

ETH zürich

ADRIAN PERRIG & TORSTEN HOEFLER

Networks and Operating Systems (252-0062-00)
Chapter 5: Memory Management

<http://support.apple.com/kb/HT5642>
 "Description: The iOS kernel has checks to validate that the user-mode pointer and length passed to the copyin and copyout functions would not result in a user-mode process being able to directly access kernel memory. The checks were not being used if the length was smaller than one page. This issue was addressed through additional validation of the arguments to copyin and copyout."

ETH zürich

Oldskool: signal ()

```
void (*signal(int sig, void (*handler)(int)))(int);
```

- Unpacking this:
 - A handler looks like `void my_handler(int);`
 - Signal takes two arguments...
 - An integer (the signal type, e.g. `SIGPIPE`)
 - A pointer to a handler function
 - ... and returns a pointer to a handler function
 - The previous handler,
- "Special" handler arguments:
 - `SIG_IGN` (ignore), `SIG_DFL` (default), `SIG_ERR` (error code)

ETH zürich

Unix signal handlers

- Signal handler can be called at *any time!*
- Executes on the current user stack
 - If process is in kernel, may need to retry current system call
 - Can also be set to run on a different (alternate) stack

⇒ User process is in *undefined* state when signal delivered

ETH zürich

Implications

- There is very little you can safely do in a signal handler!
 - Can't safely access program global or static variables
 - Some system calls are *re-entrant*, and can be called
 - Many C library calls cannot (including `_x` variants!)
 - Can sometimes execute a `longjmp` if you are careful
 - With `signal`, cannot safely change signal handlers...
- What happens if another signal arrives?

ETH zürich

Multiple signals

- If multiple signals of the *same* type are to be delivered, Unix will *discard all but one*.
- If signals of *different* types are to be delivered, Unix will deliver them *in any order*.
- Serious concurrency problem:
How to make sense of this?

ETH zürich

A better signal () POSIX sigaction ()


```
#include <signal.h>

int sigaction(int signo,
              const struct sigaction *act,
              struct sigaction *oldact);

struct sigaction {
    void (*sa_handler)(int);
    sigset_t sa_mask;
    int sa_flags;
    void (*sa_sigaction)(int, siginfo_t *, void *);
};
```


Annotations:

- New action for signal `signo`
- Previous action is returned
- Signal handler
- Signals to be blocked in this handler (cf. `fd_set`)
- More sophisticated signal handler (depending on flags)

ETH zürich  spoc.inf.ethz.ch
@spoc_eth


Signals as upcalls

- Particularly specialized (and complex) form of **Upcall**
 - Kernel RPC to user process
- Other OSes use upcalls much more heavily
 - Including Barrelfish
 - "Scheduler Activations": dispatch every process using an upcall instead of return
- Very important structuring concept for systems!

ETH zürich  spoc.inf.ethz.ch
@spoc_eth

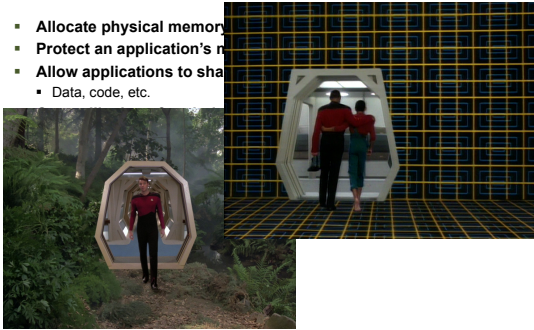
Our Small Quiz


- True or false (raise hand)
 - Mutual exclusion on a multicore can be achieved by disabling interrupts
 - Test and set can be used to achieve mutual exclusion
 - Test and set is more powerful than compare and swap
 - The CPU retries load-linked/store conditional instructions after a conflict
 - The best spinning time is 2x the context switch time
 - Priority inheritance can prevent priority inversion
 - The receiver never blocks in asynchronous IPC
 - The sender blocks in synchronous IPC if the receiver is not ready
 - A pipe file descriptor can be sent to a different process
 - Pipes do not guarantee ordering
 - Named pipes in Unix behave like files
 - A process can catch all signals with handlers
 - Signals always trigger actions at the signaled process
 - One can implement a user-level tasking library using signals
 - Signals of the same type are buffered in the kernel

ETH zürich  spoc.inf.ethz.ch
@spoc_eth

Goals of Memory Management


- Allocate physical memory
- Protect an application's memory
- Allow applications to share memory
 - Data, code, etc.



ETH zürich  spoc.inf.ethz.ch
@spoc_eth


In CASP last semester we saw:

- Assorted uses for virtual memory
- x86 paging
 - Page table format
 - Translation process
 - Translation lookaside buffers (TLBs)
 - Interaction with caches
- Performance implications
 - For application code, e.g., matrix multiply

ETH zürich  spoc.inf.ethz.ch
@spoc_eth

What's new this semester?

- Wider range of memory management hardware
 - Base/limit, segmentation
 - Inverted page tables, etc.
- How the OS uses the hardware
 - Demand paging and swapping
 - Page replacement algorithms
 - Frame allocation policies

ETH zürich  spoc.inf.ethz.ch
@spoc_eth

Terminology

- Physical address: address as seen by the memory unit
- Virtual or Logical address: address issued by the processor
 - Loads
 - Stores
 - Instruction fetches
 - Possible others (e.g., TLB fills)...

ETH zürich spezial@ethz.ch @spci_eth

Memory management

1. Allocating physical addresses to applications
2. Managing the name translation of virtual addresses to physical addresses
3. Performing access control on memory access

- Functions 2 & 3 usually involve the hardware Memory Management Unit (MMU)

ETH zürich spezial@ethz.ch @spci_eth

Simple scheme: partitioned memory

ETH zürich spezial@ethz.ch @spci_eth

Base and Limit Registers

- A pair of **base** and **limit** registers define the logical address space

The diagram shows a vertical stack of memory addresses from 0x0000000 to 0xffffffff. The top section is labeled 'Operating System' (0x0000000 to 0x1000000). Below it are three 'Process' blocks. The first process block starts at 0x5600ba0, with a 'base' register pointing to this address. The second process block ends at 0x8f20010, with a 'limit' register pointing to this address. The bottom section is labeled 0xb000000 to 0xffffffff.

ETH zürich spezial@ethz.ch @spci_eth

Issue: address binding

- Base address isn't known until load time
- Options:
 1. Compiled code must be **completely position-independent**, or
 2. **Relocation Register** maps compiled addresses dynamically to physical addresses

ETH zürich spezial@ethz.ch @spci_eth

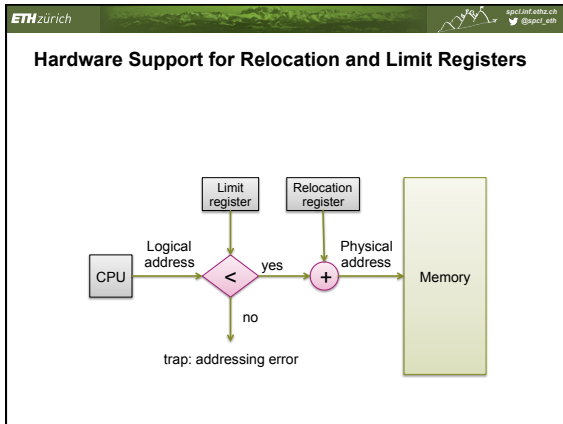
Dynamic relocation using a relocation register

The diagram illustrates the dynamic relocation process. A CPU provides a logical address of 346. This address is processed by the MMU, which also contains a relocation register with the value 14000. The MMU adds the logical address and the relocation register value (346 + 14000) to produce a physical address of 14346, which is then used to access memory.

ETH zürich spezial@ethz.ch @spci_eth

Contiguous Allocation

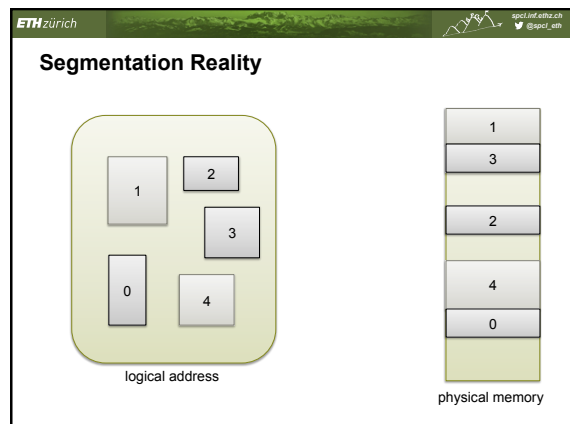
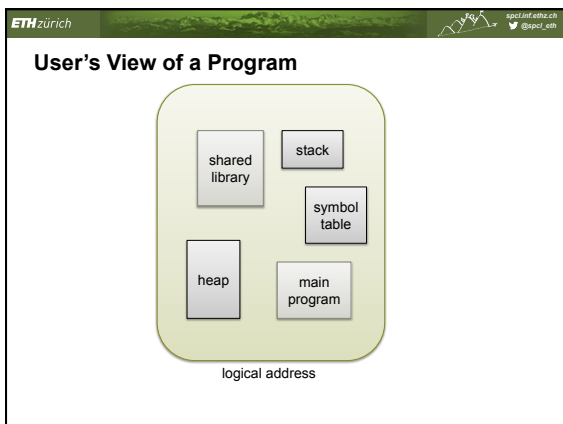
- Main memory usually into two partitions:
 - Resident OS, usually in low memory with interrupt vector
 - User processes in high memory
- Relocation registers protect user processes from
 1. each other
 2. changing operating-system code and data
- Registers:
 - **Base register** contains value of smallest physical address
 - **Limit register** contains range of logical addresses
 - each logical address must be less than the limit register
 - MMU maps logical address dynamically

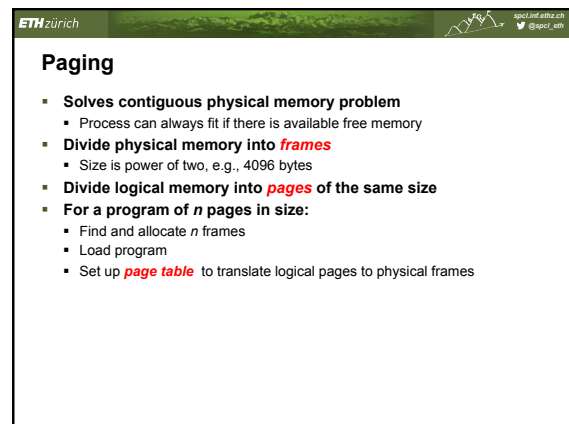
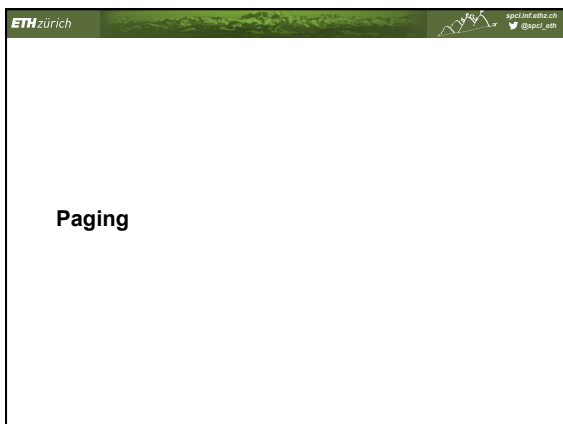
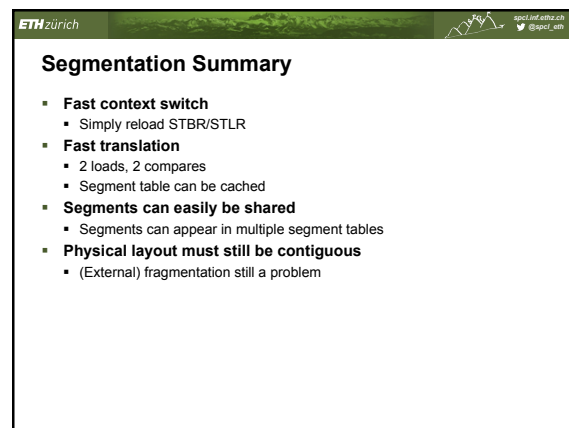
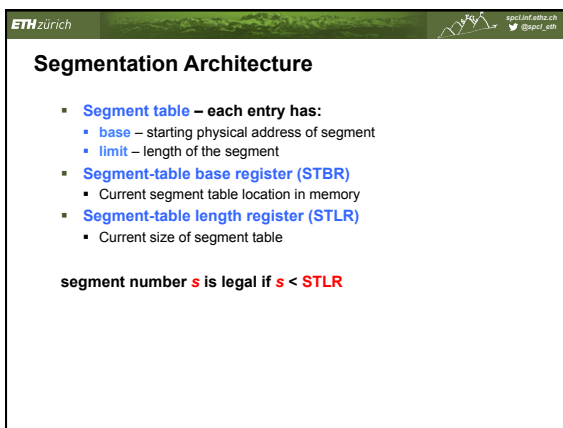
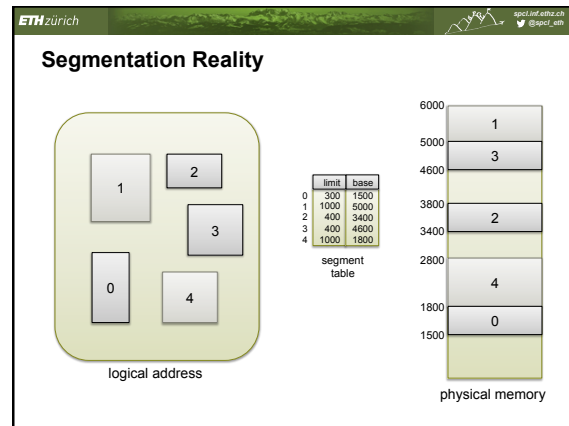
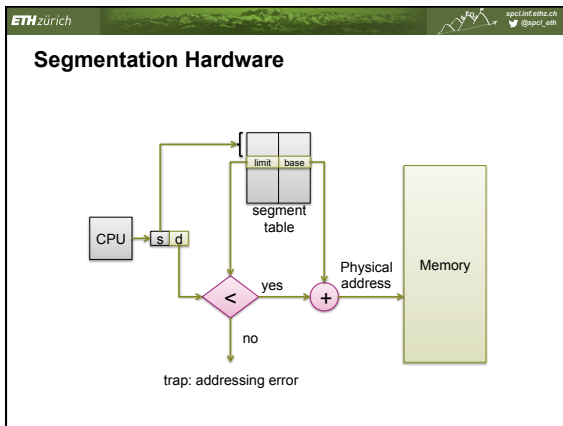


- ### Base & Limit summary
- Simple to implement (addition & compare)
 - Physical memory fragmentation
 - Only a single contiguous address range
 - How to share data between applications?
 - How to share program text (code)?
 - How to load code dynamically?
 - Total logical address space \leq physical memory

Segmentation

- ### Segmentation
- **Generalize base + limit:**
 - Physical memory divided into **segments**
 - Logical address = (segment id, offset)
 - **Segment identifier supplied by:**
 - Explicit instruction reference
 - Explicit processor segment register
 - Implicit instruction or process state





Page table jargon

- Page tables maps VPNs to PFNs
 - Page table entry = PTE
- VPN = Virtual Page Number**
 - Upper bits of virtual or logical address
- PFN = Page Frame Number**
 - Upper bits of physical or logical address
- Same number of bits (usually).

Recall: P6 Page tables (32bit)

- Pages, page directories, page tables all 4kB

x86-64 Paging

Problem: performance

- Every logical memory access needs more than two physical memory accesses
 - Load page table entry → PFN
 - Load desired location
- Performance ⇒ half as fast as with no translation
 - Solution: cache page table entries

Translating with the P6 TLB

- Partition VPN into TLBT and TLBI.
- Is the PTE for VPN cached in set TLBI?
- Yes:** Check permissions, build physical address
- No:** Read PTE (and PDE if not cached) from memory and build physical address

In fact, x86 combines segmentation and paging

- Segments do (still) have uses
 - Thread-local state
 - Sandboxing (Google NativeClient, etc.)
 - Virtual machine monitors (Xen, etc.)

Effective Access Time

- **Associative Lookup = ϵ time units**
- **Assume memory cycle time is 1 time unit**
- **Hit ratio α =**
 - % time that a page number is found in the TLB;
 - Depends on locality and TLB entries (**coverage**)

Then **Effective Access Time:**

$$EAT = (1 + \epsilon) \alpha + (2 + \epsilon)(1 - \alpha)$$

$$= 2 + \epsilon - \alpha$$

Assuming single-level page table.

Exercise: work this out for the P6 2-level table

Page Protection

Memory Protection

- **Associate protection info with each frame**
 - Actually no - with the PTE.
- **Valid-invalid bit**
 - "valid" \Rightarrow page mapping is "legal"
 - "invalid" \Rightarrow page is not part of address space, i.e., entry does not exist
- **Requesting an "invalid" address \Rightarrow "fault"**
 - A "page fault", or....

Remember the P6 Page Table Entry (PTE)

	31		12 11		9 8		7 6		5 4		3 2		1 0
Page physical base address		Avail		G	O	D	A	CD	WT	U/S	R/W	P=1	

Page base address: 20 most significant bits of physical page address (forces pages to be 4 KB aligned)

Avail: available for system programmers

G: global page (don't evict from TLB on task switch)

D: dirty (set by MMU on writes)

A: accessed (set by MMU on reads and writes)

CD: cache disabled or enabled

WT: write-through or write-back cache policy for this page

U/S: user/supervisor

R/W: read/write

P: page is present in physical memory (1) or not (0)

P6 protection bits

	31		12 11		9 8		7 6		5 4		3 2		1 0
Page physical base address		Avail		G	O	D	A	CD	WT	U/S	R/W	P=1	

Page base address: 20 most significant bits of physical page address (forces pages to be 4 KB aligned)

Avail: available for system programmers

G: global page (don't evict from TLB on task switch)

D: dirty (set by MMU on writes)

A: accessed (set by MMU on reads and writes)

CD: cache disabled or enabled

WT: write-through or write-back cache policy for this page

U/S: user/supervisor

R/W: read/write

P: page is present in physical memory (1) or not (0)

P bit can be used to trap on any access (read or write)

Protection information

- **Protection information typically includes:**
 - Readable
 - Writeable
 - Executable (can fetch to i-cache)
 - *Reference bits* used for demand paging
- **Observe: same attributes can be (and are) associated with segments as well**

ETH zürich

Page sharing

ETH zürich

Shared Pages Example

Process P₁

3
4
6
1

page table for P₁

0
1 data 1
2
3 code 1
4 code 2
5
6 code 3
7
8
9
10
11

ETH zürich

Shared Pages Example

Process P₁

3
4
6
1

page table for P₁

Process P₂

3
4
6
7

page table for P₂

0
1 data 1
2
3 code 1
4 code 2
5
6 code 3
7 data 2
8
9
10
11

ETH zürich

Shared Pages Example

Process P₁

3
4
6
1

page table for P₁

Process P₂

3
4
6
7

page table for P₂

Process P₃

3
4
6
2

page table for P₃

0
1 data 1
2 data 3
3 code 1
4 code 2
5
6 code 3
7 data 2
8
9
10
11

ETH zürich

Shared Pages

- Shared code
 - One copy of read-only code shared among processes
 - Shared code appears in same location in the logical address space of all processes
 - Data segment is not shared, different for each process
 - But still mapped at same address (so code can find it)
- Private code and data
 - Allows code to be relocated anywhere in address space

ETH zürich

Per-process protection

- Protection bits are stored in page table
 - Plenty of bits available in PTEs
- ⇒ independent of frames themselves
 - Different processes can share pages
 - Each page can have different protection to different processes
 - Many uses! E.g., debugging, communication, copy-on-write, etc.

ETH zürich

spoc.inf.ethz.ch @spoc_eth

Page Table Structures

ETH zürich

spoc.inf.ethz.ch @spoc_eth

Page table structures

- **Problem:** simple linear page table is too big
- **Solutions:**
 1. Hierarchical page tables
 2. Virtual memory page tables
 3. Hashed page tables
 4. Inverted page tables

ETH zürich

spoc.inf.ethz.ch @spoc_eth

Page table structures

- **Problem:** simple linear page table is too big
- **Solutions:**
 1. Hierarchical page tables
 2. Virtual memory page tables (VAX)
 3. Hashed page tables
 4. Inverted page tables

Saw these last Semester.

ETH zürich

spoc.inf.ethz.ch @spoc_eth

#3 Hashed Page Tables

- **VPN is hashed into table**
 - Hash bucket has chain of logical->physical page mappings
- **Hash chain is traversed to find match.**
- **Can be fast, but can be unpredictable**
- **Often used for**
 - Portability
 - Software-loaded TLBs (e.g., MIPS)

ETH zürich

spoc.inf.ethz.ch @spoc_eth

Hashed Page Table

logical address p d

hash function

hash table

q s • p r • ...

physical address r d

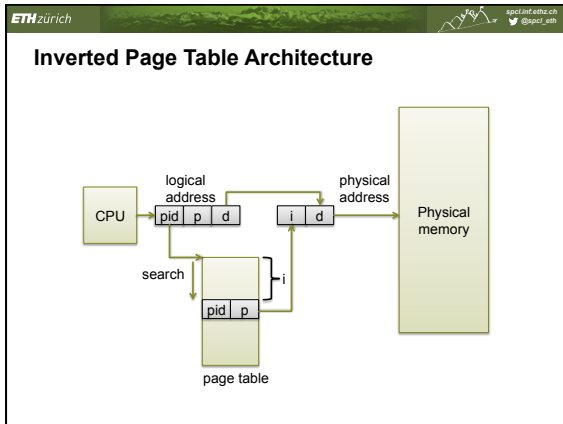
physical memory

ETH zürich

spoc.inf.ethz.ch @spoc_eth

#4 Inverted Page Table

- **One system-wide table now maps PFN -> VPN**
 - One entry for each real page of memory
 - Contains VPN, and which process owns the page
- **Bounds total size of all page information on machine**
 - Hashing used to locate an entry efficiently
- **Examples:** PowerPC, ia64, UltraSPARC



The need for more bookkeeping

- Most OSes keep their own translation info
 - Per-process hierarchical page table (Linux)
 - System wide inverted page table (Mach, MacOS)
- Why?
 - Portability
 - Tracking memory objects
 - Software virtual → physical translation
 - Physical → virtual translation

TLB shutdown

TLB management

- Recall: the TLB is a **cache**.
- Machines have many MMUs on many cores → many TLBs
- Problem: TLBs should be coherent. Why?
 - Security problem if mappings change
 - E.g., when memory is reused

TLB management

	Process ID	VPN	PPN	access
Core 1	0	0x0053	0x03	r/w
TLB:	1	0x20f8	0x12	r/w
Core 2	0	0x0053	0x03	r/w
TLB:	1	0x0001	0x05	read
Core 3	0	0x20f8	0x12	r/w
TLB:	1	0x0001	0x05	read

TLB management

	Process ID	VPN	PPN	access
Core 1	0	0x0053	0x03	r/w
TLB:	1	0x20f8	0x12	r/w
Core 2	0	0x0053	0x03	r/w
TLB:	1	0x0001	0x05	read
Core 3	0	0x20f8	0x12	r/w
TLB:	1	0x0001	0x05	read

Change to read only

ETH zürich spezialinf.ethz.ch @spezialinf

TLB management

	Process ID	VPN	PPN	access
Core 1	0	0x0053	0x03	r/w
TLB:	1	0x20f8	0x12	r/w
Core 2	0	0x0053	0x03	r/w
TLB:	1	0x0001	0x05	read
Core 3	0	0x20f8	0x12	r/w
TLB:	1	0x0001	0x05	read

Change to read only

ETH zürich spezialinf.ethz.ch @spezialinf

TLB management

	Process ID	VPN	PPN	access
Core 1	0	0x0053	0x03	r/w
TLB:	1	0x20f8	0x12	r/w
Core 2	0	0x0053	0x03	r/w
TLB:	1	0x0001	0x05	read
Core 3	0	0x20f8	0x12	r/w
TLB:	1	0x0001	0x05	read

Change to read only

Process 0 on core 1 can only continue once shutdown is complete!

ETH zürich spezialinf.ethz.ch @spezialinf

Keeping TLBs consistent

- Hardware TLB coherence**
 - Integrate TLB mgmt with cache coherence
 - Invalidate TLB entry when PTE memory changes
 - Rarely implemented
- Virtual caches**
 - Required cache flush / invalidate will take care of the TLB
 - High context switch cost!
 - ⇒ Most processors use physical caches
- Software TLB shutdown**
 - Most common
 - OS on one core notifies all other cores - Typically an IPI
 - Each core provides local invalidation
- Hardware shutdown instructions**
 - Broadcast special address access on the bus
 - Interpreted as TLB shutdown rather than cache coherence message
 - E.g., PowerPC architecture

ETH zürich spezialinf.ethz.ch @spezialinf

Tomorrow: demand paging