Front-end (generic)
  - Support program compilation, submission and debugging
- File server (generic)
  - Store data that the I/O nodes read and write



Source: Jose Moreira et al. "Designing Highly-Scalable Operating System: The Blue Gene/L Story", <http://sc06.supercomputing.org/schedule/pdf/pap178.pdf>

### Software Stack in Compute Node

- CNK controls all access to hardware, and enables bypass for application use
- User-space libraries and applications can directly access torus and tree through bypass
- As a policy, user-space code should not directly touch hardware, but there is no enforcement of that policy



Source: [http://www.research.ibm.com/bluegene/presentations/BGWS\\_05\\_SystemSoftware.ppt](http://www.research.ibm.com/bluegene/presentations/BGWS_05_SystemSoftware.ppt)

## Compute Node Kernel (CNK)

- **Lean Linux-like kernel (fits in 1MB of memory)**
  - stay out of way and let the application run
- **Performs job startup sequence on every node of a partition**
  - Creates address space for execution of compute process(es)
  - Loads code and initialized data for the executable
  - Transfers processor control to the loaded executable
- **Memory management**
  - Address spaces are flat and fixed (no paging), and fit statically into PowerPC 440 TLBs
- **No process scheduling: only one thread per processor**
- **Processor control stays within the application, unless:**
  - The application issues a system call
  - Timer interrupt is received (requested by the application code)
  - An abnormal event is detected, requiring kernel's attention

## CNK System Calls

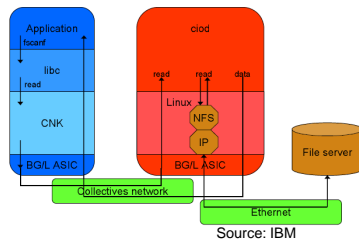
- **Compute Node Kernel supports**
  - 68 Linux system calls (file I/O, directory operations, signals, process information, time, sockets)
  - 18 CNK-specific calls (cache manipulation, SRAM and DRAM management, machine and job information, special-purpose register access)
- **System call scenarios**
  - **Simple** calls requiring little OS functionality (e.g. accessing timing register) are handled locally
  - **I/O** calls using file system infrastructure or IP stack are shipped for execution in the I/O node associated with the issuing compute node
  - **Unsupported** calls requiring infrastructure not supported in BG/L (e.g. *fork()* or *mmap()*) return immediately with error condition

## Function Shipping from CNK to CIOD

- **CIOD processes requests from**
  - Control system using socket to the service node
  - Debug server using a pipe to a local process
  - Compute nodes using the tree network

### I/O system call sequence:

- CNK trap
- Call parameters are packaged and sent to CIOD in the corresponding I/O node
- CIOD unpacks the message and reissues it to Linux kernel on I/O node
- After call completes, the results are sent back to the requesting CNK (and the application)

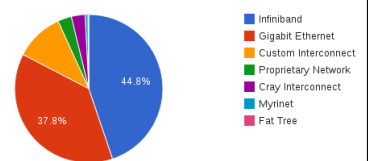


## How to communicate?

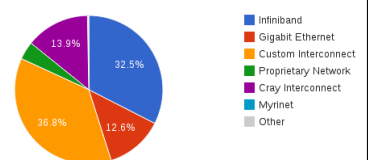
- **Communication is key in problem solving ☺**

- Not just relationships!
- Also scientific computations

Interconnect Family System Share



Interconnect Family Performance Share



Source: top500.org

## Remote Direct Memory Access

- **Remember that guy?**
  - EDR
  - 2x2x100 Gb/s → ~50 GB/s
  - Memory bandwidth: ~80 GB/s
  - 0.8 copies ☹
- **Solution:**
  - RDMA, similar to DMA
  - OS too expensive, bypass
  - Communication offloading



Want to learn to program RDMA?

## InfiniBand Overview

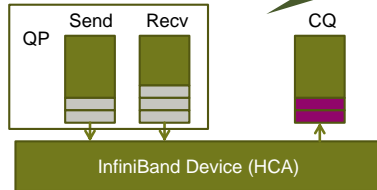
- **Components:**
  - Links/Channel adaptors
  - Switches/Routers
- **Routing is supported but rarely used, most IB networks are "LANs"**
- **Supports arbitrary topologies**
  - "Typical" topologies: fat tree, torus, islands
- **Link speed (all 4x):**
  - Single data rate (SDR): 10 Gb/s
  - Double data rate (DDR): 20 Gb/s
  - Quad data rate (QDR): 40 Gb/s
  - Fourteen data rate (FDR): 56 Gb/s
  - Enhanced data rate (EDR): 102 Gb/s

Want to find better topologies (good at group/graph theory)?

## Interaction with IB HCAs

- **Systems calls only for setup:**
  - Establish connection, register memory
- **Communication (send/recv, put, get, atomics) all in user-level!**
  - Through “verbs” interface

Want to think about a better way to interact with RDMA?



## Open Fabrics Stack

- **OFED offers a unified programming interface**
  - Cf. Sockets
  - Originated in IB verbs
  - Direct interaction with device
  - Direct memory exposure
    - Requires *page pinning* (avoid OS interference)
- **Device offers**
  - User-level driver interface
  - Memory-mapped registers

## iWARP and RoCE

- **iWARP: RDMA over TCP/IP**
  - Ups:
    - *Routable with existing infrastructure*
    - *Easily portable (filtering, etc.)*
  - Downs:
    - *Higher latency (complex TOE)*
    - *Higher complexity in NIC*
    - *TCP/IP is not designed for datacenter networks*
- **RoCE: RDMA over Converged Ethernet**
  - Data-center Ethernet!

## Student Cluster Competition

- **6 BSc students, 1 advisor, 1 cluster, 2x13 amps**
  - 8 teams, 4 continents @SC
  - 48 hours, five applications, non-stop!
  - top-class conference (>13,000 attendees)
- **Lots of fun**
  - Even more experience!
- **A Swiss team 2018?**
  - Search for “Student Cluster Challenge”
  - HPC-CH/CSCS will help

Want to become an expert in HPC?



## What to remember in 10 years!

## The Lecture’s Elevator Pitch

- **Roles:**
  - Referee, Illusionist, Glue
- **Example: processes, threads, and scheduling**
  - R: Scheduling algorithms (batch, interactive, realtime)
  - I: Resource abstractions (memory, CPU)
  - G: Syscalls, services, driver interface
- **Slicing along another dimension:**
  - Abstractions
  - Mechanisms



## The Lecture's Elevator Pitch

- **IPC and other communications**
  - A: Sockets, channels, read/write
  - M: Network devices, packets, protocols
- **Memory Protection**
  - A: Access control
  - M: Paging, protection rings, MMU
- **Paging/Segmentation**
  - A: Infinite memory, performance
  - M: Caching, TLB, replacement algorithms, tables

## The Lecture's Elevator Pitch

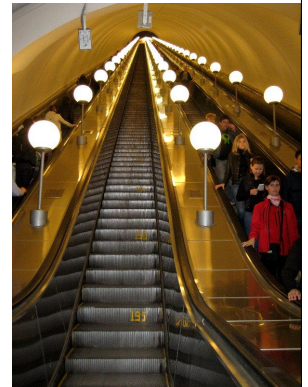
- **Naming**
  - A: (hierarchical) name spaces
  - M: DNS, name lookup, directories
- **File System**
  - A: Files, directories, links
  - M: Block allocation, inodes, tables
- **I/O**
  - A: Device services (music, pictures ☺)
  - M: Registers, PIO, interrupts, DMA

## The Lecture's Elevator Pitch

- **Reliability:**
  - A: reliable hardware (storage)
  - M: Checksums, transactions, raid 1/5
- **And everything can be virtualized!**
  - CPU, MMU, memory, devices, network
  - A: virtualized x86 CPU
  - M: paravirtualization, rewriting, hardware extensions
  - A: virtualized memory protection/management
  - M: writable pages, shadow pages, hw support, IOMMU

## The Lecture's Escalator Pitch

- **Ok, fine, it was an escalator pitch ... in Moscow**
- **Please remember all for at least 10 years!**
  - Systems principles
  - ... and how to make them **fast** ☺



## Finito

- **Thanks for being such fun to teach ☺**
  - Comments (also anonymous) are always appreciated!
- **If you are interested in parallel computing research, talk to me!**
  - Large-scale (datacenter) systems
  - Parallel computing (SMP and MPI)
  - GPUs (CUDA), FPGAs, Manycore ...
  - ... spcl-friends mailing list (subscribe on webpage)
  - ... on twitter: @spcl\_eth ☺

Hope to see you again!  
 Maybe in *Design of Parallel and High-Performance Computing* next semester ☺

Or these:  
<http://spcl.inf.ethz.ch/SeMa/>

